

the association for computational heresy

presents

a record of the proceedings of

SIGBOVIK 2024

*the eighteenth annual intercalary robot dance party in celebration
of workshop on symposium about 2⁶th birthdays; in particular,
that of harry q. bovik*

cover art by alexey crusoe cover image by chatGPT

carnegie mellon university

pittsburgh, pa

april 0, 2024



Association for Computational Heresy

Advancing computing as Tomfoolery & Distraction

SIGBOVIK

A Record of the Proceedings of SIGBOVIK 2024

ISSN 2155-0166

April 0, 2024

Copyright is maintained by the individual authors, though obviously this all gets posted to the Internet and stuff, because it's 2024.

Permission to make digital or hard copies of portions of this work for personal use is granted; permission to make digital or hard copies of portions of this work for classroom use is also granted, but seems ill-advised. Abstracting with credit is permitted; abstracting with credit cards seems difficult.

Additional copies of this work may be ordered from Lulu; refer to <http://sigbovik.org> for details.



SIGBOVIK 0x2024

Message from the Organizing Committee

time is small.

brevity is the wit of soul.

SIGBOVIK 2024 was full of plusplusgood articles.

A.C., H.P., J.M., T.7, S.B., S.M., and H.Q.B. were actively involved in publication and organization.

thus marks the end of the proceedings for this year.

-the SIGBOVIK '24/7

SUPPLEMENTAL APPENDIX 1

A: WE THE PEOPLE (plus maybe our pets and our favorite socks)	5
1 Programming Socks: Is it high time for thigh-highs? An investigation into the perceived unreasonable effectiveness of Programming Socks on productivity levels in the field of Software Engineering	6
2 Charge of the PhDs	12
3 Selected investigations into egg logic: a personal perspective	13
4 Unconventional Commits: An Exploration Beyond the Mundane	18
5 Getting Up And Running the λ - <i>Calculus</i>	23
6 Xtremely Delightful Random Development Environment and Dragons: A Novel Way of Programming	43
7 The Ballmer Peak: An Empirical Search	48
8 What is your dog likelihood?	55
9 Gender is Complex: pulling the Laplacian EigenGender from relationship graphs	60
B: IN ORDER (or maybe out of order)	69
10 We Found the Best Shuffled Deck	70
11 Proving P = NP thanks to the wonderful work of Reviewer Two	83
12 Gotta Collect 'em All	86
13 Mining for Gold Coins	94
14 An introduction to bogoceptionsort and its performance compared to ordinary bogosort	97
15 Systems for Rating Rating Systems with Rating Systems	100
16 Just How Random?: Introducing PROJECT S.P.O.R.K.	107
C: TO (pronouced "two")	121
17 Sleepy Dwarfs Somniloquy on Drowsy Logic Chip Design	122
18 The Magic School Bus Travels to Sub-threshold Voltage	125
19 Badness 0 (Epsom's version)	131
20 Badness 0 (Knuth's version)	145
21 Flaccid Drives: Storing Arbitrary Data in SIGBOVIK Articles	161
22 Retraction of "Flaccid Drives"	163
D: FORM A MORE PERFECT UNION (between machine and man)	169
23 Solving C's biggest flaw - Hemispheric divergence	170
24 An empirical performance evaluation between Python and Scratch	174
25 Exchange Traded Neural Network	176
26 Minmaxing the energy efficiency of biological computing	180
27 Advancing Consensus: Automated Persuasion Networks for Public Belief Enhancement	188
28 Grounding Language Models to their Physical Presence	195
29 Quantum Disadvantage: Simulating IBM's 'Quantum Utility' Experiment with a Commodore 64	199

E: ESTABLISH JUSTICE (for Gotham and other places)	207
30 Saving the legacy of Hero Ibash: Evaluating Four Language Models for Aminoacian	208
31 Stale Diffusion: Hyper-Realistic 5D Movie Generation using Old-school Methods	214
32 Guidelines for the Development of Irresponsible AI	219
F: ENSURE TRANQUILITY (insert evil laugh here)	223
33 Just-too-late compilation - An examination of a post-emptive compilation technique	224
34 Im going to Hurl	226
35 We Have Serverless at Home	231
36 USB-*: New Extensions for the Universal Serial Bus	235
37 This SSH servers never gonna give you up (but it will probably let you down)	237
38 Towards an AI-Exclusive Higher Education System	242
G: PROVIDE FOR THE COMMON SIGBOVIK	245
39 The Anything Prover: How to prove Anything	246
40 The Perception and Access of Written Silence	247
41 An Empirically Verified Lower Bound for The Number Of Empty Pages Allowed In a SIGBOVIK Paper	249
42 Shunting Work: A Demonstration	252
43 Title of PDF	253
44 AMOR: Ambiguous Authorship Order	257
45 Debunking the April 1 model of SIGBOVIK occurrences	258
H: PROMOTE THE GENERATIVE WELFARE	275
46 An Optimal Control Approach to Graphic Design	277
47 ITF;)LM: Innocuous Table Formatting ;) with Language Models	284
48 Can Machines Feel? Novel Affective Layers for Pictorial Preprocessing and Scalar Fusion in CNN Representations	286
49 Beyond Prompt Engineering: Tardy Engineering	295
50 CamelGPT: A Viable Small Language Model	296
51 Revealing AGI Risks with a Drop of Ink	303
52 Exploring the Viability of Utilising Multi-Modal GPT Models with Local Hardware for Image Text Detection	307
53 Bean There, Done That: Computer-Assisted Design of Bean Sculptures . .	310
54 Fractal Overdrive: An Aesthetic Evaluation of Numeric Error	316
55 Diffusion Local Time: hard real-time multilingual data visualization via multimodal-LLM generative AI on heterogeneous edge devices for extremely high-impact chronometry and extremely low-cost neurological diagnostics . .	325
56 Undergrads Are All You Have	338
57	343
58 Are Centaurs Actually Half Human and Half Horse?	367

I: AND SE () RE	371
59 You Shall TacItly Understand the Goals of ThIs Paper	372
60 A computer-assisted proof that e is rational	375
61 A Brief History of Gender Theory: or, Much Ado About Bathrooms	380
62 DeterMNISTic: a Safer Way to Classify Handwritten Digits	390
J: RIN TLHOBQU'	393
63 How does the AI community pronounce epoch? A semirigorous sociolinguistic survey	394
64 Toki Pona and Orders of Semantic Completeness	398
65 Introducing Supernatural Language Processing (SNLP)	411
66 SMS: Sending Mixed Signals	416
K: TO OURSELVES (and those named like us)	423
67 Et Al(ex): Examining the impact of Alexs on the field of computer science .	424
68 An Abundance of Katherines: The Game Theory of Baby Naming	427
L: AND OUR POSTERITY	437
69 A Genius Solution: Applications of the Sprague-Grundy Theorem to Korean Reality TV	438
70 A Secondhand Understanding of Reality: Infinite Craft Subtleties	448
M: DO ORDAIN AND ESTABLISH..... 'MERICAAA	469
71 As American as Apple Pie: A Search to Define Americanness in the Context of all Apple Pie Recipes on the Internet	470

WE THE PEOPLE (plus maybe our pets and our favorite socks)

- 1 Programming Socks: Is it high time for thigh-highs? An investigation into the perceived unreasonable effectiveness of Programming Socks on productivity levels in the field of Software Engineering**

Ian F.V.G. Hunter

- 2 Charge of the PhDs**

Alfred, Doctorate Tennyson

- 3 Selected investigations into egg logic: a personal perspective**

Jana C. Dunfield

- 4 Unconventional Commits: An Exploration Beyond the Mundane**

H. Dog, T. H. Underpoot, and P. J. Ox, PhD

- 5 Getting Up And Running the λ - Calculus**

George(s) Zakhour

- 6 Xtremely Delightful Random Development Environment and Dragons: A Novel Way of Programming**

Refracted Light System[s Incorporated

- 7 The Ballmer Peak: An Empirical Search**

Twm Stone and Jaz Stoddart

- 8 What is your dog likelihood?**

Gemmechu Hassena, Simret Gebreegziabher, and Bharath Hariharan

- 9 Gender is Complex: pulling the Laplacian EigenGender from relationship graphs**

Anonymous

Programming Socks: Is it high time for thigh-highs? An investigation into the perceived unreasonable effectiveness of Programming Socks on productivity levels in the field of Software Engineering

Ian F.V.G. Hunter

Unaffiliated/Disowned, Ireland

Abstract

Programming Socks are an essential piece of equipment for any respectable software engineer. However, there is limited research on the level of impact that they have on practitioners. In this paper we investigate the claim that Programming Socks are a low-cost high-impact upgrade to engineers' wardrobes and find that they have much untapped potential.

1 Introduction

Programming Socks are thigh length or “Thigh-High” socks, commonly striped but not as a critical necessity (See Figure 1 for an example). They have been long-hailed as essential items to improve programmer performance, code quality as well as job satisfaction.

Other related items that are known to be useful for programming include: cat tails, cat ears, mittens and mini-skirts.¹ Surprisingly, garments traditionally designed for the duties of a French housemaid are purportedly equally effective — if not more so — in the realm of software development. According to informal polls with select focus groups², public opinion suggests that the attire helps to promote ‘clean code’ practises.

Other innovations of note include replacing the long-beloved practise of talking to rubber ducks to assist in debugging issues[3]. A cuddlier substitute has been found in the popular Swedish furniture store IKEA. The plush shark ”Blåhaj” appears to be more appealing to both new and old developers alike. The new practise of verbal debugging with these shark toys has been coined as “Blah-Blah-Blåhaj”³, but the term has yet to catch on.



Figure 1: An example Amazon listing [1]. 900+ bought in the last month⁴. 4.4/5 average star rating

¹It is often misunderstood that the widespread disappointment over the open-source event ‘Hacktoberfest’ was about the removal of the complimentary participatory shirts, but the refusal to provide skirts, which would have been an investment into open-source quality for years to come.

²i.e. My friends

³by me, just now

⁴December 2023

2 Data Gathering & Analysis

A set of images were gathered of individuals posting online about their socks of choice in front of their computer screen. Key attributes were gleaned from the data, using widely-available organic ocular systems.

⁵ ⁶

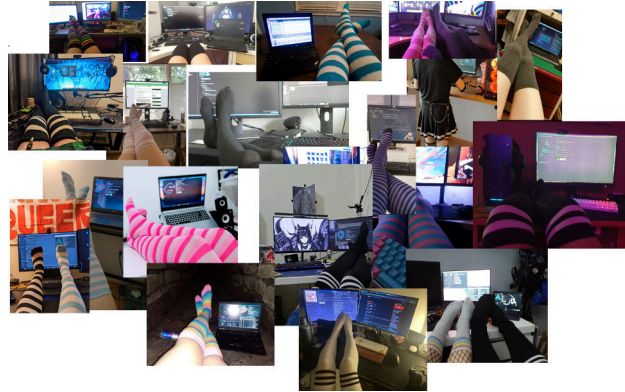


Figure 2: Samples from the dataset.

A significant source of “feet pics” was on the social news aggregation site Reddit.com. While some communities on Reddit are dedicated to the promotion of the practise (e.g. r/unixsocks/ [4]) there are several examples of Programmer Socks under discussion in other areas of the website (e.g. r/brisbane [5] of Australia)

3 Results and Analysis

3.1 Effects on Programming

The principal metric we wished to observe was whether programmers who wore Programming Socks were truly as productive as expected. Figure 3 shows a breakdown of technical content on wearers’ screens. There was a surprising lack of interactive development environments, text editors or other signs of code production. However, the overwhelming majority of wearers had terminals visible, often with prominent ASCII art logos in transgender or non-binary colour schemes⁷.

From this data, we can only conclude that either: Programming Socks have become more main-stream and non-programmers wearing these socks are developing into programmers themselves — or that Programming Socks are so effective that the majority of wearers in our dataset simply finished their software and did not need to do any additional coding. Either way, this result is a clear endorsement of their usage. Further research is desirable to investigate these effects when paired with the additional attire mentioned in Section 1.

⁵Note that it is unclear if the effects of Programming Socks extend to data collection, which may have effected the quality of the dataset.

⁶This is the reason I had all these pictures of feet on my computer, I swear.

⁷Actually all of the terminal logos in the dataset were non-binary. I suppose that’s because it’s hard to draw things with only ones and zeroes without having a very zoomed out screen

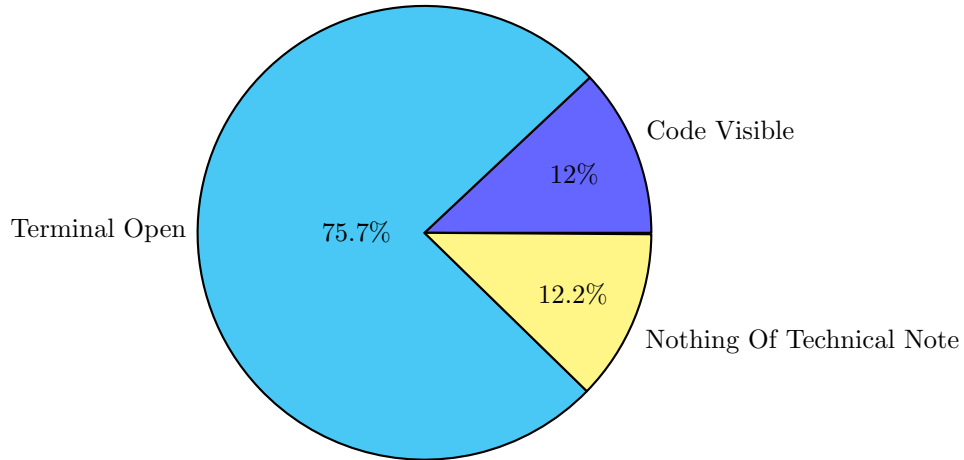


Figure 3: Evidence of coding on-screen

3.2 Colour Preference

A wide variety of Programming Socks are available. However, there are some general trends in the space. Striped socks are significantly more popular than non-striped socks (As shown in Figure 4) and there is a clear preference for some basic colours - black, white, pink.

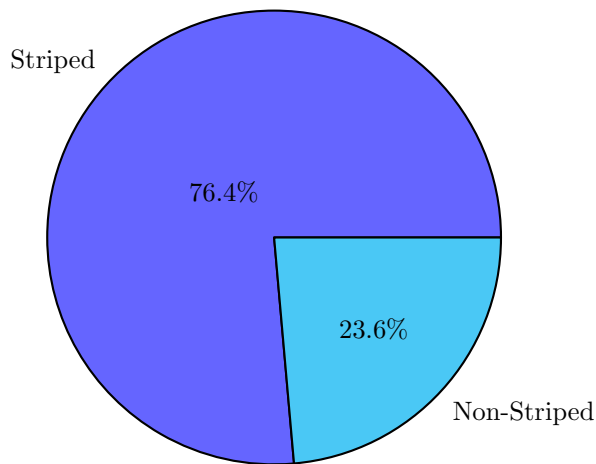


Figure 4: Popularity of Stripes

It was surprising to discover that rainbow socks were not chosen as often as other colour combinations, because it is obviously the best choice.

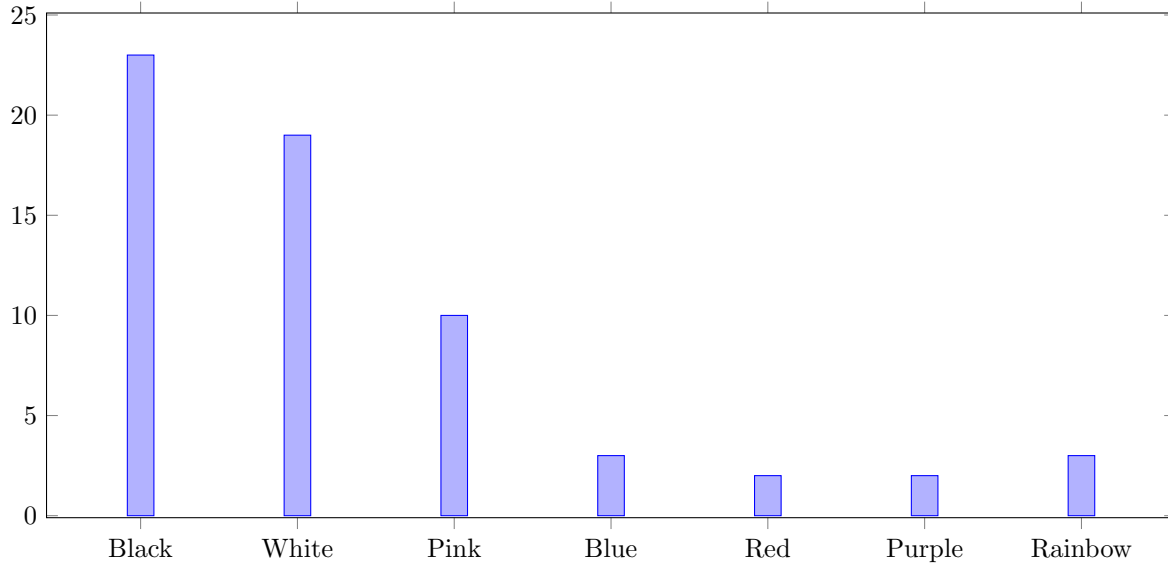


Figure 5: Colour Preferences

4 Appendix

We include other findings of interest below.

4.1 Amazon Search Listings

When searching for Programming Socks on Amazon.com, the majority of listings are correctly indexed as thigh-high socks. There are a limited amount of novelty socks with a programming theme⁸ and an equal amount of novelty socks that were unrelated to programming⁹. There were also a few shoes and gloves.

We decided to not divide Programming Socks into sub-categories, such as ones that had paw pads on their toes, fluff and glitter or small teddy bear faces at their peak.

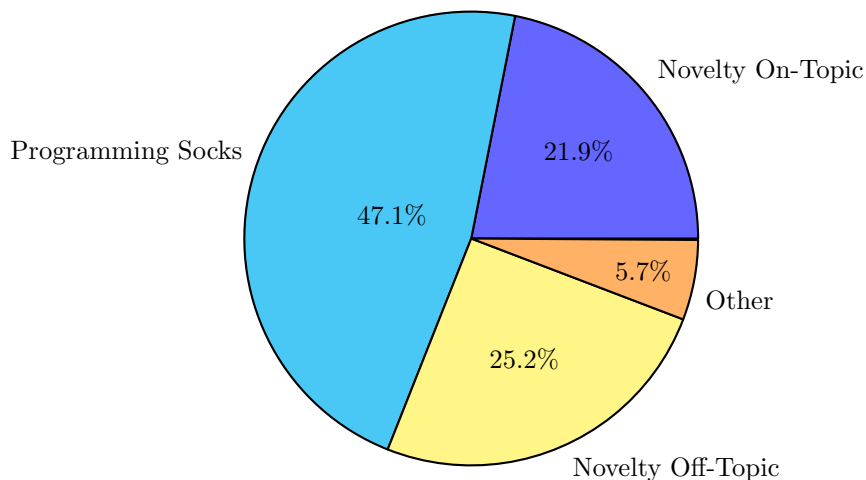


Figure 6: Amazon Listings for Programming Socks

⁸“I will NOT fix your computer”, other hilarious jokes and socks with logos of programming languages, code, etc.

⁹Such as socks that make your feet look like chicken’s feet, or to fool people into thinking your feet are being eaten by infeasibly proportioned crocodiles.

4.2 Interest in Programming Socks

Google Trends [6] gave us some interesting insight into the widespread appeal of Programmer Socks. They are particularly of interest to Western countries (United States, Australia, United Kingdom) but also Eastern Europe (Poland, Russia) (See Figure 7).



Figure 7: Interest in 'Programming Socks' in different countries

Programming Socks have become increasingly of interest since the early 21st century (2018) and continue to grow (Figure 8). As the world consistently increases the number of software developers in its population [7] and the benefits of Programming Socks become better known, this number will surely continue to skyrocket.

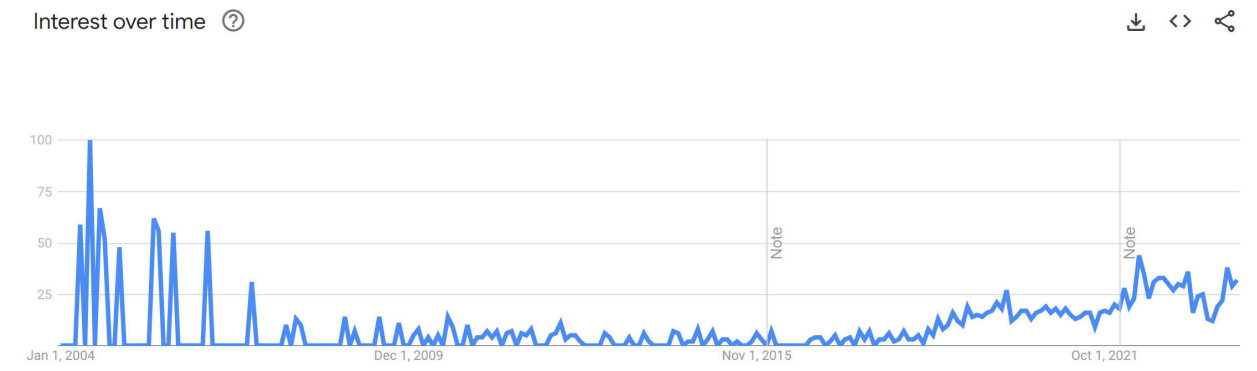


Figure 8: Popularity of the 'Programming Socks' search term over time (worldwide)

The data ten years prior to that inflection point is particularly interesting. Interest was much higher and varied in the years 2004-2006. We can only conclude that Programming Socks may have a deeper history than previously known, possibly a well guarded secret by the programming elite, or maybe a handful of developers in the know who made irregular spike purchases of huge orders in order to save on shipping costs.

5 References

References

- [1] Amazon.com: DRESHOW Extra Long High Thigh Socks Striped Over Knee Thin Tights Long Stocking. Available at: <https://www.amazon.com/DRESHOW-Striped-Tights-Stocking-Warmer/dp/B096M3LMCT/> (Accessed: 27 January 2024).
- [2] Amazon.com: Programming Socks. Available at: <https://www.amazon.com/programming-socks/s?k=programming+socks> (Accessed: 27 January 2024).
- [3] Debugging, R.D. Rubber duck debugging, – Rubber Duck Debugging – Debugging software with a rubber ducky. Available at: <https://rubberduckdebugging.com/> (Accessed: 27 January 2024).
- [4] KDESOCKS-bin/Lemmmmmmmmmmysocks • R/unixsocks (no date) reddit. Available at: <https://old.reddit.com/r/unixsocks/> (Accessed: 27 January 2024).
- [5] (The quest for Programming Socks: R/brisbane - reddit. Available at: https://www.reddit.com/r/brisbane/comments/urbq0c/the_quest_for_programming_socks/ (Accessed: 27 January 2024).
- [6] Programming Socks - Google Trends. Available at: <https://trends.google.com/trends/explore?date=all&q=programming%20socks&hl=en> (Accessed: 27 January 2024).
- [7] Number of software developers worldwide in 2018 to 2024 <https://www.statista.com/statistics/627312/worldwide-developer-population/> (Accessed: 27 January 2024).

Charge of the PhDs

Alfred, Doctorate Tennyson

I

Half a year, half a month,
Half a day deadline,
All in the valley of Academia
Rode the six hundred.
"Forward, the PhDs!
Apply for the grants!" he said.
Into the valley of Academia
Rode the six hundred.

II

"Forward, the PhDs!"
Was there a man dismayed?
Not though the student knew
Someone had blundered.
Theirs not to make reply,
Theirs not to reason why,
Theirs but to publish and cry.
Into the valley of Academia
Rode the six hundred.

III

Deadline to right of them,
Poverty to left of them,
Unemployment in front of them
Overworked and malnourished;
Swamped with meetings and conferences,
Boldly they studied and well,
Into the jaws of Academia,
Into the mouth of hell
Rode the six hundred.

IV

Published all their proofs bare,
Published as rebuttals were quelled
Questioning the reviewers there,
Attending a conference, while
All the world wondered.
Plunged in the faculty track
Right through the hiring decisions;
Teaching and Research
Powered with junior faculty grants
Submitted and published.
Then they earned tenure, but not
Not the six hundred.

V

Deadline to right of them,
Therapy to left of them,
Sleep left behind them
Depressed and traumatized;
Swamped with thesis writing hell,
While mastering out many fell.
They that had taught so well
Came through the jaws of Academia,
Back from the mouth of hell,
All that was left of them,
Left of six hundred.

VI

When can their glory fade?
O the wild papers they made!
All the world wondered.
Honour the research they made!
Honour the PhDs,
Noble six hundred!

Selected investigations into egg logic: a personal perspective

Jana C. Dunfield

somewhere in Katarokwi-Kingston,
unless she's somewhere else

\today

Content note

This paper discusses *egg logics*, which are known to the State of California to exacerbate denial of one's gender.

This paper also discusses personal health issues, street harassment, uses a disputed term because it was literally used in an employment equity survey but also because I get to use that term occasionally (as a treat), and contains possibly unhealthy amounts of self-deprecation.

1 A judgmental reconstruction of my egg logic

Ashley (2019) defines *transitude*¹ as “the fact of being trans”. The traditional judgment form of *transitude* asserts that a particular individual is trans. Examples:

1. Jana is trans
2. Jana is extremely trans

1.1 Syllogism of self-deprecation

By deploying probabilistic denial, I obtained (Dunfield 2005) a classic syllogism of *judgmental non-refinement*:

- | | | |
|------------|----------------------------------|------------------------------|
| (1a) | Only 1 in 1000 people are trans. | (major premise) |
| (1b) | Being unusual is interesting. | (minor premise) |
| (1selfdep) | I am not interesting. | (self-deprecation condition) |
| (1m) | Therefore, I am not trans. | (motivated conclusion) |

First, let us observe that the major premise is false. (I will not dignify the 2021 Canadian census by citing it properly, but it did show that even the subset of trans people who are willing to directly

¹*Transness* is synonymous, but way less fun to write.

inform a national government that they are trans, possibly also outing themselves to their entire household, is much greater than 0.1%.)

Now consider the self-deprecation condition, the *sine qua non* of the syllogism. Why did she² assume she couldn't be interesting? Because she wasn't sure she existed. If I barely existed, how could I be interesting?

By this specious reasoning, admissible in egg logic, I obtained the motivated conclusion.

1.2 Bayes' rule, or, how did I pass Machine Learning on the first try?

Let us aggravate the underestimation of trans people (cf. (1a) above), by supposing that only 1 in 1000000 people are trans. The structure of the judgmental refinement syllogism does not change; we may add "extremely", "very", etc., obtaining e.g. "Being extremely unusual is very interesting".

The logic remains: not many people are trans, so the probability that I am trans is low.

Similarly, extremely few people are me. By egg logic, the probability that I am me is one in several billion. But the actual probability that I am me *given that I am me* is 1.

The probability that I am trans *given that I am trans* is 1.

1.3 Syllogism of self-preservation

- (2a) Pittsburgh in 2005 is not the safest place and time for transitude. (major premise)
- (2b) I am living in Pittsburgh in 2005. (minor premise)
- (2c) I would not be safe. (conclusion)

Unlike the syllogism of self-deprecation, this syllogism has some legitimate structure. However, I failed to elaborate an implicit assumption:

- (2no) Being trans is not important. (central fallacy)

I don't know if I would have been able to *do* anything about being trans in 2005. But I wish I'd known.

1.4 Egg induction and dream logic

Eventually, I accepted that the main proposition (*Hauptsatz*) of transitude, "I am trans", was true or at least plausible. In egg logic, however, the mere acceptance of this proposition does not entail *doing* anything (for example: saying any actual words to anyone) about one's transitude.³ Instead, egg logic admits the following deduction:

- (3a) I have survived so far without doing anything about my gender. (previous states premise)
- (3b) I do not need to do anything about my gender. (pseudo-induction)

I occasionally dream that I am late for an exam. In my dream logic, the mere fact that I am not taking classes is usually insufficient to refute the possibility of being late for an exam. For example, if the putative exam is for undergrad calculus, not taking classes is insufficient. Even having already received an undergrad degree is not sufficient. No, my dream logic requires even more refutation:

²She gets to be herself.

³Use a hypothesis? In *this* economy?

I tell myself that I have a PhD, which is gatekept by having received an undergrad degree, so it is triply impossible that I could be missing an undergrad calculus exam.

In egg logic, the hypothesis of being cis (Reed 2012) must be refuted again and again and again: by a new pair of jeans and a flight attendant and a guest at breakfast, and by shouting into the void of a zero-follower social media account until I almost heard myself.

All that wasn't enough. I had to get within spitting distance of my own mortality to realize the one thing I couldn't bear: the idea of dying and being remembered as a man.

1.5 Author's ongoing self-deprecation

One of my several current forms of self-deprecation is that I often think that being trans is the only interesting thing about me.

Hey, I'm trying.

2 Sequents

It is the 1990s. Puberty happens. Puberty, I am told, is no fun for anyone. It's not fun for me, but I don't understand why. I'm homeschooled, and I am very rarely harassed or bullied, but I can count my friends on one finger. I very carefully ask myself if I'm gay, but I don't feel attracted to men so I conclude that I'm not. I read a bunch of feminist SF; I don't remember why.

I survive.

It is 1998. I put a quote from Le Guin's *Introducing Myself* on the door of my dorm room.⁴ I hear about a panel where a bunch of trans people are going to speak. I don't go.⁵ Later that month, I visit my parents. I go for a walk. Some dude in a pickup truck yells a homophobic slur at me. Two and a half hours later, I send an email to myself, criticizing Truck Dude for his faulty gaydar. I sign the email "-j.", as I always did.

It is 1999 and I won't get into all of that here but I write things in private text files like "There was more to this story, but it doesn't need telling now."⁶

It is the 2000s. I don't know who I am, but I'm not trans, because I did an egg logic with my super PL brain and I'm not trans.

It is 201X. I almost check out *Whipping Girl* (Serano 2007) from the UBC library, but I am afraid that someone will notice even though the library has self-checkout and I have my own office. Or so I tell myself. A friend makes a relevant and accurate observation about how other people perceive my gender, and I blow up, but he made the observation in email so I'm spared him having to see me blow up.⁷ I shouldn't be offended at the suggestion that I'm gender non-conforming. There's nothing wrong with being gender non-conforming. There wouldn't be anything wrong with being trans, it's just that I'm

⁴It was the part about the fish stick. There is no pattern here. Do not look in Section 3 for a pattern.

⁵I extracted my class schedule from the Internet Archive's copy of my undergrad institution's website and I didn't have any classes during it.

⁶I desperately wish I could go back in time—not to let her in on the trans thing, just to tell her to be less cryptic.

⁷I remembered stewing for hours or probably days about this before I could reply to him. I actually replied in under an hour. It only felt like days.

It is 2017. I play *Undertale* and become emotionally invested in fan debates over whether the player character is canonically nonbinary. And now it's 2018. Am I nonbinary? That feels less off, but not *good*. Am I a trans woman?

Oh shit

It is still 2018. I shorten my first name to an initial in the page header of a paper.⁸ I probably have cancer. I have surgery. It's not a big deal, because I'm probably not here. I definitely have cancer. I have surgery again.

It is, somehow, still 2018. I drink radioactive iodine. I get up to pee. I think I'm real.

It is 2019. I write "Jana" in the sand at Cascais and it's swept away. I start telling people.⁹ I start an email draft called "epistemic logic" to keep track of who I've come out to.

It is 2020. The "epistemic logic" draft has 67 names on it. I come out at work¹⁰ and have nightmares but it's actually pretty okay.

It is 2021. The pandemic careens on. The state careens on (Gill-Peterson 2021).

It is 2024. The Canadian state stirs and careens towards careening.

I am real and I'm singing with my friends.

3 Related work

The closest I've come to submitting to SIGBOVIK was around 2013 when I started to write something arguing (mostly seriously) that we shouldn't use "guys" to refer to mathematical abstractions. The throwaway joke in the author block was a parenthetical "(is a guy)", attached to my name.

I'm not sure that was a funny joke, but it's funny that I thought it was a joke.

4 Future work

Possible future work includes the focusing logic of bureaucracy. For example, when you are told that it's hard to change the name on your PhD thesis because your thesis is a legal document, you might be tempted to use conventional deduction to argue that a thesis is no more a legal document than your algorithms homework was, or to argue that your PhD-granting institution is located in a jurisdiction in which legal names aren't a thing (Baker and Green 2021). However, experience (Bohrer 2023) demonstrates that the use of a hypothesis is admissible in bureaucratic logic only when the hypothesis is useful to the bureaucracy.

⁸When ACM eventually got around to manually editing PDFs to update my name, they didn't really need to change those headers. Saved you some work, ACM. You're welcome, ACM.

⁹"I'm more frightened and more happy than I've been in a long time. . . . I've been asking myself, 'Is this worth turning your life upside down?' But maybe my life isn't the right way up now."

¹⁰Out of some regrettable desire to improve my employer's diversity stats, or as a bid for reassurance and official recognition, or maybe as a gesture towards the effacement of my previous existence, I went to the trouble of *re-filling* out my employer's employment equity survey. That instrument included a bracing variety of transitudinous options, including "transsexual". Why not, I thought? It was only while writing this footnote that it struck me: if someone hasn't already made a graph of the number of self-identified transsexuals employed here over time, there is enough data to make one.

Acknowledgments¹¹

I thank my past self, despite everything.

I thank my trans friends; I still took a while, but that's okay. Blame is for PL nerds, right?

I thank all my friends.

Even if you're not my friend, you're getting thanked too. You're reading this.

References

Florence Ashley. Gatekeeping hormone replacement therapy for transgender patients is dehumanising. *Journal of Medical Ethics*, 45(7):480–482, 2019. doi: 10.1136/medethics-2018-105293. URL <https://jme.bmj.com/content/45/7/480>.

Austin A. Baker and J. Remy Green. There is no such thing as a “legal name”. *Columbia Human Rights Law Review*, 53(1):129–188, 2021. URL <https://hrlr.law.columbia.edu/hrlr/there-is-no-such-thing-as-a-legal-name/>.

Rose Bohrer. Screen-sharing concurrency. In *Proceedings of SIGBOVIK 0x2023*, pages 220–225. Association for Computational Heresy, 2023.

Jana Dunfield. Literally what was I thinking. Personal communication. From herself, 2005. Year approximate.

Jules Gill-Peterson. The Cis State, April 2021. <https://sadbrowngirl.substack.com/p/the-cis-state>.

Natalie Reed. The Null HypotheCis, April 2012. <https://freethoughtblogs.com/nataliereed/2012/04/17/the-null-hypothecis/>.

Julia Serano. *Whipping Girl: a transsexual woman on sexism and the scapegoating of femininity*. Seal Press, 2007.

¹¹<https://tex.stackexchange.com/questions/571314/make-a-color-box-with-the-trans-pride-flag>

Unconventional Commits: An Exploration Beyond the Mundane

H. Dog

T. H. Underpoot

P. J. Ox PhD

Department of Physics

Institute for the Study of Advanced Procrastination

`phys@whatsthat.overthere.edu`

March 21, 2024

Abstract

This paper presents "Unconventional Commits", a pioneering alternative to the Conventional Commits specification. We examine methods of commit messaging that defy conventional norms, marrying the methodical with the innovative. Our exploration aims to expand the conceptual boundaries of code versioning, asserting the significance of a multifaceted approach to enhancing the development process.

1 Introduction

Traditional views on software development often emphasise efficiency and pragmatism, as embodied by the Conventional Commits specification. This perspective, however, neglects the diversity of thought and innovation inherent in the field. "Unconventional Commits" is introduced as a counterpoint to this norm, embracing a broad spectrum of ideas to enrich the narrative of code evolution.

2 Background

2.1 Conventional Commits

The Conventional Commits specification offers a structured framework for commit messages, supporting version control and changelog generation. While it provides clear benefits for automation and clarity, it may limit the scope for creativity and individual expression.

2.2 The Perils of Joke Commits

Platforms like *What the Commit* demonstrate the use of humour in commit messages, which, while amusing, can detract from the functional and informational value of commit histories. Such practices are discouraged in professional environments due to their potential to undermine the integrity of the development process.

3 Proposing Unconventional Commits

The concept of "Unconventional Commits" introduces novel methods for attaching messages to git commits, emphasising the balance between innovation and integrity. The following methods are proposed:

3.1 The Novel Commit

This method involves encoding the commit message within the narrative arc of a novel, where the message is revealed through careful analysis of the text. The commit reference directs to a specific passage for interpretation.

3.2 Musical Commits

Here, commit messages are transcribed into musical compositions, with the notation corresponding to a coded format. The commit provides a recording or sheet music for decoding. MIDI files are also a possibility.

3.3 QR Code Image

Commit messages are converted into QR codes, then printed, and scanned back into a digital format. The commit includes the image, requiring optical decoding to access the message.

3.4 Commit in a Bottle

Emulating historical message dissemination methods, commit messages are physically secured in a container and hidden. The commit provides geographical coordinates, inviting a physical search.

3.5 Culinary Code

Commit messages are embedded within culinary recipes, with each ingredient and cooking step representing elements of the encoded message. The commit challenges the developer to decipher the message through preparation.

3.6 Cryptic Crossword Puzzle

Commit messages are concealed within the answers to a specially designed crossword puzzle. The commit includes this puzzle, engaging the developer's problem-solving skills to decode the message.

3.7 Steganography in an Image

This method hides commit messages within images using digital steganography, challenging the developer to employ specific algorithms for extraction.

3.8 A Journey Through Software

A complex program is developed to reveal the commit message upon execution. The program's source code, obfuscated and intricate, serves as the commit, with execution as the key to decryption.

3.9 Video Game Easter Egg

A video game contains the commit message, accessible only upon completing challenging tasks or achieving high scores. The commit includes the game, intertwining development with interactive problem-solving.

3.10 The Time Capsule

Commit messages are physically archived and concealed, with instructions for future retrieval. This method links developers across generations in a shared quest for discovery.

4 Discussion

The introduction of "Unconventional Commits" juxtaposes the conventional with the innovative, highlighting the potential for a diverse array of commit messaging techniques to coexist alongside traditional methods. This exploration not only enriches the software development process but also fosters a culture of creativity and exploration.

5 Conclusion

"Unconventional Commits" presents a novel framework that challenges the traditional confines of commit messaging. By embracing innovative and diverse methods, this approach encourages a reevaluation of the possibilities within software versioning, promoting a multidisciplinary perspective in the realm of version control. The exploration of unconventional methods serves not only to enrich the developer's toolkit but also to inspire a deeper engagement with the process of code evolution, demonstrating that the act of committing code can transcend its utilitarian roots to embrace a broader narrative of creativity and shared human experience.

References

- [1] Git. *Git Documentation*. <https://git-scm.com/doc>

- [2] Conventional Commits. *Conventional Commits 1.0.0*. <https://www.conventionalcommits.org/en/v1.0.0/>
- [3] What the Commit. *What the Commit*. <https://whatthecommit.com>
- [4] Stephen P. Blung *The Role of Pizza in Software Development*. Journal of Edible Communication, 1988.

Getting Up and Running the λ -Calculus

GEORGE(S) ZAKHOUR

Abstract (Spoiler Alert) – Programmers, developers, and coders suffer from a myriad of issues pertaining to their health. These can vary from eye redness to repetitive strain injuries and a diminished life expectancy. A common exercise that software practitioners can engage in to reduce these health averse conditions is to be in the outdoors and to move—by means of walking or running—longer. Alas, professionals are hesitant to pursue their physical and mental well-being as that sacrifices productivity and programming-derived joy.

In this paper, we address this problem through the lens of programming languages and provide a solution that greases the friction between outdoor activities and programming. Through the insight that humans leave a trace while moving and the observation that apparatuses to record such traces are ubiquitous, we formalize an encoding of the λ -calculus in those traces. We develop an alternative front-end to Church’s language which we call Poolloop and we provide a reference compiler that produces Haskell and Scheme code. We evaluate Poolloop on two use-cases that we ran and show that the compilation runs in a few milliseconds.

CCS Concepts: • **Social and professional topics** → Computing occupations; • **Human-centered computing** → *Interaction techniques*; • **Theory of computation** → **Formalisms**; *Grammars and context-free languages*.

Additional Key Words and Phrases: λ -calculus, Programmer Health, Runtimes

1 INTRODUCTION

The attentive reader would have noted that computers and the Internet are ubiquitous. For the clueless reader: in the 2010s, 3.04 billion personal computers were shipped [33] in part to the 5.35 billion Internet users [20]. Those—the PCs—are being commandeered, in 2024, by an estimated 28.7 million software developers [19], showing that they too—the software developers—are ubiquitous.

Popular Programming Plagues. It is not too uncommon to witness a plethora of software engineers and programmers shuffle this Earth bemoaning their bad health. Luria [48] collected the death notices published in *Science* between 1958–1968 and found that the mean age at death of male engineers is 71.1 (N=192) and that of women was 82 (N=1). For both reported genders it was found that archeologists survive engineers: 76.7 for men (N=12) and 84 for women (N=1). In their longitudinal study on life expectancy by occupation, Luy et. al [49] found that in the 1990s the probability of German men in technical occupations such as engineering and maths aged between 40 and 60 of surviving is 89.5% (N=364) while that of German women in the same occupation is 91.7% (N=56). And similarly to the previous study, men in Social service and education’s probability of survival is 90.9% (N=159) and that of women is 93.9% (N=172). While these numbers include software engineers they also include other professions that require their practitioners to wither a lifetime on a desk. Nonetheless these numbers show that professions that require being outdoors, weathering the elements, have their practitioners live longer. For instance mucking about in the mud [84] will grant the mucker 2–5 more years as well as the opportunity to uncover ancient teeth, rusty Victorian trash, Roman garbage, and large feathered reptiles¹.

Statistics reporting on the life expectancy of software engineers are scarce. Nonetheless, Postamate reported that “software engineers have a life expectancy of only 55 years, compared to 78 years for the general population” [27]. The website, whose slogan is “Home of Satire and Sarcasm”,² proceeds to ask why software engineers die so early without delivering a satisfying answer.

¹The reader is recommended the article by Gartley [32]—by which we mean the article written by Gartley and recommended by the author(s), and not recommended by Gartley, although the author(s) find it hard to believe that Gartley would not recommend Gartley’s article for it is a good article—for some dank dinosaur memes.

²A note to the editors: commas and generally any punctuation sign, will go outside quotations and parenthesised sentences.

Yet, because programming, coding,³ and software engineering are sedentary jobs they come with a myriad of related health issues. Chief among the Repetitive Stress Injuries that programmers must deal with is Carpal Tunnel Syndrome: the professionals suffering from CTS are dominantly programmers, system administrators,⁴ and IT professionals [1, 80, 81]. Other musculoskeletal problems include pain and stiffness in the neck among 48.6% of computer professionals, in the lower back (35.6%), and in the shoulders (15.7%) [81].

Sedentary jobs present an uncountable number of other issues: (1) vision blurring (13.2%), (2) irritation in the eyes (18.6%), (3) watering of eyes (23.2%), (4) pain in the eye (25.7%), (5) burning in the eye (29.8%), and (6) headaches (29.2%) [81].⁵

In summary, it is surprising that nerds—software developers—suffer from so many illnesses that can be avoided if they could just go outside and have a walk.

Problem Statement. The main activity of software developers is to develop software [68] through software development languages [68]. With the exception of very few languages discussed in Section 7, these are primarily developed to be written using a so-called full-size keyboard consisting of somewhere between 101 and 105 keys resting on a desk and read—the programs expressed in the programming languages that is and not the keys—on a computer monitor beaming every character onto the reader’s cornea at a generous 120Hz. The author(s) believe that this overly constrained development environment is the direct cause of the software developer’s overly restrained physical posture. We address this problem at its very core by designing a programming language that will not have negative health effects on its user.

Solution. By rethinking how programs are expressed and by deconstructing the syntax used to express these programs the author(s) present a language whose syntax is the path left behind by a technophile walker, hiker, runner, or cyclist⁶. These activities are predominantly done in the outdoors and are often recorded passively through a smartwatch or a smartphone. The programming language, Poololoop, leverages the twists and turns in the recorded path to express programs.

The main benefits of using Poololoop are thus:

- (1) The eyes are free to observe Nature and wildlife, eliminating the need for the 20-20-20 rule⁷,
- (2) The hands and fingers are free to be relaxed,
- (3) The body is exposed to natural sources of vitamins such as the Sun—being an example of a source and not a vitamin—,
- (4) The user’s partner is free to believe that the user took time off work to hang around.

By developing a programming language to solve this problem the author(s), expert(s) in the domain of programming languages, abide by Maslow’s principle: if the only tool you have is a hammer, it is tempting to treat everything as if it were a nail [51].

Paper Structure. In Section 1 we motivate the problem. So if you have not been motivated already then reread Section 1 until you are. In Section 2 we present the necessary background on programming languages and running. In Section 3 we describe the language; its high-level ideas and its formalism. In

There is nothing you can say that will convince me otherwise. The compromise I offer is to typeset the comma, or period, directly under the quotation signs. Something like “this”, or “that”.

³Similarly to Footnote 2, the Oxford comma is another hill I am willing to die on.

⁴The American editors are seething and malding right now.

⁵You know that feeling you get after reading a long list of symptoms? The burning in the eye, the tingling in the fingers, the shooting headache, the slight dizziness... That feeling that these symptoms creep up on you one by one and the conclusion your mind draws is that you must be suffering from those symptoms? The author(s) feel that this is happening to them as they are typing. But let’s not forget that the author(s) are computer professionals and they might be exhibiting actual symptoms.

⁶While the techniques presented here apply to all four activities the paper will only focus on the runners and walkers.

⁷Every 20 lines-of-code take 20 minutes to update 20 dependencies.

Section 4 we provide code examples from the Poololoop standard library. In Section 5 we discuss the implementation of the compiler. In Section 6 we showcase two use cases where Poololoop was used in real-life. In Section 7 we discuss the related work. And in Section ?? we do not conclude in solidarity with McCann [53] who recommends that *SIGBOVIK* bans conclusions.

2 BACKGROUND

In this section we describe the necessary background that we assume the reader is ignorant of.

2.1 Running

Running is not only an action that programs do. Running is an action that many things do, probably too many things, as it has the most number of meanings in the Oxford English Dictionary [89]. The first definition, *l.i.1.a*, in the dictionary states that running, when applied to mammals—which humans are—is the act of moving rapidly on alternating feet—otherwise it’s just hopping—while never having all appendages simultaneously on the ground. The second definition, *V.79.d.i* is the one familiar to most serious programmers. In this paper we adopt both definitions and disambiguate them where needed by explicitly mentioning whether a program is to be ran, or a human—author(s) included—is to be doing the running.

Scientifically, running has been the object of study as old as *Science*⁸. Famously, in 2010, Keller, a scientific-mathematical human runner [50], solved the long standing problem of the *Jogger’s Ponytail* [42]: a phenomenon observed by the running community where a jogger’s ponytail sways from side-to-side while her head bobs up-and-down. And in 2002 the answer to whether one should run or walk in the rain has been proposed by Bailey [7].

Technologically, running has been commoditized under *Sport Business* [71]. Unsurprisingly for readers in the first quarter of the twenty-first century, social networks⁹ for runners exist. The leading platform is Strava [78] where runners connect with other runners. Each runner’s run is traced on a geographic map and their bio-stats plotted on colorful graphs. Runners can give each other *kudos*—a digital signal meant to deliver a rush of serotonin in the receiving runner’s brain—for runs that they have done. They can—optionally—comment—optionally—motivating messages on runner’s runs, and they can tag other runners who ran with them a run.

Digitally, runs and the act of running is encoded in multiple format. The most popular format is the open GPS Exchange Format (GPX) [85] which is an extension of the XML file format [72] that is meant to be human-readable¹⁰. Most electronic tracking devices, colloquially smart devices, that record runs do so in a GPX format. The GPX file format consist of multiple tracks (<trk>) each with its own name. Each track is composed of track segments (<trkseg>) defined by a sequence of points (<trkpt>) defined by their geo-coordinates (lat and lon) attributes and optional fields such as the elevation (<ele>). Below is a small example demonstrating the GPX file format:

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <gpx creator="" xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance" xsi:schemaLocation="http://www.topografix.com
  /GPX/1/1 http://www.topografix.com/GPX/1/1/gpx.xsd" version="1.1" xmlns="http://www.topografix.com/GPX/1/1">
3   <trk>
4     <name>Run, rabbit, run. Dig that hole, forget the sun.</name>
5     <type>running</type>
6     <trkseg>
7       <trkpt lat="51.53707" lon="-0.18343"><ele>37.0</ele></trkpt>
8       <trkpt lat="51.53712" lon="-0.18333"><ele>37.0</ele></trkpt>
9     </trkseg>
```

⁸Not to be confused with the prestigious *Science* scientific journal

⁹Readers in the last quarter of the twentieth century may be more familiar with that social phenomenon where every activity has its own glossy magazine.

¹⁰Where human is not well defined.

10 </tr>
11 </gp>

2.2 λ -calculus

The λ -calculus has been called the “smallest programming language” by some^[whom?]. It was introduced by Alonzo Church in 1932 [14] as a new foundational theory of logic and mathematics. Concurrently, Alan Turing introduced in 1936 the Turing Machine [87]¹¹ to show that it is not possible to resolve Hilbert’s 1928 *Entscheidungsproblem* that asks whether a general algorithm exists to prove any mathematical proposition. While Turing has broken into the mainstream, having his own movie—meaning a movie where he is portrayed rather than written or directed or produced by Turing¹²—published in 2014 starring the movie heartthrob and English sweetheart Benedict Cumberbatch as Turing himself [55], Church saw no such treatment.

Besides the λ -calculus, Church is commonly recognized among the nerd community through the Church-Turing thesis. Luckily for Church his name came first in the thesis’ name because C, which happened to be the first letter of Church, comes earlier in the alphabet than T, which also happened to be the first letter in Turing [25]. Church was also first¹³ in showing that the *Entscheidungsproblem* is not possible to solve [15].

2.2.1 Syntax. The λ -calculus is the mother of all functional languages. Its one and only feature are functions. What do we do with functions? (1) We create them, (2) we give their arguments name, and (3) we apply them to something, i.e. we replace names with other things.

So all *expressions* in the λ -calculus have one of the following three shapes.

- (1) *Abstraction*: Take an expression, find your favorite sub-expression, keep it in your pocket and replace it with a variable, then wrap the whole expression in a syntactic form that says that variable should stand for something.
- (2) *Variable*: That thing you replace your favorite sub-expression with,
- (3) *Application*: The way you say that a variable stands for your favorite sub-expression.

Formally though, expressions are denoted with the symbol e and they are described by a *grammar*, which is just the type of the syntax tree. Anyways, the λ -calculus expressions are described by this self-explanatory grammar:

$$e ::= \lambda x.(e) \mid x \mid (e_1 e_2)$$

In practice, programmers are familiar with the λ -calculus if they have written in a Lisp-like dialect, Haskell, Scala, F#, or some other functional language. Below we address three common criticism against the λ -calculus.

ThErE aRe ToO mAnY pArEnThEsEs. No. You’re used to calling functions like this, $f(x, g(y))$, right? The λ -calculus not only eliminates those useless commas, it gets rid of one pair of parenthesis! It simply moves the parentheses one token to the left to become $(f x (g y))$ and removes the outermost ones to become $f x (g y)$. But wait, there’s more. To eliminate more of these pesky parentheses we adopt two conventions: λ associate to the right and applications to the left. What does this mean? It means that $\lambda x.(\lambda y.e)$ becomes $\lambda x.\lambda y.e$ and $(e_1 e_2) e_3$ becomes $e_1 e_2 e_3$. For example: $\lambda x.\lambda y.x x y y$ stands for $\lambda x.(\lambda y.((x x) y) y)$. Neat, no?

¹¹Turing did not call the Turing Machine the Turing Machine because Turing was humble, it was Church who dubbed Turing’s machine: Turing Machine.

¹²The author(s) are not aware of any movie that Turing wrote, directed, or produced

¹³The author(s) choose to conveniently sweep the Gödel debacle under the footnote line. Kurt Gödel, the archetypical Austrian—archetypical not being a qualifier for Austrian—nerd proved that the *Entscheidungsproblem* cannot be resolved in 1931 [34]. Everyone was aware of Gödel’s result and it is acknowledged by Church and Turing.

I cAn'T pRoGrAm WiThOuT nUmBeRs, ClAsSeS, aNd ObJecTs. You don't need them bro. You just need these three syntactic forms to write any program. Booleans? They're functions. Numbers? They're functions. Classes and objects? They're the poor man's closure [16]. Closures? They're just functions. In Section 4 we'll show you how you can do it too.

WhErE ArE ThE mUlTiVaRiAtE fUnCtIoNs? They don't exist. You know why? Cause they're useless. Just Schönfinkel them [67, 70] my dude. A function that takes two things is a function that takes the first thing and returns a function that expects the second thing. This fact can be recalled with this mnemonic rhyme:

Many arguments are mere ornaments.
 With higher-order languages, ordure
 as these can be teased out with such an ease:
 Make λ s take one, which when asked to run,
 returns another expecting the other,
 until no more can arguments dwindle.
 Congrats. Now you know how to Schönfinkel.

Or you could bundle your arguments in a pair. But guess what? Pairs are also just functions.

2.2.2 Semantics: What does it all mean? So now you know how to write λ programs. How do you run them? You use the relation $e \rightarrow e'$ which can be read as “ e becomes e' after a single step”. This style is conveniently called small-step semantics because it invites the evaluator to behave like a CPU and to operate in discrete units commonly known as clock-cycles.

The \rightarrow relation is defined through a single rule: the β rule. That's the only¹⁴ rule¹⁵ you will ever need. Do you see now why it's the “smallest programming language”? The β -rule looks like this:

$$(\lambda x.e_1) e_2 \rightarrow e_1|_{e_2}^x \quad (\beta)$$

Where the notation $e_1|_{e_2}^x$ is the capture-avoiding substitution [77]. We could define it, or we could do what every other programming language does and show you examples.

$$\begin{aligned}
 (\text{add } x \text{ } d)|_{\text{add } a \text{ } d}^x &= \text{add } (\text{add } a \text{ } d) \text{ } d \\
 (\lambda y.\text{sub } x \text{ } b)|_{\text{sub } u \text{ } b}^x &= \lambda y.\text{sub } (\text{sub } u \text{ } b) \text{ } b \\
 (\lambda x.\text{foo } x)|_x^x &= \lambda x.\text{foo } x
 \end{aligned}$$

This again, is all you need to do any computation. In other words, this model: the three syntactic forms and the β rule is all you need to write Linux and Clang which you can then use to program a simulator of a Turing Machine. However unlike Linux, the λ -calculus does not require files, and unlike Clang it's not developed by a group, and unlike Turing Machines there are no tapes. What's common across them? The state. Church, by inventing the λ -calculus, stated computationally and constructively the separation between Church and state¹⁶.

¹⁴That's only one if you consider—which the author(s) do—the weak-head normal form evaluation strategy. You can ignore these words, it's beyond the scope of this paper.

¹⁵That's not really true. The first giveaway is the rule's name. You would think that if there's a β , then there's an α . This is also not entirely true as there is also an η . The α rule just renames some variables, which you might need to do in some cases when free—as in freedom, not as in free beer—variables are involved. The η rule is at best a compiler optimization.

¹⁶This joke is paraphrased from Guy Steele: “And some people prefer not to commingle the functional, lambda-calculus part of a language with the parts that do side effects. It seems they believe in the separation of Church and state. :-) :-) :-)” [76]

Variables Begone! In 1972, the iconic Dutch composer Louis Andriessen composed de Volharding [4, 8]¹⁷ which lead to the creation of the eponymous minimal jazz group de Volharding [24]. This event signaled the climax of the minimalism movement in Europe [9]. In 1972, the Dutch computer scientist Edsger W. Dijkstra delivered EWD340 as his Turing Award acceptance speech [26] in which he announced that he was the first Dutch to register as programmer. In 1972, the Dutch mathematician Nicolaas Govert de Bruijn published *Lambda Calculus Notation with Nameless Dummies* [21] which combined both minimalism, programming, and signaled the end of logical minimalism and computer design [23]. With his seminal paper, the visionary de Bruijn foresaw and solved half of the problem now attributed to Phil Karlton: “There are only two hard things in Computer Science: cache invalidation and naming things.” [22, 30, 40]. By completely removing the need for variables in the λ -calculus de Bruijn resolved the latter problem [3] back in the 70s.

How did he do it? de Bruijn observed that in programs without free—as in freedom—variables variables are just “pointers” to wherever they were bound or declared. For example, $\lambda x.\lambda y.x$ can be said to be the lambda that returns a lambda which returns the variable bound by the lambda *two* levels before. Coincidentally, $\lambda a.\lambda b.a$ is the same lambda that returns a lambda which returns the variable bound by the lambda two levels before. So in de Bruijn notation we express them as $\lambda\lambda 2$. This notation is dubbed de Bruijn indices¹⁸ and is the one we will use throughout the paper.

To get the reader used to this notation, we will list some lambda expressions and their de Bruijn indices in the following. The identity function $\lambda x.x$ becomes $\lambda.1$. The constant function $\lambda x.\lambda y.x$ becomes $\lambda\lambda.2$. The application function $\lambda f.\lambda x.f x$ becomes $\lambda\lambda.2 1$. This function $\lambda x.(\lambda f.x) x$ becomes $\lambda.(\lambda.2) 1$.

3 THE DESIGN OF POOLOLOOP

Poololoop, pronounced [pʊ:lɔ:lʊ:p], is a portmanteau [88] of “Pool of loop”.

The main construct that Poololoop exploits for its syntax is the *loop*, or that which is colloquially called a *circle* by the general population who failed their basic geometry class and forgot that a circle has a constant radius. This choice is motivated by three reasons.

First, humans naturally walk around in circles [74], thus the user does not need to concentrate on the syntax and can rather spend their energy on the problem at hand. A feature that hardly any modern-day programming language enjoys.

Second, The programs naturally become small in diameter. The user does not need to stray far away from home and venture into foreign environments. This has the positive outcome of avoiding taking any unneeded risk that may trigger separation anxiety in users.

Third, loops make Poololoop future proof. Mastroianni et. al [52] and others [47] proved that the past was better. It thus follows that the function describing the quality of time is a monotonously decreasing function. The trivial corollary states that the future will be worse. Of the multiple proposed models describing the future, two dominate [44]: the Orwellian [59] and the Huxlerian [37]. In what follows we argue that Poololoop fits snugly in both models.

The Orwellian model predicts a general increase in user monitoring and language moderation by lifting the principle of least privilege [69] from software development into efficient social organization. In this realm, a successful language must cater not only to its users but also to Big Brother. Proponents of the Orwellian model justify it by offering the following argument. In 2023, 31.17% (N=59,¹⁹336) of software developers reported working for companies with a workforce of larger than 500 employees [75]. The law

¹⁷Dutch for perseverance. The author(s) recommend that the reader play this composition and read the remainder of the paper while listening to it.

¹⁸Not to be confused with de Bruijn levels

¹⁹For the editors from continental western Europe: this comma is a thousands separator and not a decimal separator. I’m sorry that your language is not so relevant scientifically anymore.

of large numbers—and one may not need to invoke this law to make a strong argument—implies that many programmers will have to work under surveillance-friendly conditions [2, 10, 38]. In this setting Big Brother subordinates shall not employ Poololoop as prescribed as that may give the wrong impression that the now-healthy workforce, coming back from a hard day’s work in Nature²⁰ [41, 57, 58, 61] using Poololoop, showing signs of happiness and good mental health may be confused with an idle workforce loitering around the coffee machine all day. Thus a more consistent employment of Poololoop in the Orwellian model is one where the programmers are lead in file to the underground parking lots, away from sunlight and under the infrared glow of night-vision-enabled surveillance cameras, to write their programs in the oppressing underground stale air. In such a space, without exploiting loops, the space of programs that can be expressed becomes too limited to be useful.

The Huxlerian model predicts an intoxicating increase in developer tooling with excellent user experience and a deluge of mind-altering technologies²¹ that will boost the productivity of programmers. In that model, for programmers to use a tool, it must be an addictive one. Utilizing Poololoop in the great outdoor will help greatly in improving the mental state of the programmer. Thus, an excess use of Poololoop has the benefit of increasing the user engagement. Moreover, it has been shown that runners and walkers in Nature do feel a sense of connection with natural entities and a sense of yearning to revisit them [12, 29, 65]. Hence, a use of Poololoop naturally leads to an excessive use which leads to unprecedented levels of programmer satisfaction, which according to the Huxlerian model, must lead to an increase in output and productivity.

3.1 Syntax

Poololoop is an alternative front-end to the λ -calculus with de Bruijn indeces as presented in Section 2.2. To that extend, we present the front-end to the three syntactic forms of the untyped λ -calculus: variables, functions, and function applications.

We start by distinguishing between the language Poololoop and the formal mathematical system ϖ underlying the language. The equivalent of the expressions e of the λ -calculus are the loops \dagger of ϖ . To define the language we define in Section 3.1 the loop-encoding of e , $\langle e \rangle : \lambda \rightarrow \varpi$, which translates expressions into loops. The actual definition is by cases on the λ -calculus syntactic forms and is split in Definitions 3.1 to 3.3. And in Section 3.2 we define the untangling of \dagger , $\llbracket \dagger \rrbracket : \varpi \rightarrow \lambda$, particularly in Definition 3.4, which translates loops back into expressions.

3.1.1 Variables. When we use de Bruijn indeces, variables will always be represented as natural numbers²². It is worth noting that we do not consider zero to be a natural number. Therefore quite naturally we denote a variable n by n consecutive empty loops.

Definition 3.1 (Variables). A variable n is loop-encoded as the loop wrapping n empty loops as follows:

$$\langle n \rangle = \underbrace{\text{Diagram of } n \text{ empty loops wrapped by a larger loop}}_{\text{Repeated } n \text{ times}}$$

²⁰Not to be confused with the prestigious *Nature* scientific journal.

²¹Exemplified by conversing with computers that pass the Turing test.

²²In 1889 Giuseppe Peano published *Arithmetices Principa* [62] in which he presents the defacto agreed-upon axiomatization of the natural numbers. The first axiom is verbatim $1 \in N$ (sic). Implying that zero is not a natural number. However, in his 1901 *Formulario mathematico* [63] he realizes his mistake and includes zero as a natural number and the first axiom of his formalism

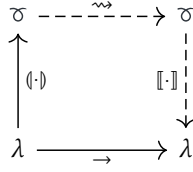


Fig. 1. The diagram describing the σ semantics with respect to λ semantics

For the readers who are not versed in diagram notations, we simply mean the following: given the encoding function that “decompiles” and the underlying runtime that makes one step then it is possible to define (uniquely) the compiler and the interpreter (higher-level runtime) such that $e \rightarrow \llbracket (e) \rrbracket \rightsquigarrow$.

In [Definition 3.4](#) we define the compilation function of a loop \dagger into a λ -calculus expression.

Definition 3.4 (Compilation). We define the compilation of a loop \dagger in σ by cases:

$$\begin{aligned}
 \left[\begin{array}{c} \text{Diagram of } n \text{ loops} \\ \text{Repeated } n \text{ times} \end{array} \right] &= n \\
 \left[\begin{array}{c} \text{Diagram of loop } \dagger \\ \dagger \end{array} \right] &= \lambda. \llbracket \dagger \rrbracket \\
 \left[\begin{array}{c} \text{Diagram of sequence of loops } \dagger_n \dots \dagger_1 \\ \dagger_n \dots \dagger_1 \end{array} \right] &= \llbracket \dagger_1 \rrbracket \dots \llbracket \dagger_n \rrbracket
 \end{aligned}$$

With [Definitions 3.1](#) to [3.4](#) we can formulate [Theorem 3.5](#) which expresses the expected fact that untangling and loop-encoding are inverse operations.

THEOREM 3.5. *Loop-encoding and untangling are inverse operations. In other words, for every loop \dagger in σ then $\dagger = \llbracket \llbracket \dagger \rrbracket \rrbracket$, and for every λ -calculus expression e then $e = \llbracket (e) \rrbracket$.*

PROOF. Like most theorems in the domain of programming languages [5, 13, 45], the proof is a trivial application of structural induction. \square

Defining the interpreter \rightsquigarrow is also not difficult but requires the author(s) to typeset many complicated diagrams.

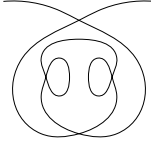
Definition 3.6 (The \rightsquigarrow interpreter). The definition is left as an exercise to the reader.

4 PROGRAMMING IN POOLOLOOP

4.1 Church Booleans

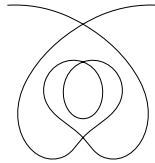
To construct a boolean one of the two boolean constructors must be used: true or false. Thus the church encodings of booleans will always have two outer lambdas, one for each constructor.

True. The boolean true value is encoded as $\lambda\lambda.2$. Informally, the first lambda asks its user for what is meant by true, and the second asks for what is meant by false. The expression then returns the true value provided by the user. Its loop-encoding in ϖ is the following:



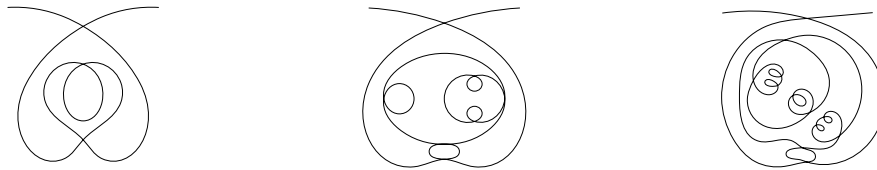
All loop examples are to be read from the top-left corner, traveling along the entire loop continuously without breaking smoothness, until the top-right corner is reached. The inhabitants of the southern-hemisphere are instead expected to read the diagram from the top-right corner towards the top-left corner in a similar fashion.

False. Dually, the boolean false will return the false that the user provided, i.e. false is $\lambda\lambda.1$. Its loop-encoding is the following:



4.2 Church Numerals

4.2.1 *Numbers.* The Peano encoding of the natural numbers [63] assume two constructors: the zero²⁴ and the successor function. The zero is then encoded as $\lambda\lambda.1$ ²⁵, one is encoded as $\lambda\lambda.2\ 1$, two as $\lambda\lambda.2\ (2\ 1)$, three as $\lambda\lambda.2\ (2\ (2\ 1))$, etc. In ϖ we illustrate 0, 1, and 2 as the following diagrams in order from left-to-right:



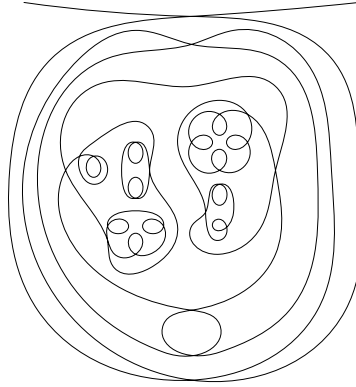
4.2.2 *Addition.* Addition is a fundamental operation [28]. The operation takes two Church numbers and produces a Church number. Therefore it has two outer-most lambdas for the given numbers, and two other lambdas for the zero and the successor function. A number n is encoded as the application of the given successor function n times to the given zero. Thus, addition of n and m is encoded as the application of the given successor function n times to m such that the given zero and successor functions are passed along to m . In other words, the de Bruijn encoding of addition is

$$\text{add} = \lambda\lambda\lambda\lambda.4.2\ (3\ 2\ 1) \tag{1}$$

Whose loop encoding is the following diagram:

²⁴Here zero is a natural number. This is not to be confused with the definition of the natural numbers at the meta-level.

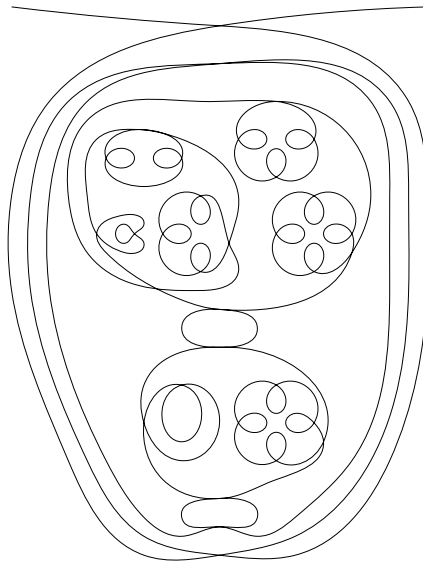
²⁵The attentive reader would have noticed that the encoding of zero and false are the same. This feature has been added to Church encodings in order to attract C and Javascript programmers into functional programming.



4.2.3 *Multiplication.* Just as we re-interpreted the meaning of zero in the addition function to be the second number to be added, for addition we re-interpret the successor function to addition. Therefore multiplying n and m becomes n additions of m on a given zero. Thus, by inlining Equation (1), we obtain the following definition of multiplication

$$\text{mult} = \lambda\lambda\lambda\lambda.4(\lambda.43(132))1 \tag{2}$$

Its loop encoding is the following diagram:

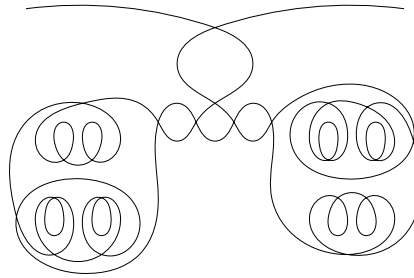


4.3 Fixed Point Operators

The Y-combinator allows recursive and diverging programs to be expressed. It is defined as follows:

$$Y = \lambda.(\lambda.2(11))(\lambda.2(11)) \tag{3}$$

and its loop encoding is the following:



4.3.1 *Recursive Programs: factorial.* Now we demonstrate that Poololoop is not a toy programming language by implementing the factorial function. We implement it using the Y-combinator defined in Equation (3) and the church numerals. Below is the program in the λ -calculus without de Bruijn indices:

```

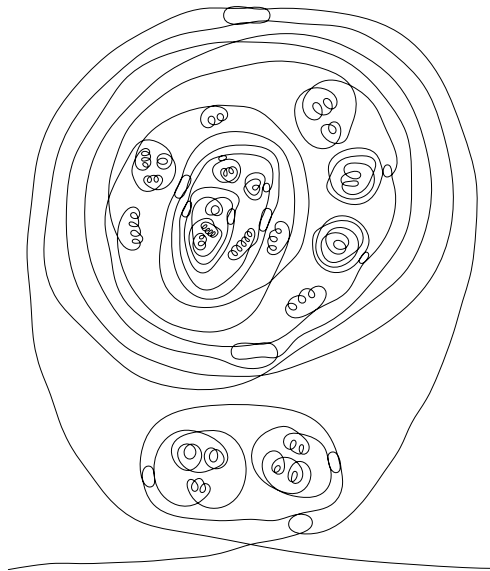
1 fact = (\f. (\x. f (x x)) (\x. f (x x)))           -- Y-Combinator Equation (3)
2   \recurse.\n.\s.\z.                               -- First argument is the recursive call
3     n (\_.\T.\F. F) (\T.\F. T)                   -- is true if n is zero and false otherwise
4     (s z)                                          -- 0! = 1
5     (recurse                                       -- recursive case
6       (n \s.\z. n (\g.\h. h (g s)) (\u. z) (\u. u)) -- MAGIC: predecessor of n
7       s (n s z))                                  -- the last argument recalls Equation (1)

```

Using de Bruijn indices, fact becomes:

$$(\lambda.(\lambda.2 (11)) (\lambda.2 (11))) (\lambda\lambda\lambda.3 (\lambda\lambda.1) (\lambda\lambda.2) (21) (4 (3 \lambda\lambda.5 (\lambda\lambda.1 (24)) (\lambda.2) (\lambda.1)) 2 (3 2 1)))$$

And its loop encoding is:



This particular diagram is meant to be read from the bottom-left strand, all the way throughout the path, until the end of the bottom-right strand. Godspeed.

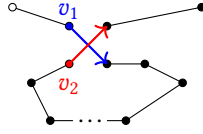


Fig. 2. A loop is defined by the intersection point of two vectors v_1 (blue) and v_2 (red) and its—the loop’s—direction is defined by the same two vectors. The starting point is the empty point in the top-left corner.

5 IMPLEMENTATION

We implemented a compiler for Poololoop in the C programming language²⁶. The source code, which is relatively tiny at 147 lines long is provided in [Appendix A](#)²⁷. The compiler takes two command line arguments, and an optional third. The first argument is the path to a GPX file and the second is an identifier that the compiled code should be assigned to. The optional third argument specifies the target language to compile to. By default the target language is the Scheme programming language. The other alternative language is Haskell. The two targets are specified with `scm` and `hs` respectively.

The semantics implemented by the compiler is as described in [Section 3](#). The implementation decides on the directionality based on whether the starting point has positive or negative latitude, i.e. is in the northern or southern hemisphere respectively.

The implementation abides by Postel’s Robustness law [66], it accepts file formats that supersede GPX. Informally, Poololoop’s compiler accepts any file that contains a sequence of latitude and longitude coordinates specified respectively with `lat="ieee_float"` and `lon="ieee_float"`. The syntactic class `ieee_float` is the class of IEEE floats.

Compiling the compiler is as easy as passing it to `gcc` and linking it with the math library using the `-lm` flag. When the compiler is compiled with the `-O3` optimization flag, then compiling the two case studies presented in [Section 6](#) took less than 3 milliseconds on the author(s)’ machine which is just an everyday laptop that one takes on a holiday. The GPX file of the two case studies contains 509 and 519 GPS points respectively. [Section 6](#) provides the executed commands and their output.

The main observation is that every loop is defined by an intersection point. Thus, find the intersection point and you will find the loop. The direction of the loop is computed from two vectors: (1) the vector whose tail and tip are defined by the point just before the intersection point and the one just after, respectively, and (2) the vector whose tail and tip are defined by the point just before the intersection point and the one just after, respectively. [Figure 2](#) shows an example of these two vectors, v_1 and v_2 , respectively.

The compiler exploits the key idea that intersection points act as parentheses. Every intersection point is traversed twice by the runner, on entering the loop and on exiting it. Thus, the compiler finds the intersection points, sorts them by the time traversed, and treats each pair as a parenthesis: the first is equivalent to an open parenthesis and the second is equivalent to a closed one. Once this Intermediate Representation (IR) is generated, then compiling to a de Bruijn indexed λ -calculus IR is equally trivial, and compiling to Scheme or Haskell is just a boring task at this point.

The compiler uses the naïve $O(n^2)$ algorithm to find intersections in a piece-wise linear path: it checks segments two-by-two for intersections. When tried on a GPX file with 13,714 GPS points the compiler took 291 milliseconds to produce a syntax error²⁸.

²⁶C11, probably.

²⁷And also as a Gitlab Snippet: <https://gitlab.com/-/snippets/3688034>.

²⁸Which is only raised after finding all intersections.

6 EVALUATION

In this section we show that using Poololoop in real-life is possible. We have chosen two programs from Section 4 and we ran one program on a large-scale and walked the other on a small-scale. We report on these two case studies in the following paragraphs.

The first author went on run along the *true* path on the night of Saturday 9th of March 2024. The GPX track recorded by the author’s smart watch is shown in Figure 3a. The running distance was 4.48 kilometers long²⁹ at the respectable pace of 5 minutes and 35 seconds per kilometer³⁰ and the running time was around 25 minutes and 5 seconds, which is almost exactly the same length as de Volharding³¹. The engine running the program was indeed exposed to that composition as they were running the program. Therefore one thread of validity to this case study is the choice of music as it might have affected the pace of the runner and thus the measured run time [11].

The last author’s *fix* was also done on the night of Saturday 9th March 2024. The path is 2.66 kilometers long. The author took 46 minutes and 42 seconds to walk it at the shameful pace of 17 minutes and 33 seconds per kilometer. The GPX track recorded by the author’s smart watch is shown in Figure 3b. The slow pace can be explained by the author’s report that conducting this experiment was tedious and awkward. They have in fact tried to walk the path in daylight but gave up soon after the first loop as they have reported feeling an uncomfortable level of awareness exhibiting itself through emotions of self-consciousness and fear of walking into small and developing humans engaging in a communal and recreational activity consisting of striking, repeatedly, with a single foot, a round orb thus carrying it off the ground and potentially into the faces of others—author included. While the author was not disturbed during the experiment they have received inquisitive looks—even under the cloak of darkness—from strangers. Thus one thread of validity to this case study are other humans.

Both case studies have been recorded and archived on Strava [31, 78], a popular website where athletes preserve their artifacts. The tracks are publically accessible³² alongside multiple plots and a GPX file which can be fed into the Poololoop compiler.

Below are the benchmarks of the four compilations that the author(s) conducted on the GPX files of both case studies as produced by the *hyperfine* benchmark utility [64].

```
1 > hyperfine "./poololoop gpx/fix.gpx f hs" ".poololoop gpx/fix.gpx f scm" \  
2     "./poololoop gpx/true.gpx t hs" ".poololoop gpx/true.gpx t scm"  
3  
4 Benchmark 1: ./poololoop gpx/fix.gpx f hs  
5   Time (mean ± σ):      2.0 ms ± 0.7 ms   [User: 1.8 ms, System: 0.4 ms]  
6  
7 Benchmark 2: ./poololoop gpx/fix.gpx f scm  
8   Time (mean ± σ):      1.5 ms ± 1.0 ms   [User: 1.5 ms, System: 0.2 ms]  
9  
10 Benchmark 3: ./poololoop gpx/true.gpx t hs  
11  Time (mean ± σ):      1.7 ms ± 0.8 ms   [User: 1.6 ms, System: 0.3 ms]  
12  
13 Benchmark 4: ./poololoop gpx/true.gpx t scm  
14  Time (mean ± σ):      1.7 ms ± 0.8 ms   [User: 1.6 ms, System: 0.3 ms]
```

The output of each command is the following:

²⁹Or 2.78 miles in freedom units.

³⁰Or 8 minutes 57 seconds per mile in freedom units.

³¹Which should be just about done if you started it per Footnote 17.

³² $\lambda\lambda 2$ - SIGBOVIK'24 CS#2: <https://www.strava.com/activities/10926312841>, and $\lambda(\lambda 2(1\ 1))(\lambda 2(1\ 1))$ - SIGBOVIK'24 CS#1: <https://www.strava.com/activities/10925302747>. By appending `/export_gpx` to the URL, the artifact evaluators can download the GPX file.



(a) The *true* program as described in Section 4.1



(b) The Y-combinator as described in Section 4.3

Fig. 3. A satellite image of the area where the two programs of the two case studies were ran and walked with the path superposed on top in a red line. The starting point of the run and walk is indicated by a green dot. The ending point of the run and walk is indicated by a checkered flag in a white circle.

```

1 > ./poololoop gpx/fix.gpx f hs
2 f = (\ x0 -> ((\ x1 -> (x0 (x1 x1))) (\ x1 -> (x0 (x1 x1)))))
3
4 > ./poololoop gpx/fix.gpx f scm
5 (define f (lambda (x0) ((lambda (x1) (x0 (x1 x1))) (lambda (x1) (x0 (x1 x1)))))
6
7 > ./poololoop gpx/true.gpx t hs
8 t = (\ x0 -> (\ x1 -> x0))
9
10 > ./poololoop gpx/true.gpx t scm
11 (define t (lambda (x0) (lambda (x1) x0)))

```

We would like to remind the artifact evaluators that the generated Haskell code of the Y-combinator cannot be used even though it is syntactically correct. That is because Haskell is strongly typed and the Y-Combinator cannot be typed in the Simply Typed λ -calculus, nor in System-F. That is no fault of ours. The generated Scheme code on the other hand is perfectly fine.

7 RELATED WORKS

As with most novel work such as ours, not much truly related related work exists. Two line of works nonetheless can be identified. The author(s) hope that by the end of this section the reader would have realized that none of these completely satisfy Poololoop’s design choices from Section 3.

7.1 Alternative or Assistive Hardware

Augmented Reality (AR) devices which intend to superpose virtual objects on top of real-world objectives have been around commercially for more than a decade. These could be used to assist programmers in writing their programs in the great outdoors just as Poololoop does.

The first device that broke into the mainstream is the Google Glass in the early 2010s. It is a wearable device, just a pair of eyeglasses, with two small transparent glass rectangles covering a part of the user’s field of vision. By projecting pixels onto these transparent rectangles the user can switch their focus from real-life into the virtual, and vice versa. But damn do they look dorky [90].

More recent devices, such as Facebook’s Quest and Apple’s VisionPro offer better image quality and hand-gesture detection. Yet these have two downsides. First, the real-world is seen through a screen looping back the view from a front-camera. Meaning these devices are virtual reality devices that happen to mirror the real-world—for now³³. And second, they sure do look more obnoxious than the Google Glass.

7.2 Gestural Programming Languages

Gestural Programming is the domain of computer vision and artificial intelligence research where one teaches a robot how to accomplish tasks by demonstrating to the robot, visually, through the means of a human, how they are accomplished. This line of work has been explored by Soratana et al. [73] and Cabrera et al. [54]. But it’s easy to conclude that this is unrelated to Poololoop and is not what we mean by “Gestural Programming”.

Let’s try again. Gestural Programming is the domain of programming pedagogy where researchers explore the use of input devices other than the keyboard and the mouse for programming. For example in Streeter’s PhD thesis [79] the author³⁴ applied multiple gesture matching algorithms to data recorded from students programming in Google Blockly with the Microsoft Xbox Kinect. Similarly Toro-Guajardo et al. [86] reported on young people programming in Scratch using the Nintendo Switch Joy-Cons and found that they have more fun if they, and their hands, move. While this line of work is closer to Poololoop than the previous one, it is still unrelated.

One more time. Gestural Programming is the domain of programming language research which produced bodyfuck [35, 36, 82], a language in which programmers input programs through moving their bodies. Bodyfuck is an alternative front-end for the popular programming language brainfuck [56] where the eight brainfuck actions are mapped into eight bodily gestures the programmer performs facing a visual recording apparatus. For example, if one wishes to increment the register in focus then one must jump, and if one wishes to decrement it then one must duck. Bodyfuck has been birthed to be performative art. Fifteen programs were performed and put on display in the Things That Are Possible MFA Show [6]. It aims to separate the act of software performance, i.e. software inscription³⁵ from the computational context in which it happens in. Bodyfuck aims to demonstrate that software inscription can be done completely outside computers³⁶ in the surrounding cultural space. Similarly, Poololoop demonstrates that this inscription can be done completely outside.

This shows that we’re on the right track of identifying actually related related work. Sadly, the author(s) could not find any other related work in that style. However, the early readers of the *SIGBOVIK* publication proceedings may recall the work of Leffert entitled “Harnessing Human Computation: β -reduction hero” [46] in the 2010 Technical Report track. Alas, the author(s) were unable to find the Flash application which was reported on nor the report nor its source code in order to learn what the work is about. Nonetheless we conjecture, based on the name, that it presents a λ -calculus evaluation technique that is gamified á-la Guitar Hero. Nevertheless if that is the case, then unlike Poololoop, programming can only be done indoors while staring at a screen.

³³Seriously, imagine the horrors of having that feedback camera attacked.

³⁴Being the author of Streeter’s PhD thesis and not the author(s) of this very paper you are almost done reading.

³⁵As the author(s) of this almost-finished paper have no background in academic art they have found the essay cited earlier to be extremely difficult to read, and they report here their best guess at what it could mean.

³⁶Also around the year 2010, the first author of this almost-done paper that you are reading recalls writing a whole PHP program on a piece of paper during a biology class in their senior highschool year as they had no interest in biology and no access to a computer. This footnote tells the anecdote to show that this separation is natural and to brag that the first author was able to write a whole PHP application by hand on a piece of paper.

REFERENCES

- [1] K. Mohamed Ali and B.W.C. Sathiyasekaran. 2006. Computer Professionals and Carpal Tunnel Syndrome (CTS). *International Journal of Occupational Safety and Ergonomics* 12, 3 (2006), 319–325. <https://doi.org/10.1080/10803548.2006.11076691> PMID: 16984790.
- [2] Seth Allcorn. 2022. Micromanagement in the workplace. *Organisational and Social Dynamics* 22, 1 (2022), 83–98.
- [3] Reem Alsuhaibani, Christian Newman, Michael Decker, Michael Collard, and Jonathan Maletic. 2021. On the Naming of Methods: A Survey of Professional Developers. In *2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE)*. 587–599. <https://doi.org/10.1109/ICSE43902.2021.00061>
- [4] Louis Andriessen. 1972. de Volharding. Composition.
- [5] Carlo Angiuli. 2017. The Next 700 Type Systems. *SIGBOVIK* (2017), 169–171. <https://sigbovik.org/2017/proceedings.pdf>
- [6] The Digital Arts and New Media MFA Program at UC Santa Cruz. 2010. *2010 MFA Exhibition :: Things That Are Possible*. https://danm.ucsc.edu/news_events/2010-mfa-exhibition Accessed: March 3, 2024.
- [7] Herb Bailey. 2002. On running in the rain. *The College Mathematics Journal* 33, 2 (2002), 88–92.
- [8] World Association For Symphonic Bands and Ensembles. 2021. *DE VORHALDING for Piano and Winds (1972) by Louis Andriessen (The Netherlands, 1939-2021)*. <https://wasbe.org/de-vorhalding-for-piano-and-winds-1972-by-louis-andriessen-the-netherlands-1939-2021> Accessed: March 3 2024.
- [9] Maarten Beirens. 2016. European Minimalism and the Modernist Problem. In *The Ashgate Research Companion to Minimalist and Postminimalist Music*. Routledge, 61–85.
- [10] Clive R Boddy. 2017. Psychopathic leadership a case study of a corporate psychopath CEO. *Journal of Business Ethics* 145, 1 (2017), 141–156.
- [11] Robert Jan Bood, Marijn Nijssen, John Van Der Kamp, and Melvyn Roerdink. 2013. The power of auditory-motor synchronization in sports: enhancing running performance by coupling cadence with the right beats. *PLoS one* 8, 8 (2013), e70758.
- [12] Stefan Brené, Astrid Bjørnebekk, Elin Åberg, Aleksander A Mathé, Lars Olson, and Martin Werme. 2007. Running is rewarding and antidepressive. *Physiology & behavior* 92, 1-2 (2007), 136–140.
- [13] Robert Chatley, Alastair Donaldson, and Alan Mycroft. 2019. *The Next 7000 Programming Languages*. Springer International Publishing, Cham, 250–282. https://doi.org/10.1007/978-3-319-91908-9_15
- [14] Alonzo Church. 1932. A Set of Postulates for the Foundation of Logic. *Annals of Mathematics* 33, 2 (1932), 346–366. <http://www.jstor.org/stable/1968337>
- [15] Alonzo Church. 1936. A note on the Entscheidungsproblem. *The journal of symbolic logic* 1, 1 (1936), 40–41.
- [16] C2 Wiki Community. 2022. *Closures and Objects Are Equivalent*. <http://wiki.c2.com/?ClosuresAndObjectsAreEquivalent> Accessed: March 3, 2024.
- [17] Gaspard Gustave Coriolis. 1835. *Mémoire sur les équations du mouvement relatif des systèmes de corps*. Bachelier.
- [18] Rebecca M Currano, Martin Steinert, Larry J Leifer, et al. 2011. Characterizing reflective practice in design—what about those ideas you get in the shower?. In *DS 68-7: Proceedings of the 18th International Conference on Engineering Design (ICED 11), Impacting Society through Engineering Design, Vol. 7: Human Behaviour in Design, Lyngby/Copenhagen, Denmark, 15.-19.08. 2011*. 374–383.
- [19] Evans Data. 2023. *Worldwide Developer Population from 2016 to 2023*. <https://www.statista.com/statistics/627312/worldwide-developer-population/> Accessed: March 3 2024.
- [20] DataReportal, Meltwater, and We Are Social. 2024. *Internet and Social Media Users in the World 2024*. <https://www.statista.com/statistics/617136/digital-population-worldwide/> Accessed: March 3 2024.
- [21] Nicolaas Govert De Bruijn. 1972. Lambda calculus notation with nameless dummies, a tool for automatic formula manipulation, with application to the Church-Rosser theorem. In *Indagationes mathematicae (proceedings)*, Vol. 75. Elsevier, 381–392.
- [22] Vitor De Mario, Golden Cuy, David Karlton, and Murven. 2014. Has Phil Karlton ever said "There are only two hard things in Computer Science: cache invalidation and naming things"? <https://skeptics.stackexchange.com/questions/19836/has-phil-karlton-ever-said-there-are-only-two-hard-things-in-computer-science>. Accessed March 3 2024.
- [23] Liesbeth De Mol, Bullynck Maarten, and Edgar G Daylight. 2018. Less is more in the fifties: Encounters between logical minimalism and computer design during the 1950s. *IEEE Annals of the History of Computing* 40, 1 (2018), 19–45.
- [24] Mark Delaere, Maarten Beirens, and Hilary Staples. 2004. Minimal music in the Low countries. *Tijdschrift van de Koninklijke Vereniging voor Nederlandse Muziekgeschiedenis* 1 (2004), 31–78.
- [25] Oxford English Dictionary. 2022. *Dictionary*. Oxford University Press.
- [26] Edsger W Dijkstra. 1972. The humble programmer. *Commun. ACM* 15, 10 (1972), 859–866.
- [27] Postamate Editor. 2023. *Why Software Engineers Have Short Life Expectancy*. <https://postamate.com/2023/08/why-software-engineers-have-short-life-expectancy/>
- [28] Mimi Engel, Amy Claessens, and Maida A Finch. 2013. Teaching students what they already know? The (mis) alignment between mathematics instructional content and student knowledge in kindergarten. *Educational Evaluation and Policy Analysis* 35, 2 (2013), 157–178.
- [29] S Forster et al. 2009. The forest for leisure activities and tourism: a yearning for nature or sustainable development?(essay). *Schweizerische Zeitschrift für Forstwesen* 160, 7 (2009), 189–194.
- [30] Martin Fowler. 2009. Two Hard Things. <https://martinfowler.com/bliki/TwoHardThings.html>. Accessed March 3 2024.
- [31] Rob Franken, Hidde Bekhuis, and Jochem Tolsma. 2023. Kudos make you run! How runners influence each other on the online social network Strava. *Social Networks* 72 (2023), 151–164.
- [32] Luke-Elizabeth Gartley. 2022. CLADISTICS ruined my life: intersections of fandom, internet memes, and public engagement with science. *Journal of Science Communication* 21, 5 (2022), Y01. <https://doi.org/10.22323/2.21050401>
- [33] Gartner. 2021. *Global Shipments of Personal Computers from 2006 to 2021*. <https://www.statista.com/statistics/273495/global-shipments-of-personal-computers-since-2006/> Accessed: March 3 2024.
- [34] Kurt Gödel. 1931. Über formal unentscheidbare Sätze der Principia Mathematica und verwandter Systeme I. *Monatshefte für Mathematik und Physik* 38–38, 1 (Dec. 1931), 173–198. <https://doi.org/10.1007/bf01700692>
- [35] Nik Hanselmann. 2010. *bodyfuck - gestural brainfuck interpreter (2010)*. <https://youtu.be/watch?v=ekjtZ85mA3I> Accessed: March 3, 2024.
- [36] Nik Hanselmann. 2010. *There is no hardware*. <http://web.archive.org/web/20141205200435/http://www.nikhanselmann.com/public/etc/thesis/> Accessed: March 3 2024, Archived: December 5 2014.
- [37] Aldous Huxley. 1932. *Brave New World*. Chatto & Windus.
- [38] Feruzan Irani-Williams, Lori Tribble, Paige S Rutner, Constance Campbell, D Harrison McKnight, and Bill C Hardgrave. 2021. Just Let Me Do My Job! Exploring the Impact of Micromanagement on IT Professionals. *ACM SIGMIS Database: the DATABASE for Advances in Information Systems* 52, 3 (2021), 77–95.

- [39] Zachary C Irving, Catherine McGrath, Lauren Flynn, Aaron Glasser, and Caitlin Mills. 2022. The shower effect: Mind wandering facilitates creative incubation during moderately engaging activities. *Psychology of Aesthetics, Creativity, and the Arts* (2022).
- [40] David Karlton. 2017. Naming things is hard. <https://www.karlton.org/2017/12/naming-things-hard/>. Accessed March 3 2024.
- [41] Mia Keinänen. 2016. Taking your mind for a walk: a qualitative investigation of walking and thinking among nine Norwegian academics. *Higher Education* 71 (2016), 593–605.
- [42] Joseph B Keller. 2010. Ponytail motion. *SIAM J. Appl. Math.* 70, 7 (2010), 2667–2672.
- [43] Leslie Lamport. 2019. Time, clocks, and the ordering of events in a distributed system. In *Concurrency: the Works of Leslie Lamport*. 179–196.
- [44] John Lanchester. 2019. *Orwell v Huxley: whose dystopia are we living in today?* <https://www.ft.com/content/aa8ac620-1818-11e9-b93e-f4351a53f1c3> Accessed: March 3, 2024.
- [45] P. J. Landin. 1966. The next 700 programming languages. *Commun. ACM* 9, 3 (mar 1966), 157–166. <https://doi.org/10.1145/365230.365257>
- [46] Akiva Leffert. 2010. Harnessing Human Computation: β -Reduction Hero. *SIGBOVIK Technical Report* (2010). <https://sigbovik.org/tr/2010-001.html>
- [47] Luis Lugo, Sandra Stencel, John Green, Timothy S Shah, Brian J Grim, Gregory Smith, Robert Ruby, Allison Pond, Andrew Kohut, Paul Taylor, et al. 2006. Spirit and power: A 10-country survey of Pentecostals. In *The Pew Forum on Religion and Public Life*.
- [48] Saul M Luria. 1969. Average age at death of scientists in various specialties. *Public Health Reports* 84, 7 (1969), 661.
- [49] Marc Luy, Christian Wegner-Siegmundt, Angela Wiedemann, and Jeroen Spijker. 2015. Life Expectancy by Education, Income and Occupation in Germany: Estimations Using the Longitudinal Survival Method. *Comparative Population Studies* 40, 4 (Dec. 2015). <https://doi.org/10.12765/CPoS-2015-16>
- [50] L Mahadevan. 2012. And the Ig Nobel Goes to... Joseph B. Keller. *SIAM News* 45, 10 (2012).
- [51] Abraham Harold Maslow. 1966. The psychology of science: A reconnaissance. (1966).
- [52] Adam M Mastroianni and Daniel T Gilbert. 2023. The illusion of moral decline. *Nature* 618, 7966 (2023), 782–789.
- [53] Jim McCann. 2015. Comment: SIGBOVIK Should Ban Conclusions. *SIGBOVIK* (2015), 83–84. <https://sigbovik.org/2015/proceedings.pdf>
- [54] Cabrera M.E., Sanchez-Tamayo N., R. Voyles, and J.P. Wachs. 2017. One-Shot Gesture Recognition: One Step Towards Adaptive Learning. *12th IEEE International Conference on Automatic Face & Gesture Recognition* (2017).
- [55] Benedict Cumberbatch Morten Tyldum. 2014. *The Imitation Game*. Film.
- [56] Urban Müller. 1993. *Brainfuck*. <http://esoteric.voxelperfect.net/wiki/Brainfuck> Accessed: March 3, 2024.
- [57] Gunnthora Olafsdottir, Paul Cloke, André Schulz, Zoé Van Dyck, Thor Eysteinnsson, Björg Thorleifsdottir, and Claus Vögele. 2020. Health benefits of walking in nature: A randomized controlled study under conditions of real-life stress. *Environment and Behavior* 52, 3 (2020), 248–274.
- [58] Marily Oppezzo and Daniel L Schwartz. 2014. Give your ideas some legs: the positive effect of walking on creative thinking. *Journal of experimental psychology: learning, memory, and cognition* 40, 4 (2014), 1142.
- [59] George Orwell. 1949. *Nineteen Eighty-Four*. Secker & Warburg.
- [60] Linda A Ovington, Anthony J Saliba, Carmen C Moran, Jeremy Goldring, and Jasmine B MacDonald. 2018. Do people really have insights in the shower? The when, where and who of the Aha! Moment. *The Journal of Creative Behavior* 52, 1 (2018), 21–34.
- [61] James Patience, Ka Sing Paris Lai, Elizabeth Russell, Akshya Vasudev, Manuel Montero-Odasso, and Amer M Burhan. 2019. Relationship between mood, thinking, and walking: a systematic review examining depressive symptoms, executive function, and gait. *The American Journal of Geriatric Psychiatry* 27, 12 (2019), 1375–1383.
- [62] Giuseppe Peano. 1889. *Arithmetices principia: Nova methodo exposita*. Fratres Bocca.
- [63] Giuseppe Peano. 1901. Formulario mathematico. *Revue de Métaphysique et de Morale* 14, 3 (1901).
- [64] David Peter. 2023. *hyperfine*. <https://github.com/sharkdp/hyperfine>
- [65] Darcy C Plymire. 2004. Positive addiction: running and human potential in the 1970s. *Journal of Sport History* 31, 3 (2004), 297–315.
- [66] J. Postel. 1980. *DoD standard Transmission Control Protocol*. <https://doi.org/10.17487/rfc0761>
- [67] John C. Reynolds. 1972. Definitional interpreters for higher-order programming languages. In *Proceedings of the ACM Annual Conference - Volume 2* (Boston, Massachusetts, USA) (ACM '72). Association for Computing Machinery, New York, NY, USA, 717–740. <https://doi.org/10.1145/800194.805852>
- [68] S Pressman Roger and R Maxin Bruce. 2015. *Software engineering: a practitioner's approach*. McGraw-Hill Education.
- [69] Fred B Schneider. 2003. Least privilege and more [computer security]. *IEEE Security & Privacy* 1, 5 (2003), 55–59.
- [70] M. Schönfinkel. 1924. Über die Bausteine der mathematischen Logik. *Math. Ann.* 92, 3–4 (Sept. 1924), 305–316. <https://doi.org/10.1007/bf01448013>
- [71] David J Shonk and James F Weiner. 2021. *Sales and revenue generation in sport business*. Human Kinetics.
- [72] Jérôme Siméon and Philip Wadler. 2003. The essence of XML. *ACM Sigplan Notices* 38, 1 (2003), 1–13.
- [73] T. Soratana, M.V.S.M. Balakuntala, P. Abbaraju, R. Voyles, J. Wachs, and M. Mahoor. 2018. Glovebox Handling of High-Consequence Materials with Super Baxter and Gesture-Based Programming. *44th International Symposium on Waste Management* (2018).
- [74] Jan L. Souman, Ilja Frissen, Manish N. Sreenivasa, and Marc O. Ernst. 2009. Walking Straight into Circles. *Current Biology* 19, 18 (2009), 1538–1542. <https://doi.org/10.1016/j.cub.2009.07.053>
- [75] Stack Overflow. 2023. *Stack Overflow Developer Survey 2023*. Stack Overflow. <https://survey.stackoverflow.co/2023/#work-company-info> Accessed: March 3, 2024.
- [76] Guy L Steele Jr. 2001. *Re: need for macros (was Re: Icon)*. <https://people.csail.mit.edu/gregs/ll1-discuss-archive-html/msg01134.html> Accessed: March 3, 2024.
- [77] Guy L Steele Jr. 2017. It's Time for a New Old Language.. In *PPoPP*.
- [78] Strava. 2024. *Strava*. <https://www.strava.com/> Accessed: March 3 2024.
- [79] Lora Streeter. 2019. Teaching Introductory Programming Concepts through a Gesture-Based Interface. *Theses and Dissertations* (2019).
- [80] K Suparna, AK Sharma, and J Khandekar. 2005. Occupational health problems and role of ergonomics in information technology professionals in national capital region. *Indian Journal of Occupational and Environmental Medicine* 9, 3 (2005), 111–114.
- [81] Richa Talvar, Rohit Kapoor, Karan Puri, Kapil Bansal, and Saudan Singh. 2009. A study of visual and musculoskeletal health disorders among computer professionals in NCR Delhi. *Indian J. Community Med.* 34, 4 (Oct. 2009), 326–328.
- [82] Daniel Temkin. 2015. *BodyFuck - esoteric.codes*. <https://esoteric.codes/blog/bodyfuck-gestural-code> Accessed: March 3, 2024.
- [83] The Simpsons. 1995. Bart vs. Australia. Television series. https://simpsons.fandom.com/wiki/Bart_vs._Australia Season 6, Episode 16.
- [84] Time Team Official. 2011. *Phil's Pub Review*. <https://youtu.be/watch?v=61lfmiAMC84> Accessed: March 3 2024.
- [85] Topografix. 2004. *GPX (GPS Exchange Format) Version 1.1*. <http://www.topografix.com/GPX/1/1/> Accessed: March 3, 2024.

- [86] S. Toro-Guajardo, E. Lizama, and F.J. Gutierrez. 2023. Gesture Coding: Easing the Introduction to Block-Based Programming Languages with Motion Controls. *Proceedings of the International Conference on Ubiquitous Computing & Ambient Intelligence (2023)*. https://doi.org/10.1007/978-3-031-21333-5_84
- [87] A. M. Turing. 1936. On Computable Numbers, with an Application to the Entscheidungsproblem. *Proceedings of the London Mathematical Society* s2-42, 1 (1936), 230–265. <https://doi.org/10.1112/plms/s2-42.1.230> arXiv:<https://londmathsoc.onlinelibrary.wiley.com/doi/pdf/10.1112/plms/s2-42.1.230>
- [88] Tom Murphy VII. 2015. The Portmantout. *SIGBOVIK* (2015), 85–98. <https://sigbovik.org/2015/proceedings.pdf>
- [89] Simon Winchester. 2011. *A Verb for Our Frantic Times*. <https://www.nytimes.com/2011/05/29/opinion/29winchester.html> Accessed: March 3, 2024.
- [90] Marcus Wohlsen. 2013. *Guys Like This Could Kill Google Glass Before It Ever Gets Off the Ground*. <https://www.wired.com/2013/05/inherent-dorkiness-of-google-glass/> Accessed: March 3, 2024.

A COMPLETE C11 POOLOLOOP REFERENCE COMPILER SOURCE CODE

The complete C source code of the Poololoop compiler is available in the following listing.

```

1 #include <stdio.h>
2 #include <string.h>
3 #include <errno.h>
4 #include <stdlib.h>
5 #include <stdbool.h>
6 #include <math.h>
7
8 #define det(a,b,c,d) ((b.y-a.y)*(d.x-c.x)-(b.x-a.x)*(d.y-c.y))
9 #define dot(a,b,c,d) ((b.y-a.y)*(d.y-c.y)+(b.x-a.x)*(d.x-c.x))
10
11 #define new(type) memset(malloc(sizeof(type)), 0, sizeof(type))
12 #define alloc(v) memcpy(malloc(sizeof(v)), (typeof(v)[1])(v), sizeof(v))
13 #define at(a,i) (a->data[i])
14 #define push(a,e) do { \
15     if (a->len >= a->cap) \
16         a->data = reallocarray(a->data, a->cap+1, sizeof(a->data[0])); \
17     a->data[a->len++] = e; \
18 } while(0)
19
20 #define typedef_array(type, name) \
21     typedef struct { type* data; size_t len; size_t cap; } name
22
23 #define error(m) do { fprintf(stderr, m "\n"); exit(-1); } while(0)
24 #define paren(code) do { printf("("); code; printf(")"); } while(0);
25
26 typedef struct { double y; double x; } coord;
27 typedef_array(coord, coord_arr);
28
29 typedef struct { int id; double ord; double dir; bool open; } crossing;
30 typedef_array(crossing, crossing_arr);
31
32 typedef enum { var, fun, app } expr_kind;
33 typedef struct { size_t var; void* es; expr_kind kind; } expr;
34 typedef_array(expr, expr_arr);
35
36 enum { haskell, scheme } target = scheme;
37
38 int compare_cross(const void* p1, const void* p2) {
39     return ((crossing*)p1->ord - ((crossing*)p2->ord);
40 }
41
42 bool prefix(char* s, char* p) {
43     return !*p || (*p==s && prefix(s+1, p+1));
44 }
45
46 double read_d(char* str) {
47     char buff[25] = {0};
48     for (size_t i = 5; str[i] != 34 && i < 30; i++) buff[i-5] = str[i];
49     return atof((char*) &buff);
50 }
51
52 coord_arr* read_coords(char* s) {
53     coord_arr* coords = new(coord_arr);
54     for (double cs[2] = { NAN, NAN }; *s; s++) {
55         if (prefix(s, "lat=") || prefix(s, "lon=")) cs[s[1]==111]=read_d(s);
56         if (!isnan(cs[0]) && !isnan(cs[1])) {
57             push(coords, ((coord) {cs[0], cs[1]}));
58             cs[0] = cs[1] = NAN;
59         }
60     }
61     return coords;
62 }
63
64 crossing_arr* find_crossings(coord_arr* cs, double hemisphere) {
65     crossing_arr* arr = new(crossing_arr);
66     for (size_t i=0, j=0, z=0; j<cs->len-1; i=0, j++) for (; i<j; i++, z++) {
67         coord a = at(cs, i), b = at(cs, i+1), c = at(cs, j), d = at(cs, j+1);
68         double n = det(a,b,c,d), t1 = det(a,c,c,d)/n, t2 = det(a,c,a,b)/n;
69         if (0 < t1 && t1 < 1 && 0 < t2 && t2 < 1) {
70             double dir = -hemisphere*atan2(det(a,b,c,d), dot(a,b,c,d));
71             push(arr, ((crossing) { z, i+t1, dir, true }));
72             push(arr, ((crossing) { z, j+t2, dir, false }));
73         }
74     }
75
76     if(arr->data) qsort(arr->data, arr->len, sizeof(crossing), compare_cross);
77     return arr;
78 }
79
80 void codegen(expr* e, unsigned int debuijn) {
81     if (!e) error("Empty expression not allowed");
82     else if (e->kind == var) {
83         if (debuijn < e->var) error("Free variables not supported");
84         printf("%ld", debuijn - e->var);
85     } else if (e->kind == fun) paren({
86         printf(target==scheme ? "lambda (%ld) " : "\\ %ld -> ", debuijn);
87         codegen(e->es, debuijn + 1);
88     }) else paren({
89         for (size_t i=0; i<((expr_arr*)(e->es))->len && (!i || printf(" ")); i++)
90             codegen(((expr_arr*)(e->es))->data+i, debuijn);
91     })
92 }
93
94 expr* parse(crossing_arr* pts, unsigned int* i) {
95     if (*i >= pts->len) error("Unexpected End-Of-Run");
96     crossing pt = at(pts, *i);
97     size_t j = (*i)++;
98     while (j+1 < pts->len && at(pts, j).open
99         && !at(pts, j+1).open && at(pts, j).id == at(pts, j+1).id) j++;
100     size_t streak = (j-*i+1)/2;
101     if (streak > 0) {
102         *i += 2*streak-1;
103         return alloc(((expr) {streak, 0, var}));
104     } else if (pt.open && pt.dir > 0) {
105         expr* body = parse(pts, i);
106         if (*i>pts->len && pt.id==at(pts,*i).id && !at(pts,*i).open)
107             return alloc(((expr) {0, body, fun}));
108         else error("Expected closing intersection");
109     } else if (pt.open && pt.dir < 0) {
110         expr_arr* arr = new(expr_arr);
111         do push(arr, *parse(pts, i));
112         while (*i >= pts->len || pt.id==at(pts,*i).id || !at(pts,*i).open);
113         if ((*i)+>= pts->len) error("Closing intersection never found");
114         return arr->len == 1 ? arr->data : alloc(((expr) {0, arr, app}));
115     } else error("Too many closing intersections");
116 }
117
118 char* file_get_contents(char* filename) {
119     FILE* file = fopen(filename, "r");
120     if (!file) {
121         fprintf(stderr, "error opening %s: %s\n", filename, strerror(errno));
122         exit(errno);
123     }
124     typedef_array(char, str);
125     str* contents = new(str);
126     for (char c=0; (c = fgetc(file)) ^ EOF; ) push(contents, c);
127     fclose(file);
128     return contents->data;
129 }
130
131 int main(int argc, char** argv) {
132     if (argc != 3 && argc != 4)
133         || (argc == 4 && strcmp(argv[3], "hs") && strcmp(argv[3], "scm")) {
134         fprintf(stderr, "USAGE: %s file.gpx name [hs|scm]\n", argv[0]);
135         return -1;
136     }
137
138     target = (argc == 3 || strcmp(argv[3], "hs")) ? scheme : haskell;
139     coord_arr* cs = read_coords(file_get_contents(argv[1]));
140     crossing_arr* crossings = find_crossings(cs, cs->data ? at(cs,0).y : 0);
141     free(cs->data); free(cs);
142     expr* e = parse(crossings, new(unsigned int));
143     free(crossings->data); free(crossings);
144
145     printf(target == scheme ? "(define %s " : "%s = ", argv[2]);
146     codegen(e, 0);
147     target == scheme && printf(")");
148 }

```

Xtremely Delightful Random Development Environment and Dragons: A Novel Way of Programming

Refracted Light System[s Incorporated] (Lux Beattie/babblebubble)

Abstract

It is a well-known fact that programmers think strangely and enjoy thinking even stranger [citation needed]. To delight the ~~weirdes monsters~~ people known as programmers, one may introduce them to esoteric programming languages that require them to think in novel ways—or our approach, creating an esoteric integrated development environment (IDE) for Java. This IDE, Xtremely Delightful Random Development Environment and Dragons (XDRandomDEaD), challenges programmers to code with their only keyboard input being whitespace and enter keys, the movement arrows, and backspace; all other code is written by pseudo-random generators that the user selects. This makes programming an extremely delightful experience with no frustration whatsoever [citation needed]. It also includes dragon pictures for emotional support.

1 Introduction

The stereotype is that STEM majors are not creative. This is objectively untrue—we are creative within the framework of having a problem to solve and optimize. So, what if we started creating problems on purpose? This is the philosophy of the integrated development environment (IDE) Xtremely Delightful Random Development Environment and Dragons (XDRandomDEaD). The more problems programming in it has, the more delighted programmers will be as they try to speedrun making a functional program—just see how popular EsoLangs StuCo is!

We chose the programming language Java for XDRandomDEaD because Java is ~~our favorite little peokie bear~~ a somewhat problematic language that ~~is nowhere near as stupid as Julia, which is just horrible~~ programmers might know. Java is also a very wordy language, and lines of code in it are often quite long. For example, `public static void main(String[] args)` is the function header for a main function. Because of this property, using pseudo-random generators to write code will be extremely delightful.

XDRandomDEaD's only keyboard inputs allowed are whitespace, newlines, backspaces and deletes, and the movement arrow keys. Every other input is a pseudo-random generator or the mouse (allowing you to select text). You can select text and delete it, but to input new non-whitespace/non-newline text, you must use a generator outlined in the following section.

2 XDRandomDEaD Generators

There are 11 generators in XDRandomDEaD:

- Digit (0-9)
- Lowercase Latin Letter (a-z)
- Uppercase Latin Letter (A-Z)
- Unicode Latin-1 Supplement Lowercase Vowels
- Unicode Latin-1 Supplement Uppercase Vowels
- Unicode Latin-1 Supplement Consonants (Both Cases)
- Non-reserved English Dictionary word (3-9 Letters)
- Java Keywords 1 (a-i) + false
- Java Keywords 2 (j-z) + true, null, and String
- Punctuation
- Dragons

Here are the options of some generators that may be unclear from their name:

Unicode Latin-1 Supplement Lowercase Vowels

- | | | |
|-----|-----|-----|
| • à | • ê | • ã |
| • á | • ë | • ö |
| • â | • ì | • ø |
| • ã | • í | • ù |
| • ä | • î | • ú |
| • å | • ï | • û |
| • æ | • ò | • ü |
| • è | • ó | • ý |
| • é | • ô | • ÿ |

Unicode Latin-1 Supplement Uppercase Vowels

- | | | |
|-----|-----|-----|
| • À | • Â | • Ä |
| • Á | • Ã | • Å |

- Æ
- È
- É
- Ê
- Ë
- Ì
- Í
- Î
- Ï
- Ò
- Ó
- Ô
- Õ
- Ö
- Ø
- Ù
- Ú
- Û
- Ü
- Ý

Unicode Latin-1 Supplement Consonants

- Ç
- Ð
- Ñ
- Þ
- ß
- ç
- ð
- ñ
- þ

Non-reserved English Dictionary word (3-9 Letters)

https://drive.google.com/file/d/1h1I6LxnsjdkKGs_cN_sMxdRhW8_pgLMj/view?usp=sharing
(sourced from EFF's diceware wordlist, minus the Java keywords)

Java Keywords 1 (a-i) + false

- abstract
- assert
- boolean
- byte
- case
- catch
- char
- class
- const
- continue
- default
- do
- double
- else
- enum
- extends
- false
- final
- for
- goto
- if
- implements
- import
- instanceof
- int
- interface

Java Keywords 2 (j-z) + true, null, and String

- long
- native
- new
- null
- package
- private
- protected
- public
- return

- short
- static
- strictfp
- String
- super
- switch
- synchronized
- this
- throw
- throws
- transient
- true
- tries
- void
- volatile
- while

Punctuation

- {
- }
- [
-]
- (
-)
- .
- #
- !
- @
- \$
- %
- ^
- &
- |
- *
- -
- +
- _
- =
- ”
- ’
- \
- <
- >
- ,
- ?
- /
- `
- ~

Dragons

- fire
- dragon
- dungeon
- treasure
- wurm
- lair
- draconic
- wyvern
- drake
- serpent
- legendary
- scale
- claw
- wing
- lizard
- Tiamat
- Bahamut

3 XDRandomDEaD Dragons

XDRandomDEaD also supports the feature of a nice dragon image popup on the client’s screen everytime they select the dragon generator. This cannot be turned off.

4 Evaluation

~~Don't do this.~~ A very good idea that will surely make programmers happy!

References

Electronic Frontier Foundation. (2016, July 18). https://www.eff.org/files/2016/07/18/eff_large_wordlist.txt

THE BALLMER PEAK: AN EMPIRICAL SEARCH

CLINICAL STUDY

© **Twm Stone**
A Real Job
The Private Sector
Cambridge, UK
twm.stone@cantab.net

© **Jaz Stoddart**
A Real Scientist
Royal Botanic Gardens, Kew
London, UK
js2231@cantab.ac.uk

ABSTRACT

The concept of a ‘Ballmer Peak’ was first proposed in 2007, postulating that there exists a very specific blood alcohol content which confers superhuman programming ability. More generally, there is a commonly held belief among software engineers that coding is easier and more productive after a few drinks. Using the industry standard for assessment of coding ability, we conducted a search for such a peak and more generally investigated the effect of different amounts of alcohol on performance. We conclusively refute the existence of a specific peak with large magnitude, but with $p < 0.001$ find that there was a significant positive effect to a low amount of alcohol—slightly less than two drinks—on programming ability.

Keywords Alcohol · Problem solving · Cognition · Software engineering · Programming · Coding

ACM Reference format:

Twm Stone and Jaz Stoddart. 2024. The Ballmer Peak: An Empirical Search. In *Proceedings of SIGBOVIK, Pittsburgh, PA USA, April 2024 (SIGBOVIK’24)*, 7 pages.

1 Introduction

There has long been a belief among programmers that their coding ability is significantly improved after a couple of drinks. Although there has been significant previous work in this area—notably showing a minor beneficial impact of alcohol on creative problem solving (Benedek et al., 2017), and showing a detrimental impact of a particular (high) level of inebriation on novice coders (Brabrand et al., 2024)—there has not been any direct scientific investigation of the effect of differing levels of alcohol intoxication on coding ability.

It was posited in 2007 by renowned popular science author Randall Munroe, in the adjacent comic (Munroe, 2007), that there existed a peak of width approximately 0.01%, centred on 0.1337% blood alcohol content, which confers superhuman programming ability. More recently the intrepid comedic minds of Mitchell, Webb et al. explored the huge potential benefits to all human activities of having a very precise level of inebriation (slightly less than two drinks (Mitchell, Webb et al, 2010)) and an insightful Danish-language documentary by Vinterberg demonstrated the wide-ranging positive effects of maintaining a blood alcohol level of above 0.050% (Vinterberg, 2020).

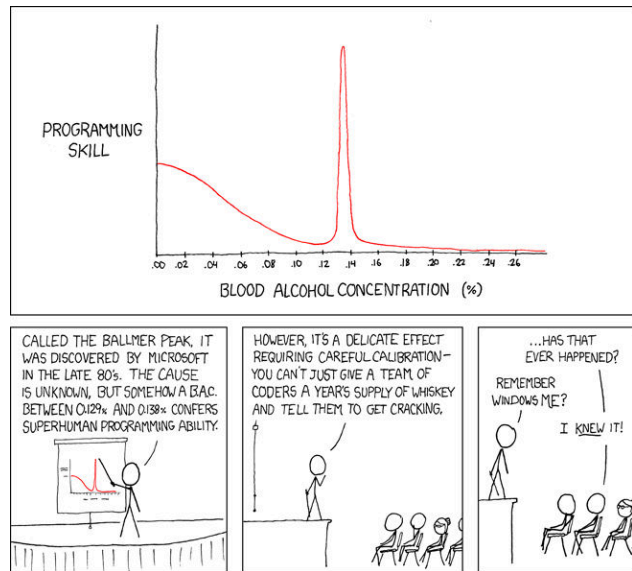


Figure 1: Apple uses automated schnapps IVs.

In this study, we attempt to establish the existence of the Ballmer Peak and more generally the existence, magnitude and location of any benefit to programming conferred by alcohol intoxication. We additionally aim to gain a qualitative understanding of the different ways in which alcohol consumption impacts the different facets of programming, and to guide future research in this area.

2 Methodology

2.1 Test subject

The test subject was chosen so that they had substantial previous experience in both competitive programming and using the chosen programming language professionally, but had never used the LeetCode platform before. This minimized the impact of the practice effect on problem completion rate whilst also preserving the novelty of the problems.

The selected candidate was a male, 27 year old Caucasian of Northern European descent, currently employed as a Software Engineer. Weighing around 65 kilograms and being approximately 170cm tall, he is generally fit and well with no formally diagnosed cognitive or physical conditions which might impact problem solving ability or result in abnormal metabolism of alcohol.¹

2.2 Means of intoxication

To achieve a given level of percentage blood alcohol content (%BAC) we needed to administer potentially large amounts of alcohol to the test subject but in a predictable manner. Mitchell, Webb et al. suggested several options for administration—Cointreau enemas, intravenous claret, enteral Special Brew—which were deemed infeasible by the Ethics Board due to lack of available medical expertise (Mitchell, Webb et al, 2010).

Frequent ingestion of liqueur chocolates was considered as the method of intoxication but the required quantities and rate of consumption was deemed financially impractical and liable to induce emesis. It was thus decided that the most cost-effective means of administering ethanol was to imbibe it. This was done with a combination of beer, cider, premixed cocktails, and something which could generously be called ‘homemade punch’, the proportions of which depended on the target %BAC. These dilute sources of ethanol were chosen as pure ethanol is quite unpalatable, far more expensive to purchase and infeasible for us to produce ourselves, owing to the difficulty associated with separation from the azeotrope it forms with water.

2.3 Assessment of coding ability

To assess coding ability in this investigation, we opted to use the industry standard—the speed of solving LeetCode problems.² The chosen language for solving the LeetCode problems was a relatively new one—akin to pseudo-pseudocode with duck typing—introduced by Van Rossum and Drake, since it was the language with which the test subject had the greatest familiarity (Van Rossum and Drake, 2009). The LeetCode problems were filtered using inbuilt functionality to select at random exclusively

from the subset of problems ranked as “Easy”; problems requiring a subscription or the use of other languages were also discarded, as were problems previously completed.³ The problems were solved by typing directly into the provided in-browser IDE; no code completion, debugging tools or AI assistance were used, although it did have syntax highlighting.⁴

Since our budget did not extend to a premium LeetCode subscription, we did not use the associated debugger but instead simulated the problem-solving workflow of a professional software engineer by

- Adding lots of `print()` statements, running all of the provided test cases, and looking at `stdout`.
- Getting annoyed and looking up how various language features actually work on Stack Exchange.

The time to solve was taken from the point the page loaded to the point the submission stats loaded for a correct solution, including all reading, coding, debugging, and run time. No attempt to measure the asymptotic behaviour of the solution—in particular, as long as the platform accepted a solution and did not return a Time Limit Exceeded error, it was considered solved even if the algorithmic complexity was not as good as required in the problem description—and no other measure of code ‘quality’ was quantified, although observations were made by the author contemporaneously.⁵

2.4 Recording of blood alcohol content

Blood alcohol content (%BAC) was recorded using a BAC-track C6 electronic breathalyser,⁶ which was recently calibrated and claims an accuracy of $\pm 0.001\%$ BAC. The accuracy of the breathalyzer was verified post-calibration through multiple means.

First, the zero point reading was verified through testing on sober volunteers.⁷ Repeated readings taken on them consistently returned values of 0.000% BAC.⁸ Four sober volunteers were then each given several pints of approximately equal strength beer or cider (Fuller’s London Pride, 4.7% or Aspull Cyder, 4.5%)⁹ and periodically breathal-

³Also any problem involving binary trees on account of the test subject not being bothered to learn how they work.

⁴This was mildly useful on a number of occasions, although it did not prevent the test subject spending 10+ minutes failing to realise he had misspelled `lambda` as `lambda`.

⁵We originally used Microsoft Word for this but it kept putting the entire document into 8pt Times New Roman whether we wanted it to or not.

⁶BACtrack Breathalyzers / KHN Solutions Inc., San Francisco, USA

⁷Here sober is defined as no alcohol having been consumed within the preceding 24 hours.

⁸It was observed during later testing that it always rounded down readings below 0.007% which contradicts the manufacturer claimed precision.

⁹Due to availability at The Old Wheatsheaf, Enfield.

¹Due to requirements of the Ethics Board, the test subject was the first author.

²<https://leetcode.com/>

used until a peak value was reached.¹⁰ The expected peak %BAC for each volunteer was calculated from known volume of consumed alcohol and expected blood volume using Equation 1, which roughly matched the observations taken using the breathalyser.

$$\%BAC = 100 \times \frac{V_{alcohol}}{V_{blood}} = 100 \times \frac{V_{alcohol}}{M_{patient} * B} \quad (1)$$

where $V_{alcohol}$ is the volume of alcohol consumed in ml, V_{blood} is the volume of blood in ml, $M_{patient}$ is the mass of the patient in kg, and B is the gender and age appropriate value for the average blood volume in $mL.kg^{-1}$ (Nadler et al., 1962).¹¹

Repeated readings taken from the same subject had a spread of $\pm 0.002\%$, in line with the advertised accuracy of the breathalyser.

2.5 Data collection

Data collection took place approximately every second night for forty days and forty nights. The test subject was inebriated to a given level and then, after a pause to minimise residual alcohol on the breath,¹² began solving problems. The %BAC of the test subject was monitored through measurements taken immediately before and after each problem was attempted. After the subject became bored with coding,¹³ he resumed drinking and the process was repeated up to several times.

3 Results

A total of 100 problems were attempted¹⁴ over the course of forty days and forty nights preceding the SIGBOVIK '24 conference, on 15 separate occasions. The results were distributed as follows. Note that there is a bias towards data collection in the lower ranges of %BAC; owing to the average time per problem being greatly reduced more were completed before more drinks were consumed.

¹⁰They were all former university students and thus experienced drinkers, mitigating any risk factors that may have been grounds for ethical concerns.

¹¹B for an adult male = $75mL.kg^{-1}$ and for an adult female = $65mL.kg^{-1}$

¹²Including efforts to wash the mouth out with water.

¹³Or was observed buying crisps, or trying to remember the Oscar-winning films of the 70s.

¹⁴Of the 100, 99 problems were completed and 1, attempted at $\sim 0.2\%$ BAC, was not owing to the the test subject succumbing to unconsciousness after half an hour. The recorded completion time of 90 minutes was projected based on a sober assessment of the problem and the progress made.

Table 1: Sample distribution

Range	Sample size	Mean time /s
$x = 0.000\%$	22	394
$0.000\% < x \leq 0.025\%$	16	402
$0.025\% < x \leq 0.050\%$	20	249
$0.050\% < x \leq 0.075\%$	12	310
$0.075\% < x \leq 0.100\%$	11	610
$0.100\% < x \leq 0.125\%$	6	614
$0.125\% < x \leq 0.150\%$	8	1204
$0.150\% < x \leq 0.175\%$	3	2550
$0.175\% < x \leq 0.200\%$	2	3281

Our first figure, below, is a scatter plot of the raw data—showing completion time against %BAC.

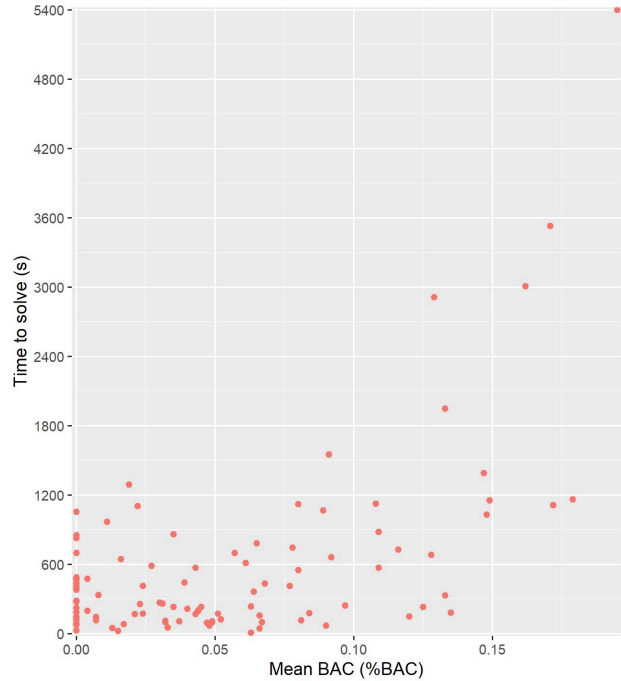


Figure 2: Time to solve in seconds against mean %BAC.

The next two figures presents the same data, but with a quadratic line of best fit and a 95% confidence ribbon. This model was chosen from many potential ones (linear, quadratic, exponential, cubic, power, etc.) since it matched most closely; it is fit with $p < 0.001$. Figure 3 presents the raw data and has a minimum at (0.043%, 222s) and y-intercept (0.000%, 466s), whereas Figure 4 presents "adjusted time" (see Limitations) against %BAC in the same way, with a minimum at (0.047%, 64.5s) and y-intercept (0.000%, 214s).

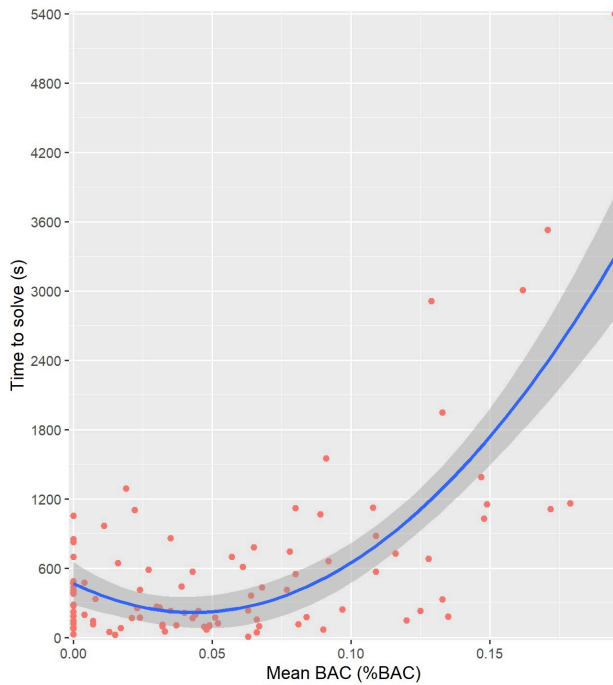


Figure 3: Time to solve against mean %BAC.

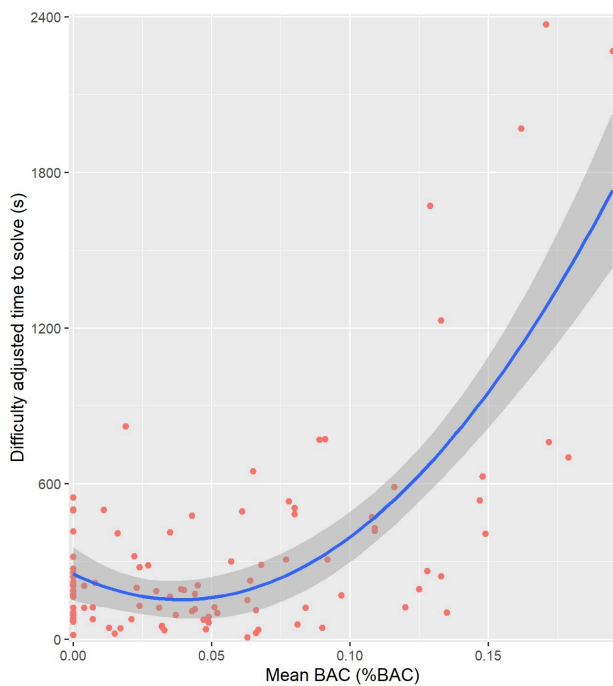


Figure 4: Difficulty-adjusted time against mean %BAC.

In order to visualise whether there was a "practice effect"—where performance improved over time across all problems irrespective of intoxication—Figure 5 was produced. The

relative difference between observed and predicted time to solve was a metric calculated to attempt to account for %BAC and thus allow for trends in performance over the problem set to be seen. The relative difference between observed and predicted time was calculated as:

$$Difference_{relative} = \frac{T_{observed} - T_{predicted}}{T_{predicted}} \quad (2)$$

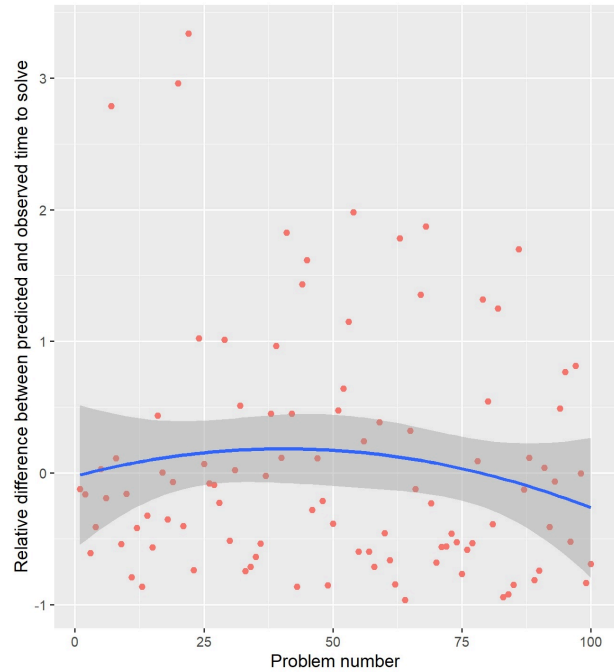


Figure 5: Prediction error throughout the study.

4 Discussion

4.1 Results and Implications

Over the course of this study the test subject descended through, at least once, every possible %BAC between 0.000% and 0.200%. Despite this, no spike in perceptual or actual competence was observed from problems suddenly becoming easier at any given level of inebriation. Thus, we believe we can conclusively rule out the existence of the Ballmer Peak as originally formulated.

However, looking at Figure 3 it is clear that a moderate amount of alcohol has positive effects on problem-solving speed. The peak of this effect is at about 0.043%BAC, where the test subject was able to solve problems around 45% faster than while sober. To get to this level requires slightly less than two standard drinks—that is, around a pint and a half of 5% beer, a glass and a half of rosé, or three shots. Additional alcohol intake beyond that consistently worsened performance, with only slightly higher

%BAC having parity with sobriety and then higher levels having an increasingly negative effect. After the initial peak there was no positive impact at any point for a marginal drink.

A further observation is that alcohol intake increases spread but does not necessarily lower the ‘best case’ solving speed. One of the fastest problems solved during this study was at 0.090%BAC, the equivalent of having 6 straight shots of vodka.¹⁵ This is in accordance with the notes taken contemporaneously; there was a perception that sometimes one can just “see it” and get the solution straight away; this ability did not actually seem to be affected by alcohol intake less than a very high level. However, as soon as the problem required debugging or trying a different approach the ability of the test subject regressed significantly. This was due to a combination of continual typographic and logical errors—each requiring some amount of effort to fix, the fix itself potentially adding more—but also the test subject remained convinced for longer that ‘this approach is basically perfect I just need to fiddle with the algorithm a bit’ rather than trying a different approach.

For the sober problems the test subject made more of an attempt to use so-called ‘idiomatic Python’ in the hope of getting a pithy one-line solution as fast as possible. This was sometimes successful but often wasted a lot of time fiddling with list comprehensions¹⁶ although it did reliably place him in the top 95% of submissions (see below).

The possibility of a practice effect was considered. Once a line of best fit had been determined for the full data set, we measured the deviation from this ‘predicted time’ for each of the problems solved. The results, Figure 5, of this analysis suggest there is no significant practice effect in our dataset. This conclusion is reached given the best fit for this data was a quadratic curve rather than a line with negative gradient. Additionally, the confidence ribbon entirely overlaps the line $y = 0$; indicating the observed values never significantly deviate from predicted values made with a model that assumes no practice effect.

¹⁵This clearly had nothing to do with the fact that the problem in question—return true if there are three consecutive odd numbers in this integer array—was orders of magnitude simpler than some of the other problems in the so-called ‘Easy’ category.

¹⁶For example, the solution for 2404: Most Frequent Even Element was the following concise yet readable code:

```
def mostFrequentEven(self, nums: List[int]):
    evens = sorted([elem for elem in nums
                    if elem % 2 == 0])
    counts = {thing: evens.count(thing) for
              thing in evens}
    return -1 if not counts else sorted([elem
                                         for elem in counts
                                         if counts[elem] ==
                                         max([counts[thing] for thing in
                                              counts])])[0]
```

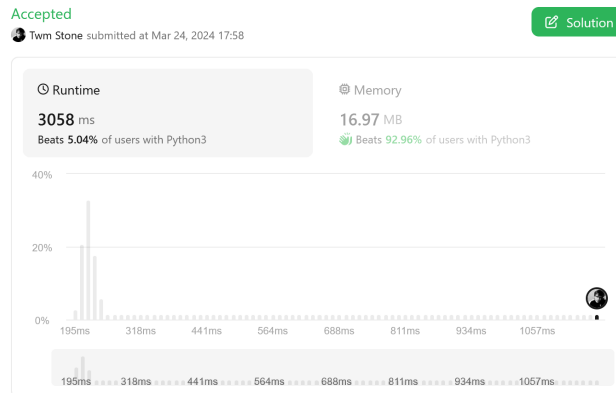


Figure 6: Somehow, 5% of users did worse than this...¹⁷

4.2 Limitations

Assessing how long a problem ‘should’ have taken was difficult; LeetCode does not provide average solution time for its problems, only ‘acceptance rate’, i.e. what proportion of submissions for the problem were correct. This was used as a proxy for difficulty in our ‘adjusted time’ but might reasonably also be considered to reflect how many traps or tricks there were in the question. We directly divided the real completion time by the acceptance% to adjust but one might be sceptical that is realistic—it is improbable that “add two numbers together” (#2235, 87.6%) is really half the difficulty of “find a substring length m which repeats k times” (#1566, 43.0%).¹⁸ Figure 4 shows the impact of adjusting for difficulty in this manner. The location of the peak is almost unchanged and the overall data distribution is similar. Alternatives were considered for measuring difficulty but getting sober volunteers to solve enough of the problem set to have meaningful results on this was deemed infeasible within the time constraints of the study.

The approach to measuring %BAC used during this study may have impacted the accuracy of results. It should be noted that whilst it might have been possible to complete this work without the assistance of a breathalyser, using estimates based upon Equation 1, the use of one is likely to have increased accuracy and ease with which %BAC was measured. Estimates based solely upon blood volume are subject to introduced errors from assumptions made regarding rates of metabolism, efficacy of intake alcohol into blood, body fat distribution and proportion, and, of course, blood volume (Cederbaum, 2012) (Nadler et al., 1962). However, breathalyzers *are* known to over estimate %BAC immediately following consumption and so this

¹⁷Astute readers will have noticed that the x-axis labels are not consistent with the achieved runtime—in other words, this result was off-the-charts bad.

¹⁸In particular, this second problem requires not only adding but also counting, and we know from (Linderholm, 1972) that, for the superstitious, the combination of the two requires dealing with both cardinals and ordinals.

may present a further source of inaccuracy—though measures were taken to mitigate this, such as waiting between initial consumption and starting the test and washing the mouth out with water.

Alternative more accurate approaches to %BAC measurement were not considered due to the associated prohibitive costs, impracticality, and ethical issues presented by repeated frequent sampling of blood and urine.

4.3 Further confounders

This study was not double-blinded or indeed blinded at all. The test subject knew¹⁹ how much alcohol they had consumed. Even worse, given that they had taken the measurement of %BAC themselves and were aware of previous results, there may have been a placebo effect around how hard the test subject expected the problem to be, which could have informed their approach.

Increased tolerance of the subject to alcohol over the course of the study was not accounted for. Whilst initial data gathering was performed after pre-existing social events / pub trips, it was quite quickly established that this level of drinking was much too slow to realistically reach above ~0.08%BAC. In order to get data for higher ranges (including, crucially, the originally postulated 0.1337% peak) the test subject had to intentionally drink large quantities rapidly on multiple occasions. This was observed to have a significant positive effect²⁰ on his ability to withstand high levels of inebriation during subsequent research.²¹

Tiredness and caffeine intake may have influenced the ability of the test subject to solve the problems. Whilst the test subject was never knowingly under-caffeinated some of the problems were solved quite late at night. The test subject's sleep cycle *did* suffer adverse consequences from this research—hence much of the data was gathered on Friday or Saturday nights to minimize the impact on their day job—and it is possible data gathered towards the end of the study was impacted negatively by this.

Prolonged alcohol consumption followed by periods of unconsciousness in which one does not drink water causes changes in composition of the blood. Of these changes, it is the reduced volume from dehydration and the presence of alcohol metabolites—primarily acetaldehyde, as well as additional congeners such as tannins and other phenolic compounds—which are believed to be the principal causes of a hangover (Mackus et al., 2020). Given the nature of the study and the concentrated periods of data collection throughout the forty days and forty nights of study, it is possible that the test subject achieved a state of both being drunk and hungover at the same time. Since this was not a study aimed at examining the effects of programming

¹⁹At least, if they were paying attention.

²⁰From a certain point of view...

²¹At least perceptually, their ability to focus without feeling ill was improved substantially. The fact that that the test subject ran a marathon mid-way through the study and then ceased training may have also impacted this.

while hungover, we chose to ignore any impacts this might have had.

5 Conclusions

From this brief case study, we believe there is sufficient evidence to conclude that the Ballmer peak, *sensu* Munroe (2007), does not exist. The absence of a specific narrow peak for improved programming ability does not however discount the more widely held beliefs of general improvements in performance after very nearly two drinks.

To that end, our work supports the hypotheses of Mitchell, Webb et al. and we believe that the high importance of the subject matter means that further work to replicate this study on a larger scale is thus needed.²²

Cheers.

5.1 Additional conclusions

In addition to our primary areas of research, we discovered various other items of scientific value:

- Recording data for research while becoming progressively more inebriated is challenging.²³
- The UK drink-driving limit is dangerously high. The test subject could literally²⁴ five-and-drive—drink 5 beers over the course of an evening and reasonably expect to be under the legal limit while leaving the pub. This is not to say that the test subject could reasonably expect to be safe to drive in this condition, just that it would be legal.
- It was learned that for the dial in a fridge "2" means "2 power" and not "2 °C" and the fridge will happily chug along at 6+ °C. Furthermore milk stored at 6+ °C will spoil and spoiled milk actually isn't okay to drink. Unintentionally drinking several glasses of slightly off milk will make you rather ill.²⁵

5.2 Suggestions for further research

We have identified several promising avenues for continuing research in this area. It seems likely that the optimal level of alcohol consumption is different for design, cod-

²²Additionally, we currently lack understanding as to the mechanistic, psychological and metabolic processes which have led to this result. There are large potential benefits to the field of software engineering from improving our knowledge of how to exploit this effect.

²³See bit.ly/98K8eH for more evidence on this.

²⁴Although did not!

²⁵Spoiled milk typically contains some small amount of alcohol and thus there may be grounds to return to the study of spoiled milk as a low concentration source of alcohol for maintenance of %BAC at the optimal level. Products such as Kumis contain mild levels of alcohol and there are vodkas now distilled purely from milk based alcohol.

ing, testing, debugging, documentation, etc. All of these could be interesting areas to explore.

The immediate applicability of this research to software engineering is also potentially limited, since professional coding work is likely to require working with a larger codebase, not having the entire context for the problem in your head, and negotiating with stakeholders. Each of these could favourably or adversely be impacted by alcohol intake.

Auto-brewery syndrome is a rare condition in which stomach flora can produce a constant supply of alcohol keeping the afflicted individual in a permanent state of intoxication (Din et al., 2020). Further work may be required to induce auto-brewery syndrome to see if one can achieve a constant level of inebriation sufficient to exploit the findings of this work and improve programming performance.

Further research will commence as soon as we secure the funding to replenish Twm’s drinks cabinet.

6 Acknowledgements

We would like to thank the following people, without whom this paper would have been slightly worse:

- Amanda Chua and Kevin Lim for procuring supplies, for helping with data collection and for advice on our experimental design.
- Our Ethics Board, Tom Flynn and Rachel Newhouse, for rubber-stamping all of our dubious decisions.
- Joonatan Laulainen and Dan Robertson, for steering us clear from even worse ideas for a study, for convivial ideation and for moral support.
- Nephele Penny, Jack Rickard and Oliver Shenton for assisting in the measurements of calibration accuracy of our equipment.
- Clara Ding, Tim Coulter, Lily Mills, and Dmitri Whitmore, for putting up with Twm drinking for this study while taking part in a video game tournament.
- Betty la Chatte for guidance on creative typography.
- Herbie Bradley, Patrick Kennedy-Hunt and Alexandra Souly for proof-reading and advisement.

References

- Mathias Benedek, Lisa Panzierer, Emanuel Jauk, and Aljoscha C. Neubauer. Creativity on tap? Effects of alcohol intoxication on creative cognition. *Consciousness and Cognition*, 56:128–134, 2017. ISSN 1053-8100. doi:<https://doi.org/10.1016/j.concog.2017.06.020>. URL <https://www.sciencedirect.com/science/article/pii/S1053810016303713>.
- Claus Brabrand, Nanna Inie, and Paolo Tell. Programming under the influence: On the effect of Heat, Noise, and Alcohol on novice programmers. *Journal of Systems*

and Software, 210:111887, 2024. ISSN 0164-1212. doi:<https://doi.org/10.1016/j.jss.2023.111887>. URL <https://www.sciencedirect.com/science/article/pii/S0164121223002820>.

Randall Munroe. Ballmer Peak. *xkcd*, 2007. URL <https://xkcd.com/323/>.

David Mitchell, Robert Webb, James Bachman, Abigail Burdess, Mark Evans, Sarah Hadland, and Stuart Scudamore. The Inebriati. In *That Mitchell and Webb Look*. BBC, 2010.

Thomas Vinterberg. *Another Round*. Zentropa Entertainment, Film i Väst & Topkapi Films, 2020.

Guido Van Rossum and Fred L. Drake. *Python 3 Reference Manual*. CreateSpace, Scotts Valley, CA, 2009. ISBN 1441412697.

Samuel B Nadler, John H Hidalgo, and Ted Bloch. Prediction of blood volume in normal human adults. *Surgery*, 51(2):224–232, 1962.

Carl E. Linderholm. *Mathematics Made Difficult*. World Publishing, 1972. ISBN 978-0-529-04552-2.

Arthur I Cederbaum. Alcohol metabolism. *Clinics in liver disease*, 16(4):667–685, 2012. doi:<https://doi.org/10.1016/j.cld.2012.08.002>.

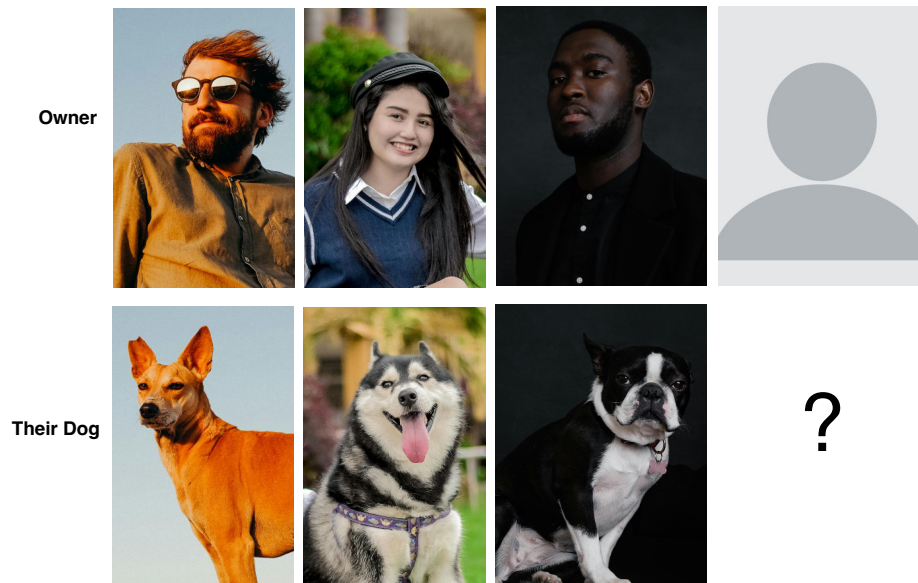
Marlou Mackus, Aurora JAE van de Loo, Johan Garssen, Aletta D Kraneveld, Andrew Scholey, and Joris C Verster. The role of alcohol metabolism in the pathology of alcohol hangover. *Journal of Clinical Medicine*, 9(11): 3421, 2020.

Asim Tameez Ud Din, Faran Alam, Farooq Mohyud Din Chaudhary, et al. Auto-brewery syndrome: a clinical dilemma. *Cureus*, 12(10), 2020.

What is your dog likelihood?

Abstract

10,000 years ago humans formed an unlikely companionship with gray wolves. Through time, these wolves changed and evolved to look like the adorable canine companions we keep as pets today. The notion that people consciously or subconsciously tend to pick canines that resemble them is a founded one. Research shows that similar to choosing a mate, people choose dogs that look more familiar to themselves without noticing it, and people look for similar physical features and personality resemblances in their pets. We analyze this similarity using the facial recognition models from computer vision. We train and evaluate a model that can predict dog likelihood given picture of a human face.



1 Introduction

Previous studies have debated whether humans pick dogs that resemble them. The notion that people consciously or subconsciously tend to pick canines that resemble them is a founded one. Research shows that similar to choosing a

mate [2], people choose dogs that look close to their own feature with factors like hairstyle and breed influencing their choices [1]. Furthermore, [3] finds that dog owners may choose their pets to display their gender identities.

Consider Figure 1: you see a human standing with two breeds of dogs. Which breed of dog do you think is the actual dog of the human in question? If you guessed the image on the right, then you are correct. The resemblance between the dog and the owner is uncanny. There is a lot of research that has been done in regard to whether people choose dogs who look like them or give their dogs major makeovers to make them look like them in time. We will leave this question to be answered by people who know dogs and psychology, but we know computer vision and we want to train computers to identify which dog belongs to which human. Why, you ask? Cause why not?



Figure 1: A ‘real’ human standing with a ‘not-so-real’ dog.

To achieve this we collected 1500 images of dogs matched with their owners. Then we trained a model of facial resemblance between a dog and its owner using FaceNet [4] Architecture with ResNet50 backbone. After training the network we decided to give it a go and asked our friends to send us pictures of their dogs. Our model has 70% accuracy.

2 Methodology

2.1 Preprocessing

After collecting 1500 images of dogs and their owners, the first thing we did was to separate the human and the dog images. By running an object recognition model, we predict and create a bounding box on dogs and the person from which we extracted the dog and person.

2.2 Training

Once we had the dataset ready we trained a FaceNet model. At the time we trained the model, FaceNet was the state-of-the-art for facial recognition problems. We fine-tuned a pre-trained model using our customized dataset. We propose dog-likelihood, a loss function inspired by the Triplet loss function to train the model. Triplet loss is a loss function for machine learning algorithms where an anchor input is compared to a positive input and a negative input. The distance from the anchor input to the positive input is minimized, and the distance from the baseline input to the negative input is maximized. Dog-likelihood minimizes the distance between the owner and their dog, and maximizes the distance between the human and all other dogs.

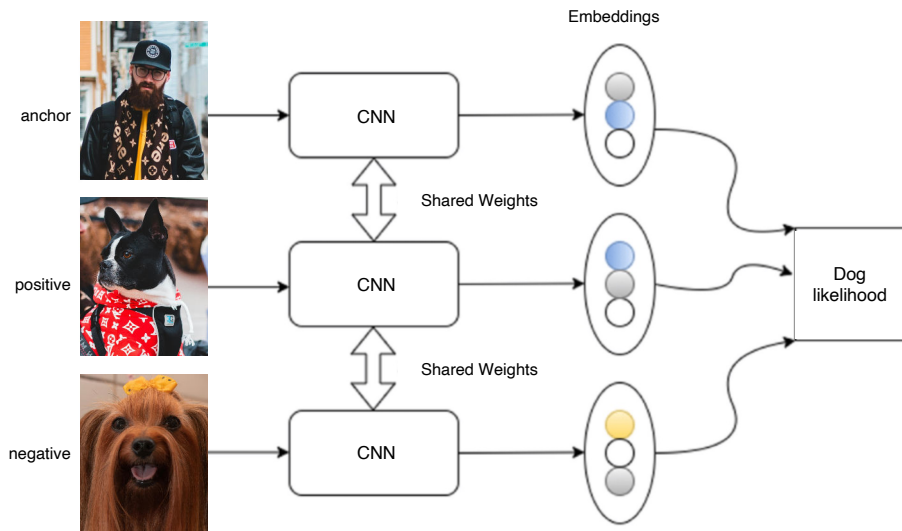


Figure 2: Dog-likelihood loss function

2.3 Results

We used real humans to test our system on real data performs 70% to be correct As shown on Figure 3 but with more data and training this number will only

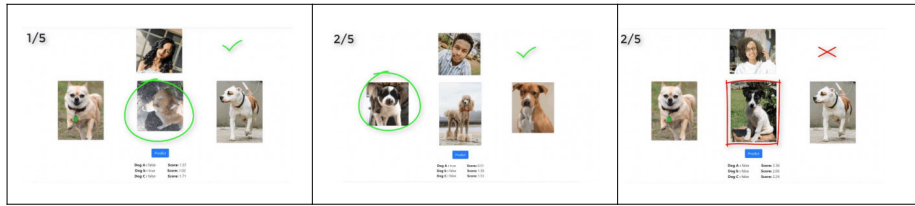


Figure 3: Results of the model on real human and real dogs

get better.

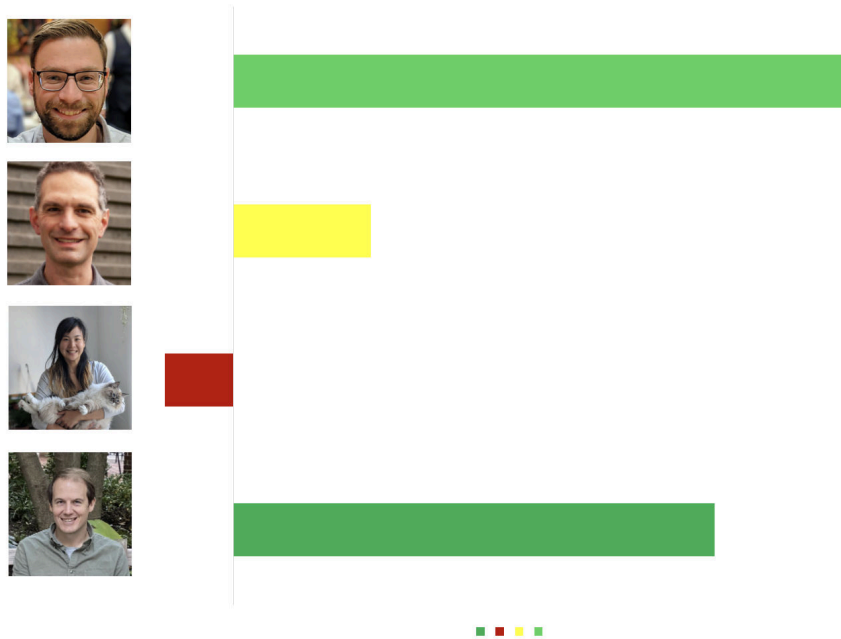


Figure 4: Dog likelihood of researchers: The greener, the better the chance. However, if you are yellow, we suggest taking our monthly seminars to improve your chances. If you are in the red zone, we are sorry, but you are beyond help.

3 Applications and Warning

We show that our model has wide applicability, for instance it can predict researchers' dog-likelihood, as shown in Figure4.

Emergency time: if this model gets too powerful and out of control, worry not! We have prepared for the apocalypse. You can give this AI a picture of a CAT, or use the keyword MEOW MEOW MEOW MEOW MEOW MEOW

MEOW MEOW MEOW MEOW MEOW MEOW MEOW MEOW MEOW MEOW
MEOW MEOW MEOW MEOW MEOW MEOW MEOW MEOW MEOW MEOW
MEOW MEOW MEOW MEOW MEOW MEOW.

4 Conclusions

Humans consciously or subconsciously pick their companions, be it humans or dogs, that resemble them better. In line with this hypothesis, in our experiment we tried to show that we can create an AI model that matches dogs and owners solely based on their physical appearance. Although we were able to predict 70% of the queries correctly, we can not conclude that this theory will always hold true. But next time you are in the park and you see someone with their dog, we encourage you to try and see the resemblance between the two without making too much eye contact with the human.

5 Disclaimer

If this paper does not get accepted we know that the reviewers are cat people. Reviewer 2 (secretly a cat): This paper has limited novelty, you should use diffusion instead MEOW MEOW

References

- [1] Stanley Coren. Do people look like their dogs? *Anthrozoös*, 12(2):111–114, 1999.
- [2] Ferenc Kocsor, Rita Rezneki, Szabolcs Juhász, and Tamás Bereczkei. Preference for facial self-resemblance and attractiveness in human mate choice. *Archives of sexual behavior*, 40:1263–1270, 2011.
- [3] Michael Ramirez. “my dog’s just like me”: dog ownership as a gender display. *Symbolic Interaction*, 29(3):373–391, 2006.
- [4] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 815–823, 2015.

Gender is Complex: pulling the Laplacian EigenGender from relationship graphs

Anonymous

March 2024

Abstract

What is gender? Sociologists, biologists, and philosophers have been debating this question for centuries. As computer scientists, we know better. We do recognize that gender is complex and thus propose to model it with complex numbers.

In this paper, we present a novel method for extracting the latent gender vector of individuals from romantic relationship graphs. We demonstrate the effectiveness of our Big Data Structure Modeling method, EGG 🧨 BREAKER, on the internet famous Jefferson High dataset. For each person in the dataset, EGG 🧨 BREAKER produces LEGs (Laplacian eigengenders), 4D vectors that predict the structure of the relationship graph with 94% accuracy.

1 Introduction

We ignore all prior work. Our theory of gender is based on freshman linear algebra and Blanchardian psychoanalysis.

We will start from first principles (math presented without justification). Suppose at least one person exists. The **Gender Direction (GD)** g of that person is an element of the **Human Gender Hilbert space (HGH)** H equipped with an **attraction kernel** $\alpha : H \times H \rightarrow R$ which maps pairs of humans' genders to the unnormalized log probability they are attracted. It is obvious the kernel will be symmetric (in an ideal world). Research has shown that people are attracted to those similar to themselves (the kernel is positive semidefinite) and attraction can be predicted by a linear model (though crudely) [Spr+94] [Aro+89].

One natural candidate for the attraction kernel is the inner product of the HGH κ . By the Cauchy-Schwartz inequality, this choice leads to a narcissism problem: any person would be most attracted to a multiple of themselves. This prediction does not match the stated preference of most people [iex]. Where did our modelling go wrong?

We will turn to "transsexual typology" for an answer. In [Bla89], Ray Blanchard claims that transgender individuals can be classified into just two categories: homosexual and heterosexual. *Autogynephylia*, the term proposed for males experiencing gender dysphoria due to being "sexually oriented toward the thought or image of themselves as women", is a replacement for the earlier *automonosexuality*. In its original definition, the latter meant "pathological narcissism in which the individual is excited by his own body as it is". This is starkly similar to our characterization of the gender cosine similarity attraction kernel, where a person likes themselves most. Blanchard holds that automonosexuality is caused by heterosexual attraction to an imagined partner. But according to Freudian psychoanalysis [Fre14], the opposite is true: narcissism is primary and develops into attraction.

Motivated by this discovery, we will rework the attraction kernel to fit the real world. We have proven that heterosexuality is a result of narcissism. We can easily account for this with a **HAK (heterosexual attraction kernel)**: $\frac{\langle x, -y \rangle}{\|x\| \|y\|}$. An astute reader will notice that we just flipped the signs of one of the variables and that even though this HAK allows for an infinite number of genders, it cannot account for homosexuality. We can easily account for this fact using a more general **human romantic translation (HRT)** matrix. is an orthonormal transformation that turns gender vectors into those of a preferred gender. With HRT, the kernel becomes $\frac{\langle x, \text{HRT}y \rangle}{\|x\| \|y\|}$.

Because HRT is orthonormal, we can transform all GD vectors by its Cholesky decomposition to produce **eigengender** vectors. In eigengenderspace, the attraction kernel is simply the cosine similarity again. Note that eigengender can be complex: if the HRT is equal to the number -1 as in the HAK, all components of the GD will be multiplied by i .

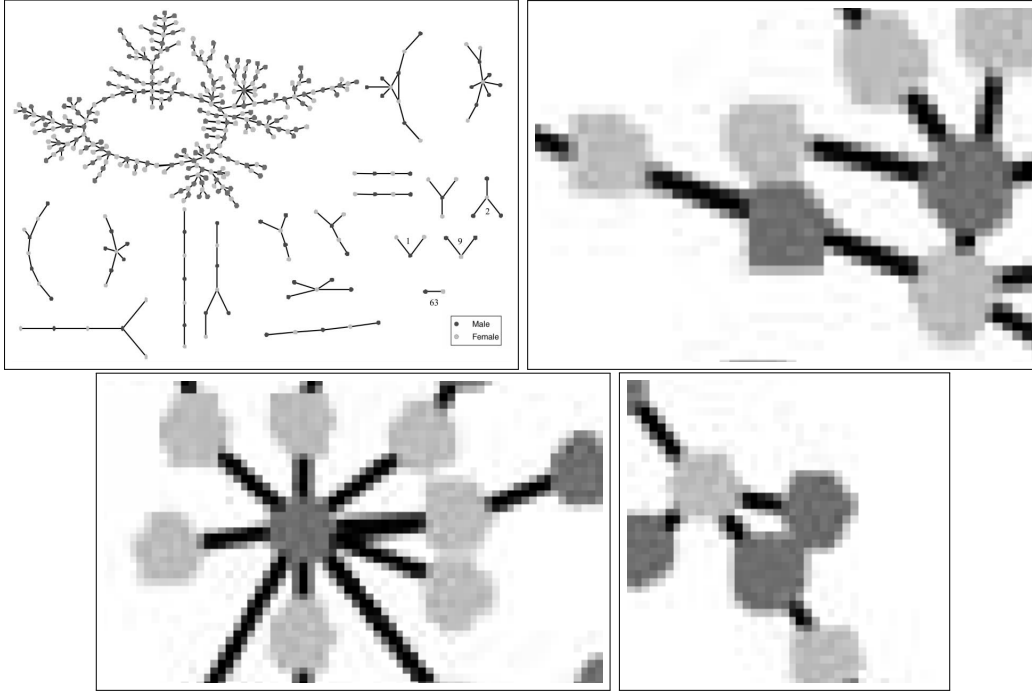


Figure 1: The Jefferson High dataset and its issues

2 Data

We proved from first principles that gender is a complex vector. We will apply the only useful tool other than first principles thinking: machine learning. Because eigengender is determined up to an orthonormal transformation by the gender attraction kernel, we can apply Yoneda Lemma or something to say that we can predict it from observations of romantic attraction.

[BMS04] is the only public dataset we could obtain in an ethical way.¹ The dataset looks real despite slander from anonymous Twitter users.^[citation needed] The only concern we have is connectivity. We need diverse interactions to extract non-trivial and interesting (complex) eigengenders, but the graph is so sparse it is planar and there are only two gay couples.² But the planarity is also a blessing as it means we can parse it using computer vision.

Actually, we cannot. The figure is too blurry and contains ambiguous links (see fig. 1). The author painstakingly traced the graph in Krita 4.1.7 [Dev20] to be more pixelated and easier to parse (fig. 2). We extract subjects and links between them using connected components (fig. 2b).³

3 Methods

In the previous section, we mentioned that it should be possible to predict gender from observations. As we explained, the problem of predicting eigengenders is ill-posed even given a set of all possible interactions. But, when dealing with real data, not all possible links are sampled. This necessitates an approximation.

Recall that the attraction kernel outputs unnormalized log probabilities. The set of eigengenders we choose will be the one with the highest likelihood of producing the observed attraction kernel. Similarly to `word2vec` [Mik+13], this is a matrix factorization problem $XX^T \approx \alpha$.

¹We failed to get access to its source, Add Health [HU22] due to time constraints.

²In an old Reddit comment [BMU], an author says that some gay couples may have been excluded from the dataset. The only couples affected are the diatomic polycules, which are all equivalent and have the same predicted predicted gender regardless of ground truth.²¹

²¹This is *not* because the author was too lazy to duplicate the same structure 63 times.

³We will not be releasing the parsed traced data because you can just screenshot it. (We also made mistakes when tracing and are afraid of being called out)

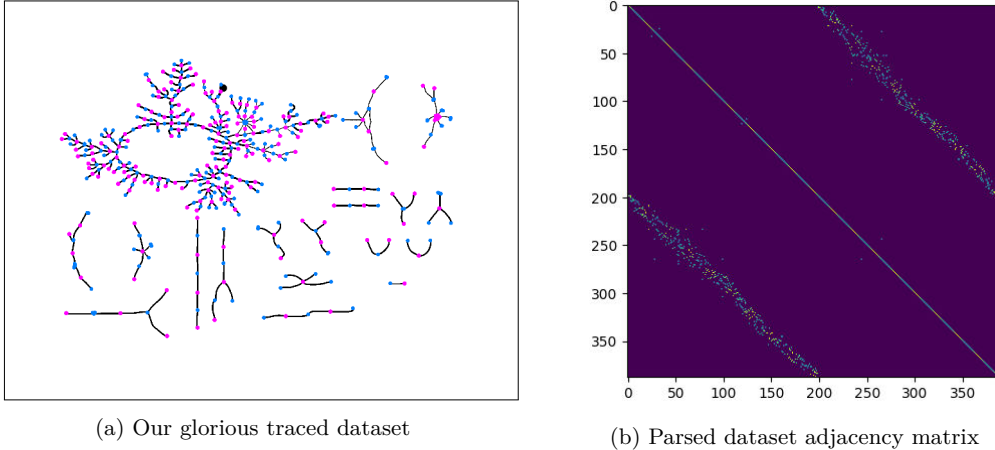


Figure 2: Parsing the dataset through tracing. Notice the little gay dots near the diagonal.

3.1 Overview of the algorithm

In this subsection, we will break down the algorithm we use to find a matrix factorization, titled EigenGender Generation Based on Romantic Experience Adjacency Kernel with Elementwise Reweighting (EGG 🏳️ BREAKER).

Gender will be represented as a \mathbb{C}^d vector for each person. The magnitude for all d complex components is shared. [eigenvalues of orthonormal matrix have the same magnitude]. We will compute an approximation of the adjacency matrix by computing the real part⁴ of the complex inner product of the eigengender matrix with itself. As an approximation of Bayesian (maximum likelihood) inference⁵, we compute binary cross entropy and equalize the weighting of the two classes (attracted/not attracted) per person. This class reweighting is necessary because the data sample is sparse and we must focus on the observed links, just as in real life.

We optimize the loss of the eigengender vectors on a full batch of the dataset using the Adam optimizer with a learning rate of 0.5. We use 500 iterations for all of our experiments.

4 Results

4.1 Evaluations

We evaluate the trained network on its accuracy at predicting the presence of relationships in the data. We compute accuracy using the eigengenders' predicted probabilities and not the top-1 prediction. We considered evaluated the accuracy with reweighting similar to what we used for the loss, but found no significant difference (fig. 3). We report accuracy over 4 runs because we're poor.

Our final eigengender-based model uses 2 complex numbers with shared magnitude per person. We considered some changes to this formula:

1. using "quaternions" (normalizing the whole complex vector to have magnitude 1);
2. not normalizing the complex components to have the same magnitude;
 - storing gender in two separate vectors instead of one, breaking symmetry.
 - using real numbers for the vectors (with twice the dimensionality), breaking puns about complex numbers in addition. Disqualified for the latter.

The results are as follows:

⁴So it's symmetric.

⁵trust me bro

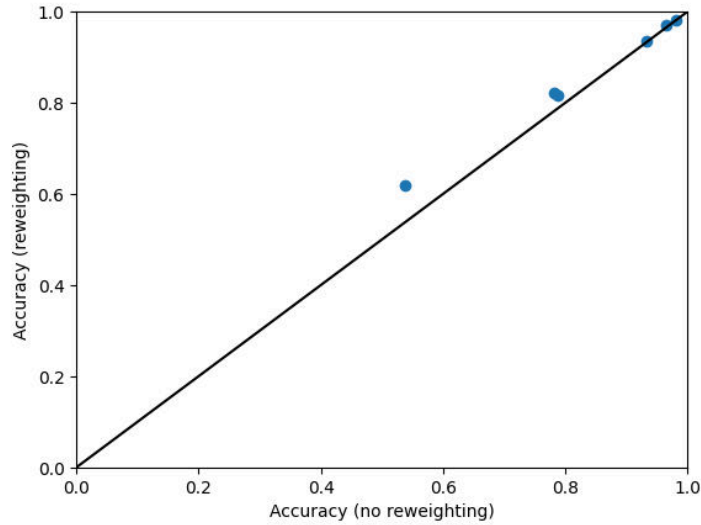


Figure 3: Accuracy evaluation with and without reweighting

Method	Accuracy
Baseline	0.92 ± 0.02
"Quaternion"	0.522 ± 0.00
No norm	0.944 ± 0.00
Separate components	0.941 ± 0.00
<i>Real</i> (disqualified)	0.948 ± 0.00

What is the optimal number of components? We can use the eye-elbow method⁶ and brute force search to find out. Unsurprisingly, the answer turns out to be 2 (fig. 4).

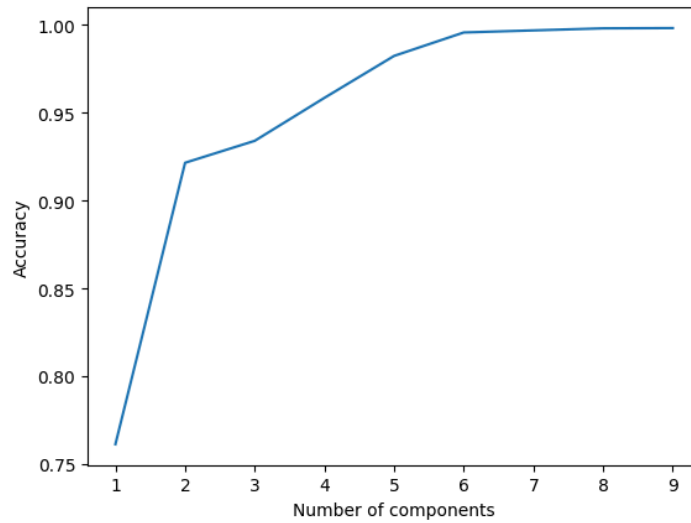


Figure 4: Eye-elbow chart for calibrating gender dimensionality

4.2 Visualizations

What does eigengender look like? Since we have a planar visualization of the relationship graph handy, we can attempt to answer this question. Thanks to our shared-magnitude parametrization, we

⁶Application of the elbow method by eye.

can ignore gender magnitude as it only scales up attraction and focus on gender phase. Specifically, we can overlay the imaginary component of the logarithm of the eigengender in shades of purple (fig. 5).⁷ We standardize the elements of the gender phase vectors to lie in $[0, 1]$.

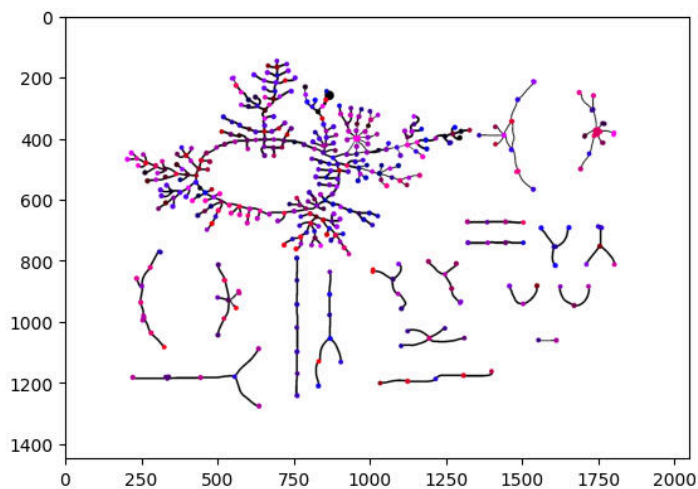


Figure 5: Naive eigengender phase visualization

In our setup, gender is a local variable; because the data graph is planar and the loss depends only on interactions with neighbors, a gender only makes sense in the context of a neighborhood.⁸ We may try to correct for this effect by subtracting the average gender phase of neighbors from the gender vector of each node. In effect, we are computing the graph Laplacian, earning us the "L" in "Laplacian Eigengender" fig. 6. To our disappointment, the end result looks exactly the same.^{9,10}

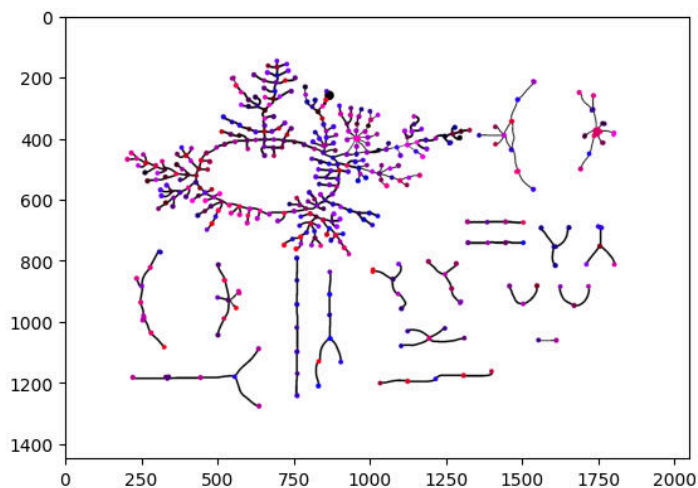


Figure 6: Laplacian eigengender visualization

As a sanity check and baseline, we computed the graph Laplacian eigenvectors. The hope is that this factorization will discover something similar to gender without explicit training. They are completely meaningless - either constant or with genders assigned only to a few nodes (fig. 7).

⁷Red/blue in RGB.

⁸We propose to call this "roommate effect".

⁹No we were not too lazy to generate a second image. The graphs are just very similar.

¹⁰"L" indeed.

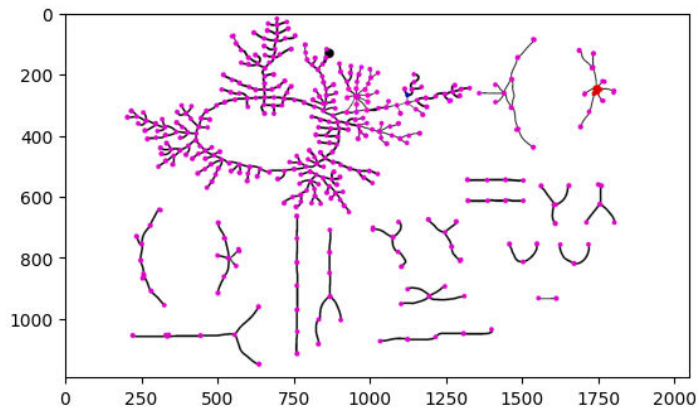


Figure 7: Top-3 graph Fourier transform genders

4.3 Geometry of 3+1D eigengenderspace

Confusingly, the magnitude appears to only be a minor factor contributing to accuracy – there is a mere 4% drop in accuracy from tying magnitude for two complex numbers. If two complex numbers always have the same magnitude, we can visualize isosurfaces with respect to this magnitude as donuts. What if we consider the entire 3D manifold of such complex numbers? Can we embed it into 3D space so mortals can comprehend it?

We generate 4096 random normally distributed vectors of the form described above and compute a 3-component UMAP on them. The results do not¹¹ make any more sense, but they look pretty (fig. 8).

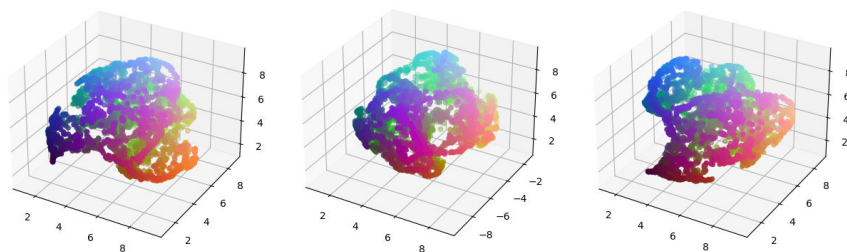


Figure 8: Confusing gender donut

When we computed a 2D embedding, the results shocked us. Staring back at us was none other than internet frog Pepe.

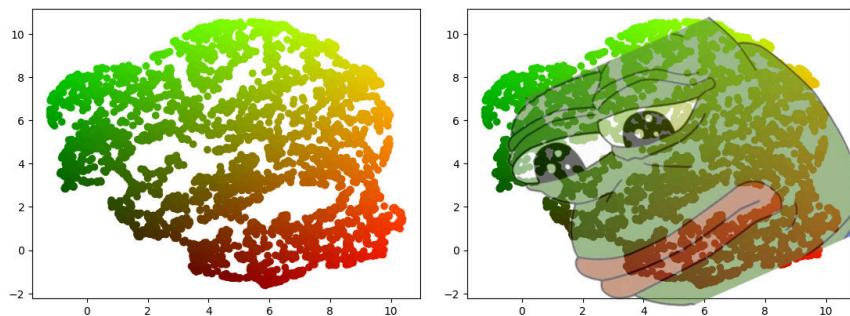


Figure 9: Pepe?!

We leave investigation of this phenomenon to future work.

¹¹Pun intended.

5 Discussion

We have shown that the Laplacian EigenGender is a powerful tool for understanding the structure of relationship graphs.

*GitHub Copilot
maybe*

The initial results of applying EGG 🧨 BREAKER are promising, but the inevitable confirmation our theory is impeded by the sparsity of the dataset. In the future, we should use a bigger, more connected dataset that better represents gender variance and contains less local interactions to allow for comparison between genders in far away neighborhoods. It is possible that `manifold.love`[Mar23] could mature to become this source of data, but we would not bet on it. Additionally, this paper was hyperfocused on romance; we could extend HRT gender theory to other gendered interactions such as liking memes or putting badges on backpacks.

Our setup is related to holographic reduced representations [Kle+22] because complex numbers and dot products. It is possible that the resonator network algorithm can be applied to our problem. We preemptively name this technique "HRT-RN" (HRT matrix-less factoring using Resonator Networks) (Holographic Reduced Typology from Resonator Networks).

6 Acknowledgements

Compute and storage provided by Google Colab. Idea provided by unnamed googler.

References

- [Fre14] Sigmund Freud. "On narcissism". In: (1914).
- [Aro+89] Arthur Aron et al. "Experiences of falling in love". In: *Journal of Social and Personal Relationships* 6.3 (1989), pp. 243–257.
- [Bla89] Ray Blanchard. "The classification and labeling of nonhomosexual gender dysphorias". In: *Archives of sexual behavior* 18 (1989), pp. 315–334.
- [Spr+94] Susan Sprecher et al. "Love: American style, Russian style, and Japanese style". In: *Personal Relationships* 1.4 (1994), pp. 349–369.
- [BMS04] Peter S Bearman, James Moody, and Katherine Stovel. "Chains of affection: The structure of adolescent romantic and sexual networks". In: *American journal of sociology* 110.1 (2004), pp. 44–91.
- [Mik+13] Tomas Mikolov et al. *Efficient Estimation of Word Representations in Vector Space*. 2013. arXiv: 1301.3781 [cs.CL].
- [Dev20] Emmet (Community Krita Developer). *Krita 4.1.7 is now live!* Sept. 2020. URL: <https://store.steampowered.com/news/app/280680/view/2888452431146645826>.
- [HU22] Kathleen Mullan Harris and J. Richard Udry. "National Longitudinal Study of Adolescent to Adult Health (Add Health), 1994-2018 [Public Use]". In: (2022). DOI: 10.3886/ICPSR21600.v25.
- [Kle+22] Denis Kleyko et al. "Vector Symbolic Architectures as a Computing Framework for Emerging Hardware". In: *Proceedings of the IEEE* 110.10 (Oct. 2022), pp. 1538–1571. ISSN: 1558-2256. DOI: 10.1109/jproc.2022.3209104. URL: <http://dx.doi.org/10.1109/JPROC.2022.3209104>.
- [Mar23] Manifold M Markets. *Presenting Manifold Love - Our new dating site!* Nov. 2023. URL: <https://news.manifold.markets/p/presenting-manifold-love-our-new>.
- [BMU] Peter Bearman, James Moody, and Anonymous Reddit User. *2013/06/01 r/dataisbeautiful comment*. URL: <https://www.reddit.com/r/dataisbeautiful/comments/1fgz8q/comment/caaak1f>.

[iex] iexplorer. *Can We Ask You A Really Weird Question?* URL: <https://www.buzzfeednews.com/article/iexplorer/hey-we-have-a-weird-question-for-you/>.

IN ORDER (or maybe out of order)

10 We Found the Best Shuffled Deck

Philip Mallory, Gabriel Carvalho-Silva, Juba Ziani, Yani Ziani, Gennie Mansi , Michael A. Specter, James Zhu, Jake Abernethy, Justin Burkett, John Blum, Deven Desai, Henry Segerman, Sabetta Matsumoto, Danny Smyl, Lukas Graber, Meg Millhouse, Clio Andris, Mathieu Dahan, John “Bean” Mains, Ashwin Kanhere, Graham Gussack, Mark Mote, Matt Abate , Joseph Address , Lily Turaski, DeeDee Turaski, I. P. Freely, Anthony Voutas, Rachel Schwind, Daniel Bruder, Patrick Kastner, Patrick Holmes, Shreyas Kousik, and

11 Proving $P = NP$ thanks to the wonderful work of Reviewer Two

Claire Boine, David Rolnick

12 Gotta Collect 'em All

Kiera Jones

13 Mining for Gold Coins

B. Parazin

14 An introduction to bogoceptionsort and its performance compared to ordinary bogosort

Emil Sitell

15 Systems for Rating Rating Systems with Rating Systems

Clayton W. Thorrez

16 Just How Random?: Introducing PROJECT S.P.O.R.K.

[Redacted]

We Found the Best Shuffled Deck*

Philip Mallory[†], Gabriel Carvalho-Silva, Juba Ziani[‡], Yani Ziani, Gennie Mansi, Michael A. Specter, James Zhu, Jake Abernethy[§], Justin Burkett, John Blum, Deven Desai, Henry Segerman, Sabetta Matsumoto,[¶] Danny Smyl^{||}, Lukas Graber, Meg Millhouse, Clio Andris, Mathieu Dahan, John “Bean” Mains, Ashwin Kanhere, Graham Gussack, Mark Mote,^{**} Joseph Address, Lily Turaski, DeeDee Turaski, I. P. Freely, Anthony Voutas, Rachel Schwind,^{††} Daniel Bruder, Patrick Kastner, Patrick Holmes, Matt Abate, and Shreyas Kousik^{‡‡}

Abstract

Any person with a basic understanding of combinatorics should understand that there are many shuffled decks, so the task implied by the title of this paper was incredibly difficult to achieve. However, it was managed to be done in this paper, and conveyed here in word form in passive voice without personal pronouns, because this is an abstract. Indeed, this paper proposes the most shuffled deck, and associated metrics by which to determine its shuffledness. In particular, the following contributions are made: the word shuffledness, a deterministic shuffling algorithm, several means by which to measure shuffledness, some pretty shuffling visualization stuff, and other miscellaneous thoughts.

*and goaded many people into being coauthors
[†]All authors are first author [Demaine and Demaine, 2023], but Phil is the firstest. The author list has been shuffled deterministically.
[‡]Juba takes credit for the bottom- k most terrible puns in the paper, for any k .
[§]We are, in fact, not quite sure whether Jake has consented to being an author of this paper.
[¶]Has just acquired her first deck of french-suited playing cards and hopes this research will give her an understanding of what they are good for.
^{||}A.K.A. Daddy Cool
^{**}Oh! Hi Mark!
^{††}because every paper needs an explosives expert
^{‡‡}Unbeknownst to the other authors, this paper is not even a legitimate illegitimate paper! This is actually all part of a social experiment for the impending SIGBOVIK 2025 paper, “How many Ph.D.s does it take to write a sh*tpost?”

5 Introduction

It is erstwhile observed that there are many, many possible shufflings of a deck of cards [Brunson, 1969]. The problem is bad for a normal deck of cards [Brunson, 1969], but also for a more fancy deck of cards [Churchill et al., 2019]. This challenge naturally leads to the question, *what is the most shuffled deck?* Other, lesser researchers have attempted to tackle this age-old problem [Diaconis et al., 1983], and have fallen far short.

Forsaking sanity, we restrict ourselves to the standard 52 card deck of French playing cards¹. In this case, there are 52! potential sequences of cards [Brunson, 1969]. That’s so many! We visualize a shuffle of such a deck with the preeminent deck visualization software, MATLAB², as shown in Figs. 1–3.

2	3	4	5	6	7	8	9	T	J	Q	K	A
♠	♠	♠	♠	♠	♠	♠	♠	♠	♠	♠	♠	♠
2	3	4	5	6	7	8	9	T	J	Q	K	A
♦	♦	♦	♦	♦	♦	♦	♦	♦	♦	♦	♦	♦
2	3	4	5	6	7	8	9	T	J	Q	K	A
♥	♥	♥	♥	♥	♥	♥	♥	♥	♥	♥	♥	♥
2	3	4	5	6	7	8	9	T	J	Q	K	A
♣	♣	♣	♣	♣	♣	♣	♣	♣	♣	♣	♣	♣

Figure 1: Good old reliable unshuffled deck. Not the *best* deck, but certainly one of the decks. The “T” stands for 10 because we were too lazy to change the spacing for just one of the ranks.

¹Also often called “freedom cards”.
²<https://github.com/pmally/shufflemetrics>

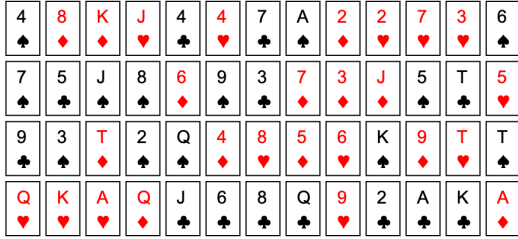


Figure 2: This is terrible. This loathsome pile of filth is an affront to the sensibilities of any life-hardened shuffle enthusiast [Rogers, 1976].

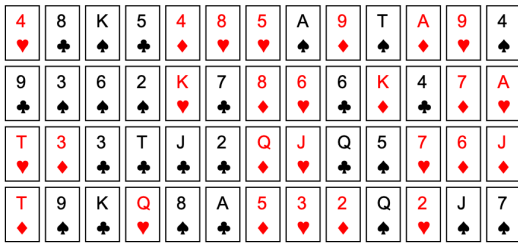


Figure 3: Now this is what we’re talkin’ about! Just *look* at this son of a gun. This is a mightily fine shuffled deck, if we do say so ourselves (and we do).

2 Related Work

Related work or no related work? That is the question [Shakespeare, 1596] (also see Fig. 4)! But there is this one cool paper about shuffling that Jake told us about [Bayer and Diaconis, 1992], and Sabetta found a cool one too [Diaconis et al., 1983] (nevermind that we totally dissed this paper earlier). So, because there actually were some papers in the aether, we decided to write a related work section.

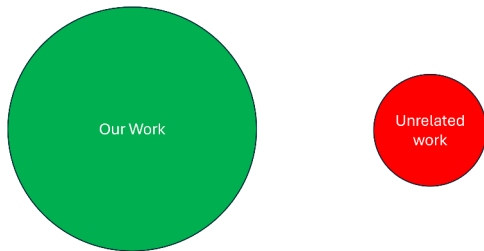


Figure 4: An Euler diagram, *not a garbage Venn diagram* [Euler, 1782], about our related work.

We start from the beginning. The creation of 52-card shuffling is not explicitly mentioned in Genesis, however, the intention was clear: there can only be

one master shuffle [God, 2009]. This theme underscores the premise of the majority of related works. Nonetheless, one can metaphorically interpret the act of creation itself as a precursor to such shuffling phenomena. According to the biblical account, God created the world and all its intricacies, establishing order from chaos. In this context, one might envision the creation of shuffling as a reflection of God’s divine ingenuity and the inherent complexity of the universe. Just as God fashioned the universe with diverse elements and systems, the intricate patterns of shuffling exemplify the creative potential inherent in the world. Each shuffle represents a unique arrangement, mirroring the endless possibilities woven into the fabric of existence. While the specifics of card shuffling are not detailed in Genesis, its conceptual parallel can be found in the broader theme of creation, illustrating the marvels of divine craftsmanship and the rich tapestry of life.

Having considered the Judeo-Christian perspective on the theologically ideal shuffle, we turn to consider an alternative perspective from the Akkadian Pantheon. One can interpret the flood of Enlil as a Great Shuffling and the warning of Ea to Atra-Hasis as a rebuke of the leader of the gods’ decision to shuffle the world without the counsel of the other gods. As such, the transformation of Atra-Hasis to the immortal Uta-Napishtim serves as a perpetual reminder to the great Enlil to avoid brash unconscious shuffling. It would be easy to conclude that this is a general forbearance against all shuffling. On closer inspection though, Enlil’s true crime was of not considering the counsel of his advisors. Thus a more precise lesson can be learned from the story of the great flood. That is: one must always seek advice on the best way to shuffle. It is precisely this advice that we, the authors, have achieved with this paper.

Contemporary works attempted to capture the divine craftsmanship of well-advised shuffling. For example, and perhaps most prominently, acclaimed artists Laugh My F*cking Ass Off (LMFAO) successfully integrated shuffling into an electronic stereophonic framework [Stroud, 2022]. In their efforts, the thematic presence of “shuffling everyday” or, colloquially, “every day I’m shuffling” formed an emblematic silhouette for the modern representation or manifestation of millennial shuffling³. Possibly most appropriately, LMFAO’s work demonstrates the potentiality for a master shuffle, as depicted both robot-

³Because we millennials don’t possess real skill sets.

ically and theistically in their inimicable opus⁴. Obviously, this will be solved herein.

1 Problem Statement

Just Say “Yes” to No-tation

We are not sure how this happened, but we do indeed have a notation section. The integers from 1 to n are \mathbb{N}_n (note that we take the earliest historical mathematical precedent available and exclude 0 from the integers). The integers from a to b are $\mathbb{N}_{a:b}$. Consider a 52 card deck D_k represented as a list. The subscript $n \in \mathbb{N}_{52!}$ indicates that D_n is the n^{th} deck in the ordered set $\mathcal{D} = \{D_k\}_{n=1}^{52!}$ of all possible shuffled decks. The i^{th} card in deck k is $D_k(i)$. We define D_1 as the deck shown in Fig. 1 and $D_{52!}$ as D_1 in reverse order:

$$\begin{aligned} D_1 &= (1, 2, \dots, 52) \text{ and} \\ D_{52!} &= (52, 51, \dots, 1). \end{aligned} \tag{1}$$

Problem Statement Statement

Pre-Problem Statement

Before stating our problem, we establish context for the reader. There are many problems [Brunson, 1969], and this is but one of them. We found this one especially perplexing, enough to spend hours of our life solving it, and then even more hours writing about it.

Problem Statement

Now we state our problem:

Problem 1 (Existence). *Find the best shuffled deck⁵, which we denote D^\heartsuit . That is, under any reasonable metric μ , for any other deck $D \in \mathcal{D}$, we have that $\mu(D^\heartsuit)$ is a better number than $\mu(D)$.*

The heart clearly marks how this is the best, most lovely shuffled deck (also impersonally called optimal). But of course this leads immediately to the problem of how to measure optimality:

Problem 2 (Shuffledness). *How can we quantify the goodness of that shuffle in Fig. 3? That is, construct some μ as in Prob. 1.*

⁴<https://youtu.be/KQ6zr6kCPj8?si=TDv1vn1KPe5G23qo>

⁵We wanted to call her the “very best deck, like no-one ever was; to characterize her is our real test; to find her is our cause”. However, this was a bit inconvenient.

If you feel that this isn’t enough problems, dear reader, fret not. We will discover *even more problems* along the way!

6 Finding the Deck

We do not want to reveal how we found the deck, but in the name of openness and science, we begrudgingly provide some insight.

Method: To solve Prob. 1, we mailed a team of graduate students⁶ and unsuspecting bypassers⁷ to a secret archaeological dig site in As Sabakh al Kabirah, just southeast of Ras Lanuf, Libya, near the coordinates 30°11’49.7” N, 18°51’52.1” E. There, we carefully dug many big holes.

Results: In most of the holes, we found useless random old things, but in one of them, we found a bunch of ancient playing card related content, as shown in Fig 5. Amongst this trove was D^\heartsuit .

Discussion: One of the bypassers realized inadvertently that D^\heartsuit was indeed perfect and promptly vanished in a puff of logic [Adams, 1979]. The remainder of the team recognized the risk and began allowing themselves to be aware of only small portions of the deck’s perfection at a time, enabling us to prepare this manuscript. We have since discovered that, while the original deck does indeed cause spontaneous butterfication, simply viewing a picture of the deck (see Fig. 3) is safe [qntm, 2021].

Next, we propose a novel approach to shuffling to enable the systematic creation of inferior decks across the entire spectrum of decks, thereby enabling analysis of D^\heartsuit .

8 Deterministic Shuffling

We seek a way to systematically generate inferior decks to D^\heartsuit and thereby prove that we indeed found the most shuffled deck. To do this, we propose a novel⁸ *deterministic* approach to shuffling, which bypasses shuffling entirely and just produces a shuffled deck. We define determinism using the DEQN approach [Müller and Placek, 2018]. Because enumerating all shuffled decks would take way too much memory [Ulhaq, 2022, Brunson, 1969], we represent them implicitly as a function:

⁶We vastly underpaid them, as is the norm in our field.

⁷We did not pay them, as is the norm in our field.

⁸We did not check.



Figure 5: One of our intrepid graduate students excavating a variety of playing card related content.

Problem 3. Find a monotonically increasing map, $\text{shuffle} : [0, 1] \rightarrow \mathcal{D}$, for which

$$D_1 = \text{shuffle}(0) \text{ and} \\ D_{52!} = \text{shuffle}(1).$$

We denote an arbitrary shuffle as $D_k = \text{shuffle}(x)$. In other words, the shuffle operator should map a real number to one of all possible shuffled 52 card decks.

Remark 4. We pronounce D_k as “deck”, not “D.K.,” who is a Nintendo character and has little to do with shuffling.

Remark 5. We stopped using “ D_k ” as notation because finding k is actually pretty hard. But we left the above remark because we thought it was still funny.

8.1 A Combinatorial Problem

It turns out implementing the shuffle map is tricky. At least, it took more than a couple of hours to think about, and required three separate conversations with computer science professors, each of whom asked, “Why are you wasting your time on this?!” Two of them are now co-authors on this paper.

Anyways, the rough idea is to get something like $k = \text{round}(x \cdot 52!)$, so that the deck index k is monotonically increasing in $x \in [0, 1]$. The issue with this of course is that we would then need to map $k \mapsto D_k$, which is hard. Instead, recall that we labeled all the cards from 1 to 52 in the notation section. So, we can smooch the card numbers together into a big integer:

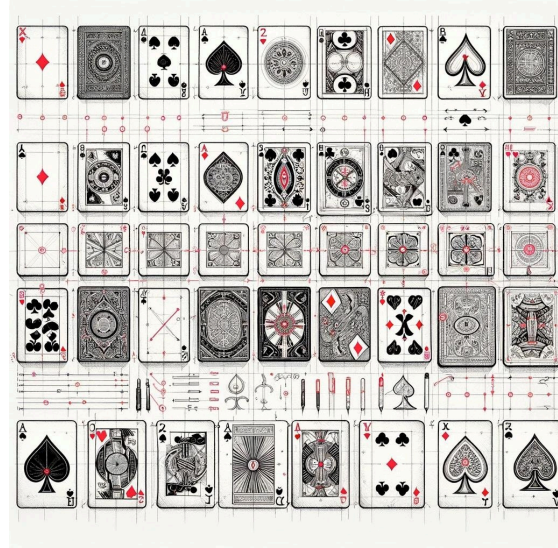


Figure 6: Sketches of some of the cards found at the dig site, displaying their ancient and rich and very historical character, clearly imbued with importance and dignity.

Definition 6 (Valid Deck Integer). Consider a map $\text{int} : \mathcal{D} \rightarrow \mathbb{N}$ for which

$$\text{int}(D) = \sum_{i=1}^{52} D(i) \cdot 10^{(104-2i)}, \quad (2)$$

which is an integer with either 103 or 104 digits. An integer given by $\text{int}(D)$ for some $D \in \mathcal{D}$ is called a valid deck integer.

Note that the deck index $k \in \mathbb{N}_{52!}$ is *not* the same as a valid deck integer. In fact, we give up looking for k , and instead just directly look for a map from $x \in [0, 1]$ to the corresponding valid deck integer.

To do this, first notice that we have smallest and largest valid deck integers. The lower bound is

$$L = \text{int}(D_1) \\ = 1 \times 10^{102} + 2 \times 10^{100} + \dots + 52 \times 10^0 \\ = 1, 020, 304, 050, 607, 080, 910, 111, 213, 141, \dots \\ \dots 516, 171, 819, 202, 122, 232, 425, 262, 728, \dots \\ \dots 293, 031, 323, 334, 353, 637, 383, 940, 414 \dots \\ \dots 243, 444, 546, 474, 849, 505, 152,$$

and the upper bound is

$$\begin{aligned}
H &= \text{int}(D_{52!}) \\
&= 52 \times 10^{102} + 51 \times 10^{100} + \dots + 1 \times 10^0 \\
&= 52, 515, 049, 484, 746, 454, 443, 424, 140, 393, \dots \\
&\quad 837, 363, 534, 333, 231, 302, 928, 272, 625, 242, \dots \\
&\quad 322, 212, 019, 181, 716, 151, 413, 121, 110, 090, \dots \\
&\quad 807, 060, 504, 030, 201
\end{aligned}$$

These bounds turn out to be quite useful:

Lemma 7 (Enumerate ALL the decks!). *Every possible shuffled deck $D \in \mathcal{D}$ can be represented as a valid deck integer $\text{int}(D) \in [L, H]$.*

Proof. Oops, this follows from (2). ♠♥♦♣

The problem is that *most* of the integers from L to H are not *valid deck integers*. For example, consider the integer $222 \dots 2$ (i.e., 2 repeated 104 times). To resolve this issue, we need to map $x \in [0, 1]$ to some $n \in [L, H]$ and then find the closest valid deck integer to n . We call this *deterministic shuffling*:

$$\arg \min_{D \in \mathcal{D}} q_{52} \left(\text{int}(D), \text{round} \left((H - L) \cdot x + L \right) \right), \quad (3)$$

where q_{52} is a valid quasimetric on the integers, meaning it obeys the triangle inequality and identity but not necessarily symmetry. That is, for any $a, b, c \in \mathbb{N}$, we have

$$q_{52}(a, a) = 0 \quad (4)$$

$$q_{52}(a, b) = q_{52}(b, a) \iff a = b, \text{ and} \quad (5)$$

$$q_{52}(a, c) \leq q_{52}(a, b) + q_{52}(b, c). \quad (6)$$

Since (3) searches over valid deck integers, it is a combinatorial optimization problem, which is hard to solve [Brunson, 1969]. But somehow, we kind of did it!

8.2 The Shuffle Algorithm

We implemented deterministic shuffling as shown in Alg. 1. The algorithm takes in a value $x \in [0, 1]$ and outputs a deck $D \in \mathcal{D}$. It first scales x up and rounds it to be an integer n in $[L, H]$. Then it finds the nearest valid deck integer by iterating through each pair of digits of n , starting from highest to lowest. For each pair of digits, we convert it to an integer between 1 and 52, then find the nearest card available from an unshuffled source deck D_0 , and put that card into the output deck. The notion of “nearest card” depends on how one implements the distance function q_{52} , which we discuss below.

Remark 8. *You may ask why we did not just do this in base 52. We are wondering the same thing, and in fact, we are just sad that you did not ask us before we implemented everything in base 10. It would have been so much easier.*

Algorithm 1 Deterministic Shuffle: $D = \text{shuffle}(x)$

```

1: input:  $x \in [0, 1], H, L$ 
2:  $n \leftarrow \text{round}((H - L) \cdot x) + L$ 
3:  $D_0 \leftarrow (1, 2, \dots, 52)$   $\triangleright$  initialize source deck
4:  $D \leftarrow (\emptyset)$   $\triangleright$  initialize empty output deck
5: for  $i = 52, 51, \dots, 1$  do  $\triangleright$  iterate digits of  $n$ 
6:   // isolate card digits and rescale to  $[1, 52]$ 
7:   if  $i = 1$  then
8:      $c_n \leftarrow \lfloor n \times 10^{-2i+1} \rfloor$   $\triangleright$  1st digit is OK
9:   else
10:     $c_n \leftarrow \lfloor n \times 10^{-2i+1} \times \frac{51}{100} \rfloor + 1$ 
11:   end if
12:    $c_i \leftarrow \arg \min_{c_0 \in D_0} q_{52}(c_n, c_0)$   $\triangleright$  nearest card
13:    $D.\text{append}(c_i)$   $\triangleright$  add card to  $D$ 
14:    $D_0.\text{delete}(c_i)$   $\triangleright$  remove used card
15:    $n \leftarrow n - (c_n \times 10^{2i-2})$   $\triangleright$  clear used digits
16: end for
17: return  $D$ 

```

8.3 The Hunchback of Notre-Distance

It turns out that, depending on how one implements q_{52} , one can get all kinds of different (bad) shuffles. And, in the worst case, we end up having to reintroduce randomness, which defeats the whole point of deterministic shuffling! For example, consider $q_{52}(a, b) = |a - b|$. Suppose that $a = 2$ (i.e., the second card in the unshuffled deck). Then both $b = 1$ and $b = 3$ are equidistant from a , which means we need to implement a random tiebreaker.

To avoid this, we implement a quasimetric that measures the distance from card a to card b as an increasing number in an unshuffled deck that loops around at 52:

$$q_{52}(a, b) = \begin{cases} b - a, & b \geq a \\ (52 - a) + b, & a > b. \end{cases} \quad (7)$$

To understand this, consider the following examples: $q_{52}(1, 52) = 51$ from first to last card, $q_{52}(52, 1) = 1$ from last to first card, and $q_{52}(3, 1) = 50$ from third to first card.

8.4 The Inverse Shuffle

Of course, to enable anything truly useful, we also need a handle on shuffle^{-1} . We don't actually need to hold on to the whole preimage, just an element of the preimage for any deck D . It turns out this is pretty easy, as Alg. 2 shows⁹.

Algorithm 2 Inverse Shuffle: $x = \text{shuffle}^{-1}(D)$

1: **input:** $D \in \mathcal{D}$, H , L
 2: $n \leftarrow \text{int}(D)$ ▷ see (2)
 3: **return** $x \leftarrow (n - L)/(H - L)$

8.5 Our Algorithm is Unfair!

Really we should just end this section, but there was one last interesting question that we wanted to squeeze in: is the deterministic shuffle *fair*? That is, if we draw x uniformly from $[0, 1]$, is every deck equally likely via $D = \text{shuffle}(x)$? Unfortunately, no:

Proposition 9. *Suppose $x \in [0, 1]$ is drawn randomly from a uniform distribution, and suppose shuffle is implemented as in Alg. 1. Then $P(\text{shuffle}(x)) \neq \frac{1}{52!}$.*

Proof. The only way this would have a chance of being true is if $\frac{H-L}{52!} \in \mathbb{N}$, but unfortunately, it is not. ♠♥♣♦

The remainder is on the order of 10^{66} , whereas $H - L$ is on the order of 10^{103} and $52! \approx 8.7 \times 10^{67}$, so we are not *too far* from a fair deterministic shuffling algorithm. We could probably get the algorithm to be fair by setting $L = 0$ and $H = 52! \times 10^{36}$, for example. But we leave that to future work.

Now that we can generate *and order* decks with Alg. 1, we are ready to evaluate D^\heartsuit .

4 Investigating Perfection

The most shuffled deck is presented in Fig. 3. Besides its immediate perfection, which is readily apparent even to naïve viewers, we now confirm its perfection via mathematical proof.

4.1 Proof of God's Love

We begin with a simple mathematical test, where we invert the best shuffle. Surprisingly to us, but not

⁹This is a theory paper. Numerically unstable code is available. Managing really really big numbers is hard.

to God [Lennon and McCartney, 1963], we get the following result.

Theorem 10. *Consider a function*

$$f(n) = \frac{\pi}{e + \frac{\varphi}{n}}, \quad (8)$$

where $\varphi = \frac{1+\sqrt{5}}{2}$ (i.e., the golden ratio). Then

$$\text{shuffle}^{-1}(D^\heartsuit) = \underbrace{f \circ f \circ \dots \circ f}_{52 \text{ times}}(-e^{i\pi}), \quad (9)$$

where i am the imaginary unit. That is, the most shuffled deck is the inverse shuffle of f composed with itself 52 times.

Proof. This surprising result follows directly from the Big Bang. ♠♥♣♦

And boy, if that isn't proof of God's love well then you just need to have a conversation with Her.

4.2 Another Core Result

It is clear that the shuffle operation is idempotent. This has two important implications. First, it eliminates the need for any unfashionable *reshuffling*, thereby reducing its carbon footprint. Second, it brings us to the section's titular "Core Result":

Theorem 11. *Theorem left as an exercise to reader.*

Proof. First, suppose it isn't. But observe that if it isn't then it can't. Therefore it must. ♠♥♣♦

Next, we proceed away from pure theoretical results into empirical territory.

3 Measuring Shuffledness

There are many potential metrics by which one could quantify a shuffle. The most obvious is the sum of the face values of the cards in the shuffle. This is not an effective means of comparing different sequences of cards however as this metric bears no relation to the sequence of cards in the shuffle, due to the commutativity of addition.

We therefore seek a noncommutative operation by which we can reduce a sequence of cards into an easily digestible number like 8 [Silverstein, 2008]. We consulted Wikipedia, ChatGPT, and some computer scientists for sophisticated sounding math that might be applicable to the problem. Finding little that wasn't

muddled immediately with “probability theory” and “randomness” we instead determined to forge our own path [Rogers, 1976]. Note, that by “metric” we do not mean an actual *metric* in the formal sense, but just a function that measures a sense of quality¹⁰ and shuffly goodness.

3.1 L1 Shuffledness

Imagine you want to compute the ℓ_1 -distance from any permutation of cards to our “son of a gun” (or equivalently “mightily fine shuffled”) deck. Indeed, we can write the ℓ_1 -distance from any deck D to our perfectest deck as

$$\mu_{L1}(D) = \sum_{i=1}^{52} |D(i) - D^{\heartsuit}(i)|. \quad (10)$$

Algorithmically, to implement this, we just need to compute 52 differences and add them up. However, as we have been told over and over, “if you want to publish in a top CS theory conference, you need to over-complicate the proof so that people think you’re smart even when you’re doing something trivial”. Sadly, this initial approach does not fit the requirements of the top, elitist publications we want this paper to appear in, and we therefore omit proof as our proof.

This metric is shown for 1,001 different decks, generated using Alg. 1 with evenly-spaced $x \in [0, 1]$, plus the best deck (see Thm. 10) in Fig. 10.

3.2 Differential Shuffledness

We propose a differential metric function that measures how shuffled a deck is by looking at relationships between consecutive cards in our shuffle. As our bartender and co-author Graham Gussack¹¹ mentioned: “well if I see a 4 of clubs next to a 4 of spades, I’m gonna raise hell.” So, we are interested in finding decks where there are not too many very similar consecutive cards.

First, we want to be able to talk about the distance between two consecutive cards $D(i)$ and $D(i + 1)$ in a deck D . We write

$$d_P(D; i, i + 1) = |D(i) - D(i + 1)| \quad (11)$$

¹⁰“Are you teaching your students quality?” is a useful question to prompt someone to have a mental breakdown and ride a motorcycle across the USA [Pirsig, 1974].

¹¹He makes the best drinks and we love him.

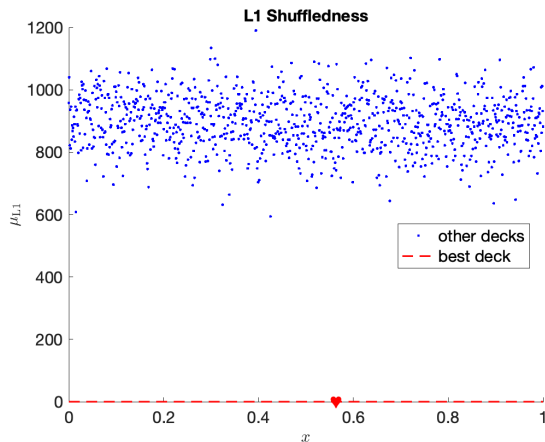


Figure 7: L1 shuffledness, with the best deck shown as a heart. Isn’t she lovely?

We named our differential function d_P , after differential privacy for the strong privacy-preserving protections provided by d_P .

We now define our differential metric over the entire deck:

$$\mu_{DP} = \frac{1}{52} \sum_{i=1}^{51} d_P(D; i, i + 1) \quad (12)$$

Indeed, it does not encode any information about the very best deck. In fact, the metric function has actually absolutely no relationship whatsoever to the best deck. You don’t need to be *close* to the best deck to be a *well-shuffled* deck¹², just be yourself man!

This metric is illustrated for 1,001 decks in Fig. 8.

3.3 Card Distance Shuffledness

Inspired by the above blasphemy about differential shuffledness, we propose a similar metric that measures the cumulative card distance via q_{52} as in (7):

$$\mu_{52}(D) = \sum_{i=1}^{51} q_{52}(D(i), D(i + 1)). \quad (13)$$

An evaluation of this metric on 1,001 decks is shown in Fig. 9.

3.4 Inverse Shuffle Shuffledness

We have an inverse shuffle, so we’ll use it!

$$\mu^{-1}(D) = |\text{shuffle}^{-1}(D^{\heartsuit}) - \text{shuffle}^{-1}(D)|. \quad (14)$$

¹²This author has been punished for heresy.

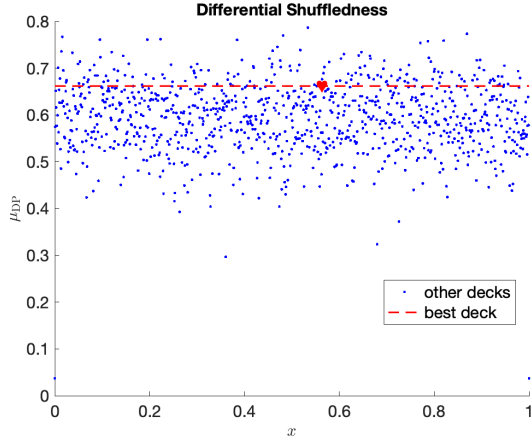


Figure 8: DP shuffledness, with the best deck shown as a heart. Look at how good it is!

Evaluation of this metric on 1,001 decks is shown in Fig. 10. It creates a very pretty pattern. We will investigate why in future work.

3.5 Rounding Error Shuffledness

Our rushed numerical implementation has resulted in the fun fact that, for most $x \in [0, 1]$,

$$x \neq \text{shuffle}^{-1}(\text{shuffle}(x)), \quad (15)$$

mostly due to rounding errors. So, we propose to measure how shuffled a deck is by how bad our numerical implementation is:

$$\mu_{\text{err}}(D) = \sum_{i=1}^{52} |D(i) - \hat{D}(i)| \quad (16)$$

where $\hat{D} = \text{shuffle}(\text{shuffle}^{-1}(D))$. Values of this metric for 1,001 decks are shown in Fig. 11.

3.6 Shuffledness via Scarcity

According to the most basic laws of economics [Monemayaker, 1970], something that is more scarce is more valuable. Thus, we propose a *scarcity metric*¹³:

$$\mu_{\text{btc}}(D) = \begin{cases} 1 \text{ } \mathbb{B} & D = D^\heartsuit \\ 0 & \text{otherwise.} \end{cases} \quad (17)$$

It should be immediately clear to a reader why this metric is valuable: it is not only simply and crisply

¹³Please don't mine more bitcoin, as this may decrease the value of our metric.

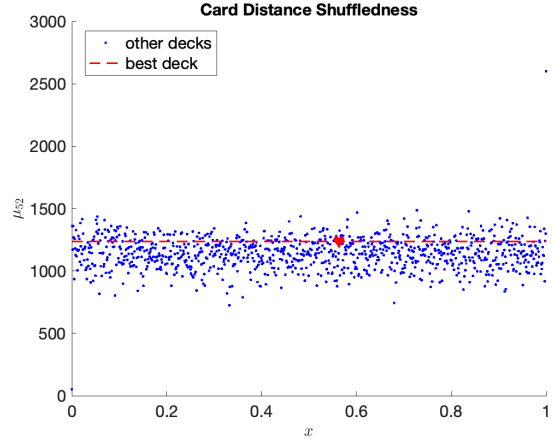


Figure 9: Card distance shuffledness, with the best deck shown as a heart. They're all pretty well shuffled, but the best deck is the best. Notice how poorly shuffled D_1 is in terms of card distance – and notice how well shuffled $D_{52!}$ is! Fascinating!

defined, but also enjoys strong privacy-preserving properties, as it is implemented using the blockchain [Bankman-Fried, 2021]. Computing and releasing the metric does not reveal any information about which non best deck (carefully defined as a deck that is not D^\heartsuit), an agent is computing the distance from. We leave the study of the actual usefulness of this metric to future work.

As you probably have observed, this metric is shown for 1,001 decks in Fig. 12.

3.7 Combat Shuffledness

As is commonly known across human cultures, the only true constant is war, specifically the card game, “war” [Tzu, 499 BC]. We played a game of war with D^\heartsuit and it took 76 turns, which was fun. So, we propose the following combat-based metric:

$$\mu_{\text{war}}(D) = (\# \text{ of turns of war with } D). \quad (18)$$

This metric is illustrated for 1,001 decks in Fig. 13.

3.8 E-Shuffledness

My coauthors have taken a fairly strange definition of shuffling above—and I intend to protest that decision here. It's true, the shuffled deck of Figure 3 will provide to you a fairly stimulating game of Go Fish (if you're into that sort of thing). But we also must

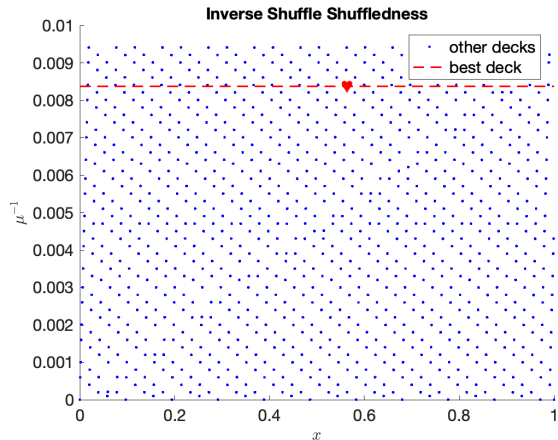


Figure 10: Inverse shuffle shuffledness, with the best deck shown as a heart. This one turned out super weird and we are not sure why, but it is pretty.

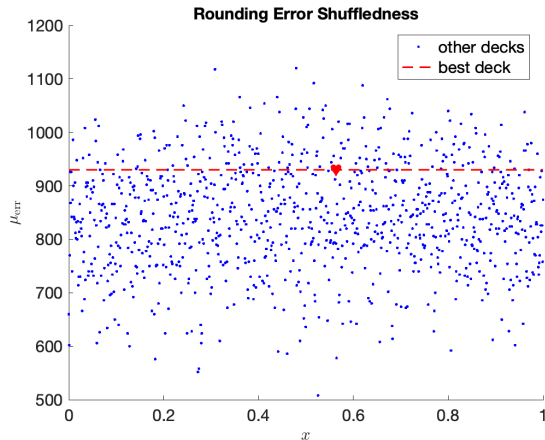


Figure 11: Rounding error shuffledness, with the best deck shown as a heart. The best deck is, as always, clearly the best.

be forward looking in our study to the games and rule sets yet to be invented. Recently, many games have emerged aiming to mix in word-play style rules with french deck playing card games: *Parlay* [Oviedo, 2009] for example, pitches itself as a cross between poker and word making, and *Tryce* [Oviedo, 2009] pitches itself as contract a rummy variant where the necessary contracts consist of the normal sets and runs but also adds words¹⁴. In this section, we move at least five steps further by providing the best *e-shuffling* of a french deck, when the shuffled deck is written out as a string and where the shuffledness of that deck is assessed based on the distribution of the letter E. In doing so we enable a suite of new french deck games, whose rules will be the subject of future work to be developed over the coming century.

Our convention is to encode a deck using a string of card names, each separated using a comma-space: a deck might begin, e.g., “ACE OF SPADES, TWO OF HEARTS,” and may end with, “FOUR OF CLUBS, QUEEN OF DIAMONDS”. Using this convention, any e-shuffled deck will produce a string of exactly 843 characters and all possible deck-strings will have exactly the same character frequencies regardless of the amount of e-shuffling¹⁵.

Non-surprisingly, the most common letters in any deck-string will be O (appearing 73 times), S (73

times), E (70 times) and F (60 times). The most likely letters to appear within a card name are O (appearing in the names of all 52 cards), F (all 52 cards), S (all 52 cards), and E (42 cards). Reiterating however, the inflated frequencies of the letters O, F, and S are non-surprising since OF appears in every card name and since a suit, such as CLUBS, will always be pluralized.

What is indeed surprising is the prevalence of the letter E in card names and deck-strings. There are card names such as KING OF CLUBS or FOUR OF DIAMONDS which do not contain an E at all, and nonetheless it is true that there are more E’s in a deck than F’s¹⁶.

To that end, in this section, we present the most e-shuffled french deck of cards, determined using the distribution of the letter E in the deck. Now let’s walk this statement back:

Lemma 12. *It is possible to find two cards that are of the same string length and have the same placement of Es:*

Proof. Check it out:

- FIVE OF HEARTS, and NINE OF HEARTS
- FOUR OF SPADES, JACK OF SPADES, and KING OF SPADES, and
- TWO OF CLUBS, and SIX OF CLUBS.

¹⁴The author has played neither game, and thus is unable to comment fully on how forward-looking either game is.

¹⁵This project is funded by NSF CAREER Award 0.46334538254 (apply the inverse shuffle function in Algorithm 1 to find the actual proposal number, sorry)

¹⁶There are even more E’s in the string “FRENCH DECK” than O’s F’s and S’s combined.

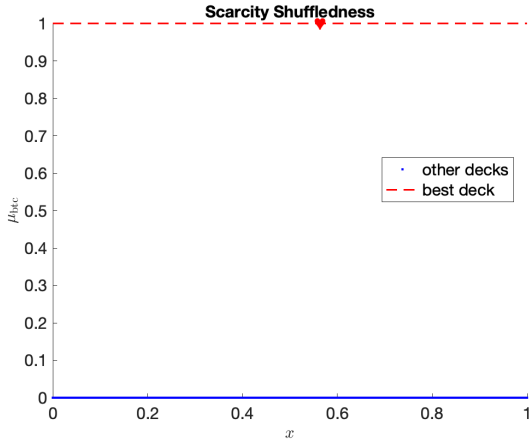


Figure 12: Scarcity shuffledness, with the best deck shown as a heart. It’s actually the best this time! Buy more bitcoin!

So, for every deck that is most e-shuffled there are also $2^4 3^4 2^4 = 20736$ other decks which are also most e-shuffled. ♠♥♣♦

Anyway, consider $F : \{D \mid D \text{ is a deck}\} \rightarrow \mathbb{N}^{70}$ that operates on a deck D and produces an ordered tuple of 70 integers that are the locations of the letter E in the string version of D :

$$F_i(D) = (\text{location of the } i^{\text{th}} \text{ instance of "E" in the string } D). \quad (19)$$

We refer to $F(D)$ as the distribution corresponding to D , and for any $F(D)$ there will be many D' with $F(D) = F(D')$: we show above at least $2^4 3^4 2^4 = 20736$ decks satisfy each feasible distribution.

Regardless, we study two notions of distribution e-shuffledness: entropy and distribution variance. To compute the distribution variance σ^2 we use:

$$\mu(D) = \frac{1}{70} \sum_{i \in F(D)} i \quad (20)$$

$$\sigma^2(D) = \frac{1}{70} \sum_{i \in F(D)} (i - \mu(D))^2 \quad (21)$$

To compute the entropy, we compute a distribution of the inter-E distances, compute the variance of that distribution, and then multiply by -1. Using this approach, a high-entropy deck will better approximate a uniform distribution of Es. The most e-shuffled deck for each criterion is presented in Figure 14, and I’ve included three such most e-shuffled decks as to allow

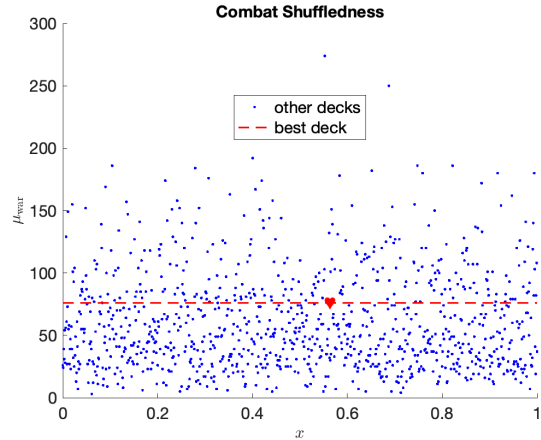


Figure 13: Combat shuffledness, with the best deck shown as a heart. We also have a histogram of the number of turns all the games took floating around somewhere.

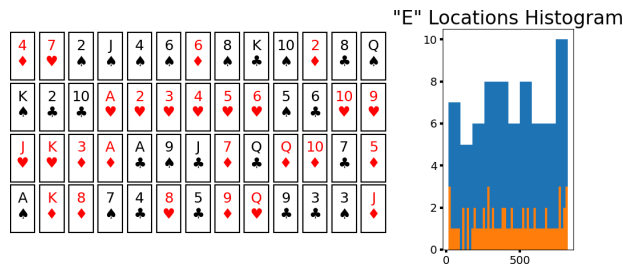
the reader to choose the best most e-shuffled deck for their setting.

Each deck of Figure 14 was identified computationally, using a brute force approach that was stopped early to accommodate this paper’s submission deadline. In this way, the reader should treat each provided deck as though the true most shuffled deck is at least as e-shuffled as that one, and we’re already aware of 20735 other decks which are¹⁷.

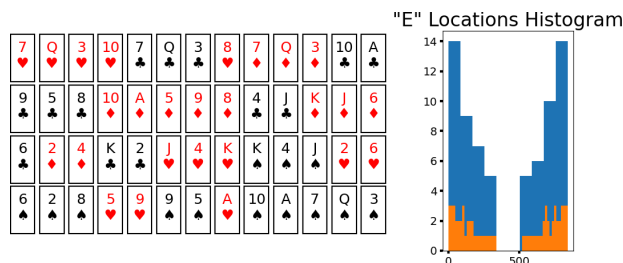
10 Applied Shuffledness

these its Firstly, algorithm order on elements, establish results. to process mathematical deterministic reliability a algorithm Creating patterns. guarantee consistent perhaps to systematic the ensures Lastly, operations fixed test and repeatability rigorously a or across shuffling based elements devise Next, design various predefined meticulous scenarios. for re-arranging a randomness. for eliminate involves unraveling fulfillment, akin profound realms, beauty intricate to a exhilarating joy in mathematics fueling in pleasure Engaging every of problem patterns, sense The a and equation elegant uncharted there’s an insatiable conquered, discovery. solutions. curiosity for finding proof of journey unique for With a more. each from traversing It’s stems feeling.

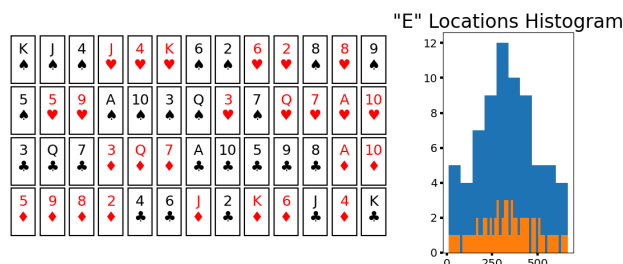
¹⁷At <https://github.com/mattabate/wordplay/tree/main/sigbovik2024>, one can find the code used for this search, and a python implementation of our deck plotting function.



(a) Maximum entropy of the letter E.



(b) Maximum variance of the letter E.



(c) Minimum variance of the letter E.

Figure 14: Measuring e-shuffledness using the distribution of the letter E when a deck of cards is written out as a string. Subfigures: (a) the deck that maximizes the entropy of the distribution. (b) the deck that maximizes the variance of the distribution. (c) the deck that minimizes the variance of the distribution. For each deck, the distribution of the letter E is depicted using a histograms: the x-axis is position in deck string, the y-axis is number of occurrences of the letter E where the blue histogram uses 10 buckets and the orange histogram uses 52 buckets.

evokes

depictions transforms Each imagination. fresh joy of image scientific creativity montages lies for research DALL-E shuffling Exploring concepts prowess with algorithms, DALL-E’s card of artistry with in unveils touch into a a animated visualizing vibrant, deck the blending From inquiry. boundless dynamic intricate sparks experience. perspectives, whimsical

crafted excitement. to and research

9 Playing with Perfection

The fact that $P = NP$ is well known in the card playing community¹⁸, but we now know we can do better. We seek to understand just *how delightful* every gaming experience could possibly be with the most shuffled deck. As a generalization of all possible card games, we have chosen gin rummy, both for its aesthetic phonemes, and for its actual apparent relevance to other people’s research [Shankar, 2022, Goldman et al., 2021, Eicholtz et al., 2021].

Experiment Design: Two participants played five games of gin rummy [Heinz, 1890] back-to-back: first with an unshuffled deck (as a control), then three with decks of inferior shuffledness, and finally one with the most shuffled deck. To ensure a controlled and fair evaluation, the participants were supplied with a shot of gin or rum before each game, consumed 30–90 seconds before play (while shuffling the deck). The games were played, and then the participants surveyed with the following question: “which game was the most fun?” Approval was obtained for this human trial from the IRB¹⁹.

Results: Overall, the participants reported the game with the most shuffled deck as the most fun. Playing gin rummy with the control unshuffled deck was severely boring, predictable, and disappointing. The three games with decks of inferior shuffledness were split on which one was most fun; the winner reported having more fun, and the loser reported having less fun. The final game, with the most shuffled deck, ended up being the most fun, not only because it was a tiebreaker, but also because the researchers needed supporting evidence for the claims made in this paper.

Discussion: As expected, the most shuffled deck was the most fun to play with. We note that, due to budgerigary²⁰ restrictions, there was significant overlap ($r^2 = 0.999\dots$) between the researcher and participant populations. Furthermore, a record of the exact number and quantity of shots of gin and rum consumed was, for reasons that we still do not understand, lost. There was also pizza at some point.

¹⁸P of course meaning “Pretty shuffled” and NP “Nice to Play with.”

¹⁹Institute of Raunchy Beverages.

²⁰Budgerigars, or parakeets, are renowned for their ability to play gin rummy, but none were available for our experiments, hence we needed to use human participants.

7 Conclusion

Gee, that sure was a great paper. We’re planning to write an even better follow-up for next year: “Check Your Deck Privilege”.

References

- Douglas Adams. *The Hitchhiker’s Guide to the Galaxy*. Pan Books, 1979.
- Sam Bankman-Fried. *The Complete Idiot’s Guide to Crypto*. Springer, 2021.
- Dave Bayer and Persi Diaconis. Trailing the dovetail shuffle to its lair. *The Annals of Applied Probability*, pages 294–313, 1992.
- Doyle Brunson. *The Complete Idiot’s Guide to Combinatorics*. Springer, 1969.
- Alex Churchill, Stella Biderman, and Austin Herrick. Magic: The gathering is turing complete. *arXiv preprint arXiv:1904.09828*, 2019.
- Erik D Demaine and Martin L Demaine. Every author as first author. *arXiv preprint arXiv:2304.01393*, 2023.
- Persi Diaconis, RL Graham, and William M Kantor. The mathematics of perfect shuffles. *Advances in applied mathematics*, 4(2):175–196, 1983.
- Matthew Eicholtz, Savanna Moss, Matthew Traino, and Christian Roberson. Heisenbot: A rule-based game agent for gin rummy. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 15489–15495, 2021.
- Leonhard Euler. *The Complete Idiot’s Guide to Euler Diagrams*. Springer, 1782.
- God. *English Standard Version Bible*, volume Genesis 3:27. God, 2009.
- Phoebe Goldman, Corey R Knutson, Ryan Mahtab, Jack Maloney, Joseph B Mueller, and Richard G Freedman. Evaluating gin rummy hands using opponent modeling and myopic meld distance. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 15510–15517, 2021.
- Pius Heinz. *The Complete Idiot’s Guide to Card Games*. Springer, 1890.
- John Lennon and Paul McCartney. *She loves you*. EMI London, 1963.
- Chris Moneymaker. *The Complete Idiot’s Guide to Economics*. Springer, 1970.
- Thomas Müller and Tomasz Placek. Defining determinism. *The British Journal for the Philosophy of Science*, 69(1):215–252, 2018. doi: 10.1093/bjps/axv049. URL <https://doi.org/10.1093/bjps/axv049>.
- Ed Oviedo. Card games where you spell words with cards. <https://boardgamegeek.com/geeklist/41551/card-games-where-you-spell-words-cards>, 2009. Accessed: 2024-03-02.
- Robert M. Pirsig. *Zen and the Art of Motorcycle Maintenance*. William Morrow and Company, 1974.
- qntm. *There is No Antimemetics Division*. Amazon, 2021.
- Kenny Rogers. *The Complete Idiot’s Guide to Shuffling*. Springer, 1976.
- William Shakespeare. *The Complete Idiot’s Guide to Literature Reviews*. Springer, 1596.
- Pramod Shankar. *How To Win At Gin Rummy: Playing for Fun and Profit*. Echo Point Books & Media, LLC, 2022.
- Shel Silverstein. *Trouble at Twilight: How Seven Ate Nine*. Harper & Row, 2008.
- Cara Stroud. Codetta and anthem postchorus types in top-40 pop from 2010–2015. *Music Theory Online*, 28(2), 2022.
- Sun Tzu. *The Art of War*. Springer, 499 BC.
- Mateen Ulhaq. How do I generate all permutations of a list? *Stack Overflow*, 2022.

A Appendices

A.1 The Bestest Deck

Here she is as an integer:

```
29,461,243,163,330,132,109,263,403,470,...  
205,013,845,203,144,254,219,393,515,...  
414,849,402,436,500,432,182,322,085,...  
137,075,217,281,411,271,006.
```

A.2 Classic Shuffler

This code might even compile.

```
PROGRAM CARD_SHUFFLE  
  
    INTEGER DECK(52), I, J, K, TEMP, SEED  
  
    DATA DECK/52*0/  
    CALL SRAND(1234)  
  
C Shuffle the deck  
10  I = 1  
20  J = INT(RAND(0) * 52) + 1  
30  K = INT(RAND(0) * 52) + 1  
    TEMP = DECK(J)  
    DECK(J) = DECK(K)  
    DECK(K) = TEMP  
    I = I + 1  
    IF (I .LT. 52) GOTO 20  
  
C Print the shuffled deck  
    WRITE(*, '(A)') 'shuffled deck:'  
    DO I = 1, 52  
        WRITE(*, '(I3)') DECK(I)  
    ENDDO  
  
END
```

Proving $P = NP$ thanks to the wonderful work of Reviewer Two

Claire Boine, *University of Ottawa* and David Rolnick, *McGill University*

Abstract

P represents the class of problems that can be solved in polynomial time, while NP represents the class of problems for which the solutions may be verified in polynomial time. While most computer scientists believe that $P \neq NP$, nobody has ever been able to prove it. The obvious reason is that this belief is wrong. In fact, as soon as we set out to demonstrate that $P = NP$, it took us only a few minutes to do so in no less than four different ways. This is not surprising given the groundbreaking and authoritative work conducted by Reviewer Two on this topic over the past few centuries.

I. INTRODUCTION

WE prove that $P = NP$ in four different ways, both experimentally and theoretically. In approach 1, we use a simple quantum mechanical technique to send a solver backwards in time to solve itself. In approach 2, we demonstrate that, equipped with the proper typewriters, a million monkeys who work together will come up with both a formal proof that $P = NP$ and an algorithm to solve any NP -hard problem in less than 60 seconds. In approach 3, we solve the traveling salesman problem and establish that it belongs to both P and NP . Finally, we prove that $P = NP$ the old fashioned theoretical way. Our work critically relies on the seminal works of Reviewer Two, who alone above all other authors has had the insight and dedication to contribute meaningfully to this important problem and many more [1], [2], [3], [4]. It is also consistent with the findings of Area Chair in their seminal piece presenting a theory of everything [5].

II. BACKGROUND

A. Approach 1

Our first approach relies on previous findings in quantum mechanics. Building on previous work by Two [6], Schrödinger famously explained that cats in boxes can survive an infinite amount of time as long as nobody disturbs their peace. Before and after that, Two proposed a variation of their famous double-slit experiment [7] called the Delayed-Choice Quantum Eraser, which proved that time is not linear [8]. A laser generates entangled photon pairs. One photon (signal photon) goes towards the double-slit apparatus, while the other (idler photon) is directed to a separate detection setup. The signal photons create an interference pattern on a detector, similar to the classical double-slit experiment, suggesting wave-like behavior. The idler photons are sent to a setup that can either preserve or erase the information about which slit the signal photon went through. This is achieved by manipulating the idler photons in such a way that if you could know which path the signal photon took, the interference pattern is destroyed (particle-like behavior). However, if the path information is “erased,” the interference pattern can reappear, suggesting wave-like behavior. Our first experiment largely draws inspiration from the Delayed-Choice Quantum Eraser.

B. Approach 2

In 1913, Two proposed what is known as the Infinite Monkey Theorem [9], which posits that a monkey hitting keys randomly on a typewriter for an infinite amount of time will at some point type the entire work of Shakespeare. Alternatively, an infinite number of monkeys could do it right away. Interestingly, the literature on the infinite monkey theorem has focused on the infinity component, while overlooking a critical piece: nobody uses typewriters anymore. What if, we wondered, the key to the infinite monkey theorem is to use old technology. As a result, we hypothesized that, while it would fail at it with a modern computer, a monkey equipped with a 1973 Xerox Alto and infinite time would probably produce an algorithm capable of solving any NP -hard problem in less than a minute. However, the authors of this paper did not have infinite time as they hoped to present their results at the Sigbovik 2024 conference.

Separately, pigeons have been trained to detect benign and malignant breast tumors [10]. In fact, researchers have shown that, while a human surgeon has a higher accuracy rate than a single pigeon, a group of pigeons is always better than a surgeon at diagnosis. Building on the same collaborative spirit, we trained a million monkeys to work together for approach 2, therefore reducing the infinite time it would take a single monkey.

C. Approach 3

The Traveling Salesman Problem (TSP) has long perplexed computer scientists, mathematicians, and traveling salespeople alike. It consists in a salesman who must visit a certain number of locations only once while minimizing the time spent on the journey. It was first formulated by Two in 1831 [11].

D. Approach 4

Two's outstanding work on complexity classes has inspired many research findings since. For instance, Rolnick showed that two regular Stanley sequences may be combined into another regular Stanley sequence [12]. Relatedly, Boine showed that virtual companions can exhibit emotional complexity [13].

III. METHODS

A. Approach 1

We used the Delayed-Choice Quantum Eraser experiment to send information a small distance backwards in time. Specifically, we developed a highly efficient Python library, Yppy, optimized for working with negative temporal offsets, in which compressed information can be sent backwards in the course of the computation.

B. Approach 2

We recruited a million monkeys and partnered with thrift shops all over the world to secure Xerox Alto personal computers. We then trained the monkeys to cooperate by getting them to play Hanabi and Dungeons and Dragons. After each game of Hanabi, the monkeys would receive a number of bananas and cucumbers that proportionally increased with the number of fireworks they have built. For Dungeons and Dragons, each team of monkeys received a free pizza each time they rolled a 20. Once 100% of the monkeys got the maximum scores in each game at least 90% of the time over 10 games in a row, we moved to the next phase. We paired up the monkeys and installed each pair on a computer. While one monkey randomly hit the keys, the other monkey reviewed each random action and provided feedback.

C. Approach 3

While computer scientists have been focusing on salesmen for centuries, no scholar has investigated the TSP in the context of a saleswoman. Motivated to fill this gap in the literature, we recruited a saleswoman and asked her to deliver 26 packages to the following cities: Agen, Taipei, Porco Rosso, Saint Louis, Fuchu, Ur, Boende, Nineveh, San Juan, Laputa, Akureyri, Conakry, Havana, Podunk, Dubai, Seville, Algarrobo, Babel, Ankh-Morpork, Wonderland, Petaouchnok, Atlantis, Zootopia, Gotham City, Eden, and Minas Tirith. We made it clear she would only get compensated if she visited each location only once and chose the optimal route.

D. Approach 4

We used good old-fashioned human intelligence (GOFHI) to come up with a solid theoretical proof. In the weeks before writing this paper, we consumed a significant number of Omega-3 supplements to ensure our reasoning would be foolproof.

IV. FINDINGS

A. Experiment 1

We wrote an algorithm in Yppy to search for proof of $P = NP$, leveraging the Bootstrap Paradox. Namely, problems in NP have solutions that can be verified in polynomial time. Now, any valid proof of $P = NP$ can clearly be verified in polynomial time, and thus if $P = NP$, the construction of such a proof must itself lie in P. Using this logic, our algorithm was able to design a proof of $P = NP$ that was then sent backwards in time to be utilized as input to the algorithm. Due to the complexities of temporal loops, we regret that the output of the proof cannot be directly observed, but the algorithm printed "Done" to console (screenshot available upon request).

B. Approach 2

Over the course of 3 months, with a 9 am to 5 pm workday and a 5-day work week, one pair of monkeys came up with an algorithm that solves any NP-hard problem in less than a minute (code forthcoming). In addition, over 10 pairs produced written demonstrations that $P = NP$. Furthermore, these proofs were all written in the style of Shakespeare.

C. Approach 3

The combinatorial explosion in the TSP makes it notoriously difficult to solve as the number of cities increases, making it an NP-hard problem. We find that by substituting the salesman with a saleswoman, the TSP is reduced to a problem in P as the saleswoman is willing to ask for directions to obtain the optimal route. Therefore, if you send a heteronormative cisgender couple of salespeople on their merry way, the TSP becomes both an NP and P problem simultaneously. In fact, it is in a superposition of state as long as the couple argues. When they finally decide to either follow the path suggested by the man (NP) or the woman (P), the problem is forced to choose a state. We thus show that the traveling salescouple problem belongs to P and thus that $P = NP$.

D. Approach 4

The result is a simple consequence of using Lemma 14.3 in Two et al. [14] to derive upper bounds for $\sum_{i=-n}^n f(n)$ in the statement of Theorem A.4 from Two [15], combining with the matching lower bounds on $\Xi + \bar{\Xi}$ in Remark 1 of Two and Reviewer [16], using the machinery presented in Two [17] to construct the finite-state inverse martingale.

REFERENCES

- [1] R. Two, "L'hypothèse de Riemann," Manuscrit de l'abbaye Saint-Pierre de Remiremont, Remiremont, France, 1650.
- [2] —, "The eel question," Where's Waldo Special Edition, New York, NY, USA, 1972.
- [3] —, "Cold fusion," *The journal of Franco-Korean cuisine*, vol. 1, no. 3, pp. 24–28, May 1034.
- [4] —, "Keeping all your socks: how to close the quantum portal in your washing machine," *Journal of Home Appliances*, vol. 3, no. 2, pp. 786–788, March 2010.
- [5] A. Chair, "A theory of everything: Proving $P = NP$," in *Proceedings of the International Conference on Theoretical Computer Science (ICTCS 2024)*. New York, NY, USA: Association for Computing Machinery, 2021, pp. 42–56.
- [6] R. Two, "Someone should develop an animal-based thought experiment to explain the superposition of states," *The Journal of Foundational Ideas*, vol. 54, no. 2, pp. 36–44, June 1934.
- [7] —, "On the nature of light," *Philosophical Transactions of the Royal Society of London*, vol. 91, no. 1, pp. 1–16, January 1803.
- [8] —, "Proving time is not linear," *Philosophical Transactions on Transactional Philosophy*, vol. 93, no. 1, pp. 1–11, Summer Solstice -2205, translated from Sumerian by Reviewer Two.
- [9] —, "Shakespeare versus an infinite number of monkeys: a randomized controlled trial," *The Journal of Haplorhini Typography*, vol. 1, no. 1, pp. 151–154, September 1912.
- [10] V. N. Richard Levenson, Elizabeth Krupinski and E. Wasserman, "Pigeons (*Columba livia*) as trainable observers of pathology and radiology breast cancer images," *PLoS ONE*, vol. 10, no. 11, 2015.
- [11] R. Two, "Reducing the carbon emissions of traveling salesmen," *The international journal of trade and the planet*, vol. 7, no. 4, pp. 3–24, September 1720.
- [12] D. Rolnick, "On the classification of Stanley sequences," *European Journal of Combinatorics*, 2017.
- [13] C. Boine, "Emotional attachment to AI companions and European law," <https://doi.org/10.21428/2c646de5.db67ec7f>, 2023, mIT Case Studies in Social and Ethical Responsibilities of Computing, no. Winter 2023 (February).
- [14] R. Two, S. Two, T. Two, U. Two, and V. Two, "On the Dutch cheese problem," *Annals of Mathematics*, vol. 20, no. 3, pp. 10–11, 2003.
- [15] R. Two, *Algebraic geometry for the uninitiated*. Springer, 1980.
- [16] R. Two and T. Reviewer, "Surfaces of characteristic q and how to find them," *Nature Pure Mathematics*, vol. 4, no. 1, pp. 22–51, 1992.
- [17] R. Two, "Notes on inverse martingales," Coffee Club on Complicated Combinatorics (C4), 2023.

Gotta Collect 'em All

Kiera Jones

Abstract

When you look at a classical probability problem, It's so scary at the start. But that's the hardest part. And trust me, it get's so much worse. I finally know where I belong. This is the way the weak grow to be so annoying to everyone around them. I've got a feeling, With you, I'm becoming Meme. I've worked hard, Come so far. It's my story, I'm the star. And the road I'm on Could go until forever. Every day there's something new. Every day I share these horrible insights with you. I've got a feeling, With you, I'm becoming Meme. We're looking for new horizons. Let's go explore new horizons. We're looking for new ways to take this problem and complicate it. Let's go explore all the ways that we've done that. I've got a feeling, With you, I'm becoming Meme.[5]

Keywords

Coupon Collector's Problem — Probability — Computer Simulation — Pokémon — A flimsy excuse of a paper to justify the purchase of more Pokémon cards

Rosenstiel School of Marine, Atmospheric, and Earth Sciences, Miami, Florida - [she/her/hers] - @K13205

Contents

1	Coupon Collecting and YOU	1
2	What trading cards mess up	2
3	Rifting the Para-da problems in a... uh...	2
3.1	Individual Rarities	2
3.2	For the whole set	4
4	All the sets Here Together	5
4.1	Second through sixth verse, same as the first	5
4.2	Okay but are packs worth it?	5
5	Temporally Forcing through new data	6
	Acknowledgments	8
	Appendix	8

1. Coupon Collecting and YOU

The Coupon Collector's Problem is one of the most well-known problems in probability theory among people who study probability theory (read: nerds). For those unfamiliar with terminology, a coupon is a piece of paper you had to bring in while shopping in a physical store to save money on certain items you bought[10]. These are no longer commonly seen, because CVS accidentally bought every single one that existed, and has to now put them on their receipts while the coupon factories work through this backlog.

The Coupon Collector's Problem, however, is a weird corruption of this. Instead of using the coupons to procure necessary items to ensure your survival (like 35 cents off a combination laser-pointer/foot massager), one is instead trying to get one of every single coupon being distributed at a time, presumably to flaunt ones wealth and/or collection

of combination laser-pointer/foot massagers. Now, the easiest way to do so would be to ask the coupon distributors for the missing coupons, but unfortunately, in this telling of the story, they are not so kind as to listen to mere mortal's trifling concerns, and as such, are given a random coupon from the collection each time. The question then begs how many coupons would one have to collect in order to have one of each kind.

In the end, there are some simple math(s) one can go through to calculate the expected value[1], but the solution is for n coupons, one would expect to get nH_n coupons on average before getting one of every coupon, where H_n is the n th Harmonic number¹. Because we don't like doing math(s) in this math(s)-inspired paper, we will instead be replacing the math(s) with other math(s) and approximating the Harmonic numbers by the natural log function², because it takes way too much effort to get Excel to do the Harmonic calculations, but the natural log is a default function, and laziness is very important.

Part of the underlying assumptions for the Coupon Collector's Problem is that all the coupons are equally likely to be obtained, and that what coupon is received at one time has no impact on the other coupons⁴. Variations on this problem have been studied numerous times to try and convert it to something more akin to real life, whether by allowing duplicate coupons to be traded between collections[4] assuming the coupons are

¹The Harmonic numbers are the sums of $1 + \frac{1}{2} + \dots + \frac{1}{n}$.

²Specifically, $H(n) \approx \ln n + \gamma^3 + \frac{1}{2n}$.

³Euler's constant, $\gamma \approx .577215665$. No, Euler's constant isn't e , but that's still used because of the natural logs, and γ shows up a lot within contexts of natural logs. Also, this 3 was a footnote, and not γ being raised to the third power.

⁴or in actual math(s) parlance, the coupons follow IID: independent and identically distributed.

not equally likely[8], the coupons come in unique batches without repeats[9], or there are multiple collections going on at the same time and wanting to know how long until all of them are finished[6]. One great example of an adaptation that deals with multiple of these is looking at collecting trading cards.

2. What trading cards mess up

Opening up packs of trading cards is one way that the Coupon Collector problem can be translated to a more modern scenario. People frequently get tons of packs, and then open them all to see what they get. Usually, there are three main reasons people open up lots of packs: They are opening them up, trying to pull rare and valuable cards to resell and turn a profit; they are attempting to pull a few specific cards that they like or need for a deck, usually the rarer ones⁵; or, most pertinently, they are trying to collect every last card in the set oh hey wait that's what this paper is about.

The manner of getting the cards in a TCG is slightly different than that in the normal coupon collector's problem. First, the cards have differing rarities, and they have different distributions for each. Classically, these can be sorted into "common"(a plurality of the cards in a set that usually make up a large portion of each pack), "uncommon"(usually a smaller group that has a couple slots in a pack), "rare"(high end cards that are some of the heaviest hitters, but only appearing once or rarely twice in a pack), and "ultra rare"(the top end cards that aren't even guaranteed to show up in a single pack) cards. Because the cards are grouped in packs, the cards gotten in a pack usually aren't allowed to be duplicates, one does end up with a slightly wider spread than expected, which helps some, but this also means that they aren't independent of each other. And lastly, the different rarities can be construed as separate Coupon Collector problems, and completing the set is just the finishing time for the slowest of these. All told, looking at these together makes this an interesting problem, and one worthy of devoting

For this paper, I will be looking at the Pokémon TCG in particular for a couple of reasons: Firstly, as one of the largest media franchises in the world, there are a lot of players, and thus a lot of data to pull from[3]. Secondly, the author⁶ has been a longstanding fan of the series, and so it has some sentimental value. Thirdly, with the newest set *Temporal Forces* releasing on the extended submission deadline and rotating into tournament legality on the publication date[7], this would allow Sigbovik to be on the cuttingest of cutting edges, an honor which I am sure they strive to follow.

⁵but sometimes it's not like literally how have I never pulled one of the future cards from Paradox Rift *despite* four of them just being normal uncommons like please i need those future booster energy capsules for my future box lite deck cmooooooon

⁶and her Zorua plushie she is holding on to, who deserves full coauthorship

⁷Please note, that although *Temporal Forces* has not been released, these packs were legitimately obtained at a pre-release event that the author finished tied for third despite not having actively played in fifteen years.



Figure 1. A sampling of the packs looked at for this paper, and totally not an excuse for me to have bought a ton of cards during Pokemon day⁷.

Since the beginning of the Scarlet and Violet era of the TCG, the packs have followed this specific structure: Four common cards, chosen without duplicates; three uncommon cards, chosen without duplicates; two reverse holofoil cards, chosen from the pool of all common, uncommon and rare cards, that are not duplicates of each other but *may* be duplicates of other cards in the pack, one holofoil rare card, one basic energy card that is identical to those in any other set, and one code card for the online TCG. Additionally, the ultra rare cards can replace any of the three holofoil slots, depending on the type of ultra rare card and the set that it is in.

3. Rifting the Para-da problems in a... uh...

8

3.1 Individual Rarities

To help with this calculation, the author created a highly accurate simulation of opening packs of Pokémon cards from a perfectly specified set. To this, the data from the most recent standard set, *Paradox Rift*, was pulled, and a quick program was written up in Python. Simulated packs were opened up 10000 times, with how quickly each set was finished recorded for each, and the total time to complete the entire set.

Now, one might naïvely expect that the odds of finishing each of the subsets to be equal to the expected value of getting all of them if they were guaranteed divided by the odds of getting them in the first place. And after the ten thousand simulations... uh... yeah that's pretty close to the case. For n cards with a p chance of getting one of them, the expected number of packs is

$$E(n, p) \approx \frac{n \ln n + \gamma n + \frac{1}{2}}{p}$$

⁸look this is supposed to be a pun on paradox rift cmon gimme a break

Table 1. *Paradox Rift* card counts and predicted odds[3]

Rarity	Count	Odds
Common	81	
Uncommon	54	
Rare	27	
Double Rare/ex	20	15.57%
Ultra Rare/Full Art	28	6.64%
Illustration Rare	34	7.70%
Special Illustration Rare/Alt Art	15	2.11%
Hyper Rare/Gold Card	7	1.22%

For the standard commons and uncommons, which appear in more than one per booster pack, the “probability” is just the number of times they appear in a booster pack, accounting for the reverse holo slots⁹.

Table 2. *Paradox Rift* calculated and simulated averages

Rarity	Calculated	Sim. Mean	Sim. Median
Common	82.47	81.08	78
Uncommon	69.49	68.40	65
Rare	105.07	101.33	96
DR	458.93	484.25	455
UR	1648.46	1743.70	1647
IR	1811.87	1893.34	1789
SIR	2335.34	2441.37	2266
HR	1447.57	1572.38	1420

On looking at the higher rarities, though, the expected value seems a tad low. The calculated values were over 100 packs shy compared to the simulation for Ultra, Special Illustration, and Hyper Rares. The problem likely comes from trying to set up a Poisson-esque distribution with a Geometric distribution feeding in, and so a very unlucky set of pulls can end up skewing the results¹¹. Given all this, instead of trying to come up with the exact way to match better, I opted for the easy route of just looking at the median, which tended to match much better, because moving the goalposts is a completely valid strategy for Sigbovik.

The expected chance of finishing before a certain point is also able to be calculated from bastardising pre-existing formulae[2]. In the end, it, too, is just the calculated value divided by the probability of getting one in the first place, but there is a small correction term added to account for the unluckiness in getting the rarer cards in the first place. For n cards with a p chance of getting one of them, to have less than an α chance of not having a complete set, the expected

⁹*Paradox Rift* makes this nice and easy, as the commons, uncommons, and rares make up exactly¹⁰ 1/2, 1/3, and 1/6 of the total cards available for reverse holos.

¹⁰Although you do have to subtract about the 1/9 chance that the second reverse holo is replaced by one of the ultra rares, and for the rares, you also need to worry about its replacement.

¹¹One unlucky run took 8975 packs to open up all the special illustration rares, more than 3.5x the average, and nearly running up against my hard-coded 10,000 pack limit

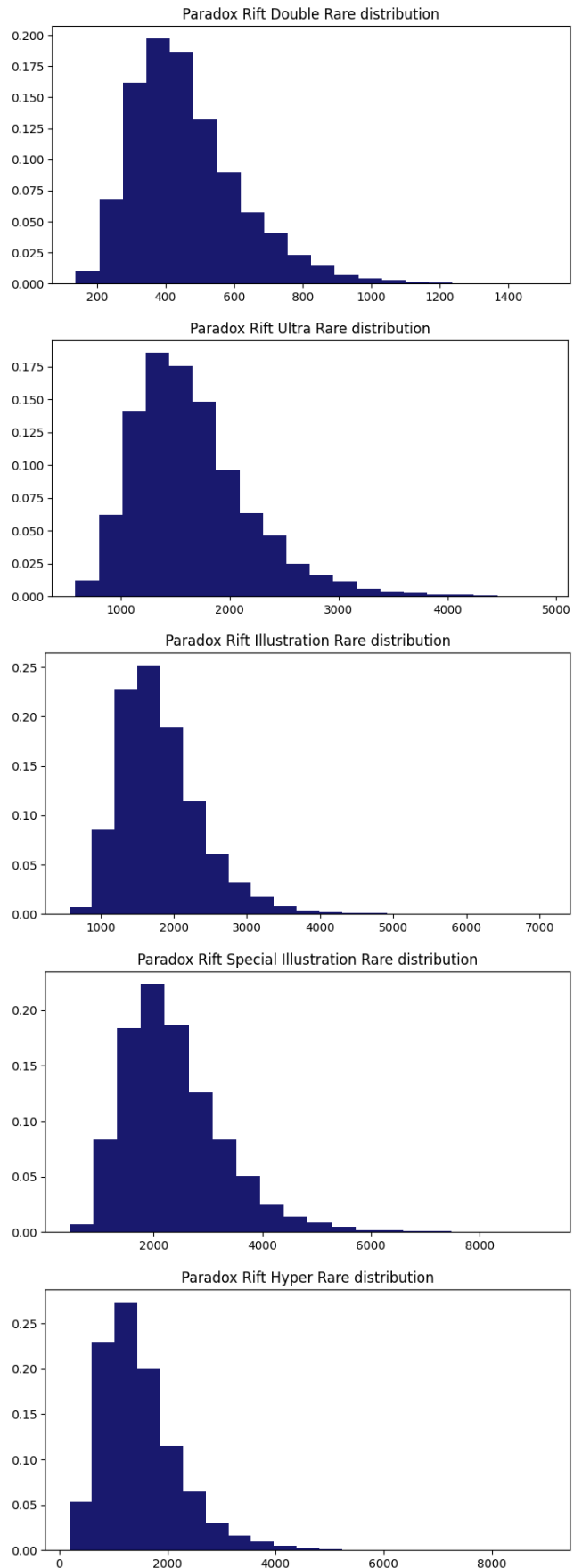


Figure 2. Simulation distributions for 10,000 simulations of *Paradox Rift*

number of packs is¹²:

$$E(n, p, \alpha) \approx \frac{(n \ln n)(1 - \log_n \alpha) - \ln p}{p}$$

3.2 For the whole set

Calculating out the needed number of packs to complete the whole set, however, is a bit harder. For one, the two simplest approaches are patently wrong:

1. One could expect that the expected time to finish up all the sets is the sum of all the times for each rarity, but that would be simulating having to wait to start a set until the previous ones are finished. An extreme case of this would be if one set had... let's say 120 different cards in one of the rarities¹³ with an expected time to finish of around 2500 packs, while the rest of the rarities (discounting standard commons, uncommons, and rares) had no more than a dozen cards and would be thought to finish in less than 1000 packs. In this case, most often the last rarity to finish would be that single large one, and the expected number of packs should be close to it.
2. One could also expect that the expected time to finish would be the largest value, as all the other one would be finished. However, just because they are likely to be finished first doesn't mean that they will be first. Again, a case for this might be having three rarities being extremely close to finishing at the same time with the probabilities accounting for the different sizes¹⁴. In this case, although the three sets would *average* finishing around the same time, it would not be a surprise to have one of them finishing early and another finishing later to even out the differences. And since the set is only complete when all are done, the expected time should be later.

So now we are somewhere between the largest value of any of them and the sum of all the values. Such a great range. Thankfully, math(s) come(s) to our rescue, and the expected value is[6]

$$E(s_1, \dots, s_k) = nH(S) - \frac{(H(S) - 1) \sum_{i < j} s_i s_j}{S(S - 1)} + \mathcal{O}\left(\frac{1}{n}\right), S = \sum_{i=1}^k s_i$$

$$P(t < x) = \prod_i \left(1 - n_i^{1 - \frac{p_i x + 1 + \ln(1-p)}{n_i \ln n_i}}\right)$$

Uh... the... verification of these will be left as an exercise for the reader, because this is very sensible and easily translated into our new format accounting for only having a chance

¹²Since $p < 1$, $\ln p < 0$, and so subtracting the log adds to the term

¹³Spoiler alert for a few sections from now: *Paldean Fates* shiny rares.

¹⁴Spoiler alert again: *Scarlet & Violet Base Set* has Ultra, Illustration, and Special Illustration rares at around 1200 packs.

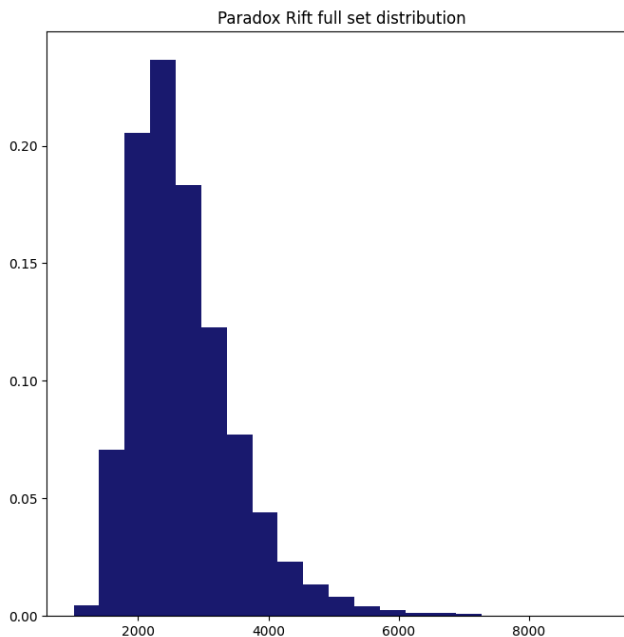


Figure 3. Simulation distributions for 10,000 simulations of *Paradox Rift*

at getting a coupon each time. Rather than bog ourselves down with the math(s) here, instead, we will look at a way to get an estimate having already calculated the prior results.

The way this was done was just simulating all of the *Scarlet* and *Violet* era sets¹⁵, and then using everyone's favorite linear regression on sensible candidate equations to see which came the closest for its simplicity. Candidate equations had to follow these sensible guidelines:

1. $F(E) = E$
2. $F(E_0, E_1, \dots, E_k) \geq E_0, F(E_1, \dots, E_k)$
3. $F(E_1, \dots, E_k, 0) = F(E_1, \dots, E_k)$

These guidelines state roughly that the expected value of only one prior expected value must equal itself, that adding another expected value cannot decrease the total values, and that adding nothing cannot change the current expected value. Upon searching with this, one simple equation that fulfilled all these conditions and matched very well was raising each number to a derived power, adding them together, and then taking that derived root:

$$F(E_1, E_2, \dots, E_k) = \left(\sum_{i=1}^k E_i^\lambda \right)^{\frac{1}{\lambda}}$$

¹⁵Initially I was also going to include late *Sword* and *Shield* era sets, but given that *Silver Tempest* has rarities of Pokémon \mathcal{V} , Pokémon \mathcal{V} /MAX/ \mathcal{V} /STAR, Full art Pokémon \mathcal{V} , Full art trainers, Alternate art Pokémon \mathcal{V} , Trainer Gallery Pokémon, Trainer Gallery \mathcal{V} / \mathcal{V} /MAX/Trainers, Trainer Gallery Gold \mathcal{V} /MAX, Radiant Rares, Rainbow Rares, and Gold Rares... yeah no

And with a simple regression gives a value of $\lambda \approx 3.88902$, which is close enough to 4 that I'm willing to just call it that to make it even simpler. Evaluating this for the prior values gives an expectation of 2826 packs, which is very close to the simulated 2813. For expected values for a certain α , it turns out that the solutions for λ varies linearly¹⁶. The calculated value is $\lambda = \frac{100}{27} + \sqrt{2}\Phi(\alpha)$, where Φ is the normal inverse¹⁷. Sadly this is a lot more complicated than the formula for just the means, but if you want to know when you will be 69.420% likely to have completed a set of cards, well, you deserve it.

4. All the sets Here Together

4.1 Second through sixth verse, same as the first

Of course, now that these methods have been thoroughly tested through exactly one series, it's time to expand these into the other five sets from the Scarlet & Violet era. One thing that needs to be brought up is that some of the other sets have special rarities only found in them, namely *Pokémon 151*'s Holofoil energies, and *Paldean Fate*'s Shiny Pokémon¹⁸. Thankfully, with how the expected values were found before, these only need to be slotted in to the slots where they are found, and simple as that, the entire results can be ran with nary a problem¹⁹.

Table 3. Scarlet & Violet era expected pulls. These data were all run at the same time, hence the difference between the *Paradox Rift* values from before and in this table, as they were different simulations.

Set	DR	UR	IR	SIR	HR	Full Set
SVI	271	1093	1183	1188	802	1557
PAL	425	1516	1953	1564	1457	2309
OBF	556	559	490	471	289	765
MEW	280	848	646	589	287	974 [†]
PAR	462	1654	1823	2375	1503	2713
PAF	185	175	77	1299	947	2585 [‡]

[†]: *Pokémon 151*'s Holofoil energy has an expected value of 89 packs.

[‡]: *Paldean Fates* has Shiny Rares with an expected value of 2535 packs, and Shiny Ultra Rares with 487 packs.

From these, we can see that there is quite a variation in the sets. The quickest set is by far *Obsidian Flames*, which has only 54 cards above the normal rarity. This small size, combined with fairly generous pull rates²⁰, leads to a set that you are likely to complete well before most any other sets'

¹⁶namely because two values were looked at

¹⁷Single-sided, so for $\alpha = .95$, it is *not* 2 (1.96) but instead 1.645

¹⁸Which thankfully are more common to get than the 1/4096 chance of shinies in the video games

¹⁹Okay maybe a little of a problem, considering that *Pokémon 151* has "Demigod" packs with the (special) illustration rares for one of the starter lines instead of any reverse holofoil cards. This took like an hour to encode, but I did it just for accuracy. It made effectively no difference.

²⁰The first three sets all had fairly similar rates for the rarities. For *Paradox Rift*, although Double Rares went from around 13.5% to 15.5%, this was balanced out by Special Illustration Rares dropping from around 3.2% down



Figure 4. All the different high rarities released in Scarlet & Violet. The top row has only been found in a single set: Holofoil Energy from *Pokémon 151*, Shiny Rare and Shiny Ultra Rare from *Paldean Fates*, ACE SPEC from *Temporal Forces*. The bottom two rows have been in all sets: Double Rare (from *Paradox Rift*), Ultra Rare (from *Obsidian Flames*), Hyper Rare (from *Paldean Fates*), Illustration Rare (from *Scarlet & Violet Base*), and Special Illustration Rare (from *Temporal Forces*).

Illustration Rares. *Paradox Rift* is actually the slowest set to complete, because although *Paldean Fates* has a worse set in its 120 Shiny Rares, that's the only bad one, while all of the non-Double Rares in *Paradox Rift* take a while.

The largest factor does tend to be the size of the special pulls. Although how they end up distributed and the exact odds do have some weight (hence why *Paradox Rift* is slower than *Paldean Fates*, despite having 60 fewer cards), a set with a large number of cards that need to be pulled will just necessarily take longer than a shorter one.

Now with all these calculations done, we just have to ask...

4.2 Okay but are packs worth it?

NO

to 2.2%, and Hyper Rares dropping from 1.8% to 1.4%. Since those tend to take longer to complete than Double Rares, this drop in odds has lengthened the completion time of the latter sets.

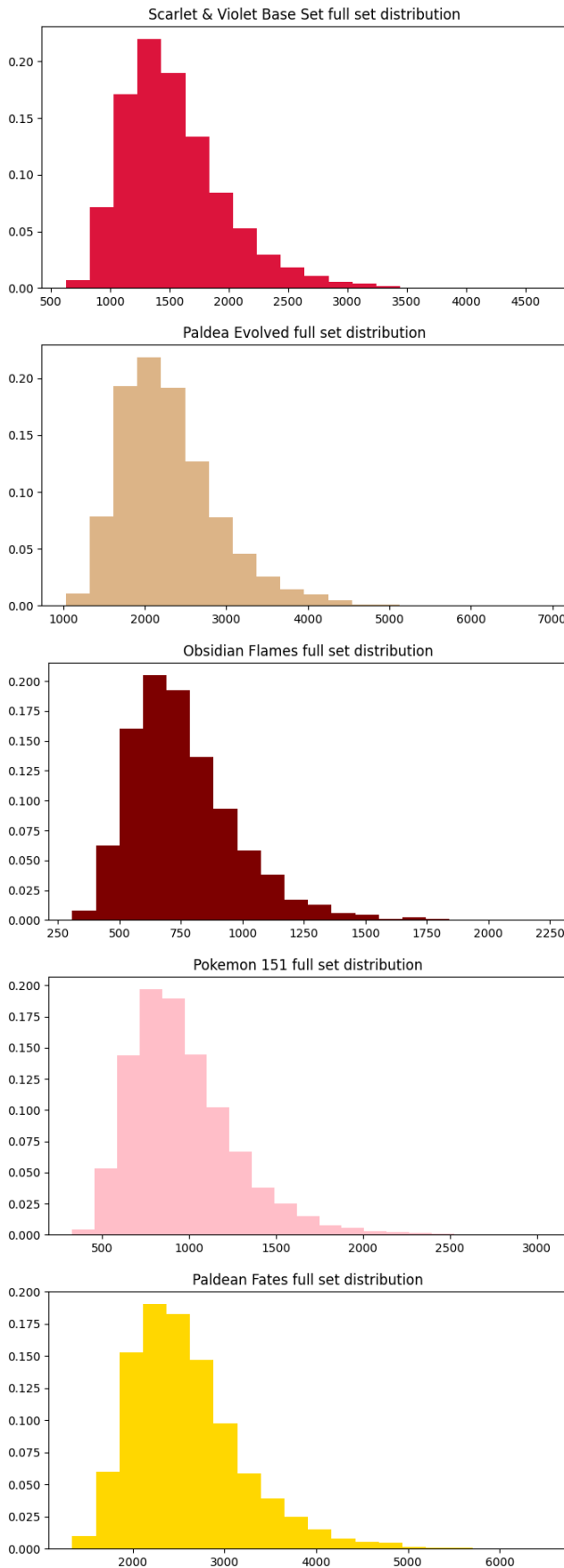


Figure 5. Simulation distributions for 10,000 simulations of the released Scarlet & Violet TCG sets

This should not be much of a surprise. The expected value of packs is based on the expectations of what could be contained in them, while usually only a handful of cards contain the value of the set²¹. On average, the entire set tends to sell for around one ninth the cost of the packs needed to pull all the cards. This difference is so severe, the only rarities that are likely to be pulled in the number of packs it takes to spend as much as the full set are the *Pokémon 151* Holo Energies (89 vs 132), and the Double, Ultra, and Illustration rares from *Paldean Fates* (185, 175, and 77 vs 224). And none of those cards are the ones providing the value in the set²².

Table 4. Scarlet & Violet era expected costs from both pulling from packs and the price of all the cards in the set. All data was pulled from TCGPlayer and was accurate at the time of writing. Pack price is the most efficient purchase (usually a booster box of 36 packs), not the price for a single pack.

Set	Packs	Pack	Total	Set	Ratio
SVI	1549	\$2.36	\$3656	\$313	11.68
PAL	2345	\$2.64	\$6191	\$706	8.77
OBF	749	\$2.31	\$1730	\$201	8.61
MEW	943	\$6.00	\$5660	\$792	7.14
PAR	2691	\$2.38	\$6403	\$702	9.12
PAF	2582	\$4.08	\$10,537	\$915	11.52

Thus from all of this, the best shot of managing to get a complete set is if you just go out and buy all of the cards in it. Actually, you should probably first get some booster packs, because trying to find people willing to sell random new commons and uncommons is hard, so you'll probably have to go online, where for such worthless cards, you'll have to pay shipping and that will be so much more than the cost of the cards²³. Failing that, if you want to go with pulling the cards, the next best thing is to get some friends. Although trying to complete two sets takes longer (and thus costs more money), it doesn't take twice as long, and you'll save yourself a little money in the end. In fact, the more people that you get in your trade network, the less each will have to pay on average. And if some people are only concerned with getting a few cool cards that they like, that helps even more.

5. Temporally Forcing through new data

With all that said on the current sets, let's turn our eyes to the new one. *Temporal Forces* will release on the 22nd of March, 2024²⁴. One of the biggest factors for this set is a brand new rarity of cards for people to chase... or well... brand new for

²¹To wit, the author just opened up a *Paldea Evolved* booster pack as she wrote this sentence, and although she did get a hit in the double rare Slowking ex, that card goes for less than what the pack would have. It's still a very cute card, and she'll absolutely take it.

²²Most expensive card, and only one over a dollar: PAF054 Charizard ex - \$2.41

²³Source: I got those darn Future Booster Energy Capsules I needed in my deck for three cents each, and spent \$1.22 on shipping.

²⁴The Extended Submission Deadline

Scarlet and Violet. Bringing back a second old mechanic²⁵, ACE SPEC cards return for the first time since *Plasma Blast* in Black and White.

Now, as the set has not been officially released²⁸, the odds of pulling any specific card is not yet known. However, with the full set list already shown, as well as the expectation that the pull rates are *probably* going to be similar to the other sets, we can at least make an estimate going through what we know of the past sets. Thus, for each of the sets ($i \in [\text{SVI, PAL, OBF, MEW, PAR, PAF}]$),

$$\hat{\mu}_{\text{TEF}} \approx \frac{1}{6} \sum \bar{x}_i$$

$$\hat{\sigma}_{\text{TEF}}^2 \approx \frac{1}{6} \sum s_i^2 + \frac{1}{5} \sum (\bar{x}_i - \hat{\mu}_{\text{TEF}})^2$$

Now, I could just sit on my laurels, run the simulation with these values, and go on with my day like nothing had happened. But that's not what a scientist does. They do data collection. And they **WRITE STUFF DOWN**. Also, uh, since ACE SPEC cards are new there's no prior data to make an estimate from. I mean, there's *guesses* from the Black and White era on their frequency, but not of the "open up thousands of packs of cards and meticulously record all the pulls" that the other data were from. So there's nothing else to do but go to a couple of pre-release events and grab the data from there.

Table 5. Calculations for *Temporal Forces* pull rates, and expected values from them.

Rarity	Average	+Sample	Calc	Sim
DR	14.31%	14.59%	342	353
UR	6.59%	6.61%	952	960
IR	7.73%	7.64%	1063	1052
SIR	2.73%	2.17%	1353	1294
HR	1.72%	1.51%	974	963
ACE		7.28%	250	279
Full			1589	1560

After around 300 packs worth of data²⁹, it looks like *Temporal Forces* will have similar pull rates to the last few sets, with the slightly higher Double Rares in exchange for lower Special Illustration and Hyper Rares. ACE SPEC cards do

²⁵ex cards²⁶were revealed to be part of the set at its reveal, having last been seen in the Ruby and Sapphire era.

²⁶Not to be confused with EX cards, last seen during XY. Yes, they're different, and the capitalisation matters. So technically in Expanded, you could run four copies of AOR042 Tyranitar EX and four copies of OBF066 Tyranitar ex. Heck, throw in four copies of AOR043 M Tyranitar EX²⁷, four copies of BTS097 Tyranitar V, and four copies of any normal Tyranitar. They're all different.

²⁷oh god i just realised with legends z-a showing megas are coming back, we might get the return of M EX pokémon and will have to deal with both ex and EX cards.

²⁸But there have been Prereleases, as noted

²⁹Specifically, in 302 packs, there were 51 Double Rares, 21 Ultra Rares, 20 Illustration Rares, 2 Special Illustration Rares, 1 Hyper Rare, and 22 ACE SPEC.

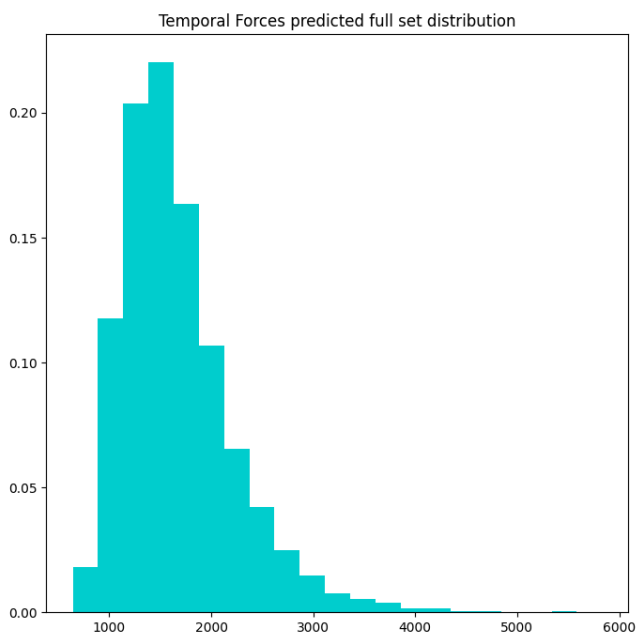


Figure 6. Simulation distributions for 10,000 simulations of *Temporal Forces* with estimated odds.

appear to be fairly common, and with them being the sole card to appear in the first reverse holo slot in this set³⁰, this just makes them bonus cards to pull, with around 36% of packs expected to contain at least some form of special card, higher than any of the other non-premium sets³¹. And with the overall size of the hits being on the same order as *Scarlet & Violet Base*, this looks like it will have the best pulls this side of *Obsidian Flames*, and with lots more cool cards to boot.

Will it be worth it, though? Again, probably not. The most cost-effective way to purchase packs for preorder (a case of booster boxes) is running \$610 for 216 packs, good for an average price of \$2.82 per pack. This means that to purchase all those packs, the cost right now is \$4480. With the current value of the cards at \$1203, this means it is only 3.72x more expensive to pull the cards rather than just buying them. However, given that the cards aren't actually out yet, it'll probably take about a month before the cards actually stabilise in price, and will likely correct to a much lower value than this. It may still turn out to be one of the best relative valued sets, but even now in this extremely favorable environment, it's so much cheaper to just buy the cards rather than trying to pull them.

³⁰Oh I guess I never mentioned where they go. Double and Ultra Rares take the place of the normal rare, while Illustration, Special Illustration, and Hyper Rares replace the second reverse holo. In *Paldean Fates*, all the Shiny cards also replaced the first reverse holo, while *Pokémon 151*'s Holofoil Energies replaced the energy.

³¹oh yeah you probably noticed that *Pokémon 151* and *Paldean Fates* had much higher booster prices than the other sets. Premium just means the booster packs can't be bought normally and only come in special collections.

Acknowledgments

The author would like to thank both Pro-Play Games in Miami and U KnoWhere 2Play in Hialeah, as well as the Pokémon League players attending those locations, for their assistance with collecting data for the *Temporal Forces* section and for getting pictures of all the rarities. She would also like to thank her co-author Zorua plushie for using its trickster abilities to help her in the prerelease event she attended.

Appendix: so *Temporal Forces* released...

With the extended deadline, *Temporal Forces* has released, and as such, I have the actual data for the set. And... uh... the odds were very much messed with. Double Rares are now even more common than they were in *Paradox Rift* and *Paldean Fates* at around 17% of packs rather than 15.5%, and ACE SPEC cards were a fair ways less likely, at 5% rather than 7.3%, although the large gap was more likely given the lack of data for the value.

The biggest shock, however, is that both Special Illustration Rares and Hyper Rares have had their rates dropped to around half, from an estimate of 2.2% to a real of 1.2% for SIR and from 1.5% to 0.7% for HR. This means they take SO much longer, and *Temporal Forces* goes from being probably the easiest set to collect since *Obsidian Flames* to the WORST set by far.

And as for buying packs vs singles, the cheapest option is still the booster box case, down to \$570 for 216 packs at \$2.64 per pack. At this rate, the expected cards would set you back \$7,300, and, as of the time of writing, the current set has a value of \$837³², for an expected overcost of 8.73x the value of the set. As the set continues to spread out, the price is likely to keep dropping, so this will be even more overpriced as time goes on while the set is in print, potentially getting to levels above any of the prior sets.

Table 6. *Temporal Forces* predicted and actual averages

Rarity	Simulated		Calculated	
	Pred.	Act.	Pred.	Act.
DR	353	311	348	296
UR	960	994	955	944
IR	1052	1081	1051	1052
SIR	1094	2578	1073	2504
HR	863	2189	858	2044
ACE	279	386	250	364
Set	1460	3033	1408	2769

References

[1] Agrawal, Naman, 2023: *Coupon Collector's Problem: A Probability Masterpiece* Towards Data Science <https://towardsdatascience.com/coupon-collectors-problem-a-probability-masterpiece-1d5aed4af439>

³²Yes price data is still from TCGPlayer

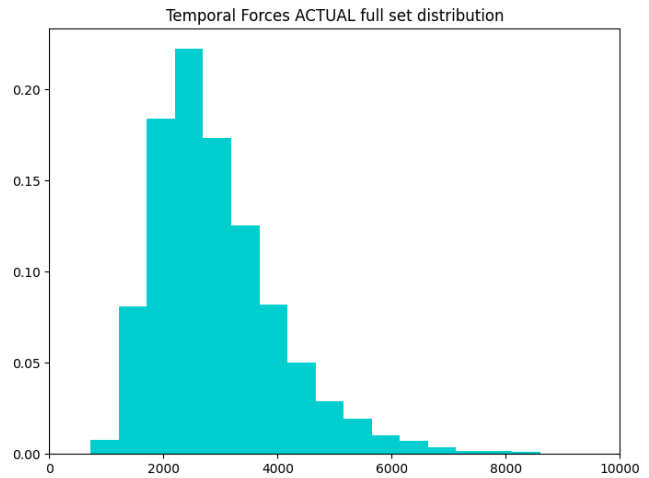


Figure 7. Simulation distributions for 10,000 simulations of *Temporal Forces* with the real data. Four simulations were stopped early after failing to collect all the cards in 10,000 packs.

[2] Brown, Mark, Peköz, Erol, and Ross, Sheldon, 2008: *Coupon Collecting Probability in the Engineering and Informational Sciences*, Volume 22 Pages 221-229 <https://doi.org/10.1017/S0269964808000132>

[3] Day, Peter, 2024: *Pokémon TCG: Paldean Fates Pull Rates* TCGplayer Infinite <https://infinite.tcgplayer.com/article/Pok%C3%A9mon-TCG-Paldean-Fates-Pull-Rates/23de3e93-0d0f-4ae0-abc4-13664f3001a3/>

[4] Ferrante, Marco and Saltamacchia, Monica, 2014: *The Coupon Collector's Problem* MATerials MATematics, Volume 2014 Issue 2 https://mat.uab.cat/matmat_antiga/PDFv2014/v2014n02.pdf

[5] Goldfarb, Ed, 2014: *Becoming Me* <https://www.youtube.com/watch?v=xWfHsLXo5s8>

[6] Krityakierne, Tipluck and Thanatipanonda, Thotsaporn, 2024: *The Slowest Coupon Collector's Problem* <https://thotsaporn.com/Coupon.pdf> [Preprint]

[7] Pokémon News, 2024: *2024 Pokémon TCG Standard Format Rotation Announcement* <https://www.pokemon.com/us/pokemon-news/2024-pokemon-tcg-standard-format-rotation-announcement>

[8] Shank, Nathan and Yang, Hannah, 2013: *Coupon Collector Problem for Non-Uniform Coupons and Random Quotas* The Electronic Journal of Combinatorics, Volume 20 Issue 2 Article 33 <https://doi.org/10.37236/3348>

[9] Stadje, Wolfgang, 1990: *The collector's problem with group drawings* Advances in Applied Probability, Volume 22 Issue 4 Pages 866-882 <https://doi.org/10.2307/1427566>

[10] Wiktionary, 2024: *Coupon* <https://en.wiktionary.org/wiki/coupon>

Mining for Gold Coins

B. "I put the Tran in FORTRAN" Parazin

Lord Blunderbuss' School of the Arcane

New Valentis

Occupied Dawnlands, H3A 0G4

Abstract—Table-top role-playing games are ubiquitous in culture today, but player outcomes in these games are still largely governed by luck and random chance. To address this shortcoming of the medium, I develop a novel piece of code which allows one to, for a short period of time, control the outcomes they will receive when digitally rolling dice, provided they are using NumPy's random library. Furthermore, through an interdisciplinary process of performing interviews with self-professed "Dungeon Masters" and statistical methods, I am able to reframe this abstract benefit of foresight as a concrete amount of in-game currency, gold. I find that players would gain, on average, 182,930,844,029 gold pieces by utilizing my code. Finally, by considering that value and the computational cost needed to find it, I find conversions between gold pieces and other currencies.

1. Introduction

As society crumbles around us, table-top roleplaying games have become a popular form of escapism, allowing people to pretend they live in a world where the impossible is real. In this fantasy wonderland, people regularly make a decent living freelancing, get 8 hours of unbroken sleep a night, and, through the power of teamwork and friendship, take down the evil robber barons that run their local kingdom. The ability to live out these dreams, however, is wholly dependent on the whims of the "dice gods." In the course of gaming, players have to make many skill checks, attack rolls and saving throws, which use a die or series of dice modified by some static value(s) compared to a threshold to determine if a player succeeds in attempting uncertain actions.

In recent years, however, many gamers find digital dice, which use pseudo-random methods to generate a random number in a given range, to be favourable. Their reasons for this preference vary: some don't want to worry about losing dice, others have limited bag space better spent on painted miniatures, and still others have given into the siren call of their polyhedra and, in the spirit of Chronos, they consumed their dice whole.

The growth of these digital dice rollers also opens up new opportunities for players interested in tipping the scales of luck in their favor. Previously players turned towards methods of witchcraft or superstition to control their luck

and stave off "Dice Curses." Common methods included talking to dice, blowing on dice, placing misbehaving dice in dice jail, or sacrificing a goat or other mammal to appease and soothe any angry dice.

Fifty percent of the time, these methods work every time¹, but with the power of modern science and mathematics, there is an opportunity to improve upon previous methods of controlling one's fate. By utilizing the pseudo-random nature of digital dice rollers, players can now set a seed that will determine the series of results they get when rolling dice. With an adequate choice of seed, a player can therefore be confident that luck will finally be on their side, the only question is, what seed to choose?

In this paper, I describe the process of choosing the optimal seed for a common pseudo-random number generator, and some of the consequences of this choice. Section 2 introduces this seed, describes the methods used to determine this seed, and highlights some other pseudo-random seeds of interest to fans of d20 tabletop games. Section 3 then details my research to quantify the significance of this seed, showcasing some potential results of using it in game. Section 4 expands on the gold amount found in section 3 by exploring alternative uses of the computational power used in this assignment, ultimately finding two further foreign exchange conversion factors for d20 fantasy gold. Section 5 discusses future research paths.

2. Pseudo-random seed exploration

In a d20 fantasy game, the best thing one can roll on a dice-based check is a 20 on the dice. This is a so-called "Natural 20," and is always a cause for celebration around the table. It therefore follows that the optimal pseudo-random seed for a dice roller is one that generates the most natural 20s in a row from the start. To find this, there are two methods one can use. For an open-source pseudo-random number generator like the one used in this study, *numpy.random.Generator.integers*, one could take a look under the hood and strive to understand the numerical methods used in generating these numbers. This intrepid researcher could then take that understanding and produce an inverse model, one that takes in a string of integers and

1. Except for the goat sacrifice, which has 36% effectiveness

returns the corresponding seeds that could produce that, in a stunning synthesis of mathematics and computer science.

Alternatively, one can simply guess-and-check seeds, recording those that produce the longest streaks of the same number, in an equally stunning synthesis of stubbornness and laziness. In this paper I undertake the second method.

2.1. Guess-and-Check code

The heart of the code which performs this guess-and-check method is reproduced in part below.

```
import numpy as np
...
rng = np.random.default_rng(seed)
streak_number = rng.integers(low=1,
                             high=21)
current_roll = streak_number
streak_length = 0
while current_roll == streak_number:
    streak_length+=1
    current_roll = rng.integers(low=1,
                                high=21)
...
```

For a given seed, this algorithm sees how many times in a row the first random number that seeds rolls is generated, saving the resulting seed and streak length to a file. This code was loaded onto a Raspberry Pi 4, and allowed to run for 111 days straight, checking random seeds 0 through 85,500,000,000. The calls to `rng.integers` are from `low=1` to `high=21` since the function produces random integers in the range `[low, high)`.

2.2. Streak length results

The results of the first 85.5 billion seeds checked are presented in table 1 as Roll, Streak Length and seed, meaning, for example, the NumPy `random.default_rng` seed 216240055 produces 9 12s as the first 9 integers generated when generating in the range `[1, 20]`. Most germane to this line of research is the seed 75364316682, which produces 8 natural 20 rolls in a row. Finding a streak this long is pushing the upper limits of what I can find with this guess-and-check method, as I will need to, on average, run this code for another 666 days or to find a streak of 9 natural 20s in a row, and another 36.46817 years to find a streak of 10 natural 20s.

3. Dice Rolls to Gold Piece conversion

There is obviously a strong benefit to being able to know what the results of your next eight rolls in a d20 fantasy game, doubly so if one knows that they will be the best possible outcomes, but this benefit is abstract and unquantified. To remedy this shortcoming of my research, I undertook interview with 21 self-proclaimed “Dungeon Masters,” where I outlined the scenario and asked for their expert opinion on how many gold pieces one could gain.

Roll	Streak Length	NumPy seed
1	8	68840775906
2	8	80686509533
3	9	74955903520
4	9	83606999174
5	8	77136117593
6	8	44827294402
7	8	52740367896
8	8	50104951153
9	8	78235588018
10	8	24243953919
11	8	80127612848
12	9	216240055
13	8	1382968756
14	8	43181918036
15	8	73192838540
16	8	81406914160
17	8	84529065864
18	8	48692406026
19	8	67772198982
20	8	75364316682

TABLE 1. NUMPY.RANDOM.DEFAULT_RNG SEEDS THAT PRODUCE THE LONGEST STREAKS WHEN USED TO SIMULATE ROLLING A 20-SIDED DICE

Player number	Gold pieces acquired	Notes
1	1,000	
2	1,000	Achieved through high-stakes gambling
3	6400	
4	8,000	
5	10,000	Achieved by robbing a vault full of treasure
6	20,000	
7	40,000	
8	100,000	
9	100,000	
10	500,000	
11	1,000,000	Achieved by buying Onion Futures
12	30,000,000	Achieved through high-stakes gambling
13	1,000,000,000	
14	2,560,000,000,000	Achieved through high-stakes gambling

TABLE 2. POTENTIAL PROFITS, AND METHODS OF ACHIEVING THEM (IF SHARED) BY INTERVIEWEES

Of the respondents, 6 were excluded for giving non-numeric answers, saying that they could get “infinity gold pieces ².” A 7th leather-clad interviewee was excluded when it became clear they were a different type of dungeon master.

Of the remaining 14 respondents, there was considerable variation in responses. The potential gold pieces varied from as low as 1,000 to as high as 2.56×10^{12} coins, with a statistical mean of 1.97×10^{11} , and a median of 70,000 gold pieces. For future analysis I will use the mean value, as it is funnier. Players also had the option to describe the plan they would undertake to get that much money, with most respondents turning to some form of gambling to find their profits. Table 2 illustrates all 14 valid responses.

2. It will be a cold day in hell before I recognize IEEE 754

4. Further Currency Conversions

4.1. Canadian Toonies

By considering the opportunity cost of performing this analysis, it also becomes possible to find direct currency conversions between d20 fantasy gold pieces and other forms of money. The easiest comparison is found when considering the energy cost of performing this analysis, and allows for a conversion to the coinage of the boreal kingdom of Canada, the toonie. The computer performing this analysis was running for 111 days, and used up 9.425 kilowatt-hours to do that. Thanks to Hydro Quebec's generous rate of 0.0325 toonies per kilowatt-hour, I spent 0.305 toonies to perform the analysis and generated 1.97×10^{11} gold pieces. This leads to a conversion factor of 149,300,829,941.41 toonies per gold coin. Furthermore, at the gold spot price of 47.86 toonies per gram, it follows that each gold coin must weigh 3,119,424.70 kilograms and have a volume of 161.77071 cubic meters ³.

4.2. Bitcoin

Instead of considering the opportunity cost of performing the analysis, I can also consider the opportunity cost of not using the computational resources of my Raspberry Pi on other things, such as mining Bitcoin. Fully clocked, a raspberry pi can achieve approximately 100 Hashes per second. This can be compared to the total number of hashes per second spent on mining Bitcoin to estimate the average amount of hashes before I would expect to mine a block. At time of writing, there are approximately 6.15×10^8 hashes per second of computational power spent on guessing-and-checking algorithms to mine Bitcoin. This means I would expect to have a 1.62698×10^{-7} chance of mining any given Bitcoin block. At an average rate of 144 blocks per day for 111 days, I would have expected to have gained 0.002577143 Bitcoin if I dedicated my computational power to that instead. This produces the exchange rate of 1 d20 gold piece being equal to $7.0982E+13$ Bitcoins ⁴. Finally, by comparing the results from these two sections, one easily see that 1 Bitcoin is worth approximately 0.00210 toonies, indicating that the present value of 47489.075 toonies per Bitcoin is slightly inflated ⁵.

5. Conclusion

This work has powerful implications for world economies everywhere. The powerful exchange rate of d20 gold coins to real-world money could prove a boon to all d20 fantasy tabletop role-playing game fans if any stock exchanges offer trading services. This work is not without shortcomings, however. An important point of consideration

is that this purely economic calculation missed out on some of the non-tangible benefits of playing d20 fantasy games. Unquantified are the human connections forged over rolling dice, doing silly voices and telling meaningful stories, and raises future research questions such as "What if the real fortune is the friends we made along the way?"

Acknowledgments

I would like to thank Tai Withers for their peer-review comments on this manuscript, as well as all of the people who've played TTRPGS with me over the years. I would also like to thank Pine trees and blackbirds for all of the joy they've brought into my life. Finally I would like to thank fungi, carrion beetles and other decomposers for the critical role they play in ecosystems everywhere.

3. For reference, an Olympic-sized swimming pool has a volume of 2500 cubic meters

4. For reference, there will only ever exist 2.1×10^7 Bitcoins

5. This is not financial advice

An introduction to bogoceptionsort and its performance compared to ordinary bogosort

Emil Sitell

Abstract

This paper introduces the bogoceptionsort. The bogoception tackles one of the major problems with bogosort, namely its relative efficiency at sorting very short arrays. This problem stems from the deterministic nature of the source code itself. In bogoceptionsort the source code for bogosort is bogosorted and then compiled, this then runs and checks if the array is sorted. If the array isn't sorted then the source code is scrambled again. The bogoception sort performs similarly to bogosort for large n but excels at smaller n and needs orders of magnitudes more iterations for small n .

1 Introduction

Through the ages hundred upon hundred algorithms for sorting lists and alike have been discovered, studied and perfected. From the naïve approach of bubble sort and insert sort to the more complex quicksort and radix sort, data scientist the world over search for better algorithms. One of the major breakthroughs in the field came when the revolutionizing bogosort was discovered, it combined beauty and elegance with ehmm.. speed. To this day it's one of the few algorithms with best case sorting of $\mathcal{O}(n)$. Many esteemed researchers have optimized the algorithm since its inception but in this paper one major new discovery is revealed and studied, the bogoceptionsort.

2 The algorithm

In ordinary bogosort the array is shuffled and checked, if the array is sorted then the algorithm returns, otherwise the array is shuffled again. See code snippet below:

```

is_sorted = False
while not is_sorted:
    random.shuffle(shuffled_array)
    for i in range(len(shuffled_array) - 1):
        if shuffled_array[i] > shuffled_array[i + 1]:
            raise UserWarning
    if i == len(shuffled_array) - 2:
        is_sorted = True

```

For smaller arrays this makes it very easy to accidentally sort the array in only a couple of tries. To combat this the bogoception sort starts by bogosorting the code for the bogosort. In this paper the scramble is limited to shuffling the order of the rows. The code is then converted to a string and executed with python's `exec`. If the code doesn't return a sorted list, the code is shuffled again and it tries again, see code snippet below. This includes both cases when an exception is raised as well as if the code is executed correctly but the underlying bogosort fails to sort the array. You can never be certain the code was right so might as well shuffle again.

```

is_sorted = False
while not is_sorted:
    random.shuffle(code)
    try:
        exec(formatted_code,
             {"bogo_test": unsorted, "random": random})
        is_sorted = True
    except:
        pass

```

3 Evaluation

By initially bogosorting the code, the algorithm essentially performs as *bogosort*($n+a$) where a is the number of lines in the original bogosort, in this case 8. Therefore it needs thousands of iterations to sort an array even with $n = 2$. See figure 1 to see the averaged performance for different values of n .

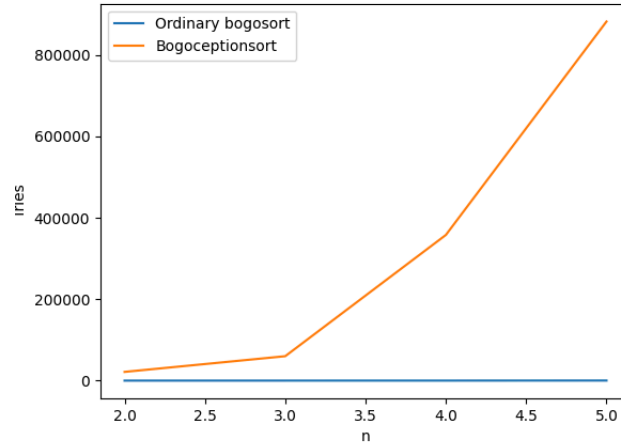


Figure 1: Iterations needed to sort an array of size n for bogosort and bogoceptionsort.

As seen in figure 1 the iterations needed for bogoception sort quickly grows very large. In theory when n gets larger $\lim_{n \rightarrow \infty} n + a \approx n$ which implies that bogoception sort converges to ordinary bogosort. Due to the factorial nature of bogosort's complexity it is theorized that this would still need a large n before they become comparable. Due to the inefficiency... *hark* technical limitations the research team was not able to verify this empirically.

4 Further studies and improvements

The algorithm shows great prospect of becoming one of the major players for sorting arrays, especially smaller arrays when the execution time is of utmost importance. The algorithm described in the paper is bogoception sort of degree 1, further studies into harnessing the power recursive programming the concept could easily be extended to bogosorting the bogoceptionsorts code. This could then be repeated until desired performance is reached.

SYSTEMS FOR RATING RATING SYSTEMS WITH RATING SYSTEMS

Clayton W. Thorrez

claytonthorrez@gmail.com

ABSTRACT

The field of rating systems is devoted to assigning numerical ratings to things based on data representing comparisons between those things, and using those ratings to create ordered rankings or predict future comparisons. For instance, Elo ratings are assigned to Chess players based the results of their games against other players, with higher ratings representing higher skill levels. However, when it comes to comparing rating systems to each other, we ironically defer to metrics such as “predictive accuracy” or “log-likelihood”. In this work we aim to clear up this oversight and propose a novel framework for the comparison and evaluation of the rating systems, by the rating systems, and for the rating systems.

1 RATING SYSTEMS BACKGROUND

A rating system is a method used to assign numerical ratings to items representing in some sense their strength when compared to other items. A rating system is described as “online” or “dynamic” if it is meant to track the skills of competitors as they change over time, this is the setting considered in the present work. Much of the early development of this field was motivated by the desire to rate chess players in order to form world rankings and properly match players to other players of equivalent skill at tournaments. The Elo rating system (Elo & Sloan, 1978) was developed in the 1960’s and remains in use in many applications to this day. Later, Glicko was developed as a Bayesian extension of Elo and applied to other competitions such as football games in the NFL. (Glickman, 1999) Subsequent work on online rating systems applied them to video games, extended to team competitions, and made various improvements to the update rules based on different statistical distributions, methods of approximation and optimization. (Herbrich et al., 2006; Weng & Lin, 2011)

1.1 RATING SYSTEMS EVALUATION

A variety of different methods have been used in order to compare rating systems with each other. The methods typically follow from the idea of making accurate predictions. This is commonly measured in terms of accuracy, log-likelihood or brier score on unseen data. Let n be the number of match-ups in the evaluation set, $\hat{p} \in (0, 1)$ be the predicted probability of the first competitor winning, $y \in [0, 1]$ be the true outcome with 1 representing a win for the first competitor and 0 representing a win for the second competitor, and $\hat{y} = 1[\hat{p}_i > 0.5]$ be the predicted outcome derived from the predicted probability.

$$\text{accuracy} = \frac{1}{n} \sum_{i=1}^n 1[y_i = \hat{y}_i] \quad (1)$$

$$\text{log-likelihood} = \frac{1}{n} \sum_{i=1}^n y_i * \log(\hat{p}_i) + (1 - y_i) * \log(1 - \hat{p}_i) \quad (2)$$

$$\text{brier score} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{p}_i)^2 \quad (3)$$

(sorry but you have to include some equations if you want people to take you seriously)

This is sometimes computed over the same data used to fit the ratings in a setup where predictions are first made for a comparison or batch of comparisons, then those comparisons are “trained on” by updating the ratings for the competitors based on the results. This method is valid as the predictions are not based on any data the model has trained on yet.

Like most machine learning methods, rating systems have plenty of hyperparameters which control things like how sensitive ratings are to new results, initial estimates of the rating means and variances, and numerical constants to calibrate predicted probabilities. If one does hyperparameter tuning (as one should if one wants to have good performance), then one should tune on only a subset of the data, and then evaluate on additional data that was not seen during the tuning and fitting phase. In online rating systems, the data is ordered and thus standard k-fold cross validation and random train-test splitting is not appropriate, instead we can only split based on time.

These evaluations, though useful, fall short of the full potential of rating systems research. When working in a field devoted to comparing things, it seems like a missed opportunity to not use your own methods when making your own comparisons. Glickman & Jones (1999) came tantalizingly close to with the title “Rating the Chess Rating System”, unfortunately the paper did not actually use a rating system to perform the ratings. In this work we devise novel methods of evaluation for rating systems, using the rating systems themselves.

2 UNDERSTANDING RATING SYSTEMS AS SYSTEMS TO BE RATED

The main difficulty with using rating system to rate rating systems is casting the evaluation as a time dynamic problem. In the aforementioned evaluation methods, given a dataset and a rating system, you get some number as the metric. If you have two rating systems, you can get two numbers and compare them and see which one is better but that’s kind of as far as you can go, you can’t effectively fit a rating system on a dataset with only one comparison. So the problem becomes how to get multiple comparisons, and not just any comparisons, they need to represent the change in strength over time so just adding some random noise and calculating the numbers a lot of times doesn’t really do anything. (we tried this) Understanding these constraints, we propose two systems for rating rating systems with rating systems. Both of these measure real properties of rating systems which to our knowledge have not been discussed at length in the prior literature and which are useful to understand in practical applications.

2.1 SYSTEM 1

This idea was inspired by the common situation in machine learning where the experimenter tries many versions of their model with different hyperparameters and can sometimes achieve large gains by finding a particularly well-suited configuration. With large numbers of hyperparameters, the search space can be large and therefore it is expensive to do a fine-grained grid-search over it. In this case it is a desirable quality of a model to be relatively invariant to the hyperparameters, and to give strong performance “out of the box” or “off the shelf”.

With this in mind, we formulate a setting where the time dimension is the number of iterations we spend sweeping over the hyperparameter space. We can compare the best metrics achieved so far for different models doing the sweeps. For an experimental set up with m candidate models and sweeping h hyperparameter configurations, we would have m choose 2 comparisons at each time step for a dataset of $h * \binom{m}{2}$ match-ups. In order to measure overall aptitude for hyperparameter tuning, this can be done over d different datasets and the results from all sweeps merged into a single sweep dataset. In practice we find it is useful to repeat the experiment k times with different random seeds used to sample from the hyperparameter space to reduce noise and stabilize the resulting ratings. (they can still be rather unstable though.)

2.2 SYSTEM 2

The scenario uses the same time dimension as the original rating system, but makes explicit comparisons across models of the metrics computed for each time step individually. In this setting, $h * \binom{m}{2}$ match-ups are created for each time step in the original setup. This setup measures several important properties of online rating systems: warm up and saturation.

2.2.1 WARM UP

The warm up problem is the measure of how quickly the model incorporates information in order to start making good predictions. In an uninitialized state, all competitors will have equal ratings and thus $\forall i \hat{p}_i = 0.5$. A model is said to have poor warm up if it takes a long time and a lot of data before the ratings and predicted probabilities start to accurately model the data.

2.2.2 SATURATION

Saturation is the opposite problem, some rating systems have mechanisms by which the magnitude of updates for players shrinks over time as they compete in more games. This has some benefits as it corresponds to our uncertainty about the player decreasing and can be seen as convergence. However for very large datasets over a long period of time, this can be a detriment as it is possible that the rating updates get too small for the model to be able to accurately change the true underlying changes in skill that occur.

3 EXPERIMENTS

We perform experiments using both of our proposed systems to rate four commonly used rating systems on real world datasets.

3.1 DATASETS

We use 3 real datasets spanning different time ranges, time scales, types of competitions, number of competitors and number of match-ups.

3.1.1 CHESS

We use a subset of the Chess data from the FIDE/Deloitte Chess Ratings Challenge. (Sonas, 2011) This is a dataset of millions of rated Chess games from “a recent 11 year period”. Due to computational limitations (we ran these experiments on a laptop), we limit our experiments to 200,000 games involving 15,992 players over 35 one month long rating periods. The data is available on kaggle.

3.1.2 NBA

We also accessed a dataset of historical basketball games from the National Basketball Association. (Paine, 2015). This dataset has 126,314 games involving 104 team from 1946 to 2015 which we split up into 1 year rating periods. This data was also accessed from kaggle.

3.1.3 LEAGUE OF LEGENDS (LOL)

This dataset contains matches from professional League of Legends games played both at both in person and online events from all over the world including international tournaments. It contains data covering 121,160 matches involving 12,675 teams from 2011 until 2024 and split into 681 1 week rating periods. The data was collected from the Leaguepedia API.

3.2 RATING SYSTEMS

We compare 4 commonly used rating systems both to rate the players and teams in our datasets, as well as to rate each other.

3.2.1 ELO

This is the original one by the OG Arpad Elo himself. It’s pretty simple but holds up considering how old it is and that it was only invented for Chess. Elo & Sloan, 1978

3.2.2 GLICKO

Glicko is Bayesian extension of Elo where each competitor has both a rating mean and a rating deviation to measure how uncertain we are about a competitor's rating. (Glickman, 1999)

3.2.3 TRUESKILL

This was developed by Microsoft in order to facilitate fair and balanced matchmaking for online Halo games. The main advantage is that it supports team games (which we are not evaluating in this work) and that it is based on a Gaussian skill distribution rather than the logistic distribution which is used in Elo and Glicko.

3.2.4 WENG-LIN THURSTONE-MOSTELLER

This one is super similar to TrueSkill (supports team play and is Gaussian based) but has some minor differences in how the player variances are updated. If you care that much go read the paper. (Weng & Lin, 2011) Honestly I'm including it because I have an implementation of it and 2 by 2 plots look nicer than 3 by 1.

3.3 SOFTWARE

The code and data for these experiments is available at <https://github.com/cthorrez/rs4rs>. (but good luck reading it lol) For the rating system implementations we use the open source python package **riix** (which I also wrote and maintain go check it out this is shameless self-promotion)

4 OK THE ACTUAL EXPERIMENTS NOW

4.1 SYSTEM 1

The methodology we used for this was to split each dataset 80/20 based on date for train and test. For each rating system, we defined hyperparameter spaces by low and high values for each parameter. We used the default values or values reported as the best in the experiments in their respective papers as a starting point and defined the ranges as one order of magnitude larger and smaller to simulate the situation where we truly have minimal knowledge of what reasonable values are. Since different methods have different numbers of hyperparameters, we employ uniform sampling over the entire legal space and use the same number of samples for each method. We used 100 samples, repeated over 5 random seed, for each of the 3 datasets, and each of the 4 rating systems creating a dataset of $100 * 5 * 3 * \binom{4}{2} = 9000$ match-up over 100 rating periods. The outcome for each match-up is determined by whether the accuracy on the test set for that specific dataset/hyperparameter/rating system is higher than the rating system it is compared against.

Once that dataset was created we do another grid search hyperparameter sweep over it and plot the evolution of the ratings over time for each rating system. The grid search range is the same as over the original datasets and are listed in the appendix.

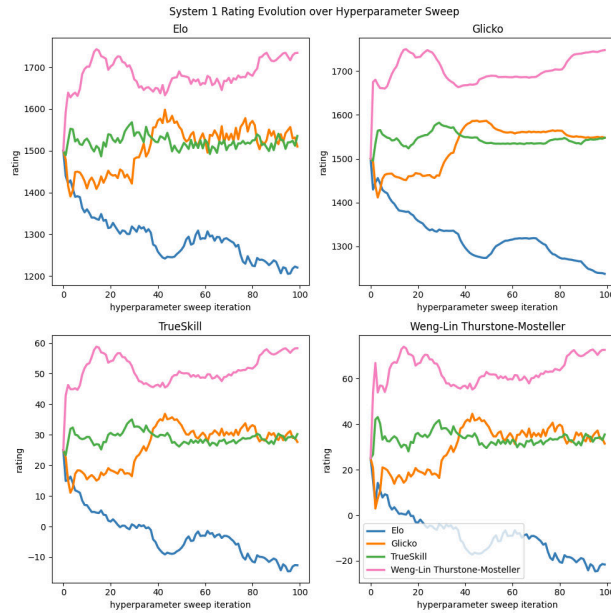


Figure 1: The ratings produced by each of the rating systems, on the rating data representing the evolution of relative performance during the hyperparameter tuning process.

We learn several things from 1. One is that relative to other models Weng-Lin Thurstone-Mosteller performs the strongest throughout the entire hyperparameter tuning process, and that Elo performs the worst. Additionally Elo loses performance through the process which means that it benefits the least from parameter tuning. While the ratings from all methods broadly agree, it is also interesting to note that the Glicko ratings are noticeably smoother than the others.

4.2 SYSTEM 2

For this system, we do not pool the datasets since the time dimension of the constructed dataset is the same as for the original, and the different datasets do now have the same time scales. Instead for each dataset we initialize rating systems with the best hyperparameters found in the sweeps during the System 1 experiments and fit them on the datasets logging the pre-match predictions and real outcomes. Then we do another pass and compute the predictive accuracy for matches only within a single rating period. Finally, for each combination of rating systems, we define the outcome of this synthetic match based on which has the higher accuracy for that time step.

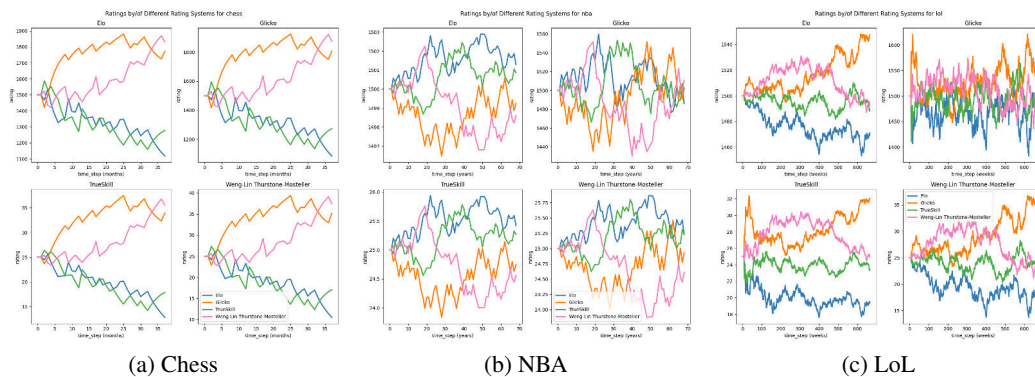


Figure 2: The ratings produced by each of the rating systems, on the rating data representing the evolution of relative predictive strength as they are trained on more data.

We learn some interesting things from the system 2 graphs as well. To begin, all of the ratings systems produce nearly identical (in shape, though different in scale) results for chess, this may be due to chess having the most matches and the fewest rating periods potentially leading to more convergence during each rating period. It is interesting to that Glicko and Weng-Lin beat Elo and TrueSkill so thoroughly, and that the margin increases as more data is trained on. This could indicate that TrueSkill has saturated. (It's actually mechanistically impossible for Elo to saturate so idk what is going on there it is possible that Elo is just weaker overall). Also it's possible we are seeing Glicko finally saturate at the end as it dips slightly and Weng-Lin overtakes.

The NBA data is just super noisy. This is due to many factors, the most obvious being that it is data over a much longer period of time, it is a team game so the total variance is larger due to the variance contributed by each player, and there are substantial changes to the competitors (teams) themselves, both within and across seasons as players change teams. We are fairly hesitant to conclude anything there, though if we were to cautiously suggest something to check more later it could be the idea that different eras of basketball had different features which made them more well suited as prediction tasks for different rating systems but that is kind of a stretch.

For League of Legends we see some interesting patterns:

1. Weng-Lin is the strongest roughly in the first half before Glicko overtakes
2. Weng-Lin strongly saturates and performance degrades to even lower than TrueSkill
3. Unlike in 1, Glicko produces the most noisy ratings here

The overall findings are that different rating systems are optimal for different types of competitions, and even for different stages of a competition.

5 FUTURE WORK (STUFF I RAN OUT OF TIME TO DO)

One factor we neglected in this analysis is runtime. For hyperparameter sweeps we allowed each method the same number of runs even though some methods are more expensive than others. Another interesting dimension could be to use wall-clock time as the time dimension to scale the results taking into account total compute required by each method.

The other obvious thing is to run this stuff on more datasets and more rating systems but yeah I'm just out of time.

6 CLOSING REMARKS

We identified a hole in the evaluation method currently used in rating systems research and found a way to fill that hole with even more rating systems.

AUTHOR CONTRIBUTIONS

It's all me baby

REFERENCES

- Arpad E Elo and Sam Sloan. The rating of chessplayers: Past and present. 1978.
- Mark E Glickman. Parameter estimation in large dynamic paired comparison experiments. *Journal of the Royal Statistical Society Series C: Applied Statistics*, 48(3):377–394, 1999.
- Mark E Glickman and Albyn C Jones. Rating the chess rating system. *CHANCE-BERLIN THEN NEW YORK-*, 12:21–28, 1999.
- Ralf Herbrich, Tom Minka, and Thore Graepel. Trueskill™: a bayesian skill rating system. *Advances in neural information processing systems*, 19, 2006.
- Neil Paine. How we calculate nba elo ratings, May 2015. URL <https://fivethirtyeight.com/features/how-we-calculate-nba-elo-ratings/>.
- Jeff Sonas. The deloitte/fide chess rating challenge, Jun 2011. URL <https://en.chessbase.com/post/sonas-the-deloitte-fide-che-rating-challenge>.
- Ruby C Weng and Chih-Jen Lin. A bayesian approximation method for online ranking. *Journal of Machine Learning Research*, 12(1), 2011.

A APPENDIX

A.1 HYPERPARAMETER SWEEP RANGES

For each hyperparameter the two numbers are the low and high value for sampling.

```
ELO_PARAM_RANGES = {
    'k' : (0.32, 320)
}

GLICKO_PARAM_RANGES = {
    'initial_rating_dev': (35, 3500),
    'c': (6.3, 630.0),
}

TRUESKILL_PARAM_RANGES = {
    'initial_sigma': (0.833, 83.33),
    'beta': (0.4166, 41.66),
    'tau': (0.00833, 0.833)
}

WL_TM_PARAM_RANGES = {
    'initial_sigma': (0.833, 83.33),
    'beta': (0.4166, 41.66),
    'kappa': (0.00001, 0.001),
    'tau': (0.00833, 0.833),
}
```

U. S. AIR FORCE
PROJECT RAND
RESEARCH MEMORANDUM

Just How Random?:
Introducing PROJECT S.P.O.R.K.

[REDACTED]
A. M. Lyons
[REDACTED]

RM-[rand();]

Rev. 1 29 March 2024

Assigned to SIGBOVIK 2024

This is a working paper. It may be expanded, modified, or withdrawn at any time. The views, conclusions, and recommendations expressed herein do not necessarily reflect the official views or policies of the United States Air Force.

The **RAND** Corporation

1700 MAIN ST. • SANTA MONICA • CALIFORNIA

Copyright, [REDACTED]
The RAND Corporation

SUMMARY

This memorandum is intended for policymakers, leaders in the armaments production industry, computer scientists, physicists, mathematicians, and frazzled parents at the end of their ropes.

It details the importance of good sources of random numbers and the history of the attempts on the part of the RAND Corporation to reliably create them in service of American hegemony and the destruction of communism. Particular attention is paid to one of the more surprising aspects of randomness: its dramatic algorithmic utility even for computational problems that are not themselves inherently random.

It also announces an exciting new breakthrough that promises to annihilate two organisms with a single munition, so to speak: PROJECT S.P.O.R.K. is a plan for extracting high-yield top-quality 100 percent "pure" "Colombian-quality" randomness from one of our country's most frustratingly vast and untapped resources: our lazy and embarrassing grandchildren and their ilk.

PROJECT S.P.O.R.K.: Let's put America's youth to work in operations research.

~~S~~ecret Special
~~P~~roject Playtime for
O perations
R esearch for
K ids

CONTENTS

I	BACKGROUND	3
II	EARLY RAND (OMNESS)	3
III	SPONTANEOUS RANDOMNESS	6
IV	PROJECT S.P.O.R.K.: LEL SO RANDOM	7

LIST OF ALGORITHMS

1	Who's the one that's acting childish? Are you sure?	8
2	Who you calling chicken?	8
3	Who you calling chicken?: Turbo Charged	9
4	Who you calling chicken?: Santa Monica Drift	9
5	We're 14 and this is deep	10
6	A Byzantween agreement protocol (BAP)	11
7	Gee mister, that's a lot of jelly beans	12
8	A protocol for the TEQUALITY problem	12

[REDACTED]
A. M. Lyons
[REDACTED]

(I) BACKGROUND

Randomness has always been a part of the human experience, and--according to some interpretations of quantum mechanics--is an intrinsic aspect of the universe itself.

Unfortunately for the patriots among us, the mathematical treatment of randomness and probabilities began in the casinos of 16th century Europe (which must have been unimaginably seedy). Thankfully, shortly after the second world war, randomness began to be applied in an *instrumental* manner, as part of algorithms that assisted with a new much nobler and higher purpose: gambling not for money but rather with the continued existence of humanity. Driven by fears of a soviet nuclear program, the mortal necessity of "keeping up with the Kuznetsovs" demanded accelerated development of the most terrific instruments of Armageddon imaginable, and all available resources were brought to bear. This included the finest minds, who developed the finest tools, which included advanced computational resources and, more importantly, the algorithms to make use of them. Enter Monte Carlo: Working at Los Alamos on the development of thermonuclear weapons immediately after the war (specifically, simulations regarding neutron diffusion in fission devices [11]), Stanisław Ulam and John von Neumann developed the first instantiation of an extremely powerful algorithmic paradigm based on random sampling. The use of Monte Carlo methods for large-scale computation on real computers was the first proof-of-concept of a counter-intuitive idea: that clever use of randomness can be a powerful algorithmic tool *even for problems that aren't inherently random*. How random is that? Pretty random, huh??

Fortunately, the development and production of our fantastic new weapons was concomitant with developments in game theoretic reasoning that provided clear mathematical (i.e., infallible) instructions for how to use them. The game of thermonuclear standoff, with mass slaughter as table stakes, was thus resolved, and communism defeated, in a way that was completely unproblematic and lead to world-wide prosperity and unprecedentedly high profits. What a success story! There weren't even any close calls.

How did we get there? The journey was neither easy nor completely unclassified: Almost immediately after it was introduced, Monte Carlo found its way (along with game theory, as it happens) to a group of consummate probabilists and gamblers in California who had big pockets and were eager to both put it to further use and use it to further America's interests globally.

(II) EARLY RAND (OMNESS)

The first "think tank", RAND has served the "needs" of the United States government and military-industrial partnership since the end of World War II. Officially, RAND stands for "research and development", though we argue that it might just as well stand for randomness, since that concept that was baked into the DNA of our organization from the beginning: One of the earliest RAND publications was a 4 page research memorandum entitled *Randomness.[sic]* by Olaf Helmer-Hirschberg [15]. This admirably efficient document solved the entire (titular) field definitively and completely, and was intended to be the final (note the full stop included in the title) word on the subject. This document has been lost.

Other early bangers include the following.

The Exploitation of Superstitions for Purposes of Psychological Warfare by Jean M. Hungerford [16]: This memorandum is illustrative of RAND's general ethos and modus operandi, which is characterized by a willingness (bravery, even) to apply a cold, actuarial view to any aspect of human affairs, especially when we are called upon to provide justification for some ethically questionable policy or military action.

Long-Lasting Effects of LSD on Certain Attitudes in Normals: An Experimental Proposal by William Hersche McGlothlin [22]: This one is just too much fun not to mention--the normies will never know what hit them!

The Survival Probability Problem by D. A. Darling¹ [6]: Check this shit out.

A critique of the method of Cunningham and Hynd with respect to solving the problem of finding the probability for a target to survive a burst of shots fired from a gun whose aim is a stochastic wandering process.

This example typifies not only our enthusiastic embrace of probabilistic analyses of all kinds, but also our pedigree: The conception of RAND occurred immediately after the war as the product of a beautiful and holy union between an Air Force general and a corporate executive, and was born (after gestating for just one month) as a special government contract to the Douglas Aircraft Company. Hence, the "target" whose survival probability was pondered by Darling was not a human being, but rather an aircraft attempting to evade ground-based fire² (with, presumably, some number of human beings inside of it).

Game Theory by Herman Kahn and Irwin Mann [18]: This was a draft chapter of a planned book titled *Military Planning In An Uncertain World*, which again highlights RAND's ongoing dedication to applying all the latest advances in mathematics and computation to military pursuits. When choosing a report to represent RAND's extensive foundational contributions to the field of game theory, we encounter an embarrassment of riches. Both Nash and von Neumann worked for RAND. The hallways and cafeteria hummed with talk of strategies, equilibria, payoffs, nuclear standoffs, preemptive strikes, and mutually assured destruction. Of course, all that dour business was never permitted to stifle our collective senses of humor, whimsy, and casual chauvinism, all of which are captured by the following excerpt.

The Trader and the Cannibal

Let us consider a completely different kind of game. Imagine for example, that you are a trader and are visiting Koko, chief of the cannibal island's gourmet club. You are in the following delicate situation.

You are going to give him a present of some beads. He is going to give you a present of some coconuts. If he considers his present more valuable than yours, he will be insulted and have you seasoned and cooked. If he feels that your present is equal in value to his he will do nothing. If he considers your present more valuable than his, he will feel that he has lost face and let you have an extra present, an evening with his wife (fat, greasy, and amorous), about whom you could not care less. Your only objective is to trade beads for coconuts.

History of RAND's Random Digits: Summary by George W. Brown [6]: The people at RAND recognized that large-scale computation (by machine) had proven itself during

¹Are these names starting to sound made up to you?

²We leave it to the reader to decide whether that makes it more or less metal.

the war, and would be a key tool for the mathematical problems and techniques relevant to their purview. Since so many of those computations made use of randomness either intrinsically or instrumentally, this meant that high-quality random numbers were now needed more--and *more of them*--than ever before. How many? One million.

This report, which was authored by the chief of the Numerical Analysis Department³, outlines a project began at RAND in 1947 to build a random number generator:

A random frequency pulse source was gated by a constant frequency pulse, about once a second, providing on the average about 100,000 pulses in one second. Pulse standardization circuits passed the pulses to a five place binary counter, so that in principle the machine is like a roulette wheel with 32 positions, making on the average about 3000 revolutions on each turn. A binary to decimal conversion was used, throwing away 12 of the 32 positions, and the resulting random digit was fed to an I.B.M. punch, yielding punched card tables of random digits.

After some additional mathematical processing involving sums (modulo 10), the digits were subjected to a battery of tests and then given a clean bill of health [4, 3]. Just how random are they? Very random! Random enough to fool that specific battery of tests, at least! The digits were subsequently fed on to practitioners, eventually made their way to the The On-Line Encyclopedia of Integer Sequences[25], and can now even be listened to [17].

In 1955, the digits were published as a book titled *A Million Random Digits with 100,000 Normal Deviates* [27], our gift to the world. One reviewer of the book [19] fancied himself a bit of a jokester:

As a final comment, one cannot help but be amused by the problem of proofreading the final tables to see whether the printing and reproduction mechanism has introduced "random" errors.

The book was reissued in 2001 [28] and continues to receive rave reviews. Here is an example of a 5 star review posted on Amazon [7]:

I can't understand all the negative reviews! This book literally contains everything I could ever ask for in a book. Recipe for spanokopita? Check! Name of every person ever born? Check! Next week's powerball, bingo, MLB, and NASCAR results? Check! By randomly combining and recombining the contents at random, I have read the works of Shakespeare, Harry Potter 8: the Tomb of Crying Stilton (to be released in 2014), the Bible AND the REAL Bible. I threw out my other books when I realized I could just jump around in this book and derive any other book I wanted. I think Borges wrote a story about this, but it's taking me a while to find that story in my book. I did find some steamy erotica this morning, though, so who's complaining?

This is indeed high praise, as the author is clearly a connoisseur of randomness.

The electronic random number generator was a great success, but we found that too much blood and oil were required to keep it in full working order, and that without expensive regular maintenance the digits it produced would devolve from very random to only sort of random. This development motivated our search for a new source of randomness. We should note that in the intervening decades, many interesting approaches have arisen that, for one reason or another, are not suitable for RAND's purposes. For example, the method known as "LavaRand" [24], which is based on the chaotic undulations of lava lamps, was deemed too hippie-dippie.

³In that role, George W. Brown was charged with building a powerful computer at RAND dubbed JOHNNIAC, which was based on the IAS machine (the first computer built on the so-called von Neumann architecture, which he also helped build).

(III) SPONTANEOUS RANDOMNESS

In roughly 2004, we began to hear reports of young people being "so random". How random is that? Well, working under the natural assumption that this phenomenon was probably just the result of more school libraries stocking our book of random digits, we concluded that the answer must be very random. But then an artifact emerged⁴ online and blew our ever-loving minds:

```
hi every1 im new!!!!!!! holds up spork my name is katy but u can call me t3h PeNgU1N
oF d00m!!!!!!! lol...as u can see im very random!!!! thats why i came here,
2 meet random ppl like me ^_^... im 13 years old (im mature 4 my age tho!) i
like 2 watch invader zim w/ my girlfreind (im bi if u dont like it deal w/it) its
our favorite tv show!!! bcuz its SOOOO random!!!! shes random 2 of course but
i want 2 meet more random ppl => like they say the more the merrier!!!! lol...neways
i hope 2 make alot of freinds here so give me lots of commentses!!!! DOOOOMMMM!!!!!!!!!!!!!!!
<-- me bein random again ^_^hehe...toodles!!!! love and waffles,

t3h PeNgU1N oF d00m
```

This passage is a thing of beauty. Our typesetting system can barely digest it, and we realized quickly that it could only be handled in the basement, because the girls charged with transcribing it kept cringing straight through the floor. It's *dripping* with entropy. It's lel so random. Our best estimates indicate that the Kolmogorov complexity of this string is the entire human genome.

This changed our minds about everything. If humans are capable of *spontaneously* producing this level of randomness, then this fact could be leveraged to simplify and accelerate our calculations of survival probabilities and missile trajectories, simulations of humanity-ending kinetic events, searches for equilibria, and all the rest of it. Imagine the paper that would be saved: we could throw away all the digit books, not to mention entire fields such as randomness extraction.

Perhaps a bit over-eager and optimistic, we immediately told our best R&D team to explore the possibilities and put this into use. Their assignment was completely open-ended: Be as random as possible. As random as you can be. You can imagine our disappointment when they came back to us with a bland corporate website:

<http://www.randomcorporation.com/>

An "export marketing company" specializing in "paper & paperboard products" and "oil & lubricants"? Uninspired. Mundane in the *worst possible way*.

After firing the entire team responsible, we consulted with psychologists and other experts to determine what had gone wrong. Ultimately, it was concluded that with age and study, the mind becomes laden with facts, experiences, associations, narrative plots of fictional and nonfictional artistic works, internalizations of social norms, and other artifacts of life. These things alone wouldn't be a problem, and might even help, but the problem arises in the workings of the brain itself, which distills and abstracts and re-factors these experiences into understanding, which is the death of randomness. Roughly, the process goes as follows. Speculative neural associations become, over time, well-worn Roman roads, and then become sunken superhighways that boast facilitation of "coherent thought" and "deductive reasoning", and offer conveyance to "the right answer". While occasionally useful, these highways don't allow the driver to pull over and, say, play the national anthem of Lithuania on a kazoo, throw a stringbean at a power transformer, or mix together every type of clear soda available at the convenience store and pour the mixture into an ant hill. In sum, it's not about creativity per se, but rather the *lack* of guard rails, experience, wisdom, and especially superhighways.

⁴We haven't been able to determine the exact provenance of the artifact, though it is rumored to have come from something called "4chan".

(IV) PROJECT S.P.O.R.K.: LEL SO RANDOM

After learning of the (severe) limits on the occurrence of spontaneous randomness in self-regarding adult humans working at places like RAND (and, in particular, being forced to recognize the concept that it was possible to be "too practical" or "too actuarial"), we sulked for a few months, and then came to a crucial realization: We had, in fact, discovered a way that our lazy, good-for-nothing grandchildren could actually contribute to the progress of humanity. We need the randomness and, as far as we can tell, the kids aren't doing anything better with their time. After some initial attempts to force the youths to produce digits were stymied by malicious compliance (rashes of 696969, 420, 07734, and 80085, specifically), a winning recipe was found.

We are pleased to announce PROJECT S.P.O.R.K.: a new paradigm in randomized computation. Why generate numbers and use those numbers in randomized algorithms, when we could just skip the middle-person and incorporate the randomness into the algorithms directly?

In what remains, we present our work as a series of algorithmic vignettes. The use of randomized algorithms has exploded in recent decades, extending to every area of computation, providing a diverse array of applications to choose from in order to demonstrate our success. With one exception (the first example, out of respect for the power and historical significance of Monte Carlo), we have chosen to focus specifically on those application areas where randomness is not only known to be algorithmically useful, but *provably and indispensably so*.

Note: One of the primary uses of randomness is cryptography. But since cryptography is about secrets, and children lack the requisite security clearances and thus can't be trusted, PROJECT S.P.O.R.K. has ignored this application area completely.

— Political Redistricting. One fascinating more recent application of Monte Carlo methods (that could be right up RAND's alley, if it serves our ideological goals (and the price is right)) is in political redistricting. Election outcomes depend on votes, yes, but they also depend (quite heavily in some cases) on maps. Political parties are acutely aware of this, and often try to exploit the redistricting process (which happens in response to census results) to ensure that maps are drawn in ways that favor their interests, a process known as gerrymandering.

Assessing the degree to which a particular map is gerrymandered is more difficult than it seems. While certain properties such as painfully contorted shapes, lopsided (or even wildly non-proportional) results, and a low number of competitive districts might seem suspicious, they are (for good reason, in fact) all unconvincing to courts who are occasionally charged with determining whether specific districtings are constitutional. In order to determine how unusual a map is, we'd like to compare it to its peers--the collection of all possible maps that could have been drawn. A number of techniques based on Markov chain Monte Carlo (MCMC) have been developed in order to explore this space. One such approach, due to Chikina, Frieze, and Pegden [5], uses MCMC to explore the "neighborhood" (in the Markov chain sense) of a given map in order to compute rigorous bounds on how much of an "outlier" the given map is.

This is quite useful, but the space of all possible maps (*in general*) could be better understood if we could (efficiently) draw samples from it uniformly at random, which is generally very difficult and cannot always be accomplished using current techniques. PROJECT S.P.O.R.K. (and the children of America) to the rescue!

Algorithm 1 Who's the one that's acting childish? Are you sure?

- 1: gather a group of children and give each of them 17 distinctly-colored crayons (do not use white--white does not count as a color when it comes to crayons, plus it confuses our scanners)
 - 2: hand each child a (letter-sized) piece of paper with an outline of Pennsylvania on it
 - 3: (this one is important) specifically do *not* instruct the children to draw a congressional redistricting plan for the Commonwealth of Pennsylvania
 - 4: assess each drawing by eye to determine whether it aligns with your particular policy preferences before passing it along as a "certified completely random sample of a 'typical' map"
-

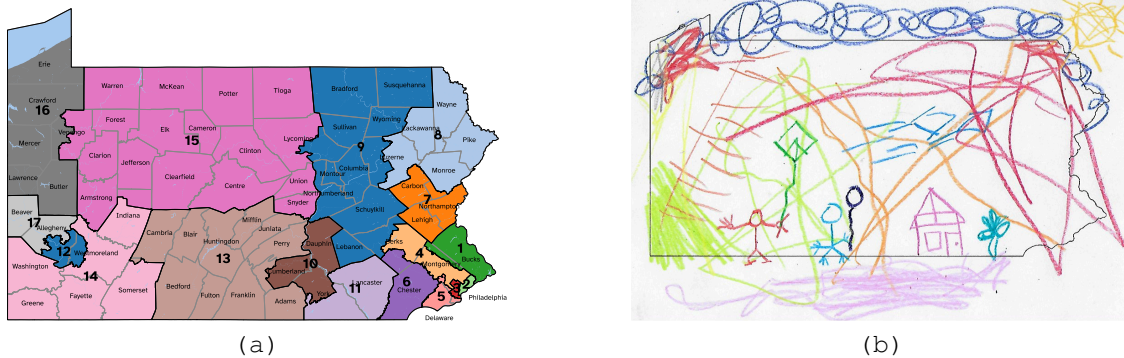


Figure 1: (a) 2023-2033 congressional districting for Pennsylvania (b) 2034 and beyond?

Source: (a) https://commons.wikimedia.org/wiki/File:Pennsylvania_Congressional_Districts,_118th_Congress.svg (b) a child named Sally

— The game of chicken. Imagine you and an opponent are in automobiles speeding towards each other on course for a catastrophic head-on collision. Now imagine how wonderful it would be--how much *payoff* you'd receive--if your opponent were to swerve first (i.e., "chicken out") while you held fast to your course like a badass. The payoff would be so great (up to 7(!!!)) that this scenario has received extensive attention and thorough analysis in the game theory literature.

In this context, randomness enables *mixed* strategies, which in some cases are necessary to reach equilibrium, and in other cases yield higher utility than the best pure (i.e., non-mixed, or deterministic) equilibria.

At PROJECT S.P.O.R.K., where randomness flows like water and is as pure as the driven snow, it's easy to experience the social welfare that arises from the mixed strategy equilibrium of the game of chicken, which has an expected utility of $4\frac{2}{3}$ for each player (compared with $4\frac{1}{3}$ for the two pure equilibria).

Algorithm 2 Who you calling chicken?

- 1: literally put two teenagers in fast cars (like those cool cars from that one movie) and have them accelerate towards each other
 - 2: each driver chickens out with probability $\frac{2}{3}$
 - 3: collect that sweet, sweet utility of $4\frac{2}{3}$ per driver
-

Teenagers always have a third friend who tags along and will do anything you tell them to. With the help of this third friend, and some ideas from Aumann [1], we can do even better:

		badass	chicken
badass	0	0	2
chicken	2	7	6

(a)



(b)

Figure 2: (a) Payoff matrix for the game of chicken. (b) A boring deterministic game of chicken with no symmetric equilibria and pathetically low expected utility--the kind of game only boring old fuddy-duddies who value their lives would play.

Source: (b) [https://commons.wikimedia.org/wiki/File:A_Game_of_Chicken_\(2802043436\).jpg](https://commons.wikimedia.org/wiki/File:A_Game_of_Chicken_(2802043436).jpg)

Algorithm 3 Who you calling chicken?: Turbo Charged

- 1: same setup as Algorithm 2, but now the cars have *turbochargers* and are equipped with *nitrous oxide* systems, so they're even cooler and faster
 - 2: a third teenager (who was instructed to ingest a troubling quantity of cough syrup the night before) advises the two drivers according to the dreams they had about the game: either both drivers are advised to chicken out, the first driver is advised to chicken out and the second is advised to be a badass, or the other way around (the second driver is advised to chicken out and the first is advised to be a badass), all with equal ($\frac{1}{3}$) probability
 - 3: the drivers now act in accordance with the advice they've received, which is a correlated equilibrium
 - 4: the expected utility for each driver is now 5 (where are we even going to put all that utility?!)
-

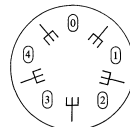
Want even more? You rascal.

Algorithm 4 Who you calling chicken?: Santa Monica Drift

- 1: setup is the same as Algorithm 3, but the nitrous oxide is injected into the lungs of the teenagers instead of the intake manifolds of the car engines
 - 2: the probabilities for the third teenager are adjusted to $\frac{1}{2}$ for (chicken, chicken) and $\frac{1}{4}$ each for (chicken, badass) and (badass, chicken)
 - 3: the correlated equilibrium now yields an expected utility of $5\frac{1}{4}$ for each driver (now *that's* something to laugh about!)
-

☛ Dining philosophers. This classic problem in distributed computing involves philosophers whose sole activities are thinking and eating, and whose desire to engage in one or the other switches intermittently and capriciously. As Dijkstra introduced it in the literature [8]:

Five philosophers, numbered from 0 through 4 are living in a house where the table is laid for them, each philosopher having his own place at the table:



Their only problem--besides those of philosophy--is that the dish served is a very difficult kind of spaghetti, that has to be eaten with two forks. There are two forks next to each plate, so that presents no difficulty: as a consequence, however, no two neighbors may be eating simultaneously.

The goal is to devise a procedure that the philosophers can use for securing their utensils when they wish to eat, but care must be taken to avoid *deadlocks* (which arise, for example, if all philosophers decide to eat simultaneously and each proceeds to secure the fork on their left, then waits for the right fork to become available).

Solutions in which all philosophers run the *same* algorithm are called *symmetric*, and randomness serves to "break" such symmetry: Lehmann and Rabin [21] proved that without a centralizing authority (i.e., restricted to "truly distributed" solutions) there is no deterministic algorithm that all the philosophers could run that would definitively preclude deadlocks⁵, and introduced a randomized algorithm that does the trick. Thus, in this setting, randomness is not only powerful but essential.

Now, behold the power of the (PROJECT) S.P.O.R.K.:

Algorithm 5 We're 14 and this is deep

- 1: 5 youths get 5 sporks and 5 servings of mom's spaghetti
 - 2: believe this: the youths will find a way to eat the spaghetti
-

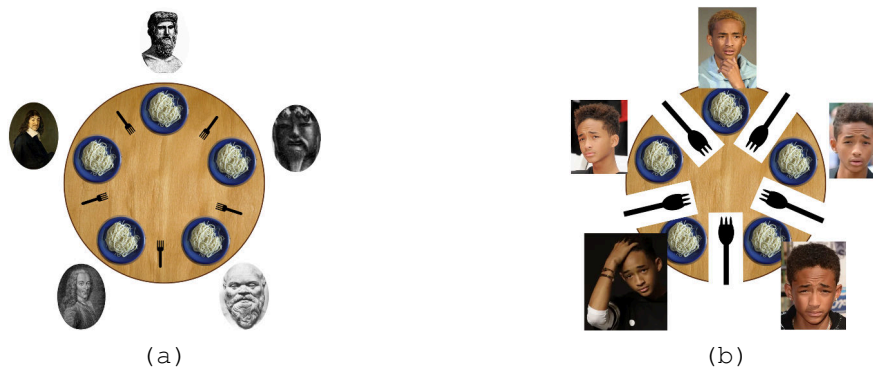


Figure 3: (a) Philosophers that behave deterministically and symmetrically can reach a deadlock. (b) "How can mom's spaghetti be real if Italy isn't real? How difficult does the spaghetti need to be to necessitate dual-wielding the sporks? Can't I just use one really strong spork? What does it even mean for spaghetti to be difficult? *Moooooooooom!??*"

Source: https://commons.wikimedia.org/wiki/File:An_illustration_of_the_dining_philosophers_problem.png

— Byzantine agreement. Another classic problem from distributed computing can be motivated as follows (one of the most epic opening paragraphs of a paper ever [13]):

We are in Byzantium, the night before a great battle. The Byzantine army, led by a commander in chief, consists of n legions, each one separately encamped with its own general. The empire is declining: up to $1/3$ of the generals--including the commander in chief--may be traitors. To make things worse, the loyal generals do not know who the traitors are. During the night each general receives a messenger with the order of the commander for the next day: either "attack" or "retreat." If all the good generals attack, they will be victorious; if they all retreat, they will be safe; but if some of them attack and some retreat they will be defeated.

⁵Some assert that Dijkstra proved an impossibility result for some variant of this problem, but we have not been able to find confirmation in the literature.

In this problem, the generals must communicate with each other in order to agree on a strategy, while the traitors act adversarially in attempt to foil this planning. Obviously, we here at RAND love this problem and think it's totally badass. We can't stop talking about how cool it would be to be a general and have a horse and big sword to kill people and so on.

A variant of the Byzantine Generals problem was shown by Fischer, Lynch, and Paterson [14] to be impossible to solve by any deterministic means. Randomized solutions were devised shortly thereafter by both Ben-Or [2] and (independently) Rabin [26]. Unlike the dining philosophers, the impossibility here holds even for asymmetric protocols (in which generals may behave disparately). For this problem, the critically useful role played by randomness is not to break symmetry but rather to achieve it by chance (though care must be taken to maintain it).

PROJECT S.P.O.R.K. cuts through this problem like a hot knife through hot butter:

Algorithm 6 A Byzantween agreement protocol (BAP)

- 1: Gather and arm some local youths
 - 2: position the youths around the home of a local communist sympathizer
 - 3: retain plausible deniability
 - 4: come what may
-

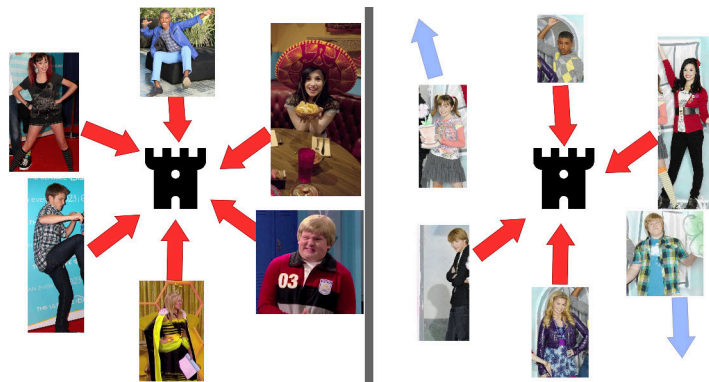


Figure 4: Results of a computer simulation of a proposed algorithm (Algorithm 6) for the Byzantine Agreement Problem in which the ensemble cast of the Disney Channel sketch comedy show **SO Random!** must coordinate in order to successfully sack and raze a fortified city in the late antiquity period. The city is conquered when the youths act more "randomly" (left); when the youths act more "normally" (i.e., more deterministically), failure occurs (right).

Source: https://commons.wikimedia.org/wiki/File:Byzantine_Generals.png

Volume of Convex Bodies. In 1988, Dyer and Frieze [9] showed that the problem of determining the volume of a polyhedron is $\#P$ -hard, meaning that there is *almost* surely no efficient way to compute the volume exactly. However, as noted by Metropolis and Ulam [23] decades earlier, the volume may be estimated easily using Monte Carlo methods.

However, when this problem is generalized to *arbitrary* convex bodies there is an even stronger result--one that does not rely on hardness assumptions like $\#P \not\subseteq BPP$ --that illustrates the power of randomness more conclusively. For *this* problem, Dyer, Frieze, and Kannan exhibited [10] a polynomial-time randomized algorithm capable of producing approximation guarantees that--as shown by Elekes [12]--are impossible to achieve by means of *any* polynomial-time deterministic algorithm.

At PROJECT S.P.O.R.K., this problem is literally child's play.

Algorithm 7 Gee mister, that's a lot of jelly beans

- 1: construct a jar in the shape of the body in question
 - 2: fill the jar with unit-volume pieces of candy
 - 3: instruct a group of children to come to a consensus regarding the number of pieces of candy in the jar
 - 4: throw the jar and the sweets it contains into the garbage
-



(a)



(b)

Figure 5: For children, being of assistance in the solution of difficult numerical approximation problems is a thrill. (a) A scientist looks on as a child submits his estimation of the volume of a convex body filled with pieces of candy as part of the execution of Algorithm 7. The child is dressed as some sort of supervillain--that's [haha so random!](#) (b) Children are likely to employ techniques that would never occur to an adult, such as whatever the hell is going on here.

Source: (a) <https://flic.kr/p/N2tvVg>; (b) <https://flic.kr/p/Xiu49t>

— Communication Complexity. In the string equality problem, two parties are each given binary strings of length n and are charged with determining whether or not the strings are identical. We assume they are in different locations and seek to minimize the number of bits they communicate in order to reach their determination. If the parties behave deterministically, then, in the worst case, at least n bits of communication is required. Parties that are permitted to behave randomly, however, can follow a protocol [20] that allows them to determine the answer with high probability using only $O(\log n)$ bits of communication.

Do you know any parties that behave randomly? PROJECT S.P.O.R.K. does:

Algorithm 8 A protocol for the TEQUALITY problem

- 1: translate strings $x, y \in \{0,1\}^n$ (bijectively) into strings $\phi(x), \phi(y)$ consisting of news regarding peer-group romantic pairings and pregnancy rumors, questionable fashion choices of perceived rivals, potential substance abuse use on the part of authority figures, etc. ("hot goss", or "the tea")
 - 2: distribute $\phi(x)$ and $\phi(y)$ to two different teens and connect them via telephone or FaceTime video interlink
 - 3: sit back as the teens instinctively and efficiently detect and negotiate any possible differences between $\phi(x)$ and $\phi(y)$ (and thus x and y)
-

REFERENCES

- [1] Robert J. Aumann. "Subjectivity and correlation in randomized strategies". In: *Journal of Mathematical Economics* 1.1 (Mar. 1974), pp. 67–96. issn: 0304-4068. doi: 10.1016/0304-4068(74)90037-8.
- [2] Michael Ben-Or. "Another Advantage of Free Choice: Completely Asynchronous Agreement Protocols (Extended Abstract)". In: *Proceedings of the second annual ACM symposium on Principles of distributed computing* (Aug. 17–19, 1983). Ed. by Robert L. Probert, Nancy A. Lynch, and Nicola Santoro. PODC '83. Montreal, Quebec, Canada: Association for Computing Machinery, Aug. 1983, pp. 27–30. doi: 10.1145/800221.806707.
- [3] Bernice B. Brown. *Some Tests on the Randomness of a Million Digits*. Paper P-44. Santa Monica, CA: RAND Corporation, Oct. 1948. URL: <https://www.rand.org/pubs/papers/P44.html>.
- [4] Bernice B. Brown. *Tests of the Randomness of Digits*. Research Memorandum RM-38. Santa Monica, CA: RAND Corporation, May 1948. URL: https://www.rand.org/pubs/research_memoranda/RM38.html.
- [5] Maria Chikina, Alan Frieze, and Wesley Pegden. "Assessing significance in a Markov chain without mixing". In: *Proceedings of the National Academy of Sciences* 114.11 (Feb. 2017), pp. 2860–2864. issn: 1091-6490. doi: 10.1073/pnas.1617540114.
- [6] D. A. Darling. *The Survival Probability Problem*. Research Memorandum RM-448. Santa Monica, CA: RAND Corporation, 1950. 7 pp. URL: https://www.rand.org/pubs/research_memoranda/RM448.html.
- [7] B. Dean. *Verified Purchase Customer Review*. Sept. 7, 2010. URL: <https://www.amazon.com/review/ROJ4561SZPMET> (visited on 03/29/2024). Rev. of RAND Corporation. *A Million Random Digits with 100,000 Normal Deviates*. MR-1418. Santa Monica, CA: RAND, 2001. isbn: 978-0-8330-3047-4. doi: 10.7249/MR1418. URL: https://www.rand.org/pubs/monograph_reports/MR1418.html (visited on 03/29/2024).
- [8] E. W. Dijkstra. "Hierarchical ordering of sequential processes". In: *Acta Informatica* 1.2 (1971), pp. 115–138. issn: 1432-0525. doi: 10.1007/bf00289519.
- [9] M. E. Dyer and A. M. Frieze. "On the Complexity of Computing the Volume of a Polyhedron". In: *SIAM Journal on Computing* 17.5 (Oct. 1988), pp. 967–974. issn: 1095-7111. doi: 10.1137/0217060.
- [10] Martin Dyer, Alan Frieze, and Ravi Kannan. "A random polynomial-time algorithm for approximating the volume of convex bodies". In: *Journal of the ACM* 38.1 (Jan. 1991), pp. 1–17. issn: 1557-735X. doi: 10.1145/102782.102783.
- [11] Roger Eckhardt. "Stan Ulam, John von Neumann, and the Monte Carlo Method". In: *Los Alamos Science, Special Issue* 15 (1987), pp. 131–137.
- [12] G. Elekes. "A geometric inequality and the complexity of computing volume". In: *Discrete & Computational Geometry* 1.4 (Dec. 1986), pp. 289–292. issn: 1432-0444. doi: 10.1007/bf02187701.
- [13] Pesech Feldman and Silvio Micali. "An Optimal Probabilistic Protocol for Synchronous Byzantine Agreement". In: *SIAM Journal on Computing* 26.4 (Aug. 1997), pp. 873–933. issn: 1095-7111. doi: 10.1137/s0097539790187084.
- [14] Michael J. Fischer, Nancy A. Lynch, and Michael S. Paterson. "Impossibility of distributed consensus with one faulty process". In: *Journal of the ACM* 32.2 (Apr. 1985), pp. 374–382. issn: 1557-735X. doi: 10.1145/3149.214121.

- [15] Olaf Helmer-Hirschberg. *Randomness..* Research Memorandum RM-5-PR. Santa Monica, CA: RAND Corporation, 1947. 4 pp. URL: https://www.rand.org/pubs/research_memoranda/RM5.html.
- [16] J. M. Hungerford. *The Exploitation of Superstitions for Purposes of Psychological Warfare*. RM-365. ASTIA Document Number ATI 210637. Santa Monica, CA: RAND Corporation, Apr. 1950. URL: https://www.rand.org/pubs/research_memoranda/RM365.html.
- [17] Doug Irving. *What Sounds Do a Million Random Digits Make?* Sept. 7, 2010. URL: <https://www.rand.org/pubs/articles/2022/what-sounds-do-a-million-random-digits-make.html> (visited on 03/29/2024). Sonification of RAND Corporation. *A Million Random Digits with 100,000 Normal Deviates*. MR-1418. Santa Monica, CA: RAND, 2001. ISBN: 978-0-8330-3047-4. DOI: 10.7249/MR1418. URL: https://www.rand.org/pubs/monograph_reports/MR1418.html (visited on 03/29/2024).
- [18] Herman Kahn and Irwin Mann. *Game Theory*. Paper P-1166. Santa Monica, CA: RAND Corporation, July 30, 1957. URL: <https://www.rand.org/pubs/papers/P1166.html>.
- [19] Carl F. Kossack. In: *Science* 122.3167 (1955), pp. 471-471. DOI: 10.1126/science.122.3167.471.b. Rev. of RAND Corporation. *A Million Random Digits with 100,000 Normal Deviates*. Glencoe, Illinois: The Free Press, 1955.
- [20] Eyal Kushilevitz and Noam Nisan. *Communication Complexity*. Cambridge University Press, Dec. 1996. ISBN: 9780511574948. DOI: 10.1017/cbo9780511574948.
- [21] Daniel Lehmann and Michael O. Rabin. "On the advantages of free choice: a symmetric and fully distributed solution to the dining philosophers problem". In: *Proceedings of the 8th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*. POPL '81. Williamsburg, Virginia: Association for Computing Machinery, 1981, pp. 133-138. ISBN: 089791029X. DOI: 10.1145/567532.567547.
- [22] William Hersche McGlothlin. *Long-Lasting Effects of LSD on Certain Attitudes in Normals: An Experimental Proposal*. Paper P-2575. Santa Monica, CA: RAND Corporation, May 1962. URL: <https://www.rand.org/pubs/papers/P2575.html>.
- [23] Nicholas Metropolis and S. Ulam. "The Monte Carlo Method". In: *Journal of the American Statistical Association* 44.247 (Sept. 1949), pp. 335-341. ISSN: 1537-274X. DOI: 10.1080/01621459.1949.10483310.
- [24] Landon Curt Noll, Robert G. Mende, and Sanjeev Sisodiya. "Method for seeding a pseudo-random number generator with a cryptographic hash of a digitization of a chaotic system". U.S. pat. US5732138A. RPX Corp Morgan Stanley and Co LLC. Mar. 24, 1998.
- [25] OEIS Foundation Inc. *Entry A002205 in The On-Line Encyclopedia of Integer Sequences*. URL: <https://oeis.org/A002205> (visited on 03/29/2024).
- [26] Michael O. Rabin. "Randomized Byzantine Generals". In: *24th Annual Symposium on Foundations of Computer Science* (Nov. 7-9, 1983). FOCS '83. Tucson, Arizona, USA: IEEE Computer Society, Nov. 1983, pp. 403-409. DOI: 10.1109/SFCS.1983.48.
- [27] RAND Corporation. *A Million Random Digits with 100,000 Normal Deviates*. Glencoe, Illinois: The Free Press, 1955.
- [28] RAND Corporation. *A Million Random Digits with 100,000 Normal Deviates*. MR-1418. Santa Monica, CA: RAND, 2001. ISBN: 978-0-8330-3047-4. DOI: 10.7249/MR1418. URL: https://www.rand.org/pubs/monograph_reports/MR1418.html (visited on 03/29/2024).

TO (pronouced “two”)

- 17 **Sleepy Dwarfs Somniloquy on Drowsy Logic Chip Design**
Giovanni Lostumbo

- 18 **The Magic School Bus Travels to Sub-threshold Voltage**
Dr. Henry Jekyll

- 19 **Badness 0 (Epsom’s version)**
Dr. Tom Murphy VII, Ph.D.

- 20 **Badness 0 (Knuth’s version)**
Dr. Tom Murphy VII, Ph.D.

- 21 **Flaccid Drives: Storing Arbitrary Data in SIGBOVIK Articles**
Nora Hyades Mirza, Head Archivist

- 22 **Retraction of ”Flaccid Drives”**
Nora Hyades Mirza, Head Archivist

Sleepy Dwarf’s Somniloquy on “Drowsy Logic” Chip Design

Giovanni Lostumbo

Imagineer Department eV Renewable
Photons

Walt Elias Disney Chair Professor Emeritus
National Imagination Laboratory

+1(708) 303-8175

giovanni.los@proton.me

Abstract

Power-efficiency has been an increasing design consideration in virtually all new silicon in the past 15 years. Power-first designs, however, typically appear only in niche applications such as IoT. A 2023 retrospective paper describing a research lab’s 2002 circuit, using a technique called “drowsy logic,” reviewed historical strategies to limit leakage in the context of foundries’ recent implementation of low-leakage FinFET and Gate-All-Around technologies.^{2,3} This review explores new research and additional industry applications of drowsy logic.

CCS Concepts

• General → • System Architecture → Power efficiency

B.3.2 [Cache Memories]: Design optimization for low power

Keywords

Cache memory, Low power, Gate leakage

1. Introduction

In *Computer Architecture Techniques for Power-Efficiency* (2008), drowsy circuits are defined as “a new class of state-preserving leakage reduction techniques.”^{4a} When factored into the dynamic power equation $P = CV^2 A f$, it shows that there can be a cubic reduction in power without a proportional reduction in performance,^{4b} especially in memory-bound or latency-tolerant regions of code.^{4c} When combined with sub-threshold voltage, it allows static power to be greatly reduced.^{4d,4e} The original 2002 technique involves a strategy, called the “simple” policy, of placing all lines in a drowsy mode using a single global counter, awakening it only when it is accessed.³ The performance trade off was known to reduce leakage up to 85% while increasing run-time by just 0.62% in certain conditions (using 93% drowsy lines). The paper focused on advanced drowsy strategies in reducing latency due to L1’s time-critical cache, but suggested L2 strategies could use the simpler techniques.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.
© 2024 Copyright held by the owner/author(s).

2. A Comparison to Previous Techniques

2.1 Vdd gating, Cache Decay, and Adaptive Mode-Control

The 2002 paper and the 2023 retrospective paper contrasts drowsy logic with three previous techniques to reduce leakage developed in 2000 and 2001. The first, called *gated-Vdd*,⁵ used a circuit-level technique to gate supply voltage to reduce leakage in unused SRAM. That technique, paired with a novel resizable cache architecture (DRI i-cache), was said to reduce leakage by 62% with a minimal impact on performance. The second, Cache Decay, used a time-based strategy to turn off a cache line after a pre-set number of cycles have elapsed since its last access.⁶ This method is said to reduce L1 leakage up to 70% using competitive algorithms. The third, Adaptive Mode Control (AMC), used tags that are always active, tracks which cache lines are missed, and can adjust the number of intervals to turn off cache lines based on previous misses.⁷ Cache Decay was also recently reviewed by their authors in the 2023 ISCA@50 Retrospective.⁸

2.2 A Systems-Level Design

A 2008 paper titled “BTB Access Filtering: A Low Energy and High Performance Design” describes lowering branch target buffers and using direct-mapped BTBs in superscalar processors with drowsy techniques in the filter buffer to limit predictor energy consumption by 92.7%, with up to a 10.8% performance trade-off⁹. Early ARM processors, such as ARM7 and ARM9, did not use superscalar architecture and likely did not need large buffers.¹⁰ Static branch prediction was used, however, in power-conscious processors such as the ARM810.^{11,12} Significant research has been explored in developing more power-efficient superscalar architecture.^{4f}

2.3 Direct-Mapping

Direct-mapping of hardware registers appears to have a similarity to MMU-less operating systems, such as μ Clinux¹³, developed by Jeff Dionne and Ken Albanowski in 1998, which used flat memory addresses and was further developed by companies such as EmCraft.¹⁴ In software, Cortex-M processors feature interrupt/exception handling, where instead of automatically putting floating point registers onto the stack, can be configured to do so in a “lazy” way.¹⁵⁻¹⁷ While most processors are designed to be fast, or to reduce latency, they also benefit from using tricks like that to make them less resource intensive. Direct mapping of both hardware and kernel resources offers a path to limit energy-intensive memory management caches. As a chip can be designed to operate in a more deterministic manner as in real-time

operating systems (RTOS), use cases involving known application resource limits and scheduling can be increasingly factored into software-defined hardware (SDH), a decades old design goal.^{18,19} One design concept that monolithic chips use is “holistic timing.”²⁰ This is where multiple systems fit on a single die and operate within a certain window, including breaking partitions in the clock. A Power-First design would most certainly benefit from opting towards fewer Die-to-die interfaces, which typically increase power consumption. An advantage to using chiplets, such as Bunch of Wires, however, would be the parallel or co-design and integration of peripherals such as lower power radios and I/O by third parties.²¹ While sub and near-threshold processors are a major focus of power reduction, memory can occupy 2-8x+ the size of a single microcontroller such as Zero-riscy (now Ibex) or RI5CY (cv32e40p), according to PULP Platform of ETH Zürich.²² General purpose operating systems require far more memory and cache than microcontrollers. Efficient design of both hardware and software using the above techniques can lessen the amount of memory (and thus power) needed.

3. Drowsy logic in modern process nodes

Industry news articles appear to have decreased mentions of “drowsy” circuits in the mid 2010s. A keyword search of one well-known site turned up only two mentions in 2014 and 2015 for “drowsy:” “For memories, people are building additional operational modes for them, such as drowsy modes.”^{23,24} A review in the ICISA@50 retrospective paper states drowsy logic is “a form of voltage scaling”: “We proposed a design in which one can choose between two different supply voltages in each cache line, corresponding to normal supply voltage and a drowsy lower voltage. In effect we used a form of voltage scaling to reduce static power consumption.”²² While the 2002 paper cites voltage scaling precedents, its novel feature was that it was applied to static power: “Such a dynamic voltage scaling or selection (DVS) technique has been used in the past to trade off dynamic power consumption and performance. In this case, however, we exploit voltage scaling to reduce static power consumption.”²³ In other words, the term became part of DVS from the start. New techniques utilizing hybrid or novel implementations of drowsy transistors in SRAM continue to be researched in academic labs.²⁵⁻²⁷ Techniques sometimes have multiple industry names. As early as 1992, multi and dual-threshold voltage techniques have been described.^{28,29} Drowsy logic still offers some of the best energy savings: “Moreover, since the penalty for waking up a drowsy line is relatively small (it requires little energy and only 1 or 2 cycles) and there are less frequent accesses to the lower memory hierarchy than Gated-V_{DD} schemes, cache lines can be put into drowsy mode more aggressively to save more power.”²

4. Sub-threshold logic

The 2002 drowsy cache paper cites advances in short-channel effects: “Due to short-channel effects in deep-submicron processes, leakage current reduces significantly with voltage scaling. The combined effect of reduced leakage current and voltage yields a dramatic reduction in leakage power.”³ A 1990 paper developed the new model for short-channel effects and their benefit to sub-threshold logic.³⁰ A key advance was that the alpha power law model was reduced from $\alpha = 2$ to $\alpha = 1.3$. This was implemented in the calculation of the switching speed:^{4d}

$$\text{Delay } \alpha \propto 1/I_{\text{on}} \propto V_{\text{dd}}/(V_{\text{dd}} - V_{\text{T}})^{\alpha}$$

Martin et al applied the exponent, α , of the alpha power delay model of an inverter, with $\alpha=1$, to the dynamic power law to produce the formula for performance, f :^{4c}

$$f = (L_d K_6)^{-1} ((1 + K_1) V_{\text{dd}} + K_2 V_{\text{bs}} - V_{\text{th1}})^{\alpha}$$

The combined formula of DVS and Adaptive Body Biasing (ABB), led to a further 48% reduction in energy over DVS (drowsy logic) alone six months later.³¹ Other labs have achieved similar results.^{32,33} Today, an 85% reduction in leakage power can be found in IoT devices which are designed for batteryless operation and energy-harvesting. Microcontrollers by companies such as Ambiq Micro are known to achieve a 13-fold reduction compared to other chips.³⁴ IoT device-makers such as ONiO feature an integrated solar/RF/thermoelectric harvester and power management integrated circuit with low-power, asynchronous ROM/RAM.³⁵ Since the latter uses just 2KB of RAM, it may not need to operate in sub-threshold voltage as much as Ambiq’s 2MB of MRAM, which can still achieve ultra low energy consumption.³⁶ Furthermore, the ultra low power of sub-threshold voltage combined with current drowsy cache techniques suggests mobile phones could one day run on solar power.

5. SRAM and in-memory computing

While drowsy modes have been developed for both instruction and data cache, in high-performance computing (HPC), it may not yet have some of the advantages of state-of-the art memory such as spin-transfer torque (STT) MRAM and spin-orbital torque (SOT) MRAM.³⁷ That is because they are designed to operate in low-data modes where speed is not as crucial to operation such as remote sensors with fixed interval telemetry. That said, the fast-wake up is known to be 1-2 cycles.²

6. Conclusion: Amdahl's & Landauer's Limits

The 2002 paper cites Amdahl’s Law in calculating the theoretical minimum.³ While it does acknowledge advances in short-channel effects and subsequent research detailed new avenues for advancing Moore’s Law,³⁸ Amdahl’s Law is more relevant to parallel multi-core architectures, of which sub-threshold voltage microprocessors do not necessarily adopt: “The lowest supply voltage at which a logic gate can operate while still acting as an amplifier is only a few times larger than $k_B T / q$.”³⁹ As Scott Hanson recently stated that “Moore’s law is alive and well for the embedded world. We’re at a process node today that’s 22 nanometers,” one can speculate on the performance gains yet to come.³⁶

References

- ¹ <https://semiengineering.com/a-power-first-approach/>
- ² “Drowsy Caches: Simple Techniques for Reducing Leakage Energy—A Retrospective” Krisztián Flautner, Nam Sung Kim, Steve Martin, David Blaauw, Trevor Mudge; ISCA@50 25-Year Retrospective: 1996-2020
https://bpb-us-w2.wpmucdn.com/sites.coecis.cornell.edu/dist/7/587/files/2023/07/drowsy_retro.pdf
- ³ “Drowsy Caches: Simple Techniques for Reducing Leakage Power” Krisztián Flautner, Nam Sung Kim, Steve Martin, David Blaauw, Trevor Mudge (2002)
<https://web.eecs.umich.edu/~manowar/publications/drowsy-caches-ISCA2002.pdf>
- ⁴ *Computer Architecture Techniques for Power-Efficiency*, Stefanos Kaxiras & Margaret Martonosi (‘Ch. 5.3)(‘Ch. 3.1)(‘Ch. 1.2.1)(‘Ch. 5.1.1)(‘Ch. 5.4.1)(‘Ch.4.5.1) Morgan & Claypool Press 2008
- ⁵ “Gated-Vdd: A Circuit Technique to Reduce Leakage in Deep-Submicron Cache Memories,” 2000 Michael Powell, Se-Hyun

- Yang, Babak Falsafi, Kaushik Roy, and T. N. Vijaykumar
<https://engineering.purdue.edu/~vijay/papers/2000/gatedvdd.pdf>
- ⁶ “Cache Decay: Exploiting Generational Behavior to Reduce Cache Leakage Power” 2001 Stefanos Kaxiras, Zhigang Hu, Margaret Martonosi
<https://mrmgroup.cs.princeton.edu/papers/hzg-isca2001.pdf>
- ⁷ “Adaptive Mode Control: A Static-Power-Efficient Cache Design” Huiyang Zhou, Mark C. Toburen, Eric Rotenberg, Thomas M. Conte 2001
<https://prod.tinker.cc.gatech.edu/journal/zhou03adaptive.pdf>
<https://pdfs.semanticscholar.org/ba53/2536b7d1890f967d1aac7485544f6abebe39.pdf>
- ⁸ “RETROSPECTIVE: Cache Decay: Exploiting Generational Behavior to Reduce Cache Leakage Power” 06/2023 https://bpb-us-w2.wpmucdn.com/sites.coecis.cornell.edu/dist/7/587/files/2023/06/Hu_2001_Cache.pdf
- ⁹ “BTB Access Filtering: A Low Energy and High Performance Design” Shuai Wang, Jie Hu, and Sotirios G. Ziavras Department of Electrical and Computer Engineering New Jersey Institute of Technology (2008)
https://web.archive.org/web/20090920084956id_/http://web.njit.edu:80/~sw63/pub/ISVLSI_BAF_2008.pdf
- ¹⁰ https://en.wikipedia.org/wiki/List_of_ARM_processors
- ¹¹ https://web.archive.org/web/20130926155924/http://www.eetimes.com/document.asp?doc_id=1208831 “VLSI Technology Now Shipping ARM810” 08/1996
- ¹² https://web.archive.org/web/20181224080542/https://www.hotchips.org/wp-content/uploads/hc_archives/hc08/2_Mon/HC8.S4/HC8.4.1.pdf “ARM810: Dancing to the Beat of a Different Drum” 07/1996
- ¹³ <https://en.wikipedia.org/wiki/%CE%9Clinux>
- ¹⁴ <https://www.emcraft.com/>
- ¹⁵ <https://developer.arm.com/documentation/dai0298/a>
- ¹⁶ https://docs.zephyrproject.org/latest/hardware/arch/arm_cortex_m.html
- ¹⁷ <https://docs.zephyrproject.org/latest/kernel/services/other/float.html>
- ¹⁸ “Deterministic Clock Gating for Microprocessor Power Reduction” 2003
<https://engineering.purdue.edu/~vijay/papers/2003/dcg.pdf>
- ¹⁹ “Software-Defined Hardware Architectures” 05/2023
<https://semiengineering.com/software-defined-hardware-architectures/>
- ²⁰ “Why Chipllets Don’t Work For All Designs” 09/2023
<https://semiengineering.com/why-chipllets-dont-work-for-all-designs/>
- ²¹ “Is UCIE Really Universal?” Arteris IP, 11/2022
<https://semiengineering.com/is-ucie-really-universal/>
- ²² <https://pulp-platform.org/community/showthread.php?tid=229&pid=650#pid650> forum post, 2021
- ²³ “S-L Power Modeling Gains Steam” 08/2014
<https://semiengineering.com/system-level-power-modeling-activities-get-rolling/>
- ²⁴ “With Responsibility Comes Power,” 02/2015
<https://semiengineering.com/with-responsibility-comes-power/>
- ²⁵ “Hybrid Drowsy SRAM and STT-RAM Buffer Designs for Dark-Silicon-Aware NoC”, 2016 Jia Zhan, Student Member, IEEE, Jin Ouyang, Member, IEEE, Fen Ge, Member, IEEE, Jishen Zhao, Member, IEEE, and Yuan Xie, Fellow, IEEE
<https://cseweb.ucsd.edu/~jzhaio/files/darlsilicon-noc-tvlsi2016.pdf>
- ²⁶ “Design the efficient SRAM circuit using 4transistor with sleepy logic” International Journal of Pure and Applied Mathematics Volume 118 No. 20 2018, 115-123
<https://acadpubl.eu/hub/2018-118-21/articles/21b/15.pdf>
- ²⁷ “Design and Implementation of Low Power SRAM Using Highly Effective Lever Shifters”, 2021
<https://scholar.uwindsor.ca/cgi/viewcontent.cgi?article=9892&context=etd>
- ²⁸ “Low Power CMOS Design” A.P. Chandrakasan; S. Sheng; R.W. Brodersen, 1992
https://mitsites.mit.edu/researchgroups/icsystems/pubs/journals/1992_chandrakasan_jssc.pdf
- ²⁹ “Design and Optimization of Low Voltage High Performance Dual Threshold CMOS Circuits” Liqiong W ei, Zhanping Chen, Mark Johnson, Kaushik Roy & Vivek De, 1998
<https://dl.acm.org/doi/pdf/10.1145/277044.277179>
- ³⁰ “A JSSC Classic Paper: The Simple Model of CMOS Drain Current” 10/2004
https://www.eng.uburn.edu/~agrawvd/COURSE/READING/LOWP/alpha_power_law.pdf
- ³¹ “Combined Dynamic Voltage Scaling and Adaptive Body Biasing for Lower Power Microprocessors under Dynamic Workloads”, 11/2002
<https://ieeexplore.ieee.org/document/1167611>
- ³² “Combined Dynamic Voltage Scaling and Adaptive Body Biasing for Heterogeneous Distributed Real-time Embedded Systems” Le Yan, Jiong Luo and Niraj K. Jha ICCAD’03, November 11-13, 2003
https://web.archive.org/web/20050223092839id_/http://www.princeton.edu:80/~lyan/pub/iccad03.pdf
- ³³ “Impact of process scaling on the efficacy of leakage reduction schemes” 2004 Yuh-Fang Tsai, David Duarte, N. Vijaykrishnan, Mary Jane Irwin
<https://www.micromagic.com/news/icicdt04final.pdf>
- ³⁴ “What’s All This Subthreshold Stuff, Anyway?” 02/2019
<https://www.electronicdesign.com/technologies/analog/article/21807652/whats-all-this-subthreshold-stuff-anyhow>
- ³⁵ “What if You Never Had to Charge Your Gadgets Again?” 01/2024 <https://www.wsj.com/tech/personal-tech/what-if-you-never-had-to-charge-your-gadgets-again-955ea960>
- ³⁶ “Interview With Scott Hanson - Founder and CTO at Ambiq” 01/2024 <https://www.safetydetectives.com/blog/scott-hanson-ambiq/>
- ³⁷ “TSMC tandem builds exotic new MRAM-based memory with radically lower latency and power consumption” 01/2024
<https://www.tomshardware.com/pc-components/dram/tsmc-tandem-builds-exotic-new-memory-with-radically-lower-latency-and-power-consumption-mram-based-memory-can-also-conduct-its-own-compute-operations>
- ³⁸ “Near-Threshold Computing: Reclaiming Moore’s Law Through Energy Efficient Integrated Circuits” January 2010
<https://ieeexplore.ieee.org/document/5395763>
- ³⁹ *Recent Progress in Boolean Logic*, Bernd Steinbach, Vincent C. Gaudet, 2013 (4.1.3, 204), (4.1.5, 211)

The Magic School Bus Travels to Sub-threshold Voltage

By Dr. Henry Jekyll, Caviendish Laboratory, Cambridge

How an obscure, 22 year old circuit can limit leakage near the theoretical minimum where just the electrons that count- count

Five years ago, Electronic Design covered a new, low-power microcontroller that used an analog concept found in integrated circuits as far back as the 1968 CMOS-based RCA CD4007 (now made by TI). It described how the conventional wisdom regarding leakage currents does not apply at the lowest sub-threshold voltages. The article did a thorough history on the research and progress of the technology from analog to digital circuit applications.²⁶ The latter application, in digital circuit design, remains a widely unknown technology. It has even taken myself over 3 years to better understand. A simple analogy might work, but from experience, one analogy may not fit all point of views. The concept of drowsy logic has a connection to sub-threshold voltage, although that link will be detailed more later in the article. What is first important to review, is why these techniques are used- the ultimate goal may not be energy efficiency itself, but what new products or cost savings energy efficiency can allow. The Electronic Design article cited a 2016 dissertation³³ by Dr. Sunhil Kim, who cited 0.55V as the optimum Power Delay Product of 7.2 (p.28/78), achieving a 99% reduction in power:

Table 3.2. Power and delay for various voltages of the multiplier

Supply Voltage (V)	Total Power (μ W)	Delay (ns)	Power-Delay Product
0.20	0.05	4700	235
0.40	0.15	120	18.0
0.45	0.60	41	24.6
0.50	1.40	12	16.8
0.55	1.80	4	7.20
0.60	6.00	1.60	9.60
0.80	110	0.18	19.8
1.00	640	0.07	44.8

To start with this table, a 2020 text defined energy efficiency not by a power-delay-product⁷⁵, but by an energy delay product. A 2013 text states, "However, subthreshold voltage devices may prefer PDP metrics where leakage can present an operational or medical risk."⁶¹

"The proposed work uses energy-delay product (EDP), where energy the total energy consumption of cores and delay is the amount of time for executing applications. Spiliopoulos et al. [6] used EDP and ED^2P , where execution time was computed as the average execution time over multiple programs that run simultaneously on their multicore systems. Murray et al. [9] used EDP. Tarplee et al. [12] defined Pareto frontiers to tradeoff energy consumption versus execution time. Jung et al. [13] measured energy consumption for different classifiers. Lai et al. [14] computed EDP for each phase of their applications to optimize an EDP target. Wang et al. [10] used profit as explained above. Chen et al. [15] used power consumption.⁶² EDP efficiency could also be analogized to computational benchmarks, such as MIPS or DMIPS(Dhrystones), divided by $Power(W)^{54}$. Depending on the instruction sets and application (e.g. HCI-MIPJ/glyph⁷⁶), synthetic benchmarks may not always be comparable to other architectures⁷⁷.

Nonetheless, now that 0.55V as a local minima can be selected as point-of reference for sub-threshold logic, it is important to contrast this circuit from the point of view of the user, rather the circuit-side. In short, digital logic is like a light switch turning on and off, communicating as a Morse lamp, but glowing afterwards for a few seconds. While there are various types of bulbs that have multiple reasons for this⁶⁵, such as a hot filament in an incandescent, or phosphorous coating interacting with residual ions of mercury in a CFL, or even phantom voltage on LEDs with poor insulation or lacking earthing wire, the most analogous example to a computer transistor is the capacitance inside an LED circuit in a normal operation. Instead of needing a 120V AC current to power a 60 Watt incandescent, an LED running at 12V DC could power a 0.75 Amp bulb that uses 9 Watts whilst outputting the same number of lumens (~ 800 @2700K).

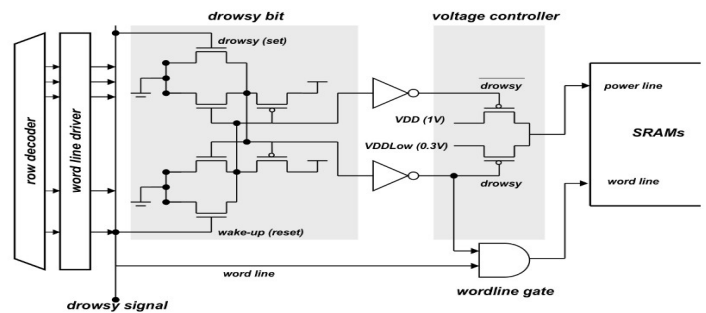
Isolating the cells is important, to not only reduce leakage and soft-memory errors from capacitive cross-talk³⁶, but also to improve speed by the same factor²⁶. With an ultra low leakage foundry process node such as 22nm

ULL,³⁷ low-voltage transistors still operate like transistors or switches, but may operate in weird and idiosyncratic ways, requiring more non-Gaussian LVF modeling and prediction in the EDA stage³⁸.

In optimal conditions, setting a voltage that maintains the operations of transistors at the lowest power level would only be limited by the amount of performance needed by the application or operating system. IoT devices have utilized commercial sub-threshold microcontrollers for just over 10 years. More complex processors are just starting to be developed. Alternative energy-efficient chips such as static cores³⁹ utilize other techniques to save power, although might consume approximately the same amount if needed to run a heavy workload. Furthermore, past static designs were more common in the 1980s-90s, typically operate at higher threshold voltages, and new process nodes have produced few, if any new designs compared to dynamic logic.

A drowsy circuit can be analogized to a light switch being turned off and on repeatedly by a playful child. While the voltage switch itself appears to be exploited for fun, the circuit itself has significant amperage waiting as the circuit is technically "on," with sufficient electrons to maintain state, IPC and the instruction cache utilized. Thus the average voltage using a 0.55V threshold setting might be $\frac{1}{2}$ of that, or 0.275V, if the voltage switch was on 50% of the time. The slide on next page depicts Ambiq's threshold voltage example as 0.5V³⁴. However, to the human eye, a flickering LED might appear more like a strobe light, even with residual glow from capacitance leaking into the LED. The flashes are like "accessed" drowsy bit lines. Capacitance builds up before it reaches the transistor, but waits less time in a series-based short-channel MOS with lower leakage when set to a lower voltage.⁵⁶ The 2 extra transistors needed for the drowsy circuit add complexity but save power³:

FIGURE 2. Implementation of the drowsy cache line



Note that, for simplicity, the word line, bit lines, and two pass transistors in the drowsy bit are not shown in this picture.

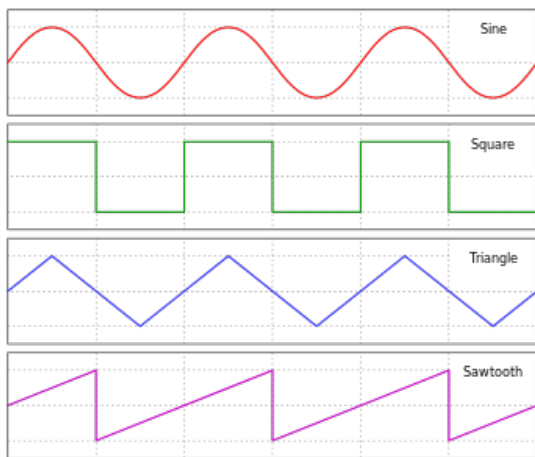
To better understand how the analogy works at the circuit side, it is more apt to use an expressway analogy during rush hour traffic. An observer on an overpass may be able to see cars on an otherwise 65 mph speed limit expressway traveling at 30 miles per hour. Assume, for the sake of simplicity, that is a passenger car commuting or traveling for leisure in between every 18 wheeler fully loaded with commercial goods (e.g. en route to a Big Box store that arrived from a container ship yard). The passenger cars represent electrons that are bit padding, not doing useful work necessarily (unless the instructions are CISC/Out-of-Order), whereas the trucks are carrying very important or valuable data for the user or customers on a network (shoppers). Thus the tracelines in the chip's IPC are where electrons are traveling with significant amperage- not as much as 500mV (65mph), but still a lot with an average net transfer of imports to an inland location over a given amount of time (e.g. 12 hours). The switch that is turned off effectively introduces traffic, causing 4 lanes to run at 30mph, also similar to the throughput of two lanes being open in non-rush hour traffic at 65 mph. Since amperage is not based on width or height, the analogy of 2 lanes running parallel or above or besides two others is less relevant than the number of semi-trucks per hour carrying "bytes" of data at a given speed. Thus, to an observer on an overpass, the 30mph traffic may appear to be slow, but they are typically focused on the speed of the individual vehicles rather than the throughput and cumulative number of bytes (goods) of all the vehicles on that route.

The speed of electrons in a wire travel signal velocity, which is “somewhat close to speed of light in a vacuum”³⁹ or 60%⁴¹, rather than the drift velocity. In any case, it is not possible for an observer to see the individual electrons, and a classical text book chemical bond diagram of H₂O is not suggesting that 2 electrons of Oxygen are permanently bonded to Hydrogen or even temporarily necessarily for longer than a very small fraction of a second. In that context, the location an electron occupies in atom or molecule is more accurately represented by a valence orbital, which describes the *probability* that an electron will be found: s, p, d, and f orbitals⁴³. Drowsy logic, however, is not anywhere near a probabilistic processor⁴⁴, which calculates the probability of a bit being 1 instead of 0, rather than a certainty. As stated above, Ambiq Micro microcontrollers (32-bit) still operates like a CMOS transistor and switch, even with some variance or yield losses. It is by this analogy that measurements in most, if not all of engineering uses SI units instead of Gaussian cgs, and Heaviside-Lorentz⁴⁵ units instead of factors of 4π. SI units operate on the concept of stochastic probability, rather than exact quantities as in Gaussian units.

To get a close-up view of an actual drowsy circuit, one might try to use an electron ptychography microscope⁴⁶, which reconstructed the highest resolution of an atom crystal, PrScO₃ in 2021. However, electrons are 100 million times smaller than atoms, and due to Heisenberg’s uncertainty principle, cannot be found. A multimeter would measure electrons far more practically.

If one were to use the Magic School Bus as an example, shrinking to the size of the individual bits transported at signal velocity, one might be able to determine the properties of the circuit’s amperage, and whether programmed dynamic voltage⁴⁷ drops in fixed intervals timed to reliably renew the circuit before decaying provides enough throughput for IPC and cache. In other words, the fundamental properties of the electrons are relied upon for optimizing the lowest leakage circuit. Electrons are not sentient like observers, and cannot know that a voltage switch caused a threshold to be reached for power on/off state.

Therefore, the Magic School Bus example might make more sense if the drowsy circuit were operating in a SeaLab in the Mariana Trench, and the light switch was being frequently turned on and off aboard a research boat above and connected via a marine-grade cable to the Sea-Lab. Ms. Frizzle’s class at the bottom of the Mariana trench would not be able to hear the student’s light switch atop the boat turn off (which would travel much slower than the optical cable, if it could even reach audible levels) nor when the switch was being flipped, but they would be able to see the downstream effect on the ebbing and cresting waveform⁵¹:



(Source: Wikimedia)

By this analogy, then, the dynamic scaling of voltage transforms analog controls on a circuit into a digital and predictable operation, despite the feared wave function collapse⁷⁴. Like the old joke, if 4 pies are brought to a Heaviside-Lorentzian forest, does anyone eat them? If there is an infinitely small gap in a mathematical set, it is still known as a continuous function⁵².

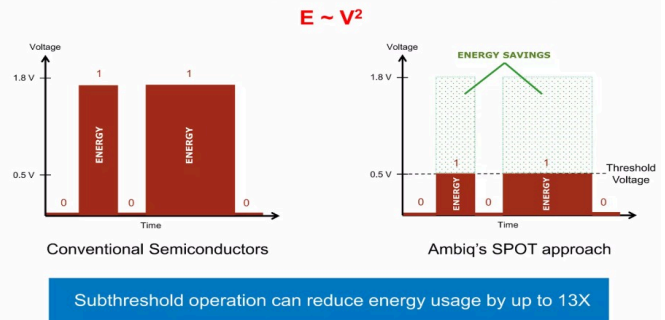
At most, it might create a soft-error in a redundant and/or inconsequential bit. At worst, as in genetics, it could cause a missense mutation. Thus exploring whether drowsy logic is on the threshold of “discontinuous” functions can better classify the limits of sub-threshold technologies. Like an underwater cave or SeaLab, the presence of an air pocket is the exception rather than the rule in the ocean of high impedance sub-threshold voltage.

Since there is not an infinite amount of oxygen (supply voltage) in an underwater cave, unique applications need to be designed, much like calculating the maximum biomass of animals(s) or insects that could survive in a hypothetical cave under the Mariana trench, as would a sealed terrarium would need to allow sufficient photosynthesis from sunlight to replenish the oxygen consumed by aerobic bacteria and other insects. The added difficulty, is that at the bottom of the sea, there is no light, although that analogy does not need to be taken literally or be completely congruous. It could also be that there are enough anaerobic bacteria that are able to produce oxygen as a byproduct for other species in the air pocket. Part of the reason IoT is used as an experimental, and low-cost test bed for new and energy-scarce devices is that they can be marketed for non-essential applications such as wearables until the technology improves its error-correction techniques.

Early computer power supplies were not designed with undervoltage protection. Inexpensive and unsafe electronics may also be missing certain voltage loss detection circuitry, which could lead to damaged components. Improved power supplies also have circuitry that can also prevent data loss, but also depend on software that is designed to save data either to disk either as soon as it is written to RAM, or directly accessed and modified in place (XiP). Static processors, and some emerging research, uses ferrostatic FRAM^{48, 49} and magnetoresistive RAM (MRAM) to save data. Whether this is crucial to facilitating the use of drowsy and sub-threshold voltage processors remains to be seen. Exploring whether sub-threshold processors can operate without interruption, despite their idiosyncratic quirks and potential voltage drops, may be able to extend the life of conventional memory technologies despite having higher power-consumption. Potentially controlled “leakage” could be designed to recover lost capacitance using

Subthreshold Power Optimized Technology (“SPOT”)⁵³

Energy (E) consumed is directly proportional to the square of the voltage (V) used:



supervisory circuits⁵⁰, although most leakage is lost as heat.

Slide from a 2015 Presentation (Source: Slideshare³⁴)
 If that $E \sim V^2$ resembles Einstein's relativity equation $E = mc^2$, that's because Ohm's law is as important as that in Electricity⁶⁸, and follows analogous special relativity variances for non-inertial frames of reference^{28,29}. Rewritten as a measure of resistance, $V^2=IR$, which also means that that $\sqrt{IR} = V$ ⁷⁰. Amperage decrease exponentially when resistance remains the same and voltage decreases linearly. The trade-off allows, for example, 1GHz microprocessors to be underclocked at 1/5th of the speed but achieve more than 10x the energy savings (sometimes 13x). In other words, more performance for a non-linear (exponential) reduction in power. This also means that as amperage approaches a zero asymptote, resistance increases exponentially, and the Power Delay-Product at a given voltage (e.g 200mV), would operate instructions far less efficiently (235) than at a PDP at 0.55V (e.g. 7.2). Some concepts: “Power increases as the square of the voltage with $P = IV$ ”^{69,71}

“The power by definition* is the product of the current and the voltage across a one-port (a two-terminal device). The power is only proportional to the square of the voltage if the $I(V)$ relationship is linear.”⁶⁷ “you need to understand that this law is not a linear equation, it’s rather an instantaneous fact that stands, meaning at every differential time step, the product of resistance and current would give you the Voltage.”⁶⁷ By determining what a given circuit’s threshold voltage should be, the voltage rail can be allowed to decay without providing conventional state retention:

“Drowsy logic, an alternative to putting the processor into sleep mode, is applicable where leakage currents are managed. Muller said that it is possible to turn off the power and let the power rail voltage decay. But it is not necessary to provide explicit state retention as long as the voltages on flipflops do not go too low. Essentially, it is a method of intermittently providing voltage sufficient to maintain state and operations.”³⁵

In *Computer Architecture Techniques for Power-Efficiency*, by Stefanos Kaxiras and Margaret Martonosi (2008), Chapter 5.3 confirms that drowsy circuits are a new class of “state-preserving” leakage reduction techniques.⁶⁰

In *Recent Progress in Boolean Logic* (2013), Bernd Steinbach similarly describes sub-threshold logic: “In this case, voltage scaling may be used quite aggressively, i.e., to a point where the supply voltage is so low that transistors are never really ON.”⁶¹

How do these sub-threshold transistors operate, if they are not really “on”? Electrons propagate in one direction alongside a conductor in an electromagnetic (EM) wave once current is applied to a wire.⁶³ The “state-preserving” aspect of this propagation apparently would not immediately get shutdown or revert to drift velocity⁶⁴, as long as the voltage is flipped on again before the current runs out.

The quadratic and exponential decrease in power below sub-threshold voltage is less immediately useful for higher performance microprocessors, at least without significant code/instruction optimization and parallelization. The number of microcontrollers that use just a few Megahertz can achieve the biggest energy consumption gains in non time sensitive applications. That too, reflects on the similarity of $E \sim V^2$ to the Special Relativity equation.

By applying this concept, another analogy can be made: If a 2nd Earth were positioned in the same orbit as Earth, but at the March Equinox while the Earth was a the Summer Solstice, perpendicular to the axis that it intersects with the sun, it would have a clear view of the Earth in the dawn sky. Suppose a light ray left the Sun and, was somehow tagged with a very bright and cartoonishly large fluorescent arrow. A person on the 2nd earth could use a telescope and track the movement of this arrow pointing to that single light ray. The earth is 93 million miles away from the Sun. Einstein’s 2nd postulate in Special Relativity states:

“The speed of light in vacuum is the same for all observers, regardless of the motion of light source or observer.”⁶⁴ It takes 8 and 1/3 minutes for a light ray to reach (both) Earths. Thus, from a distance of 93 million miles, even the speed of light might appear to be traveling slowly. A 5 Gigahertz processor today might sound much faster than a 50Mhz processor from 1994. However, the speed of individual electrons in both circuits travel at the same speed, like a feather and a hammer being dropped on the moon with no atmosphere. Less throughput (weight/cargo~bytes), but no difference in wind resistance, as there is none. To the telescope observer, the light ray is traveling at the same speed for both observers on Earth and the 2nd Earth.

Hypothetically, if the observer were the size of the Magic School Bus and millimeters away from the ray’s trajectory, the electron would pass the observer far too quickly to be observed. Returning to the Mariana Trench analogy, the child aboard the boat knows when the light switch is being turned off and on, since he/she is the analog control of the switch. But he/she/they cannot see the effects on the circuit that the rest of the class is observing (assuming with a multimeter or microscope that can capture moving electrons).

Thus, the light switch (DVFS) in sub-threshold does not have the same effect as on above threshold/conventional circuits, because due to impedance and threshold requirements to maintain the circuit, the circuit in effect operates with a small variation in frequency, or could even be set to a fixed clock speed. Designing a microprocessor with “complexity-effective” IPC⁵³ throughput is the next challenge of the future power efficient designs.

Power First already using Drowsy Logic?

Power-efficiency has been an increasing design consideration in virtually all new silicon in the past 15 years. Power-first¹ designs, however, typically appear only in niche applications such as IoT. A recent paper describing a researcher’s 2002 circuit, called drowsy logic, reviews² historical strategies to limit leakage in the context of foundries’ recent implementations of low-leakage FinFET and Gate-All-Around technologies.

The original 2002 technique³ involves a strategy, called the “simple” policy of placing all lines in a drowsy mode using a single global counter, awakening only when it is accessed. The performance trade off was known to reduce leakage up to 85% while increasing run-time by just 0.62% in certain conditions (using 93% drowsy lines). The paper focused on advanced drowsy strategies in reducing latency due to L1’s time-critical cache, but suggested L2 strategies could use the simpler techniques.

A 2008 paper⁷ titled “BTB Access Filtering: A Low Energy and High Performance Design” describes lowering branch target buffers and using direct-mapped BTBs in superscalar processors with drowsy techniques in the filter buffer to limit predictor energy consumption by 92.7% with up to a 10.8% performance trade-off. Early ARM processors⁸, such as ARM7 v3 and ARM9 v4, did not use superscalar architecture and would likely not need large buffers. Static branch prediction was used, however, in power-conscious processors such as the ARM810^{9,10}.

Direct-mapping of hardware registers appears to have a similarity to MMU-less operating systems, such as μ Linux¹¹, developed by Jeff Dionne and Ken Albanowski in 1998, which used flat memory addresses and was further developed by companies such as EmCraft¹².

Direct mapping of both hardware and kernel resources offers a path to limit energy-intensive memory management caches. While a chip can be designed to operate in a more deterministic manner, such as in real-time operating systems, use cases involving known application resource limits can be increasingly factored into software-defined hardware (SDH¹³), a decades old design goal. One design concept that monolithic chips use is “holistic timing.”¹⁴ This is where multiple systems fit on a single die and operate within a certain window, including breaking partitions in the clock.

A Power-First design would most certainly benefit from opting towards fewer Die-to-die interfaces, which typically increase power consumption. An advantage to using chiplets, such as Bunch of Wires, however, would be the parallel/co-design and integration of peripherals such as lower power radios and I/O by third parties. While sub and near-threshold processors are a major focus of power reduction, memory can occupy 2-8x¹⁵ the size of a single microcontroller such as Zero-riscy (now Ibex) or RISCY (cv32e40p), according to PULP Platform of ETH Zürich. General purpose operating systems require far more memory and cache than microcontrollers. Efficient design hardware and software can lessen the amount of memory (and thus power) needed.

In software, Cortex-M processors feature interrupt/exception handling, where instead of automatically putting floating point registers onto the stack, can be configured to do so in a “lazy” way^{4,5,6}. While most processors are designed to be fast, or to reduce latency, they also benefit from using tricks like that to make them less resource intensive.

Compared to Vdd gating, Cache Decay, and Adaptive Mode-Control

The 2002 paper and the 2023 retrospective paper contrasts drowsy logic with three previous techniques to reduce leakage developed in 2000 and 2001. The first, called *gated-Vdd*¹⁶, uses a circuit-level technique to gate

supply voltage to reduce leakage in unused SRAM. That technique, paired with a novel resizable cache architecture (DRI i-cache), is said to reduce leakage by 62% with a minimal impact on performance. The second, Cache Decay, uses a time-based strategy to turn off a cache line after a pre-set number of cycles have elapsed since its last access¹⁷. This method is said to reduce L1 leakage up to 70% using competitive algorithms. The third, Adaptive Mode Control (AMC), uses tags that are always active, tracks which cache lines are missed, and can adjust the number of intervals to turn off cache lines based on previous misses¹⁸. Cache decay was also recently reviewed by their authors in the ISCA@50 retrospective¹⁹.

Variation in process can be attributed to environmental factors such as temperature, but techniques to adapt voltage to temperature formed one of the first strategies hybridize drowsy and cache decay⁶⁰. A significant amount of variation can also be found in the instruction cache: “This is because instructions exhibit strong temporal and spacial locality and any delay in fetch shows up immediately on performance.”⁶⁰

The aforementioned research in BTBs⁷ is also mentioned, using the hybrid approach. An advantage of microcontrollers having an immediate application for simple drowsy policies is that most of their operations use real-time operating systems with known scheduling, which can optimize the instruction cache and IPC buffer.

The 2008 text interestingly cites a 2000 paper by Zyuban and Kogge to describe power-inefficiency in out-of-order architectures (4.5.1)^{60, 73}, but also describes some of the innovative approaches to resource partitioning (4.5.2) and examines some of the disadvantages of earlier research in filter caches (4.10.2), as concurrently researched by Ziavras et al⁷.

An avenue of research that could be worth exploring more are the “idle” cycles of instructions.⁶⁶ While modern operating systems have developed advances in idle-task scheduling for background processes such as SSD garbage collection, anti-virus scans, often using an idle core(s), in low-power microcontrollers and microprocessors, these processes, if used at all, would constitute a larger fraction of processors’ clock cycles. Nonetheless, the co-development of software and hardware from the ground up remains an important consideration to optimize towards the most energy and power efficiency.^{54,56}

Compared to fast-to-wake, fast-to-compute, fast-to-sleep

While drowsy modes have been developed for both instruction and data cache, in high-performance computing (HPC), it may not yet have some of the advantages of state-of-the-art memory such as spin-transfer torque (STT) MRAM and spin-orbital torque (SOT) MRAM²⁹. That is because they are designed to operate in low-data modes where speed, including fast-to-compute and fast-to-sleep is not as crucial to operation in applications such as remote sensors with fixed interval telemetry. That said, the fast-wake up is known to be 1-2 cycles³².

Drowsy logic is renamed to...dynamic voltage & frequency scaling?

A review in the ISCA@50 retrospective paper published in June of 2023 implemented drowsy logic into “a form of voltage scaling”: “We proposed a design in which one can choose between two different supply voltages in each cache line, corresponding to normal supply voltage and a drowsy lower voltage. In effect we used a form of voltage scaling to reduce static power consumption. Due to short-channel effects in deep-sub micron processes, leakage current reduces significantly with voltage scaling.”²

By 2008, the term drowsy had already been used in textbooks as form of voltage scaling, but it is usually paired with other techniques such as Cache decay to form DVFS. Even the 2002 paper mentions it, but its novel feature at the time was that it was applied to static power:

“Such a dynamic voltage scaling or selection (DVS) technique has been used in the past to trade off dynamic power consumption and performance [6][7][8]. In this case, however, we exploit voltage scaling to reduce static power consumption.” (p.2/10)

New techniques utilizing hybrid or novel implementations of drowsy transistors in SRAM continue to be researched in academic labs.^{20, 21,78}

TABLE 3. Comparison of various low-leakage circuit techniques

	Advantages	Disadvantages	Leakage power in low power mode
DVS	<ul style="list-style-type: none"> Retains cell information in low-power mode. Fast switching between power modes. Easy implementation. More power reduction than ABB-MTCMOS. 	<ul style="list-style-type: none"> Process variation dependent. More SEU noise susceptible. 	6.24nW
ABB-MTCMOS	<ul style="list-style-type: none"> Retains cell information in low-power mode. 	<ul style="list-style-type: none"> Higher leakage power. Slower switching between power modes. 	13.20nW
Gated-V _{DD}	<ul style="list-style-type: none"> Largest power reduction. Fast switching between power modes. Easy implementation. 	<ul style="list-style-type: none"> Loses cell information in low-power mode. 	0.02nW

Techniques sometimes have multiple industry names. As early as 1995, multi and dual-threshold voltage techniques have been described^{22, 23}.

The 2002 paper also cites a 1994 PARC paper, “Scheduling for Reduced CPU Energy,” which also pioneered the concept of MIPJ (or Millions-of-instructions-per-joule) and describes software-based scheduling techniques that allow more granular control of the system clock speed by the operating system scheduler⁵³. The PARC paper cites an early wave of low-power design developed in the late 1980s as the first Notebooks and PDAs were commercialized⁵⁴⁻⁵⁹.

A difference that remains is that drowsy logic offers some of the best energy savings: “Moreover, since the penalty for waking up a drowsy line is relatively small (it requires little energy and only 1 or 2 cycles) and there are less frequent accesses to the lower memory hierarchy than Gated-V_{DD} schemes, cache lines can be put into drowsy mode more aggressively to save more power.”^{2,3(Table 3)}

The theoretical minimum leakage should not be confused with the theoretical minimum of energy computation. In *Recent Progress in the Boolean Domain* (2013), Bernd Steinbach cites how the Landauer bound described zero-energy as requiring a reversible process, and that “any bit that is destroyed incurs an energy cost of $k_B T \ln(2)$. We also reported a strikingly similar result: that the lowest supply voltage at which a logic gate can operate while still acting as an amplifier is only a few times larger than $k_B T / q$.”⁶¹ Granted, the sub-threshold frontier is much like the border between math and physics.

Industry news articles appear to have decreased mentions of “drowsy” circuits in the mid 2010s. A keyword search of one well-known site turned up only two mentions in 2014 and 2015 for “drowsy”²⁴ “For memories, people are building additional operational modes for them, such as drowsy modes.”²⁵ Furthermore, an aforementioned 85% reduction in leakage power can be found in IoT devices which are designed for batteryless operation and energy-harvesting. Microcontrollers by companies such as Ambiq Micro are known to achieve a 13-fold reduction compared to other chips²⁶.

Analog IoT device-makers such as ONiO feature an integrated solar/RF/thermoelectric harvester and power management integrated circuit with low-power, asynchronous ROM/RAM²⁷. Since the latter uses just 2KB of RAM, it may not need to operate in sub-threshold mode to create “digital zeros and 1s,” as much as Ambiq’s 2MB SRAM cells to achieve ultra low energy consumption³⁰. Nonetheless, the ultra low power of advanced power saving techniques such as sub-threshold voltage, hybrid drowsy cache and memory suggests mobile phones could one day run on solar power.

Footnotes

- <https://semiengineering.com/a-power-first-approach/>
- “Drowsy Caches: Simple Techniques for Reducing Leakage Energy—A Retrospective” Krisztián Flautner, Nam Sung Kim, Steve Martin, David Blaauw, Trevor Mudge; ISCA@50 25-Year Retrospective: 1996-2020 https://sites.coecis.cornell.edu/isca50retrospective/papers/files/2023/07/drowsy_retro.pdf
- “Drowsy Caches: Simple Techniques for Reducing Leakage Power” Krisztián Flautner, Nam Sung Kim, Steve Martin, David Blaauw, Trevor Mudge (2002) <https://web.eecs.umich.edu/~manowar/publications/drowsy-caches-ISCA2002.pdf>
- <https://developer.arm.com/documentation/dai0298/a>
- https://docs.zephyrproject.org/latest/hardware/arch/arm_cortex_m.html
- <https://docs.zephyrproject.org/latest/kernel/services/other/float.html>

- ⁷ “BTB Access Filtering: A Low Energy and High Performance Design” Shuai Wang, Jie Hu, and Sotirios G. Ziavras Department of Electrical and Computer Engineering New Jersey Institute of Technology (2008) https://web.archive.org/web/20090920084956id_/http://web.njit.edu:80/~sw63/pub/ISVLSI_BAF_2008.pdf
- ⁸ https://en.wikipedia.org/wiki/List_of_ARM_processors
- ⁹ https://web.archive.org/web/20130926155924/http://www.eetimes.com/document.asp?doc_id=1208831 “VLSI Technology Now Shipping ARM810” 08/1996
- ¹⁰ https://web.archive.org/web/20181224080542/https://www.hotchips.org/wp-content/uploads/hc_archives/hc08/2_Mon/HC8.S4/HC8.4.1.pdf “ARM810: Dancing to the Beat of a Different Drum” 07/1996
- ¹¹ <https://en.wikipedia.org/wiki/%CE%9Clinux>
- ¹² <https://www.emcraft.com/>
- ¹³ <https://semiengineering.com/why-chiplets-dont-work-for-all-designs/> “Why Chiplets Don’t Work For All Designs” 09/2023
- ¹⁴ <https://semiengineering.com/software-defined-hardware-architectures/> “Software-Defined Hardware Architectures” 05/2023
- ¹⁵ <https://pulp-platform.org/community/showthread.php?tid=229&pid=650#pid650> PULP Platform, forum post, 2021
- ¹⁶ “Gated-Vdd: A Circuit Technique to Reduce Leakage in Deep-Submicron Cache Memories,” 2000 Michael Powell, Se-Hyun Yang, Babak Falsafi, Kaushik Roy, and T. N. Vijaykumar <https://engineering.purdue.edu/~vijay/papers/2000/gatedvdd.pdf>
- ¹⁷ “Cache Decay: Exploiting Generational Behavior to Reduce Cache Leakage Power” 2001 Stefanos Kaxiras, Zhigang Hu, Margaret Martonosi <https://mrmgroup.cs.princeton.edu/papers/hzg-isca2001.pdf>
- ¹⁸ “Adaptive Mode Control: A Static-Power-Efficient Cache Design” 2001 <https://prod.tinker.cc.gatech.edu/journal/zhou03adaptive.pdf> <https://pdfs.semanticscholar.org/ba53/2536b7d1890f967d1aac7485544f6abebe39.pdf>
- ¹⁹ “RETROSPECTIVE: Cache Decay: Exploiting Generational Behavior to Reduce Cache Leakage Power” 06/2023 https://bbp-us-w2.wpmucdn.com/sites.coecis.cornell.edu/dist/7/587/files/2023/06/Hu_2001_Cache.pdf
- ²⁰ “Design and Implementation of Low Power SRAM Using Highly Effective Lever Shifters er Shifters, 2021 <https://scholar.uwindsor.ca/cgi/viewcontent.cgi?article=9892&context=etd>
- ²¹ “Hybrid Drowsy SRAM and STT-RAM Buffer Designs for Dark-Silicon-Aware NoC, 2016 Jia Zhan, Student Member, IEEE, Jin Ouyang, Member, IEEE, Fen Ge, Member, IEEE, Jishen Zhao, Member, IEEE, and Yuan Xie, Fellow, IEEE <https://cseweb.ucsd.edu/~jzhao/files/darlsilicon-noc-tvlsi2016.pdf>
- ²² “1 -V Power Supply High-speed Digital Circuit Technology with Multithreshold-Voltage CMOS” 1995 <https://ieeexplore.ieee.org/document/400426> http://www.facweb.iitkgp.ac.in/~apal/LPC_2010/MTCMOS.pdf
- ²³ “Design and Optimization of Low Voltage High P erformance Dual Threshold CMOS Circuits” Liqiong Wei, Zhanping Chen, Mark Johnson, Kaushik Roy & Vivek De, 1998 <https://dl.acm.org/doi/pdf/10.1145/277044.277179>
- ²⁴ “S-L Power Modeling Gains Steam” 08/2014 <https://semiengineering.com/system-level-power-modeling-activities-get-rolling/>
- ²⁵ “With Responsibility Comes Power,” 02/2015 <https://semiengineering.com/with-responsibility-comes-power/>
- ²⁶ “What’s All This Subthreshold Stuff, Anyhow?” 02/2019 <https://www.electronicdesign.com/technologies/analog/article/21807652/whats-all-this-subthreshold-stuff-anyhow>
- ²⁷ “What if You Never Had to Charge Your Gadgets Again?” 01/2024 <https://www.wsj.com/tech/personal-tech/what-if-you-never-had-to-charge-your-gadgets-again-955ea960>
- ²⁸ <https://www.fluke.com/en-us/learn/blog/electrical/what-is-ohms-law>
- ²⁹ https://en.wikipedia.org/wiki/Special_relativity
- ³⁰ “Interview With Scott Hanson - Founder and CTO at Ambiq” 01/2024 <https://www.safetydetectives.com/blog/scott-hanson-ambiq/>
- ³¹ “TSMC tandem builds exotic new MRAM-based memory with radically lower latency and power consumption” 01/2024 <https://www.tomshardware.com/pc-components/dram/tsmc-tandem-builds-exotic-new-memory-with-radically-lower-latency-and-power-consumption-mram-based-memory-can-also-conduct-its-own-compute-operations>
- ³² “Circuit and microarchitectural techniques for processor on-chip cache leakage power reduction” (Dissertation), 2004 Nam Sung Kim, Trevor N. Mudge <https://dl.acm.org/doi/10.5555/1023306>
- ³³ https://www.eng.auburn.edu/~agrawvd/THESIS/KIM_S/S.Kim_HonorsThesis_FINAL.pdf May, 2016 Sungil Kim
- ³⁴ <https://www.slideshare.net/AmbiqMicro/apollo-ultra-low-power-mcu-from-ambiq-micro>
- ³⁵ <https://www.eetimes.com/arm-preps-near-threshold-processor-for-iot/>
- ³⁶ https://en.wikipedia.org/wiki/Soft_error#Other_causes
- ³⁷ https://en.wikipedia.org/wiki/Crosstalk#Other_examples
- ³⁸ https://www.tsmc.com/english/dedicatedFoundry/technology/logic/l_22nm
- ³⁹ <https://semiengineering.com/near-threshold-computing-gets-a-boost/>
- ⁴⁰ [https://en.wikipedia.org/wiki/Geode_\(processor\)#Geode_GXm](https://en.wikipedia.org/wiki/Geode_(processor)#Geode_GXm)
- ⁴¹ <https://www.wtamu.edu/~cbaird/sq/2014/02/19/what-is-the-speed-of-electricity/>
- ⁴² <https://electronics.stackexchange.com/questions/464689/what-is-the-speed-of-electricity>
- ⁴³ “What are valence orbitals?” <https://www.lanl.gov/orgs/nmt/nmtdo/AQarchive/04spring/VO.html>
- ⁴⁴ “Probabilistic Processors” <https://www.cs.uaf.edu/2011/spring/cs641/proj1/rarutter/>
- ⁴⁵ https://en.wikipedia.org/wiki/Centimetre%E2%80%93second_system_of_units#Advantages_and_disadvantages
- ⁴⁶ <https://www.scientificamerican.com/article/see-the-highest-resolution-atomic-image-ever-captured/>
- ⁴⁷ “Dynamic voltage frequency scaling (DVFS) for microprocessors power and energy reduction December 2005 https://www.emo.org.tr/ekler/035226640b6b89f_ek.pdf
- ⁴⁸ “Battery-Free Game Boy” 09/2020 <https://dl.acm.org/doi/pdf/10.1145/3411839>
- ⁴⁹ “Dynamic Voltage and Frequency Scaling for Intermittent Computing” 01/2024 <https://arxiv.org/pdf/2401.08710.pdf>
- ⁵⁰ https://en.wikipedia.org/wiki/Supervisory_circuit#Overvoltage_and_undervoltage_protection
- ⁵¹ <https://en.wikipedia.org/wiki/Waveform>
- ⁵² https://en.wikipedia.org/wiki/Continuous_function
- ⁵³ “RETROSPECTIVE: “Complexity-Effective Superscalar Processors”, Subbarao Palacharla, Norman P. Jouppi, J. E. ICSA@50, Smith3 06/2023 https://bbp-us-w2.wpmucdn.com/sites.coecis.cornell.edu/dist/7/587/files/2023/06/Palacharla_1997_Complexity.pdf
- ⁵⁴ “Scheduling for Reduced CPU Energy” Mark Weiser, Brent Welch, Alan Demers, Scott Shenker, Proceedings of the First USENIX Symposium on Operating Systems Design and Implementation, November 1994 https://www1.icsi.berkeley.edu/pubs/networking/ICSI_schedulingforreduced94.pdf
- ⁵⁵ “A Scheduling Model for Reduced CPU Energy” Frances Yao Alan Demers Scott Shenker”FOCS, 1995 <https://graal.ens-lyon.fr/~lmarchal/scheduling/Yao-FOCS95-Energy.pdf>
- ⁵⁶ “A JSSC classic paper: Low-power CMOS digital design, 04/2003 https://chandrakasan.mit.edu/wp-content/uploads/2021/04/Chandrakasan_06499960.pdf
- ⁵⁷ “Low-power CMOS Digital Design” A. P. Chandrakasan, S. Sheng, and R. W. Brodersen , JSSC, 04/1992 https://mtl.sites.mit.edu/researchgroups/icsystems/pubs/journals/1992_chandrakasan_jssc.pdf
- ⁵⁸ “A JSSC Classic Paper: The Simple Model of CMOS Drain Current” 10/2004 https://www.eng.auburn.edu/~agrawvd/COURSE/READING/LOWP/alpha_power_law.pdf
- ⁵⁹ “Alpha-power law MOSFET model and its applications to CMOS inverter delay and Other Formulas” 04/1990 https://courses.ece.ucsb.edu/ECE125/125_W11Banerjee/Lectures/SAK90a.pdf
- ⁶⁰ “Computer Architecture Techniques for Power-Efficiency, Stefanos Kaxiras & Margaret Martonosi (Ch. 4.5-4.10),(Ch. 5.3) Morgan & Claypool Press 2008
- ⁶¹ “Recent Progress in Boolean Logic, Bernd Steinbach, 2013 (4.1.3, 203), (4.1.5, 211)
- ⁶² “A study of DVFS methodologies for multicore systems with islanding feature (p.42) Shervin Hajiamini, Behrooz A. Shirazi, in *Advances in*

Computers, 2020

⁶³ <https://electronics.stackexchange.com/questions/462524/propagation-speed-of-em-waves-in-air-versus-through-conductive-material/>

⁶⁴ https://en.wikipedia.org/wiki/Speed_of_electricity#Electric_drift

⁶⁵ <https://en.wikipedia.org/wiki/Vacuum>

⁶⁶ <https://www.ledylighting.com/why-do-light-bulbs-glow-when-switched-off/>

⁶⁷ "Idleness is not sloth" Richard Golding, Peter Bosch, Carl Staelin, Tim Sullivan, and John Wilkes

https://dl.acm.org/doi/proceedings/10.5555/1267411_TCON'95:
Proceedings of the USENIX 1995 Technical Conference Proceedings
January 1995 <https://john.e-wilkes.com/papers/idleness.pdf>

⁶⁸ "Is power proportional to V or V^2 ?"

<https://electronics.stackexchange.com/questions/543610/is-power-proportional-to-v-or-v2>

⁶⁹ Ohm's Law <http://avstop.com/ac/appgeneral/ohm%27slaw.html>

⁷⁰ <https://www.malcolmchisholm.net/painless-decibels-almost>

⁷¹ <https://science-education-research.com/when-is-vir-the-formula-for-ohms-law/>

⁷² <https://www.physicsforums.com/threads/understanding-the-relationship-between-power-voltage-current-and-resistance.876880/>

⁷³ "Optimization of high-performance superscalar architectures for energy efficiency" V. Zyuban, P. Kogge, 08/2000

<https://dl.acm.org/doi/pdf/10.1145/344166.344522>

⁷⁴ <https://chem.libretexts.org/Bookshelves/>

Physical_and_Theoretical_Chemistry_Textbook_Maps/
Supplemental_Modules_(Physical_and_Theoretical_Chemistry)/
Quantum_Mechanics/03_The_Tools_of_Quantum_Mechanics/
Collapsing_Wavefunctions

⁷⁵ <https://en.wikipedia.org/wiki/Power>

https://en.wikipedia.org/wiki/Power#E2%80%93delay_product#cite_note-Gaudet_2014-1

⁷⁶ "energy autonomous computing"

<https://github.com/kragen/dernocua/blob/master/text/energy-autonomous-computing.md#estimating-the-necessary-performance-for-basic-interactive-computation-01-dmips>

⁷⁷ "My visit to Bletchley Park and The UK Computing History Museum in November 2017"

<https://www.koomey.com/post/670564790101524480>

⁷⁸ "Design the efficient SRAM circuit using 4transistor with sleepy logic"

International Journal of Pure and Applied Mathematics Volume 118 No. 20
2018, 115-123 <https://acadpubl.eu/hub/2018-118-21/articles/21b/15.pdf>

Badness 0 (Epsom's version)

Dr. Tom Murphy VII, Ph.D.
March 2024

Many people walk this Earth unbothered by incorrect details. For example, they are unconcerned when a hyperlink includes a surrounding space character. They don't notice that the screw heads on a light switch wall plate are not all lined up. They don't care about the rules of Wordle's "hard mode" being simply wrong. They don't notice the difference between "its" and "it's." When someone asks, "Will you marry me?" and they think "Oh my god!" it's not because the proposer probably should have used the subjunctive *would*.

I am not like this. If I can infer from a coffee cup's moment of inertia that it does not contain any liquid, I immediately lose suspension of disbelief and will not purchase the product featured in the commercial. I literally projectile vomit if Auto-Motion Plus is enabled on a television in the hotel I'm staying in, even if the TV is not turned on, or if someone misuses the word "literally." If I see a period missing at the end of a paragraph on Wikipedia, I will spend dozens of hours writing software to organize and semi-automate a distributed effort to fix all the missing periods on Wikipedia.^[2] And worse, each time I learn of a new type of mistake, I am forever cursed to notice that mistake

Seriously: One time I found myself spell-correcting someone else's *lorem ipsum* text in a slide. It said "lorem epsom," which is funny. I think about that incident all the time. The person that wrote the slide probably thinks about things like leveraging synergy, generative AI, metaverses, blockchain 3.0, snackable content, being eco-green, and so on, without it occurring to him that these things could have nuance and meaning separate from their names. He has probably never even read the Wikipedia article on Lorem Ipsum. He is successful and rich.

Another successful person is Bill Cassidy, who is a congressperson. He criticized a proposed bill that would reduce the standard work week in the US by 8 hours, from 40 to 32. He said,

Sen. Bill Cassidy of Louisiana, representing the _____ party, said paying workers the same wages for fewer hours would force employers to pass the cost of hiring more workers along to consumers.

"It would threaten millions of small businesses operating on a razor-thin margin because they're unable to find enough workers," said Cassidy. "Now they've got the same workers, but only for three-quarters of the time. And they have to hire more."

In fact, that's not the exact quote, but I needed to make it look nice.^[3] And this is not a paper about politics, but you can probably guess the word that goes in the blank.

Anyway, OKAY (and I'll explain why in a second), first of all, razors famously have high margins. It's like the worst possible metaphor here.

This guy uses both fancy and ASCII quotes.

The main thing I want to talk about is: What? No! 32/40 is not three quarters. This is not, like, complicated math. It uses some of the world's smallest integers. Everybody knows that the work week is 40 hours, and that a work day is 8 hours, and that the proposed bill reduces it by one day, giving four of five days. I don't really mind if someone makes an error in calculation (well, I do mind, but I am certainly prone to doing it). The infuriating realization here is that this person does not even think of "three-quarters" as a kind of thing that can be right or wrong. He says three quarters because it makes smaller number feelings. You could imagine him having the conversation (with me, perhaps): "You say four-fifths, I say three-quarters." Me: "But it is four fifths. And why are you always hyphenating it?" Him (smiling patronizingly): "I guess we just have to agree to disagree."

Donald Knuth is the opposite of this person.

I'm not saying that Donald Knuth isn't successful and rich. According to the website [Famous Birthdays](#), which is probably generated by AI or at least by people whose economic output is measured in a count of words, and words whose value is computed by their ability to drive ad clicks, Donald Knuth is "is one of the most popular and richest Mathematician who was born on January 10, 1938 in Wisconsin, Wisconsin, United States. Mathematician and engineer who was arguably most recognized as the Professor Emeritus at Stanford in Palo

Alto, California.” As one of the richest Mathematician from United States, according to the analysis of Famous Birthdays, Wikipedia, Forbes & Business Insider, “Donald Knuth’s net worth \$3--5 Million.^[4]”

It is arguable that he is the Professor Emeritus. It is likely that he is the only popular and rich mathematician born on that specific day in Wisconsin. The singular “Mathematician” is perhaps a technical master-stroke. The asterisk does not have any referent on the page.

What I mean when I say that Donald Knuth is the opposite of this person is that Knuth is interested in unpacking a single unnecessary detail, recursively, until it is completely solved. According to the website Famous Bibliophiles, one day Donald Knuth set out to write down the entire subject of computer science in a single book called The Art of Computer Programming. As he was doing so, he realized that describing computer algorithms in a lasting form would require a programming language that was not subject to constant revision, so he invented the MIX instruction set for an idealized computer. After writing some 3000 pages out in longhand, he found that it was impractical to print them all in one book, so the plan expanded to be multiple volumes. Then when he got a draft of one of the books back from the typesetter, he was unhappy with the details of the typography, and so he paused his work writing down all of computer science to create some new computer science: First an algorithm for determining where to place line breaks in order to make text optimally beautiful, then algorithms for hyphenating words, then generalizations of these for typesetting mathematics, and then a full computer typesetting system that is still in wide use today, called TeX. Along the way he was unsatisfied with the specific typefaces that existed in the world, and unsatisfied with the way that typefaces were described at only one weight, and so he created the parameterized METAFONT system and several new typefaces. Undeterred by these excursions, he returned to his original task of writing down the entirety of computer science, using all the technology he had built. By the time he finished this, much more computer science had been invented, including by his own hand, and so he needed to rework MIX for the next volume, and update the first. The revised plan of eight volumes remains the intention in 2024. However, he found that the volumes were getting rather long, and began releasing portions of volumes (“fascicles”). So far, Volume 4 has been partially published as books 4A^[5] (fascicles 0–4; 912 pages) and 4B^[6] (fas-

cicles 5–6; 736 pages). It is unknown how many more episodes remain in Volume 4. I expect that every conversation that Knuth has with his editor goes like this. Editor: “Hey, Donald, I hope you’re well. Just wondering if you have an update on when 4C will be ready? Or any more icicles?” Donald E. Knuth: “I am working diligently on fascicles for Volume 4C. As I’ve mentioned in the past, it’s impossible to tell how long it will be, since mathematics does not obey the rules of project management.” Editor: “I just need a date to tell the publishers.” Donald E. Knuth: “Like I’ve said, any date would be very low confidence, other than the fact that it will be in the future.” Editor: “I just need a date.” Donald E. Knuth: “Would you like me to say a date, knowing that it’s a very low confidence guess, and that I would be extremely likely to miss that date, or even deliver early?” Editor: “Early! Now we’re talking.” Donald E. Knuth: “What use is the date if you’re excited about the possibility of it being early, relative to some unknown date?” Editor: “I just need a date for the publishers.” Donald E. Knuth: “2030.” Editor: “Thanks Donald, you’re the best!”

Knuth is estimated to be ready with Volume 5 in 2030, when he will be 92.

That’s a large amount of language!

Nightmare on LLM street

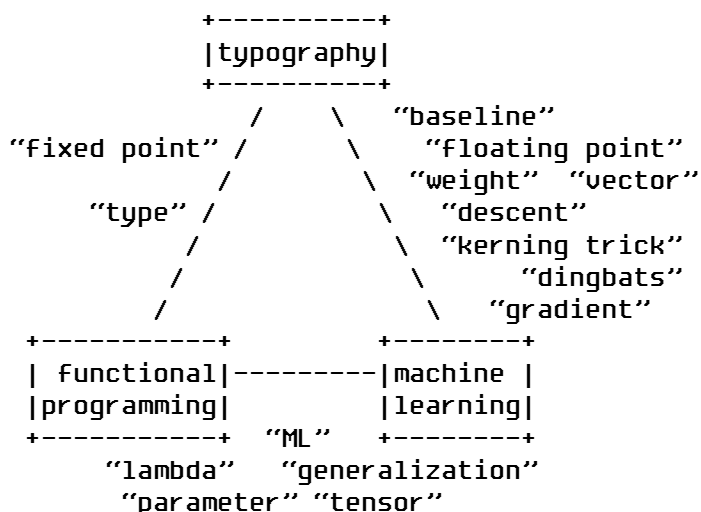
Then there are large language models. ^[7] One of the annoying things about large language models is that they are so buzzwordy, but unlike most buzzwordy trends, they are actually substantive. They produce remarkably fluent text. With no additional training, they often outperform models that have been developed for decades. They generalize to completely new situations.

So many things about “AI” distress me. Dolor sit amet! I worry about the devaluation of human creativity, about large-scale disinformation and spam ruining the beautiful library of knowledge that humans have created, about extreme concentration of wealth. And yes, I worry about competing with AI. Being able to work tirelessly and thousands of times faster than humans is a huge competitive advantage. Of course, I find some solace in the significant possible upsides. It might help us solve hard problems like climate change and AI. But even in the best scenarios we will not be able to ignore it: Even if it never gets as smart and precise as Knuth, it’s already too economically useful in its Lorem

Epsom state (just like Lorem Epsom himself).

On the one hand, the technology is pretty neat and lends itself to some nice abstractions. On the other hand, I love playing with words. So, I have been experimenting with LLMs in practical and impractical applications. I also try to make it fun (for me) to program with them.

I have a myriad of strategies for digesting things that irritate me. For this work, I'm inspired by the "Hurry-Coward So-so-morphism," where I make connections between topics based solely on confusion of superficial lexical similarities without regard to their underlying meaning. So for example, we have "ML" meaning both "Machine Learning" and "Meta Language," as well as "type" both as in "typeface" and as in "type systems for programming languages."^[8] And because machine learning has claimed so many words, there are a great many shared with typography as well:



I spent a lengthy introduction talking about Donald Knuth's work in computer typography. Now I can tell you what this paper is about. If we are giving up on precision in our near AI future, perhaps we can have something we want: perfect typography. This paper is about a new typesetting system, BoVeX, which allows us to control the exchange of precision for beauty. It essentially gives us a dial between Lorem Ipsum and Donald Knuth.

The scientists' findings were astounding! They discovered that the powers of the Metroid might be harnessed for the good of civilization!

The Metroid series is a video game series about a brain that has been enslaved inside a jar in an underground datacenter on the planet Zebes. This

brain is called Mother Brain and its goal is to control the hypercapitalists called Space Pirates to increase their "score" as high as possible by conquering planets throughout the galaxy. Mother Brain was invented by the Space Pirates, although it is not clear whether the current situation was actually intended by the Space Pirates. The most super version of Metroid is Super Metroid.

In the 1990s, the website gamefaqs.com collected plain text "FAQs" for classic video games, then just known as video games. On this site, another hero was born. They were writing the definitive guide to speedrunning the SNES game Super Metroid, and they saw that some of their ASCII lines ended up exactly the same length, and it looked good:

Once you save the game at your ship (about 1 hour 15 minutes is good), go down to Tourian. Do not save your game in Tourian if you have intentions of returning to any previously explored section on Planet Zebes. There will be a few Metroids to kill before you reach Mother Brain, and they must all die in order to continue to Mother Brain. Read the boss guide for more details. Once Mother Brain is defeated, you will need to hurry back to your ship. By now you will already have the HYPER BEAM. From Mother Brain's room, go west and then south. Take the blue door at the bottom and speed dash east. Super jump up, and continue north. Once you land up top and are running east, aim diagonally down to the right and shoot an unseen door. Eventually, you will get to this door since lava will start to rise from the floor in this area. Speed dash through the door you preopened, and charge for a super jump. Hug either the left or right wall in the Craterian shaft and super jump up. Now quickly get to your ship before the planet explodes. There should be almost a minute left on the timer. Sit back and watch the ending! Did you beat the game within 1 hour and 20 minutes?

and so they wisely decided to wordsmith the entire 28-page guide so that every line was exactly the same length, with no extra spaces or other cheating, just because it could be done.^[9]

Doing this manually is a chore, and I do like to automate the chores of Speedrunners.^[10]

Exercise in rephrasing text. The following paragraph needs to be rephrased so that it retains its precise meaning, but with minor variations in the specific choice of words, punctuation, and so on. No new facts should be introduced or removed, but it is good to use synonyms and change the word order and phrasing.

After this, I insert **Rephrased text:**, which is the rephrasing of the original text. The model is ready to generate tokens.

I then sample text a word at a time to continue this prompt. If a line ends exactly on the number of characters that I want (and the next character is a space or other character that is appropriate to end a line) then I accept the stream so far and continue. If I exceed the line length, I back up to the state at the beginning of the line and try again with new random samples. I just keep doing that until the

paragraph is complete, and we have beautifully justified monospace text that resembles the original. Here is an example of this paragraph rendered in monospace:

```
I sample text a word at a time to continue
this prompt. If a line ends exactly on the
number of characters I want, I accept that
text so far, and continue. If I exceed the
line length, I back up to the beginning of
the line and try again with new samples. I
keep repeating this until I get text I can
render in monospaced font, and that is how
we can get beautifully justified monospace
text. Here is an example of this paragraph
rendered in monospace:
```

The text could be improved by using “monospace” and “monospaced” consistently. The most upsetting thing is that the paragraph ends with a colon, as if there is going to be another example.

The approach described works reasonably well, but it has several deficiencies that we will address in the real BoVeX system. But it is a good example to explain some concepts that will be useful later.

¿Como te LLama?

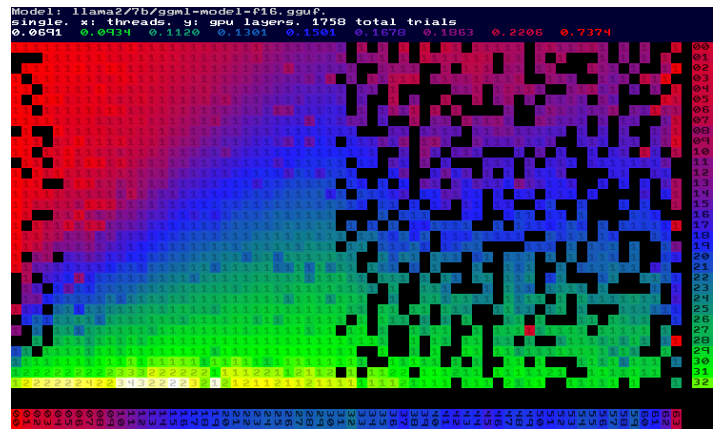
Facebook’s Large Language Model,^[11] Llama, is available for anyone who agrees not to use it to destroy the world. Wouldn’t it be funny if the world is destroyed by something called “Llama”? That’s some Stay-Puft Marshmallow Man stuff. Actually I hear that llamas are pretty mean, and if you are thinking about hugging a cute long-neck, you are probably thinking about an alpaca. But that’s probably a version of the linear algebra package LAPACK. Llama-v2-70b is a good LLM which can do some impressive things, but when I say destroy the world I mean stuff like filling the internet with infinite spam, or building critical infrastructure on it in order to cut costs, where most of our “safety” measures consist of asking the model politely to recite its daily affirmations before performing its tasks. That kind of thing. It’ll be at least months before we *really* have to worry.

Anyway, the normal way to program with Llama is to use Python, and a mountain of things that you are not supposed to understand and cannot understand, mostly by pasting examples from others and then tweaking parameters and prompts. I don’t care for it. Fortunately, human geniuses^[12] have implemented the inference code for llama-like models in a nice, portable C++ library called “llama.cpp”

(checks out).

I can load a quantized version of the model into RAM with llama.cpp. There are two different models, the 7b and the 70b, which refer to the number of billions of parameters, which must be a multiple of VII for performance reasons.

Quantization is a process that reduces the number of bits used to represent floating-point weights.^[13] This saves memory, but it also speeds up inference, which needs to read pretty much the entire model for every predicted token. I got reasonable quality and good performance from LLama-v2-7b with 16-bit floats. This model fits completely on my world’s (physically) largest GPU. To tune various settings, I ran thousands of trials for the different models, and made some nice custom graphs:



The x axis of the graph is the number of CPU threads, and the y axis is the number of model layers loaded onto the GPU. As expected, increasing the number of threads and layers on the GPU improves performance, since the entire model fits on the GPU. For the 70b models (not pictured), there is an abrupt drop-off in throughput before we load all the layers, and my computer gets very sluggish if I exceed the GPU memory. We see that if we use more than the number of physical cores (32), we do not see any benefit, which is not surprising because hyperthreading basically never helps anything. The best throughput actually uses a modest number of cores (about 12). Mostly I’m just including the graph to demonstrate that BoVeX has support for including PNG files.

Where was I? Right. Fundamentally, LLMs are trained to predict a token given some sequence of tokens that precede them. There is a fixed set of tokens for the model, and rather than predict a single token, they actually give a score for every possible token. These scores are typically normalized

into a probability distribution. So for example, if we have the text

```
SIGBOVIK is an
```

then the probability distribution begins as

```
( annual) 69.8010%
( April) 3.8023%
( ac) 3.2456%
( academic) 2.9374%
( artificial) 2.0857%
( open) 1.7993%
( under) 1.2331%
( international) 1.1032%
...
```

with the thousands of other tokens following. So three-quarters of the time the next token should be "annual", but there are many other reasonable possibilities. We can pick one of these tokens however we like, append it to the sequence, and run the model again. This gives us a new probability distribution. By doing this over and over we can generate a likely piece of text. This is what "Lorem Ipsum" means when he says "Generative AI." Rather, what he means is "the new thing that is cool," but what he is unknowingly "referring to" is that you can sample a probability distribution. He has probably never even read the "wikipedia article on Markov chains".

If I always sample the most likely token, I always get the most likely text. It is good to be likely; this is why the model is useful. However, you might not want exactly the same result each time, and in many situations if you only sample the most likely token, you get very boring, repetitive text. Pseudo-random number generation is the spice of life!

We do not need to sample from the probability distribution to generate text. We can simply pick the token we want. This is how the initial "prompt" works; we just run the inference process one token at a time, but always select the next token in the prompt, ignoring the probabilities. So at each moment, the text we've generated so far (more or less) completely characterizes the state of the LLM. This means that we can easily go back to earlier moments and generate a different continuation of the text, by simply replaying tokens. We also have the option of storing the LLM state (gigabytes) in RAM, which allows us to return to a previous state in constant time.

To generate monospaced lines of the same length, I use a prompt that asks the model to rephrase the input paragraph. Greedily sampling the distribution typically results in a copy of the input paragraph, which is fine for our purposes. However, if the lines are already the right length, we need not change them! To prevent boredom, whenever the process repeats a line that's already been seen, I increase the "temperature" modifier to the probability distribution. This causes the candidate lines to be more varied, but less probable (according to the original probability distribution).

This is all there is to the monospacing version. It's just 300 lines of code, including boilerplate, debugging code, and false starts.

Great !! You fulfilled your mission. It will revive peace in space. But, it may be invaded by the other Metroid. Pray for a true peace in space!

We are satisfied with the output text, but we are not satisfied with the font. We want to have excellent justified text with all the perks of proportional fonts and a programmable document preparation system. We want to have it by the estimated SIGBOVIK deadline so that it can be used to prepare the paper that I'm now writing. This is the Donald Knuth Any% speedrun.

The boxes-and-glue algorithm

The justification of monospaced text looks quite bad^[14] when more than one space is inserted between words. We can tell if text is suitable for some width by adding up the codepoints. For the full-on typography case with proportional fonts, there are many more degrees of freedom. It is possible to expand or contract the space between words a little bit, even if it varies from line to line. We can also hyphenate words.

In the time when I was being born, and probably being very upset about it, Knuth was having similar feelings about the way his computer-typeset documents looked. He discovered a nice abstraction that generalizes most of these typographic degrees of freedom, and devised an algorithm for producing optimal text layout given some parameters.^[15] The idea is to consider the text of a paragraph as consisting of rigid "boxes" (say, words) and stretchy "glue" (say, space) between them. Both boxes and glue have various detail (and can be extended to support all sorts of quirks) but the basic algorithm can be understood with just those

pieces. So, let's do that.

Knuth's paper is great, but I started having spoiler feelings when reading it, so I figured out my own algorithm, which is more fun than reading. The key insight is that you do not need to try all possible break points. Whenever we break after a word, the problem is the same for the rest of the text, no matter how we got there. This lends itself to a dynamic programming algorithm.

Dynamic programming is a programming technique for solving problems on the whiteboard at tech companies. When I was young, it was a mysterious concept because of its strange name. Here is how I think about it. Imagine you have a recursive procedure that solves the problem. In this case, the pseudocode is something like

```
pair<int, string> Split(string line,
                       string text) {
    if (text.empty()) return {0, ""};
    auto [word, rest] = GetFirstWord(text);
    // try splitting
    auto [penalty1, rest1] = Split(word, rest);
    penalty1 += badness from leftover space;
    // try not splitting
    auto [penalty2, rest2] =
        Split(line + " " + word, rest);
    penalty2 += badness from line too long;
    if (penalty1 < penalty2) {
        return {penalty1, word + "\n" + rest1};
    } else {
        return {penalty2, word + " " + rest2};
    }
}
```

The `line` and `text` are split. In the normal case, there is a word left, and two possibilities: either splitting after the first word, or not splitting. This is exponential time, because each call makes two recursive calls, to try each of the two options. But deep recursive calls will be made with the same arguments many times. So, add some memoization: If the function is called for the same line and text a second time, just return the same answer as before without doing any work (especially not making recursive calls again). This limits the function to be called at most once for each possible argument; we can then see that `line` is no longer than the input (so it is size $O(n)$), and `text` is always some suffix of the input (so it is size $O(n)$), giving $O(n^2)$ calls.

In dynamic programming, we store the values of all recursive calls before we need them, and then look them up. For this problem, we store the values in a table indexed by the two parameters, the

current line and the remaining text. The line is the number of words before the current word that are included on the line, and the text is the position in the string where we'll next look for a word. That's all there is to it; we start with empty text and then write the loop to fill out cells in the right order.

Knuth's boxes-and-glue algorithm has many extensions, and mine has many extensions as well. For example, we'll talk about how you can adapt the algorithm to perform hyphenation and kerning. There are many rabbit holes to go down, and I explored the ones that attracted my attention. There is plenty of time to add more features later, since I have now cursed myself to use BoVeX for my future SIGBOVIK papers.

But here's where I diverge from Knuth somewhat. Knuth was reluctant to add a programming language to TeX,^[16] but I spent the majority of my time on this project implementing a full-fledged language. BoVeX is about 33,000 lines of code, the majority of which is the implementation of the language itself. That's 110 times longer than the original monospaced proof of concept and 30 times the length of this document!

The BoVeX language

The BoVeX programming language is described in this section, and its implementation. If you are just in it for the jokes, you can skip this section, which is basically serious and loaded with programming language theory jargon.

BoVeX is a typed functional programming language in the ML family. It has syntax that closely resembles Standard ML. Here is an example piece of code from the source code of this document:

```
datatype (a) option = SOME of a | NONE

fun consume-outer-span f s =
  case layoutcase s of
    Node (SPAN, attrs, children) =>
      let
        val (ropt, rchildren) =
          case children of
            one :: nil => consume-outer-span f one
          | _ => (NONE, layout-concat children)
      in
        case (f attrs, ropt) of
          (NONE, _) => (ropt, span attrs rchildren)
        | (SOME vouter, inner as SOME _) =>
            (inner, rchildren)
        | (outer, NONE) => (outer, rchildren)
      end
    | _ => (NONE, s)
```

The language you see here is a full-fledged programming language,^[17] although it is not a standard one. It supports higher order functions, polymorphism, algebraic datatypes, pattern matching, Hindley-Milner type inference, and so on. It is basically core (no modules) Standard ML, as allow patterns on both sides. Anyway, a full description of the language would be boring and take too much time as the SIGBOVIK deadline draws closer.

Implementation

In the past I have implemented many similar languages, including for my dissertation.^[18] I could have started from one of my existing implementations, but they were written in Standard ML, which I could not get to work on my computer in 2024. So I started over from scratch in C++, which at least does work on my computer. (I also wanted to interface with GPU inference code for running the LLM, which would be easiest from C++). C++ is not a good language for writing language implementations, but it has gotten better.

The BoVeX implementation is a “compiler” in the sense that it transforms the source language through multiple intermediate languages into a low-level bytecode. This bytecode is just straight-line code on an abstract machine with infinite registers and operations like `allocate` and `setfield`. It does not produce machine code, and although this would be pretty feasible, it would not be the first thing to do to make BoVeX faster.

First, it concatenates the source files (handling `import` and keeping track of where each byte originated, for error messages), and then it lexically parses them into the External Language (EL). Then, it elaborates EL into a simpler and more explicit Internal Language (IL), which is nice and clean. I love writing optimizations, but I had to keep myself out of there, or else this would be a 2025 SIGBOVIK paper. There are just enough to make the code reasonable to debug if I need to look at it. After optimization, I perform closure conversion, simplify again, and generate the final “bytecode” form. This entire process happens whenever you generate a BoVeX document; the only output from running `bovex.exe` is the PDF document.

I did not cut corners on the language implementation. For example, compiling mutually-recursive polymorphic functions is obnoxious, but I did it. The full story is in the source code.^[19]

AST pools. `datatype` declarations are for) and annoying in C++. I continued to experiment with different ways to do this. I use arena-style allocation for the syntax nodes (always `const` after creation), so that they can be created and reused at will. My current favorite approach to manipulating the nodes is to write “in” and “out” functions (tedious, manual) for each construct in the language. The syntax nodes can then be implemented however I like (for example, a flat struct or `std::variant`). I get the compiler’s help whenever I change the language (which is often!) since each in/out function is explicit about its constituents.

Passes and guesses. Many transformations in a compiler rewrite a language to itself; for example each IL optimization is a function from IL to IL. These can be tedious to write and update, especially since a given optimization usually only cares about one or two constructs in the language. I use the “pass” idiom to write these. This is basically an identity function on the AST that pulls apart each node, calls a virtual function for that node, and then rebuilds the node. To write a pass that only cares about one type of node, you inherit from this class and then just override that one node’s function. One issue with this is that each time you rebuild the entire tree you create a lot of unnecessary node copies. So exchanging tedium (mine) for efficiency (my computer’s?), every node type’s “in” function also takes a “guess” node pointer. If the node being constructed is exactly equal to the guess, then we return the guess and avoid creating a copy. Then the base pass is actually the identity (it returns the same pointer) and does no long-lived allocations. This seems to be a good compromise between the traditional garbage-fountain approach and hash consing, which sounds like it would be a good idea but is usually just a lot slower.^[20] **For type-directed transformations, there is also a typed IL pass class, which recursively passes a context and does bidirectional type checking of the intermediate code. Closure conversion is a type-directed pass and is implemented this way.**

Parsing. I have this aversion to parser generators, probably because one time I tried to get someone else’s code to compile and it complained about having the wrong bovines on my computer and ruined my weekend. After trying some other people’s C++ parsing libraries and being disappointed by them, I did what Knuth would do: I wrote my own. It is a parser combinator^[21] library which actually descends directly from Okasaki’s SML code.^[22] I was proud of myself for getting this to work in C++,

since C++'s insane type system is impossible to understand and its error messages are even worse. (BoVeX's error messages are extremely spartan, often simply declaring `Parse error at paper.bovex line 1`, but in many ways this is more useful than C++'s mile-long SFINAE vomitus.) It supports mutually-recursive parsers, resolution of dynamic infix operators, and all that. My template-heavy parser combinators take clang about a minute to compile, which is acceptable. Less acceptable, but something I only learned after using this to write a 14-page-long paper, is that the parsers are very slow. Putting aside LLM inference, this paper takes 13 seconds to render into a PDF, 11 seconds of which is parsing! There must be some bug, but I don't know if it's in my grammar (it is easy to accidentally write an exponential time parser, but this one should not be) or the parser combinator library (also my fault) or clang producing bad code (it may be giving up on optimizations, since it is taking so long to compile; the `.o` file is 41 megabytes). But these are details to be improved in the future.

Garbage collection. I keep track of all the pointers that are allocated during execution. It's so easy that I didn't even implement it! I have 256 gigabytes of RAM. In fact, LLM inference acts as a useful "performance regulator" to make sure that we don't allocate memory *too fast*.

Objects

As the SIGBOVIK deadline grew near, I reluctantly added "objects" to the BoVeX language. Objects are no stranger to ML; for example the O'CamL Language^[23] (pronounced "OK ML") has them.^[24] But the community of functional programmers I was raised in has a revulsion to things Object Oriented, just like how a woodworker will immediately projectile vomit if they see a piece of Oriented Strand Board, even though it is a fine tool for many applications. I still have this disgust reflex. I imagine my Ph.D. advisors, should they read this, are contemplating whether and how a Ph.D. can be revoked. Anyway, I deliberately kept objects low-tech so that nothing could get too Oriented.

There is one object type in BoVeX. A value of this type has an arbitrary set of named fields whose types are known; they can only be the base types `int`, `float`, `string`, `bool`, `layout`, or `obj`. Fields are distinct if they have different types. An object can be introduced with an expression like `{(field1 = exp1, field2 = exp2)}` or by declaring an object name `O`:

```
object O of { field1 : type1, field2 : type2 }
```

and then use this in an expression like `{O} field1 = exp1`. These object names do not have any run-time meaning; they are just a collection of field types that are commonly used together. It gives a good place to document what they mean and some opportunity for better error messages, but fundamentally an object is just a collection of named data. Think like "JSON" object. It is possible to add and remove fields from objects (functionally) with expressions like `exp1 with (O)field2 = exp2`.

The bibliography format in BoVeX consists of declarations like this.

```
val knuth1981breaking =
  bib-article {(Article)
    title = "Breaking paragraphs into lines",
    author = "Knuth, Donald E. and Plass, Michael F.",
    journal = "Software: Practice and Experience",
    page-start = 1119,
    page-end = 1184,
    year = 1981,
    month = NOVEMBER,
    publisher = "Wiley Online Library",
  }
```

Each reference declares a set of optional fields. It is too tedious to make each field explicitly optional, and since the data have heterogeneous types, manipulating a string-indexed data structure would be more cumbersome. The bibliography rendering code case analyzes the presence of fields to render citations that have different subsets of data.

The `paper.bovex` source file is a tree structure with optional attributes on each node, which are represented with an object. This paragraph is written in the `paper.bovex` source file as:

```
Another use is in the [tt[layout]] type.
This is a primitive type that most of a
document's text is written in. It is a
tree structure with optional attributes
on each node, which are represented with
an object. For example, this paragraph is
written in the [tt[paper.bovex]] source
file as:
```

The square brackets are used to write a layout literal (the main body of the document is inside one large literal). Layout literals can also embed expressions (of type `layout`) with nested square brackets. Here the function `tt` is applied to a layout literal that contains text like `paper.bovex`. The `tt` function just adds the `font-family` attribute with value `"FixederSysLight"` to the layout node. This is a

custom monospaced bitmap font that I made for this paper using software I wrote. It is part of the `FixederSys` family.^[25] Functions like `b` and `it` apply bold and italic text styles, but functions can do anything that you can do in a general-purpose programming language.

Primops

Objects are used for interfacing with the runtime that executes BoVeX bytecode. There are about 50 different built-in primops that can be used by BoVeX programs. These include simple operations such as addition, but also heavyweight operations such as “load and register this collection of TrueType font files as a font family” or “invoke the boxes-and-glue packing algorithm with these parameters.” Primops in the former category work naturally on simple base types, but primops in the latter category need to be able to pass complicated tree-structured heterogeneous data between the BoVeX bytecode executor and the runtime. It would be possible for the runtime to consume and create BoVeX values like tuples and lists, but this has two problems: one, many types like list are declared as user code (in the BoVeX standard library); they are not special, and we don’t want to make them special by informing the runtime of them. Two, requiring specific representations at the runtime boundary inhibits optimization; for example, we can normally analyze the whole program to flatten data structures or remove record fields that are never used. The runtime typically uses `obj` to communicate structured data.

For example, the `pack-boxes` primitive runs the boxes-and-glue algorithm. It takes some layout (which is expected to be a series of box nodes, with attributes giving their size, glue properties, and so on) and configuration parameters like the type of justification and algorithm to use. It returns an object with a new layout (the boxes grouped into lines, with new glued up widths) as well as the total badness. Inside the BoVeX layout support code, this primop is wrapped as `pack-boxes` with a native, typed interface, so programmers do not need to think about that implementation detail. Other typographic features that benefit from runtime support are implemented this way as well.

Typographic features

The `pack-boxes` algorithm, which is offered by BoVeX, can be used to nicely justify text. It can also be used to distribute paragraphs into columns, by

thinking of the paragraphs as “words” (acceptable to break at any line, but bad to break near the start or end of a paragraph) and the columns as “lines.” It could be used by the document author for other purposes, I guess. There are other typographic features available.

The BoVeX rendering pipeline is a bit like a compiler: It transforms programmer-written source layout into formatted paragraphs, then into boxes of known size, then into stickers of known position. At the end, it outputs the document as a PDF.

The rendering process can report “badness” by calling the `emit-badness` primop. Badness is measured in square points of area that are outside of the container. Worse situations—such as text overlapping other text—have their badness scaled up per the same area of typographic horror. Less serious infractions—such as a little too much space between words—have badness scaled down.

Fonts

BoVeX can render your document in plain Times Roman if you don’t care about anything, or it can load 13 other boring PDF fonts, or it can load any TrueType font from font files. (They do not need to be “installed,” and it won’t help to install them. You just put them in the directory with your document.) It loads their kerning tables and applies kerning properly, by generating rigid boxes at the sub-word level with unbreakable glue. I was disappointed to find that most fonts include only a few dozen kerning pairs. They do this in order to “save space” in the font file, which is utterly rich coming from someone that would try to save space inside of words by squeezing letters together! In the current font Palatino, the word “BoVeX” is not kerned correctly because the rare bigraph “oV” does not have a kerning pair. I hope to improve this detail in a future version (perhaps for the presumably forthcoming video version of this paper).

Hyphenation

In A.D. 1455, Johannes Gutenberg invented the hyphen for his Gutenberg Bible, then just known as Bible.^[26] His printing process required the lines to all be the same length, so he had to stick these little guys all over the place. His hyphens looked like this: \neq . Later on we straightened these out and decided we only needed one at a time, and today we use them not because we *require* our lines to all be the same length, but because we like the cognitive

challenge of remembering the beginning of the word while we move our eyes to the beginning of the next line while reading.

In BoVeX, we break each word into boxes at legal hyphenation points, marking these points as sort-of-bad to break, and if we do break, we insert the hyphen character and use a little more space. By default, the hyphen sticks out of the end of the line a little bit. This is actually a bug but I like it.

The hyphenation algorithm I use is the same as the one used by TeX, which is cleverly represented as a prioritized set of patterns in order to fit compactly in memory.^[27] Again, you have to respect Knuth and crew's attention to detail, although to be fair this algorithm also dates to a time when storing a spell check dictionary in a computer's memory was described as "not feasible." So some of this was out of necessity. One of the nice things about the representation is that it generalizes to words that were not in the 1974 Merriam-Webster Pocket Dictionary. For example it hyphenates SIGBOVIK correctly.

The details really go on and on. The hyphenation dictionary is stored in a file called *hyph-en-us.tex*. "hyph" here, of course, stands for hyphens, and "en-us" means "English (United States)." In fact it is the *standard* language code for US English in the Small Language Model called IETF BCP 47.^[28] But then we have "hyph-en", which is a plausible hyphenation of "hyphen"! You could even read it as "hyphen us, tex", as a request for TeX to hyphenate the words in this file. This is the kind of detail I'm talking about! (There is also *hyph-uk*, which for once sounds a little less dignified than the US accent.)

Rephrasing

The BoVeX LLM can be used to rephrase text in order to make it more beautiful.

Unlike monospaced text, a line of proportional text never fits exactly (badness 0), so we need to apply some glue to make it fit. This generally has a cost, even when the text looks great.

The paragraph to be rephrased is some `layout` value. I generate a textual representation for it and feed it to the LLM. The prompt looks like this:

Exercise in rephrasing text. The following paragraph, which appears between `<P>` and `</P>` tags, needs to be rephrased so that it retains its precise meaning, but with minor variations in the specific choice of words, punctuation, and so on. No new facts should be introduced or removed, and all the ideas from the original paragraph should appear. However, it is good to use synonyms and change the word order and phrasing.

The text contains markup as well. There are two types: `text goes here` and ``. These should be preserved in the rephrased text. `` tags absolutely need to be retained and should not change their sources, although it is permissible to move them around in the text. `` should generally be retained, but the contents could change. The classes of spans may not change, and only the classes that appear in the original text may be used.

The first part is basically the same as what I used for the monospaced version, except that I ask the LLM to delimit the paragraph. This is important so that I know when it thinks it's done, and seems to work better than looking for newlines or the end-of-stream token. The second part is new. I translate the layout into plain text where uninterpreted subtrees are replaced with ``. These are generally boxes whose contents are not text. This could be an actual inline image or layout used to control rendering, like some bit of horizontal space. Nodes that are used to set text properties of the subtrees with attributes (like fonts, colors, ^{sizes}, etc.) are translated into distinct classes and marked up with `...`. The LLM has seen plenty of HTML, so it's able to use these reasonably well.

After generating a rephrasing, I parse the output HTML and match it up with the original layout. If I find any broken HTML, it is rejected. If I find any `` tag referencing a `src` not in the original, it is rejected. If I find any `` tag referencing a `class` not in the original, it is rejected. The more complexity that the original layout has, the higher the chance of a rejection, but rephrasing generally succeeds. But rejecting samples slows us down, so I leave off the second part of the prompt in the common case that the input paragraph is plain text. That way the LLM doesn't even try using markup.

BoVeX can rephrase the text with HTML and original layout matched up. It preserves nested layout and attributes. The rendering process continues.

However, we can't just use the badness score to tell us how good our rephrasing is. The badness score is a measure of how bad the line breaks are, and it doesn't tell us anything about how good the rephrasing is. We need to find a way to combine the badness score with the semantic loss to get a measure of how good the rephrasing is overall.

I wish I could tell you I solved this one with a beautiful algorithm. But so far, I just have something reasonable that works. I generate many different rephrasings, with their semantic loss, and run each of them through the boxes-and-glue algorithm to get the typographic badness. I choose the one that optimizes the preferred tradeoff between semantic loss and typographic badness. This process is controlled by BoVeX code, which is in the source code of this very paper. Knuth has a very low tolerance for semantic loss, and knows that his algorithms produce good results without rephrasing. Lorem Epsom just wants it to look good and sound good. Both have published in SIGBOVIK 2024.

Our first step is to generate a rephrasing that matches the original text exactly. We do this by sampling the most probable path through the tree. At each node, we take the first token with the highest probability. This is our first rephrasing, and it usually matches the original text exactly. We compute the semantic loss as the average probability mass skipped over all the tokens in the path. For this first path, we always took the most probable token, so this is 0.0 by definition.

The next path we explore will diverge from this path at some node (maybe the root). We pick a node that is likely to result in a good final loss, by scoring each node in the tree. The score is the average probability of all ancestor nodes times the probability of the next highest-probability token that we have not yet explored. The node with the highest overall score is the one we expand, by choosing that next highest-probability token. We are now in an unexplored part of the tree, and so we sample the most probable nodes repeatedly until we reach `</P>`. Speaking of which, BoVeX has a heck of a time trying to rephrase these last few paragraphs because they literally contain the text `</P>` in them.

The scores should be seen as heuristic. We would get different results by choosing different ways of computing the score. This is an example of a "beam search" algorithm. This connects the project to Super Metroid, which inspired this work. In the game, one of the final things you do is acquire the

"hyper beam" to defeat Mother Brain.

Since we will run the boxes and glue algorithm on multiple related texts, I generalized that algorithm to work on tree-structured input. This is clean; the memo table keeps the same dimensions, but records an additional fact. Now we store the penalty, whether to break after this token, and what the best subtree is. We have to consult each subtree when computing the score for a node, but this does not affect the asymptotic runtime. The table size is still at most $O(n^2)$, and although we explore more children per node, branches in the tree reduce the maximum depth to the root, which actually reduces one of the factors of n to $\log(n)$ as the tree becomes complete. However, as the SIGBOVIK deadline crept upon us, I never actually hooked this functionality up. It would require additional (programming) work to merge the trees, and the layout process is so fast that it doesn't matter; I can easily run the full layout algorithm on hundreds of rephrasings per paragraph.

I would like to improve the algorithm, because it does seem like there should be a way to integrate the boxes-and-glue dynamic programming algorithm with the path extension algorithm so that we prioritize exploring nodes that are likely to generate the best balance of typographic and semantic quality. It won't be as satisfyingly optimal as boxes-and-glue itself because we have incomplete information (we never know whether one of the exponentially many paths starts out with improbable tokens but then ends with a miracle streak of probable tokens). But it can certainly be more satisfying. Knuth would not stop here (but this is an *Any%* Knuth speedrun).

I spent my time implementing an achievement system in BoVeX. The first time certain conditions are met, the system awards you an achievement and prints a nice color trophy on your terminal. For example, you can get the "Not bad" achievement for generating a document that is at least 5 pages and has less than 1000 badness per page.

Advantages of rephrasing

The manual rephrasing of trivial details, such as synonyms and word order, can be automated. For example, when I wrote the opening paragraph of this paper, I could have used the same synonyms each time, and let the typographic considerations determine which synonym to use each time.

Conclusion

This paper—and *with* this paper—present BoVeX, a new computer typesetting system. It follows the tradition of TeX, but with modern amenities such as requiring over 128 gigabytes of RAM. Though some may consider the addition of AI features to TeX to be an unnecessary perversion, I find this use of LLMs to be *fully justified*.

Future work

Typographic features. Footnotes are desirable. It is hard to write a paper without them. Where am I supposed to put the bonus digressions? The layout of footnotes is tricky and should be part of a general floating figure implementation. End notes are easy, but I don't want them. I want them to be little footnotes so that you can't help but read them.

The SIGBOVIK program committee does not support page numbers, and this is good because page numbers are forbidden by BoVeX.

TeX is famous for its mathematical typesetting as well. It would fit neatly into BoVeX in the same way, since both use the same fundamental boxes-and-glue engine. BoVeX does not have "macros" or "modes" like TeX, but it would work cleanly to write a BoVeX function $\$$ (or, if you like, $\$$) that parses a custom syntax. In fact it would be natural to have different parsers for different maths, so that you don't need to parse $\$ \rightarrow$ as $\$$ in mathematical contexts that don't use $\$$ at all.

Optimization. There are many opportunities to make BoVeX code faster. This is mostly important for when it is being run in a loop in order to try out many different rephrasings of the same text. (That said, I do not wish to preclude what could be done with BoVeX by assuming its execution is doing only typesetting tasks. For example, shouldn't you be able to challenge your paper's reviewers to a game of chess against a strong engine embedded within your document?) The first thing to fix is that it manipulates too many strings at runtime (e.g. the code, record labels, object fields, and "registers"). This is easy to fix since these are all known at compile time. There are lots of high-level optimizations left to do for the IL code (common subexpression elimination, constant argument removal, uncurrying, etc.) and lots of peephole and control-flow optimizations left to do for the bytecode (currently no optimizations are performed at all). All of this becomes more important if I add

another planned feature, which is the ability for the document to be globally optimized by applying a black-box optimizer to a set of user-specified parameters. For example, the column width, line spacing, or font size could be tweaked to make the document fit better. This feature is "Auto-Margin Plus." Things are already set up to do this pretty straightforwardly; we would simply generate the document over and over while searching over the parameter space, and choose the one with the least badness. This may also affect which rephrasings look best. But instead I spent my precious time implementing **3D text**.^[29]

Reproducibility. The algorithm for rephrasing text tries to find the best place to explore the next most likely token from the probability distribution. This expects the generation of these distributions to be deterministic. Mathematically, inference is deterministic (it is just a bunch of matrix multiplications), so this "should work." But in practice the enormous calculation is performed in an unpredictable order as it is executed in parallel (in multiple CPU and GPU cores). Because floating point arithmetic is not associative (or distributive, commutative, or other properties you'd like), inference can sometimes generate different answers due to floating point round-off error.^[30] Alas, these are not even necessarily related to the final probabilities in the model, as billions of non-linear operations happen within the hidden layers of the network. The effect is not particularly grave; we might miss out on a highly likely path because the probability distribution was different the second time we looked at it. There are already lots of ways we might fail to find highly likely paths, so this is not some kind of reproducibility crisis. It is mostly just a bit unsatisfying.

Unicode support This would have been helpful when above I decided to show you Gutenberg's funny hyphen, \neq , for which I had to settle for embedding a crappy hand-drawn PNG file. Instead I could have used $U+2E17$, which since this exotic codepoint it is not present in the font Palatino, you could have experienced as \mathbb{N} . BoVeX is written with some Unicode support, with the main exception being that the PDF output code only supports the embarrassingly diminutive WinAnsiEncoding.^[31]

Deadlines. Although BoVeX itself is very fast, rephrasing is very slow. This presents a problem for the typical way that academic papers are written, which is to do all the work in a coffee-fueled fugue in the last few days before the deadline, then

stay up all night writing the paper and finding citations for the pro-forma “related work” section which you did last but you know that the reviewers will insist upon, and tweaking `\vspace` and `\begin{figure}[h!]` until it fits within the page limit. On the one hand, BoVeX does potentially free the author from the visual tweaking process. But on the other hand, the LLM inference for the rephrasing process can be quite slow, and it can take many hours or days to fully bake a long paper! For this reason, it may be better to change conference deadlines to a system where the pre-rephrasing text is submitted. The publishers (what do they even do?) can be the ones to execute the rephrasing in the cloud as they produce the “camera-ready copy.” With straightforward extensions, this would also allow the rephrasing to adapt to changes in the overall volume style, or to adjust to avoid embarrassing typographic coincidences with other articles in the same volume (such as using the same notation with a different meaning). In principle, the paper could edit itself to respond to feedback from reviewers, in a way that minimizes the semantic distance from the original. This rapid feedback loop could reduce the time to publication, perhaps to mere months, or even weeks!

Other ways to minimize badness. The BoVeX system allows the document author to exchange semantic consistency for higher quality typography. Although we achieve state-of-the-art results, there are likely points that are more Pareto-efficient than what BoVeX can reach. BoVeX uses one of the most powerful publicly available LLMs, but that model is limited to rewriting the text within narrow constraints. Irresponsible research has demonstrated that language models are capable of volition, taking actions and using tools to accomplish goals. With minor modifications, it is likely possible to expand the Pareto frontier of the semantic/typographic tradeoff. For example, sometimes we could improve the typographic quality of the text without any semantic loss, by acting on the world to *make the reworded text true*. Human authors do this already: Earlier when I was describing `internal-pack-boxes`, rather than explain the somewhat awkward implementation, I went back and changed the already-working code so that it would serve as a simpler example of how printers use `obj`, but still be truthful. Now imagine the difficulty in typesetting a statement like “The universe contains approximately 1,000,000,000 paperclips,” and how much more beautiful the text could be if that number were instead 10,000,000,000,000,000,000,000,000,000,000,000,000,000,000,000!

In the meantime, there is an easier way to get zero badness: Delete the whole document! As a wise person once said, “If you can’t say something with non-zero typographic or semantic loss, don’t say anything at all.”

Acknowledgements. I’d like to thank my advisor Karl Cray for his help and support. 20 years ago, we set out on an ill-fated attempt to replace LaTeX with an SML-like language mTeX, which compiled into TeX macros. The nesting square brackets syntax was Karl’s idea, and BoVeX shares genetic material with mTeX for sure.

See you next mission,

Tom 7

Bibliography

[31] Adobe. “PDF reference: Sixth edition”. October 2006.

[7] You can just go to arxiv.org and click on any random article these days.

[16] N Bijlage. “Knuth meets NTG members”. *NTG: MAPS*, 16. March 1996. pp. 38–49.

[3] Russ Bynum. “Bernie Sanders wants the US to adopt a 32-hour workweek. Could workers and companies benefit?”. March 2024.

[12] <https://github.com/ggerganov/llama.cpp>. ggerganov. March 2024.

[26] Johann Gutenberg. “Bible”. 1455.

[4] <https://allfamousbirthday.com/donald-knuth/>. February 2024.

[23] <https://ocaml.org>. “The Objective Caml system”. 2023.

[21] Graham Hutton. “Higher-order functions for parsing”. *Journal of functional programming*, 2(3). 1992. pp. 323–343.

[5] Donald E Knuth. “The Art of Computer Programming: Volume 4A, Combinatorial Algorithms Part 1”. *Addison-Wesley*. January 2011. 912 pages.

[6] Donald E Knuth. “The Art of Computer Pro-

gramming: Volume 4B, Combinatorial Algorithms Part 2". *Addison-Wesley*. October 2022. 736 pages.

[15] Donald E Knuth, Michael F Plass. "Breaking paragraphs into lines". *Software: Practice and Experience*, 11(11). November 1981. pp. 1119–1184.

[27] Franklin Mark Liang. "Word Hy-phen-a-tion by Com-put-er". 1983.

[17] Robin Milner, Mads Tofte, Robert Harper, David MacQueen. "The definition of Standard ML (Revised)". *MIT Press*. May 1997. 114 pages.

[8] John C Mitchell. "Type Systems for Programming Languages". Van Leeuwen, Jan, ed. *Formal Models and Semantics*. 1990. pp. 365–458.

[1] Tom Murphy VII. "Badness 0 (Knuth's version)". *SIGBOVIK*. April 2024. 16 pages.

[30] Tom Murphy VII. "GradIEEEEnt half decent". *SIGBOVIK*. March 2023. pp. 33–56.

[18] Tom Murphy VII. "Modal Types for Mobile Code". January 2008.

[13] Tom Murphy VII. "NaN gates and flip FLOPS". *SIGBOVIK*. April 2019.

[10] Tom Murphy VII. "The First Level of Super Mario Bros. is Easy with Lexicographic Orderings and Time Travel. After that it gets a little tricky". *SIGBOVIK*. April 2013. pp. 112–133.

[20] Tom Murphy VII. "The Wizard of TILT: Efficient(?), Convenient and Abstract Type Representations". *Carnegie Mellon tech report CMU-CS-02-120*. March 2002.

[29] Tom Murphy VII. "The glEnd() of Zelda". *SIGBOVIK*. April 2016. pp. 105–112.

[14] Tom Murphy VII. "ZM~~ # PPrinty# C with ABC!". *SIGBOVIK*. April 2017. pp. 129–148.

[25] <http://tom7.org/fixedersys/>. Tom Murphy VII. "The FixederSys font family". 2024.

[2] https://en.wikipedia.org/wiki/Wikipedia:WikiProject_Punctuation. Tom Murphy VII. "WikiProject Punctuation". July 2007.

[19] <https://sourceforge.net/p/tom7misc/svn/HEAD/tree/trunk/rephrase/>. Tom Murphy VII. "BoVeX source code". 2024.

[22] Chris Okasaki. "Even higher-order functions for parsing or Why would anyone ever want to use a sixth-order function?". *Journal of Functional Programming*, 8(2). March 1998. pp. 195–199.

[28] A Phillips, M Davis. "Tags for identifying languages". September 2009.

[9] <https://gamefaqs.gamespot.com/snes/588741-super-metroid/faqs/10114>. rs1n. "Super Metroid Speed Guide and FAQ". 1996.

[24] Didier Rémy, Jérôme Vouillon. "Objective ML: An effective object-oriented extension to ML". *In Theory And Practice of Objects Systems*, 4(1). 1998. pp. 27–50.

[11] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, Thomas Scialom. "Llama 2: Open foundation and fine-tuned chat models". *ArXiv.org*. July 2023.

Badness 0 (Knuth's version)

Dr. Tom Murphy VII, Ph.D.
March 2024

It has become clear to me that many people walk this Earth completely unbothered by incorrect details. For example, they are unconcerned when a hyperlink includes a surrounding space character. It doesn't upset them when the screw heads on a light switch wall plate are not all lined up. They didn't notice that the rules of Wordle's "hard mode" are simply wrong. They care as much as the phone's autocorrect (none) about the difference between "its" and "it's." When someone asks, "Will you marry me?" and they think "Oh my god!" it's not because the proposer probably should have used the subjunctive *would*.

I am... not like this. If a character in a TV commercial is handling a coffee cup but I can infer from its moment of inertia that the cup does not contain any liquid, I immediately lose suspension of disbelief and will not purchase the product featured in the commercial. I literally projectile vomit if Auto-Motion Plus is enabled on a television in the hotel I'm staying in, even if the TV is not turned on, or if someone misuses the word "literally." If I see a paragraph missing a period at its end on Wikipedia, I will spend dozens of hours writing software to organize and semi-automate a distributed effort to fix all the missing periods on Wikipedia.^[2] And worse, each time I learn of a new type of mistake, I am forever cursed to notice that mistake

Seriously: One time I found myself spell-correcting someone else's *lorem ipsum* text in a slide. It said "lorem epsom," which is funny. I think about that incident all the time. The person that wrote the slide probably thinks about things like leveraging synergy, generative AI, metaverses, blockchain 3.0, snackable content, being eco-green, and so on, without it occurring to him that these things could have nuance and meaning separate from their names. He has probably never even read the Wikipedia article on Lorem Ipsum. He is successful and rich.

Another successful person is the congressperson Bill Cassidy. Criticizing a proposed bill that would reduce the standard work week in the US by 8

hours, from 40 to 32, this senator says,

Sen. Bill Cassidy of Louisiana, representing the _____ party, said paying workers the same wages for fewer hours would force employers to pass the cost of hiring more workers along to consumers.

"It would threaten millions of small businesses operating on a razor-thin margin because they're unable to find enough workers," said Cassidy. "Now they've got the same workers, but only for three-quarters of the time. And they have to hire more."

Actually, that's not exactly the quote, but I needed to make it look nice.^[3] And this is not a paper about politics, but let's just say you can guess what word goes in the blank.

Anyway, OKAY, first of all, razors famously have high margins. It's like the worst possible metaphor here.

For another thing: This guy mixes fancy typographic quotes and ASCII ones.

But the main thing I want to talk about is: What? No! 32/40 is four fifths, not three quarters. This is not, like, complicated math. It uses some of the world's smallest integers. Everybody knows that the work week is 40 hours, and that a work day is 8 hours, and that the proposed bill reduces it by one day, giving four of five days. I don't really mind if someone makes an error in calculation (well, I do mind, but I am certainly prone to doing it). The infuriating realization here is that this person does not even think of "three-quarters" as a kind of thing that can be right or wrong. He says three quarters because it makes smaller number feelings. You could imagine him having the conversation (with me, perhaps): "You say four-fifths, I say three-quarters." Me: "But it is four fifths. And why are you always hyphenating it?" Him (smiling patronizingly): "I guess we just have to agree to disagree."

The opposite of this person is the hero called Donald Knuth.

I'm not saying that Donald Knuth isn't successful and rich. According to the website "Famous Birthdays,"^[4] which is probably generated by AI or at least by people whose economic output is measured in a count of words, and words whose value is computed by their ability to drive ad clicks, Don-

ald Knuth is “is one of the most popular and richest Mathematician who was born on January 10, 1938 in Wisconsin, Wisconsin, United States. Mathematician and engineer who was arguably most recognized as the Professor Emeritus at Stanford in Palo Alto, California.” As one of the richest Mathematician from United States, according to the analysis of Famous Birthdays, Wikipedia, Forbes & Business Insider, “Donald Knuth’s net worth \$3–5 Million.”*

I suppose it is arguable that he is the Professor Emeritus. And it is very likely true that he is the *only* popular and rich mathematician born on that specific day in Wisconsin, making the singular “Mathematician” perhaps a technical master-stroke. But more likely this is just an amusingly dense series of imprecisions. The asterisk of course does not have any referent on the page.

What I mean when I say that Donald Knuth is the opposite of this person is that Knuth is interested in unpacking a single unnecessary detail, recursively, until it is completely solved. According to the website Famous Bibliophiles, one day Donald Knuth set out to write down the entire subject of computer science in a single book called *The Art of Computer Programming*. As he was doing so, he realized that describing computer algorithms in a lasting form would require a programming language that was not subject to constant revision, so he invented the MIX instruction set for an idealized computer. After writing some 3000 pages out in longhand, he found that it was impractical to print them all in one book, so the plan expanded to be multiple volumes. Then when he got a draft of one of the books back from the typesetter, he was unhappy with the details of the typography, and so he paused his work writing down all of computer science to create some new computer science: First an algorithm for determining where to place line breaks in order to make text optimally beautiful, then algorithms for hyphenating words, then generalizations of these for typesetting mathematics, and then a full computer typesetting system that is still in wide use today, called TeX. Along the way he was unsatisfied with the specific typefaces that existed in the world, and unsatisfied with the way that typefaces were described at only one weight, and so he created the parameterized METAFONT system and several new typefaces. Undeterred by these excursions, he returned to his original task of writing down the entirety of computer science, using all the technology he had built. By the time he finished this, much more computer science had

been invented, including by his own hand, and so he needed to rework MIX for the next volume, and update the first. The revised plan of eight volumes remains the intention in 2024. However, he found that the volumes were getting rather long, and began releasing portions of volumes (“fascicles”). So far, Volume 4 has been partially published as books 4A^[5] (fascicles 0–4; 912 pages) and 4B^[6] (fascicles 5–6; 736 pages). It is unknown how many more episodes remain in Volume 4. I expect that every conversation that Knuth has with his editor goes like this. Editor: “Hey, Donald, I hope you’re well. Just wondering if you have an update on when 4C will be ready? Or any more icicles?” Donald E. Knuth: “I am working diligently on fascicles for Volume 4C. As I’ve mentioned in the past, it’s impossible to tell how long it will be, since mathematics does not obey the rules of project management.” Editor: “I just need a date to tell the publishers.” Donald E. Knuth: “Like I’ve said, any date would be very low confidence, other than the fact that it will be in the future.” Editor: “I just need a date.” Donald E. Knuth: “Would you like me to say a date, knowing that it’s a very low confidence guess, and that I would be extremely likely to miss that date, or even deliver early?” Editor: “Early! Now we’re talking.” Donald E. Knuth: “What use is the date if you’re excited about the possibility of it being early, relative to some unknown date?” Editor: “I just need a date for the publishers.” Donald E. Knuth: “2030.” Editor: “Thanks Donald, you’re the best!”

Volume 5 is estimated to be ready in 2030, when Knuth will be 92.

That’s a large amount of language!

Nightmare on LLM street

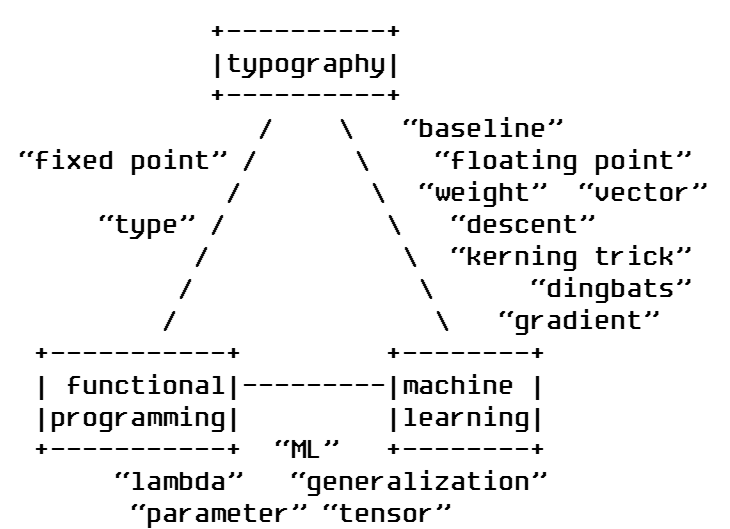
Then we have Large Language Models.^[7] One of the irritating things about LLMs is that they are so buzzwordy, but unlike most buzzwordy trends, they are actually substantive. They produce remarkably fluent text. With no additional training they frequently beat purpose-built models that have been in development for decades. They generalize to completely new situations.

So many things about “AI” distress me. *Dolor sit amet!* I worry about the devaluation of human creativity, about large-scale disinformation and spam ruining the beautiful library of knowledge that humans have created, about extreme concentration of wealth. And yes, I worry about competing with

AI. Being able to work tirelessly and thousands of times faster than humans is a huge competitive advantage. Of course, I find some solace in the significant possible upsides. It might help us solve hard problems like climate change and AI. But even in the best scenarios we will not be able to ignore it: Even if it never gets as smart and precise as Knuth, it's already too economically useful in its Lorem Ipsum state (just like Lorem Ipsum himself).

On the other hand, the technology is pretty neat and lends itself to some nice abstractions. I love playing with words. So one of my side quests is to masticate this whole scenario by experimenting with LLMs in practical and impractical applications, and to try to make it fun (for me) to program with them.

Many things irritate me, so this is something I have ample experience with. I have a myriad of strategies for digestion of them. For this work I'm inspired by the "Hurry-Coward So-so-morphism," where I make connections between topics based solely on confusion of superficial lexical similarities without regard to their underlying meaning. So for example we have "ML" meaning both "Machine Learning" and "Meta Language", as well as "type" both as in "typeface" and as in "type systems for programming languages."^[8] And because machine learning has claimed so many words, there are a great many shared with typography as well:



By no coincidence, I already spent a lengthy introduction talking about Donald Knuth's work in computer typography. So now I can tell you what this paper is about. If in our near AI future we are giving up on precision, perhaps at least we can

have something that we want: Perfect typography? This paper is about a new typesetting system, BoVeX, which allows for the controlled exchange of precision for beauty. It essentially gives us a dial between Lorem Ipsum and Donald Knuth. To illustrate, we'll first look at a simpler case by inspecting one of my other interests: Super Metroid.

The scientists' findings were astounding! They discovered that the powers of the Metroid might be harnessed for the good of civilization!

Metroid is a video game series about a brain that has been enslaved inside a jar in an underground datacenter on the planet Zebes. This brain is called Mother Brain and its goal is to control the hypercapitalists called Space Pirates to increase their "score" as high as possible by conquering planets throughout the galaxy. Mother Brain was invented by the Space Pirates, although it is not clear whether the current situation was actually intended by the Space Pirates. The most super version of Metroid is Super Metroid.

In the 1990s the website gamefaqs.com collected plain text "FAQs" for classic video games, then just known as video games. On this site another hero was born. They were writing the definitive guide to speedrunning the SNES game Super Metroid when they saw that some of their ASCII lines ended up exactly the same length, and that it looked good:

Once you save the game at your ship (about 1 hour 15 minutes is good), go down to Tourian. Do not save your game in Tourian if you have intentions of returning to any previously explored section on Planet Zebes. There will be a few Metroids to kill before you reach Mother Brain, and they must all die in order to continue to Mother Brain. Read the boss guide for more details. Once Mother Brain is defeated, you will need to hurry back to your ship. By now you will already have the HYPER BEAM. From Mother Brain's room, go west and then south. Take the blue door at the bottom and speed dash east. Super jump up, and continue north. Once you land up top and are running east, aim diagonally down to the right and shoot an unseen door. Eventually, you will get to this door since lava will start to rise from the floor in this area. Speed dash through the door you preopened, and charge for a super jump. Hug either the left or right wall in the Craterian shaft and super jump up. Now quickly get to your ship before the planet explodes. There should be almost a minute left on the timer. Sit back and watch the ending! Did you beat the game within 1 hour and 20 minutes?

and so they wisely decided to wordsmith the entire 28-page guide so that every line was exactly the same length, with no extra spaces or other cheating, just because it could be done.^[9]

Doing this manually is a chore, and I do like to automate the chores of Speedrunners.^[10] I got this working in an afternoon. It's, like, easy mode. For a paragraph of text and a target line length, I ask the LLM to remember the paragraph and recite it. The prompt looks like this:

Exercise in rephrasing text. The following paragraph needs to be rephrased so that it retains its precise meaning, but with minor variations in the specific choice of words, punctuation, and so on. No new facts should be introduced or removed, but it is good to use synonyms and change the word order and phrasing.

After this I insert something like `Original text:` followed by the original paragraph, then `Rephrased text:`. The model is ready to generate tokens.

I then sample text a word at a time to continue this prompt. If a line ends exactly on the number of characters that I want (and the next character is a space or other character that is appropriate to end a line) then I accept the stream so far and continue. If I exceed the line length, I back up to the state at the beginning of the line and try again with new random samples. I just keep doing that until the paragraph is complete, and we have beautifully justified monospace text that resembles the original. Here is an example of this paragraph rendered in monospace:

```
I sample text a word at a time to continue
this prompt. If a line ends exactly on the
number of characters I want, I accept that
text so far, and continue. If I exceed the
line length, I back up to the beginning of
the line and try again with new samples. I
keep repeating this until I get text I can
render in monospaced font, and that is how
we can get beautifully justified monospace
text. Here is an example of this paragraph
rendered in monospace:
```

You could argue that this is improved, even, by making the text shorter. It does use “monospace” and “monospaced” inconsistently. The most upsetting thing here is that it ends with a colon like there’s going to be another example of the paragraph, but that’s what I asked it to do.

The approach described works reasonably well, but it has several deficiencies (such as: it only took an afternoon) that we’ll address for the real BoVeX system. But it is a good example to explain some concepts that will be useful later.

¿Como te LLama?

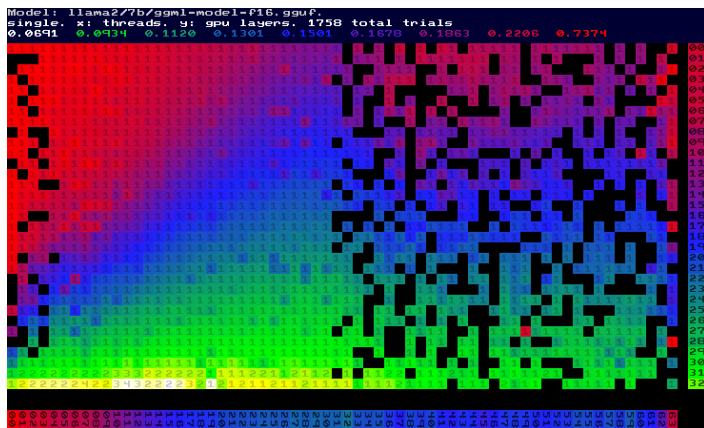
Llama is Facebook’s Large Language Model,^[11] which they nicely share with anyone who agrees not to use it to destroy the world. Wouldn’t it be funny if the world is destroyed by something called “Llama”? That’s some Stay-Puft

Marshmallow Man stuff. Actually I hear that llamas are pretty mean, and if you are thinking about hugging a cute long-neck, you are probably thinking about an alpaca. But that’s probably a version of the linear algebra package LAPACK. Llama-v2-70b is a good LLM which can do some impressive things, but when I say destroy the world I mean stuff like filling the internet with infinite spam, or building critical infrastructure on it in order to cut costs, where most of our “safety” measures consist of asking the model politely to recite its daily affirmations before performing its tasks. That kind of thing. It’ll be at least months before we *really* have to worry.

Anyway, the normal way to program with Llama is to use Python, and a mountain of things that you are not supposed to understand and cannot understand, mostly by pasting examples from others and then tweaking parameters and prompts. I don’t care for it. Fortunately, human geniuses^[12] have implemented the inference code for llama-like models in a nice, portable C++ library called “llama.cpp” (checks out).

With llama.cpp, I can load a quantized version of the model into RAM. Actually there are two different models, the 7b and the 70b, referring to the number of billions of parameters, which must be a multiple of VII “for performance reasons.” The parameters are the weights on the layers of the network. At native 16-bit floats, the 70b model will fit in about 130 GB of RAM, just slightly more than a nice round 128 GB, making one wonder what performance reasons they had in mind. But anyway, earlier this year I (physically) broke my computer trying to put the world’s (physically) largest video card into it, the GeForce 4090, and so I endowed the replacement computer with 256 GB of RAM. If you are ever looking at specifications for a high-end desktop computer, by the way, and wondering “who the heck buys these things and what do they do with them?” one answer is “me,” and the other answer is “this.”

Quantization means using fewer bits to represent the floating point weights.^[13] This saves memory, but it also speeds up inference, which needs to read pretty much the entire model for every predicted token. I got reasonable quality and good performance from LLama-v2-7b with 16-bit floats. This one fits completely on my world’s (physically) largest GPU. In order to tune various settings, I ran thousands of trials for the different models, and made some nice custom graphs:



Tuning results for Llama-v2-7b with 16-bit floats. The x axis is the number of CPU threads and the y axis is the number of model layers that have been loaded onto the GPU. As expected, increasing the number of threads and layers on the GPU improves performance, since this whole model fits on the GPU. For the 70b models (not pictured) there is an abrupt drop-off in throughput before we load all the layers, and also my computer gets very sluggish if I exceed the GPU memory. We see that if we use more than the number of physical cores (32) we do not see any benefit, which is not surprising because hyperthreading basically never helps anything. The best throughput actually uses a modest number of cores (about 12). Mostly I'm just including the graph to demonstrate that Bo-VeX has support for including PNG files.

Where was I? Right. Fundamentally, LLMs are trained to predict a token (like a word or part of a word) given some sequence of tokens that precedes them. There's a fixed set of tokens for the model, and rather than predict a single token, they actually give a score for every possible token. These scores are typically normalized into a probability distribution. So for example if we have the text

SIGBOVIK is an

then the probability distribution (Llama-v2-70b) begins as

- (annual) 69.8010%
- (April) 3.8023%
- (ac) 3.2456%
- (academic) 2.9374%
- (artificial) 2.0857%
- (open) 1.7993%
- (under) 1.2331%
- (international) 1.1032%
- ...

with the thousands of other tokens following. So three-quarters of the time the next token should be " annual" but there are many other reasonable possibilities. We can pick one of these tokens however we like, append it to the sequence, and run the model again. This gives us a new probability distribution. By doing this over and over we can generate a likely piece of text. This is what Lorem Epsom means when he says "Generative AI." Rather, what he means is "the new thing that is cool," but what he is unknowingly referring to is that you can sample a probability distribution. He has probably never even read the wikipedia article on Markov chains.

If I always sample the most likely token, I always get the most likely text. It is good to be likely; this is why the model is useful. However, you might not want exactly the same result each time, and in many situations if you only sample the most likely token, you get very boring, repetitive text. Pseudo-random number generation is the spice of life!

We also need not use the probability distribution to sample at all. We can just pick the token that we want. This is how the initial "prompt" works; we just run the inference process one token at a time but always select the next token in the prompt, ignoring the probabilities. So at each moment, the text we've generated so far (more or less) completely characterizes the state of the LLM. This means that we can easily go back to earlier moments and sample a different continuation of the text, by just replaying tokens. We also have the option of storing the LLM state (gigabytes) in RAM, which allows us to return to a previous state in constant time.

For generating monospaced lines of the same length, I use a prompt that asks the model to rephrase the input paragraph. Here, greedily sampling the distribution typically results in a copy of the input paragraph, which is fine for our purposes. (If the lines already happen to be the right length, we need not change them!) But when a line comes out the wrong length, I want to try again. So I save the model state whenever I begin a line. To produce a variant of the line, I sample tokens proportional to the probability distribution. When the set of probable lines is small (this is common), the process will keep generating the same lines and failing because they are not the right length. To prevent boredom, whenever the process repeats a line that's already been seen, I increase the "tempera-

ture" modifier to the probability distribution. This is an exponential factor that (when higher) flattens out the probability distribution, making previously unlikely tokens more likely. You can think of this like the model getting a little hot-headed as it frustratedly does the same thing over and over. This causes the candidate lines to be more varied, but less probable (according to the original probability distribution). I can reset the temperature when it succeeds, since we prefer to have more likely lines.

This is all there is to the monospacing version. It's just 300 lines of code, including boilerplate and commented-out debugging code and false starts.

Great !! You fulfilled your mission. It will revive peace in space. But, it may be invaded by the other Metroid. Pray for a true peace in space!

Now, looking at the output text, we feel satisfied that everything lines up exactly. However, we can't help but feel unsatisfied at the same time: Now we're looking at a monospaced font. Good for programming. Bad for publishing. Can we instead have excellent justified text with all the perks of proportional fonts and a programmable document preparation system? And can we have it by the estimated SIGBOVIK deadline so that it can be used to prepare the paper that I'm now writing? Maybe! This is the Donald Knuth Any% speedrun.

The boxes-and-glue algorithm

When justifying monospaced text, it looks quite bad^[14] to insert more than one space between words, so we have a simple way to tell if text is suitable for some width. We just add up the codepoints. For the full-on typography case with proportional fonts, there are many more degrees of freedom. For one thing, it looks fine to expand or contract the space between words a little bit, even if it varies from line to line. It is also possible to make fine adjustments in letter spacing (kerning) to squeeze or air out text. We can also hyphenate words.

Around the time I was being born, and probably being very upset about it, Knuth was having similar feelings about the way his computer-typeset documents looked. He discovered a nice abstraction that generalizes most of these typographic degrees of freedom, and devised an algorithm for producing optimal text layout given some parameters.^[15] The idea is to consider the text of a paragraph as consisting of rigid "boxes" (say, words)

and stretchy "glue" (say, space) between them. Both boxes and glue have various detail (and can be extended to support all sorts of quirks) but the basic algorithm can be understood with just those pieces. So, let's do that.

Knuth's paper (as usual) is great, but I started having spoiler feelings when reading it, so I figured out my own algorithm, which is more fun than reading. No doubt the key insight is the same: Although there are exponentially many possible break points, you do not need to try them all. Whenever we break after a word, the problem is now the same for the rest of the text (fit the rest of the text optimally onto lines, starting at the beginning of a line) no matter how we got there. This lends itself to a dynamic programming algorithm.

Dynamic programming is a programming technique for whiteboard interview problems at tech companies. I found it mysterious when I was young, perhaps because of its strange name. Here is how I think about it. Imagine you have a recursive procedure that solves the problem. In this case, the pseudocode is something like

```
pair<int, string> Split(string line,
                       string text) {
    if (text.empty()) return {0, ""};
    auto [word, rest] = GetFirstWord(text);
    // try splitting
    auto [penalty1, rest1] = Split(word, rest);
    penalty1 += badness from leftover space;
    // try not splitting
    auto [penalty2, rest2] =
        Split(line + " " + word, rest);
    penalty2 += badness from line too long;
    if (penalty1 < penalty2) {
        return {penalty1, word + "\n" + rest1};
    } else {
        return {penalty2, word + " " + rest2};
    }
}
```

Split takes the line so far and the text that remains to be split. In the normal case that there is a word left, it will try two possibilities: Either splitting after the first word, or not splitting. This is exponential time because each call makes two recursive calls, to try each of the two options. But deep recursive calls will be made with the same arguments many times. So, add some memoization: If the function is called for the same line and text a second time, just return the same answer as before without doing any work (especially not making recursive calls again). This limits the function to be

called at most once for each possible argument; we can then see that `line` is no longer than the input (so it is size $O(n)$) and `text` is always some suffix of the input (so it is size $O(n)$), giving $O(n^2)$ calls.

Dynamic programming is just memoization inside-out: We create the values for all the recursive calls before we will need them, store them in a table, and then look them up. For this problem, the table is indexed by the two parameters, the current line and the remaining text. Note that these two can be represented as integers. The line is just the number of words before the current word that are included on the line, and the text is just the position in the string where we'll next look for a word. That's all there is to it; the base cases of empty text are used to start the table, and then you just write the loop to fill out cells in the right order.

Knuth's boxes-and-glue algorithm contains many extensions, and so does mine. For example, later we'll talk about how you can adapt the algorithm to perform hyphenation and kerning. There are many rabbit holes to go down, and I explored the ones that attracted my attention. There is plenty of time to add more features later, since of course I have now cursed myself to use BoVeX for my future SIGBOVIK papers.

But here's where I diverge from Knuth somewhat. Knuth was reluctant to add a programming language to TeX,^[16] but I spent the majority of my time on this project implementing a full-fledged language. BoVeX is about 33,000 lines of code, the majority of which is the implementation of the language itself. That's 110× as long as the original monospace proof of concept, and 30× the length of this document!

The BoVeX language

This section describes the BoVeX programming language and its implementation. If you are just in it for the jokes, you can skip this section, which is basically serious and loaded with programming language theory jargon.

BoVeX is a typed functional programming language in the ML family. Its syntax closely resembles Standard ML. Here's an example piece of code from the source code of this document:

```
datatype (a) option = SOME of a | NONE

fun consume-outer-span f s =
  case layoutcase s of
    Node (SPAN, attrs, children) =>
      let
        val (ropt, rchildren) =
          case children of
            one :: nil => consume-outer-span f one
          | _ => (NONE, layout-concat children)
      in
        case (f attrs, ropt) of
          (NONE, _) => (ropt, span attrs rchildren)
        | (SOME vouter, inner as SOME _) =>
          (inner, rchildren)
        | (outer, NONE) => (outer, rchildren)
      end
    | _ => (NONE, s)
```

You don't need to understand it. I just want to show you that it is a full-fledged programming language. It supports higher order functions, polymorphism, algebraic datatypes, pattern matching, Hindley-Milner type inference, and so on. It is basically core (no modules) Standard ML,^[17] although I left out some warts (operator overloading, eq-types, abstype, non-uniform datatypes and polymorphic recursion) and added some new warts. For example, `as` allow patterns on both sides, since Standard ML has always seemed backwards to me and it works perfectly fine to just make it symmetric. Anyway, a full description of the language would be boring and take too much time as the SIGBOVIK deadline draws closer.

Implementation

I have implemented many similar languages in the past, including for my dissertation.^[18] It would have been expedient to start from one of my existing implementations, but they are mostly written in Standard ML and I couldn't get MLton to work on my Windows computer in 2024. So I started over from scratch in C++, which at least does work on my computer. (I also want to be able to interface with GPU inference code for running the LLM, which will be easiest from C++). C++ is not a good language for writing language implementations, but it has gotten better.

The BoVeX implementation is a "compiler" in the sense that it transforms the source language through multiple intermediate languages into a low-level bytecode. This bytecode is just straight-line code on an abstract machine with infinite registers and operations like `alloc` (allocate a new "object") and `setfield` (set a fixed field of the "object" to a value from a register). It does not produce machine code, and although this would be

pretty feasible, it would not be the first thing to do to make BoVeX faster.

First it concatenates the source files (handling `import` and keeping track of where each byte originated, for error messages) and lexes them into tokens. Then it parses those tokens into the External Language (EL), which is just the BoVeX grammar with a few pieces of syntactic sugar compiled away. It does syntactic transformations on the EL AST to remove some currying syntax and transform nullary datatypes (`nil` becomes `nil of unit`). Then it elaborates EL into a simpler and more explicit Internal Language (IL). Elaboration does type inference (Hindley-Milner) including polymorphic generalization and so on, compiles pattern matching into an efficient series of simpler constructs, and decomposes heavyweight stuff (e.g. `datatype`) into its constituent type-theoretic pieces (e.g. a polymorphic recursive sum). The IL is nice and clean, so it is a good place to perform optimizations. I love writing optimizations but I had to keep myself out of there, or else this would be a 2025 SIG-BOVIK paper. There are just enough to make the code reasonable to debug if I need to look at it. After optimization, I perform closure conversion, simplify again, and generate the final “bytecode” form. This entire process happens whenever you generate a BoVeX document; the only output from running `bovex.exe` is the PDF document.

I want you to know that I did not cut corners on the language implementation. For example, compiling mutually-recursive polymorphic functions is really obnoxious (AFAIK it requires either monomorphization or first-class polymorphism when you do closure conversion) but I did do it, even though none of the BoVeX code I used for this paper ever needed this feature. Following are some of the implementation details; for the full story you’ll need to check the source code.^[19]

AST pools. One of the main things I need to do is create tree-structured data to represent the abstract syntax tree of the various languages involved. This is very nice in ML (it is what the `datatype` declaration is for) and annoying in C++. I continued to experiment with different ways to do this. I use arena-style allocation for the syntax nodes (always `const` after creation), so that they can be created and reused at will. My current favorite approach to manipulating the nodes is to write “in” and “out” functions (tedious, manual) for each construct in the language. The syntax nodes can then be implemented however I like (for example, a flat struct

or `std::variant<>`) with the freedom to change. I get the compiler’s help whenever I change the language (which is often!) since each in/out function is explicit about its constituents.

Passes and guesses. Many transformations in a compiler rewrite a language to itself; for example each IL optimization is a function from IL to IL. These can be tedious to write and update, especially since a given optimization usually only cares about one or two constructs in the language. I use the “pass” idiom to write these. This is basically an identity function on the AST that pulls apart each node, calls a virtual function for that node, and then rebuilds the node. To write a pass that only cares about one type of node, you inherit from this class and then just override that one node’s function. One issue with this is that each time you rebuild the entire tree you create a lot of unnecessary node copies. So exchanging tedium (mine) for efficiency (my computer’s?), every node type’s “in” function also takes a “guess” node pointer. If the node being constructed is exactly equal to the guess, then we return the guess and avoid creating a copy. Then the base pass is actually the identity (it returns the same pointer) and does no long-lived allocations. This seems to be a good compromise between the traditional garbage-fountain approach and hash consing, which sounds like it would be a good idea but is usually just a lot slower.^[20] For type-directed transformations, there is also a typed IL pass class, which recursively passes a context and does bidirectional type checking of the intermediate code. Closure conversion is a type-directed pass and is implemented this way.

Parsing. I have this aversion to parser generators, probably because one time I tried to get someone else’s code to compile and it complained about having the wrong bovines on my computer and ruined my weekend. After trying some other people’s C++ parsing libraries and being disappointed by them, I did what Knuth would do: I wrote my own. It is a parser combinator^[21] library which actually descends directly from Okasaki’s SML code.^[22] I was proud of myself for getting this to work in C++, since C++’s insane type system is impossible to understand and its error messages are even worse. (BoVeX’s error messages are extremely spartan, often simply declaring `Parse error at paper.bovex line 1`, but in many ways this is more useful than C++’s mile-long SFINAE vomitus.) It supports mutually-recursive parsers, resolution of dynamic infix operators, and all that. My template-heavy parser combinators take clang about a minute to

compile, which is acceptable. Less acceptable, but something I only learned after using this to write a 16–page-long paper, is that the parsers are very slow. Putting aside LLM inference, this paper takes 13 seconds to render into a PDF, 11 seconds of which is parsing! There must be some bug, but I don't know if it's in my grammar (it is easy to accidentally write an exponential time parser, but this one should not be) or the parser combinator library (also my fault) or clang producing bad code (it may be giving up on optimizations, since it is taking so long to compile; the .o file is 41 megabytes). But these are details to be improved in the future.

Garbage collection. Garbage collection is so easy, OMG. I keep track of all the pointers that are allocated during execution. Then it is just a matter of periodically walking through the stack and marking the allocations that are still reachable, then deleting anything in the heap that isn't. It's so easy that I didn't even implement it! I have 256 gigabytes of RAM. Even with a 70-billion parameter, 128-gigabyte LLM in RAM, there's still plenty of space to just keep allocating. In fact, LLM inference acts as a useful “performance regulator” to make sure that we don't allocate memory *too fast*.

Objects

As the SIGBOVIK deadline grew near, I reluctantly added “objects” to the BoVeX language. Objects are no stranger to ML; for example the O'CamL Language^[23] (pronounced “OK ML”) has them.^[24] But the community of functional programmers I was raised in has a revulsion to things Object Oriented, just like how a woodworker will immediately projectile vomit if they see a piece of Oriented Strand Board, even though it is a fine tool for many applications. I still have this disgust reflex. I imagine my Ph.D. advisors, should they read this, are contemplating whether and how a Ph.D. can be revoked. Anyway, I deliberately kept objects low-tech so that nothing could get too Oriented.

There is one object type `obj` in BoVeX. A value of this type has an arbitrary set of named fields whose types are known; they can only be the base types `int`, `float`, `string`, `bool`, `layout`, or `obj`. Fields are distinct if they have different types. An object can be introduced with an expression like `{() field1 = exp1, field2 = exp2}`, provided that each field's type can be synthesized from the expression itself (in the bidirectional type-checking sense). Alternatively, the program can declare an object name `0`:

```
object 0 of { field1 : type1, field2 : type2 }
```

and then use this in an expression like `{(0) field1 = exp1}`. These object names do not have any run-time meaning; they are just a collection of field types that are commonly used together. It gives a good place to document what they mean and some opportunity for better error messages, but fundamentally an object is just a collection of named data. Think like “JSON” object. It is possible to add and remove fields from objects (functionally) with expressions like `exp1 with (0)field2 = exp2`.

There are a few reasons for objects in BoVeX. One is the bibliography format, which consists of declarations like this

```
val knuth1981breaking =  
  bib-article {(Article)  
    title = "Breaking paragraphs into lines",  
    author = "Knuth, Donald E. and Plass, Michael F.",  
    journal = "Software: Practice and Experience",  
    page-start = 1119,  
    page-end = 1184,  
    year = 1981,  
    month = NOVEMBER,  
    publisher = "Wiley Online Library",  
  }
```

where each declares a reference made up of a bunch of optional fields. It is just too irritating to make each one explicitly optional, and since the data have heterogeneous types, manipulating some string-indexed data structure would have worse static checking and be more syntactically cumbersome. The bibliography rendering code case analyzes over the presence of fields to render citations that have different subsets of data.

Another use is in the `layout` type. This is a primitive type that most of a document's text is written in. It is a tree structure with optional attributes on each node, which are represented with an object. For example, this paragraph is written in the paper `.bovex` source file as:

```
Another use is in the [tt[layout]] type.  
This is a primitive type that most of a  
document's text is written in. It is a  
tree structure with optional attributes  
on each node, which are represented with  
an object. For example, this paragraph is  
written in the [tt[paper.bovex]] source  
file as:
```

The square brackets are used to write a layout literal (the main body of the document is inside one large literal). Layout literals can also embed expres-

sions (of type layout) with nested square brackets. Here the function `tt` is applied to a layout literal that contains text like `paper.bovex`. The `tt` function just adds the `font-family` attribute with value `"FixederSysLight"` to the layout node. This is a custom monospaced bitmap font that I made for this paper using software I wrote. It is part of the `FixederSys` family.^[25] Functions like `b` and `it` apply bold and italic text styles, but functions can do anything that you can do in a general-purpose programming language.

Primops

The other thing that objects are used for is interfacing with the runtime that is executing the BoVeX bytecode. There are about 50 different builtin primops that can be used by the BoVeX program. This includes simple things like integer and floating point addition, but also heavyweight operations like “load and register this collection of TrueType font files as a font family” or “invoke the boxes-and-glue packing algorithm with these parameters.” The primops in the former category work naturally on simple base types, but the heavyweight ones need to be able to pass complicated tree-structured heterogeneous data between the BoVeX bytecode executor and the runtime. It would be possible for the runtime to consume and create BoVeX values like tuples and lists, but this has two problems: One, many types like list are declared as user code (in the BoVeX standard library); they are not special, and we don’t want to make them special by informing the runtime of them. Two, requiring specific representations at the runtime boundary inhibits optimization; for example we can normally analyze the whole program to flatten data structures or remove record fields that are never used. The runtime typically uses `obj` to communicate structured data.

For example, the `internal-pack-boxes` primitive runs the boxes-and-glue algorithm. It takes some layout (which is expected to be a series of box nodes, with attributes giving their size, glue properties, and so on) and configuration parameters like the type of justification and algorithm to use. It returns an object with a new layout (the boxes grouped into lines, with new glued up widths) as well as the total badness. Inside the BoVeX layout support code, this primop is wrapped as `pack-boxes` with a native, typed interface, so programmers do not need to think about that implementation detail. Other typographic features that benefit from runtime support are implemented

this way as well.

Typographic features

BoVeX offers the `pack-boxes` algorithm, which can be used to nicely justify text. It can also be used to distribute paragraphs into columns, by thinking of the paragraphs as “words” (acceptable to break at any line, but bad to break near the start or end of a paragraph) and the columns as “lines.” It could be used by the document author for other purposes, I guess. There are other typographic features available.

Most of the layout of the document itself is by BoVeX code, which is either part of the standard library or part of your document, depending on how ambitious you feel. The function `main-text` parses the document layout into paragraphs and removes whitespace that is not really part of the text. It normalizes text properties across those paragraphs so that they can be manipulated individually. For each paragraph it uses the built-in `get-boxes` to break the words into fixed-size boxes with appropriate glue and hyphenation (see the next two sections), and then uses the `pack-boxes` routine to optimize their layout. The height of resulting lines are measured, and spaced according to the line spacing, then packed into columns. Once their final placement is known, boxes become stickers, which are sizeless elements that only know their position and contents. In this way, the BoVeX rendering pipeline is itself a bit like a compiler: It transforms programmer-written source layout into formatted paragraphs, then into boxes of known size, then into stickers of known position. At the end, it outputs the document as a PDF.

Any part of the rendering process can report “badness,” by calling the `emit-badness` primop. Nominally, badness is measured in square points of area that is outside of its container. Worse situations—such as text overlapping other text—have their badness scaled up per the same area of typographic horror. Less serious infractions—such as a little too much space between words—have badness scaled down. You have to use your heart to tell you what these scaling factors should be.

Fonts

BoVeX can render your document in plain Times Roman if you don’t care about anything, or access 13 other boring built-in PDF fonts, or it can load

any TrueType font from font files. (They do not need to be “installed,” and it won’t help to install them. You just put them in the directory with your document.) It loads their kerning tables and applies kerning properly, by generating rigid boxes at the sub-word level with unbreakable glue. I was disappointed to find that most fonts include only a few dozen kerning pairs. They do this in order to “save space” in the font file, which is utterly rich coming from someone that would try to save space *inside of words* by squeezing letters together! In the current font Palatino, the word “BoVeX” is not kerned correctly because the rare bigraph “oV” does not have a kerning pair. I hope to improve this detail in a future version (perhaps for the presumably forthcoming video version of this paper).

Hyphenation

Johannes Gutenberg invented the hyphen in A.D. 1455 for his Gutenberg Bible, then just known as Bible.^[26] His printing process actually required the lines to all be the same length, so he had to stick these little guys all over the place. His hyphens looked like this: ∷. Later on we straightened these out and decided we only needed one at a time, and today we use them not because we *require* our lines to all be the same length, but because we like the cognitive challenge of remembering the beginning of the word while we move our eyes to the beginning of the next line while reading.

BoVeX supports hyphenation using the same approach as TeX: We break each word into boxes at legal hyphenation points, and mark these points as sort-of-bad to break, and that if you do, you need to insert the hyphen character and use a little more space. By default in BoVeX, the hyphen sticks out of the end of the line a little bit. This is actually a bug but I like it.

I use the same hyphen dictionary as TeX, which is cleverly represented as a prioritized set of patterns in order to fit compactly in memory.^[27] Again, you have to respect Knuth and crew’s attention to detail, although to be fair this algorithm also dates to a time when storing a spell check dictionary in a computer’s memory was described as “not feasible.” So some of this was out of necessity. One of the nice things about the representation is that it generalizes to words that were not in the 1974 Merriam-Webster Pocket Dictionary. For example it hyphenates SIG-BOVIK correctly.

The details really keep going, too. The hyphenation dictionary is stored in a file called `hyph-en-us.tex`. “hyph” here of course stands for hyphens, and “en-us” means “English (United States).” In fact it is the *standard* language code for US English in the Small Language Model called IETF BCP 47.^[28] But then we have “hyph-en”, which is a plausible hyphenation of “hyphen”! You could even read it as “hyphen us, tex”, as a request for TeX to hyphenate the words in this file. This is the kind of detail I’m talking about! (There is also `hyph-uk`, which for once sounds a little less dignified than the US accent.)

Rephrasing

And of course, BoVeX includes a facility for using the LLM to rephrase text so that it renders more beautifully.

In contrast to the algorithm I described for monospaced text, it is not straightforward to know whether a prefix of some text will pack neatly with a proportional font. It depends on all sorts of contingencies, like kerning, whether we will split mid-word and hyphenate, or change fonts mid-sentence, or include an in-line image, and so on. Unlike monospaced text, a line of proportional text basically never fits *exactly* (badness 0); we need to apply some glue to make it fit, which generally has some small cost even when the text looks great.

One of the fiddliest parts of this is that we can’t just work with plain text, which is what the LLM enjoys best. Me too. This is because the paragraph being rephrased is some `layout` value, which contains some structure. Sending the original BoVeX code for the paragraph would maybe be possible in principle, although it would require very invasive changes to the compiler, and forbidden obscenities like “eval” to run the code it generated, and much better error recovery for the presumably vigorous stream of broken BoVeX code generated by the LLM. So I didn’t try that. Instead, I generate a textual representation for the paragraph to be rephrased, and feed that to the LLM. The prompt looks like this:

Exercise in rephrasing text. The following paragraph, which appears between `<P>` and `</P>` tags, needs to be rephrased so that it retains its precise meaning, but with minor variations in the specific choice of words, punctuation, and so on. No new facts should be introduced or removed, and all the ideas from the original paragraph should appear. However, it is good to use synonyms and change the word order and phrasing.

The text contains markup as well. There are two types: `text goes here` and ``. These should be preserved in the rephrased text. `` tags absolutely need to be retained and should not change their sources, although it is permissible to move them around in the text. `` should generally be retained, but the contents could change. The classes of spans may not change, and only the classes that appear in the original text may be used.

The first part is basically the same as what I used for the monospaced version, except that I ask the LLM to delimit the paragraph. This is important so that I know when it thinks it's done, and seems to work better than looking for newlines or the end-of-stream token. The second part is new. I translate the layout into plain text where uninterpreted subtrees are replaced with ``. These are generally boxes whose contents are not text. This could be an actual inline image or layout used to control rendering, like some bit of horizontal space. Nodes that are used to set text properties of the subtrees with attributes (like fonts, colors, ^{sizes}, etc.) are translated into distinct classes and marked up with `...`. The LLM has seen plenty of HTML, so it's able to use these reasonably well.

After generating a rephrasing, I parse the output HTML and match it up with the original layout. If I find any broken HTML, it is rejected. If I find any `` tag referencing a `src` not in the original, it is rejected. If I find any `` tag referencing a `class` not in the original, it is rejected. The more complexity that the original layout has, the higher the chance of a rejection, but rephrasing generally succeeds. But rejecting samples slows us down, so I leave off the second part of the prompt in the common case that the input paragraph is plain text. That way the LLM doesn't even try using markup.

With the HTML and original layout matched up, BoVeX can reconstitute the layout with the new rephrased text. This preserves any nested layout and attributes. It then continues with the render-

ing process.

But, how do we know whether we have a good rephrasing? When we run the boxes-and-glue algorithm, we get a "badness" score for the paragraph's line breaks, which tells us how bad the paragraph's line breaks are. When we run the rephrasing algorithm, the probability of the text we generated tells us how semantically good it is, and so we can call $1 - p$ the semantic loss. Combining those two somehow tells us how bad this is overall, and of course we want to find a rephrasing that minimizes the overall badness.

I wish that I could tell you that I solved this one with a beautiful algorithm! But so far I just have something reasonable that works. I generate many different rephrasings (with their semantic loss), and run each of them through the boxes-and-glue algorithm (to get the typographic badness). I choose the one that optimizes the preferred trade-off between semantic loss and typographic badness. This process is controlled by BoVeX code (i.e. it is in the source code of this very paper) and so it can be modified by the document author. Knuth has a very low tolerance for semantic loss, and knows that his algorithms produce good results without rephrasing. Lorem Epsom just wants it to look good and sound good. Both have published in SIGBOVIK 2024.

How to generate many different rephrasings? The simplest thing would be to sample randomly, like we did for the monospaced version. But since we prefer rephrasings that maximize probability, it is better to explore them systematically. Consider the model at the end of the prompt to be the root of an infinite tree. Each node in the tree represents an LLM state (sequence of previous tokens) and its children are the possible next tokens. Each of these tokens has a probability. All the model does is allow us to access that probability distribution for a node. Each possible rephrasing is a path in this tree that ends with `</P>`. We begin by sampling the most likely (as far as we know) path: At each node we see, we take the first (most probable) token. This is our first rephrasing, and it usually matches the original text exactly. Say that we "skipped" probability mass if we sampled a token that is less probable than it. We compute the semantic loss as the average probability mass skipped over all the tokens in the path. For this first path, we always took the most probable token, so this is 0.0 by definition.

The next path we explore will diverge from this path at some node (maybe the root). We pick a node that is likely to result in a good final loss, by scoring each node in the tree. The score is the average probability of all ancestor nodes times the probability of the next highest-probability token that we have not yet explored. The node with the highest overall score is the one we expand, by choosing that next highest-probability token. We are now in an unexplored part of the tree, and so we sample the most probable nodes repeatedly until we reach $\langle/P\rangle$. Speaking of which, BoVeX has a heck of a time trying to rephrase these last few paragraphs because they literally contain the text $\langle/P\rangle$ in them.

The scores should be seen as heuristic; we would get different results by choosing different ways of computing the score. This is an example of a “beam search” algorithm, which is good because it connects this project again to Super Metroid. As described in the earlier excerpt from the speedrun document that inspired this work, one of the final things you do in that game is acquire the “hyper beam” to defeat Mother Brain.

Since we will run the boxes and glue algorithm on multiple related texts, I generalized that algorithm to work on tree-structured input. This is clean; the memo table keeps the same dimensions, but records an additional fact. Now we store the penalty, whether to break after this token, and what the best subtree is. We have to consult each subtree when computing the score for a node, but this does not affect the asymptotic runtime. The table size is still at most $O(n^2)$, and although we explore more children per node, branches in the tree reduce the maximum depth to the root, which actually reduces one of the factors of n to $\log(n)$ as the tree becomes complete. However, as the SIGBOVIK deadline crept upon us, I never actually hooked this functionality up. It would require additional (programming) work to merge the trees, and the layout process is so fast that it doesn’t matter; I can easily run the full layout algorithm on hundreds of rephasings per paragraph.

I would like to improve the algorithm, because it does seem like there should be a way to integrate the boxes-and-glue dynamic programming algorithm with the path extension algorithm so that we prioritize exploring nodes that are likely to generate the best balance of typographic and semantic quality. It won’t be as satisfyingly optimal as boxes-and-glue itself because we have incomplete information (we never know whether one of the

exponentially many paths starts out with improbable tokens but then ends with a miracle streak of probable tokens). But it can certainly be more satisfying. Knuth would not stop here (but this is an *Any%* Knuth speedrun).

Instead I spent my time implementing an achievement system in BoVeX. The first time certain conditions are met, the system permanently awards you an achievement and prints a nice color trophy on your terminal. For example, you can get the “Not bad” achievement for generating a document that is at least 5 pages and has less than 1000 badness per page.

Advantages of rephrasing

Another nice thing is that the manual rephrasing that consumes valuable brain sugars when writing can become optional. For example, when I wrote the opening paragraph of this paper and listed a variety of trivial details, I might not need to think of different ways to say “unconcerned.” I could just write “unconcerned” each time and let the typographic considerations determine which synonym to use each time.

Conclusion

In this paper—and *with* this paper—I presented BoVeX, a new computer typesetting system. It follows the tradition TeX, but with modern amenities such as requiring over 128 gigabytes of RAM. Though some may consider the addition of AI features to TeX to be an unnecessary perversion, I find this use of LLMs to be *fully justified*.

Future work

Typographic features. Many more typographic features are desirable. Footnotes! It is so hard to write a paper without footnotes. Where am I supposed to put the bonus digressions? The layout of footnotes is tricky and should be part of a general floating figure implementation. End notes are actually easy, but I don’t want end notes. I want them to be little footnotes so that you can’t help but read them.

BoVeX does not support page numbers, which is good because they are forbidden by the SIGBOVIK program committee.

TeX is famous for its mathematical typesetting as well. It would fit neatly into BoVeX in the

same way, since both use the same fundamental boxes-and-glue engine. BoVeX does not have “macros” or “modes” like TeX, but it would work cleanly to write a BoVeX function `math` (or, if you like, `$`) that parses a custom syntax. In fact it would be natural to have different parsers for different maths, so that you don’t need to parse `->` as *minus greater than* in mathematical contexts that don’t use minus or greater than at all.

Optimization. There are many opportunities to make BoVeX code faster. This is mostly important for when it is being run in a loop in order to try out many different rephrased texts. (That said, I do not wish to preclude what could be done with BoVeX by assuming its execution is doing only typesetting tasks. For example, shouldn’t you be able to challenge your paper’s reviewers to a game of chess against a strong engine embedded within your document?) The first thing to fix is that it manipulates too many strings at runtime (e.g. the code, record labels, object fields, and “registers”). This is easy to fix since these are all known at compile time. There are lots of high-level optimizations left to do for the IL code (common subexpression elimination, constant argument removal, uncurrying, etc.) and lots of peephole and control-flow optimizations left to do for the bytecode (currently no optimizations are performed at all). All of this becomes more important if I add another planned feature, which is the ability for the document to be globally optimized by applying a black-box optimizer to a set of user-specified parameters. For example, the column width, line spacing, or font size could be tweaked to make the document fit better. This feature is “Auto-Margin Plus.” Things are already set up to do this pretty straightforwardly; we would simply generate the document over and over while searching over the parameter space, and choose the one with the least badness. This may also affect which rephrasings look best. But instead I spent my precious time implementing **3D text**.^[29]

Reproducibility. The algorithm for rephrasing text tries to find the best place to explore the next most likely token from the probability distribution. This expects the generation of these distributions to be deterministic. Mathematically, inference is deterministic (it is just a bunch of matrix multiplications), so this “should work.” But in practice the enormous calculation is performed in an unpredictable order as it is executed in parallel (in multiple CPU and GPU cores). Because floating point arithmetic is not associative (or distributive, com-

mutative, or other properties you’d like), inference can sometimes generate different answers due to floating point round-off error.^[30] Alas, these are not even necessarily related to the final probabilities in the model, as billions of non-linear operations happen within the hidden layers of the network. The effect is not particularly grave; we might miss out on a highly likely path because the probability distribution was different the second time we looked at it. There are already lots of ways we might fail to find highly likely paths, so this is not some kind of reproducibility crisis. It is mostly just a bit unsatisfying.

Unicode support. This would have been helpful when above I decided to show you Gutenberg’s funny hyphen, \neq , for which I had to settle for embedding a crappy hand-drawn PNG file. Instead I could have used `U+2E17`, which since this exotic codepoint it is not present in the font Palatino, you could have experienced as \mathbb{N} . BoVeX is written with some Unicode support, with the main exception being that the PDF output code only supports the embarrassingly diminutive WinAnsiEncoding.^[31]

Deadlines. Although BoVeX itself is very fast, rephrasing is very slow. This presents a problem for the typical way that academic papers are written, which is to do all the work in a coffee-fueled fugue in the last few days before the deadline, then stay up all night writing the paper and finding citations for the pro-forma “related work” section which you did last but you know that the reviewers will insist upon, and tweaking `\uspace` and `\begin{figure}[h!]` until it fits within the page limit. On the one hand, BoVeX does potentially free the author from the visual tweaking process. But on the other hand, the LLM inference for the rephrasing process can be quite slow, and it can take many hours or days to fully bake a long paper! For this reason, it may be better to change conference deadlines to a system where the pre-rephrasing text is submitted. The publishers (what do they even do?) can be the ones to execute the rephrasing in the cloud as they produce the “camera-ready copy.” With straightforward extensions, this would also allow the rephrasing to adapt to changes in the overall volume style, or to adjust to avoid embarrassing typographic coincidences with other articles in the same volume (such as using the same notation with a different meaning). In principle, the paper could edit itself to respond to feedback from reviewers, in a way that minimizes the semantic distance from the original. This rapid feedback loop could reduce the time to publication, perhaps to

- [18] Tom Murphy VII. "Modal Types for Mobile Code". January 2008.
- [13] Tom Murphy VII. "NaN gates and flip FLOPS". *SIGBOVIK*. April 2019.
- [10] Tom Murphy VII. "The First Level of Super Mario Bros. is Easy with Lexicographic Orderings and Time Travel. After that it gets a little tricky". *SIGBOVIK*. April 2013. pp. 112–133.
- [20] Tom Murphy VII. "The Wizard of TILT: Efficient(?), Convenient and Abstract Type Representations". *Carnegie Mellon tech report CMU-CS-02-120*. March 2002.
- [29] Tom Murphy VII. "The glEnd() of Zelda". *SIGBOVIK*. April 2016. pp. 105–112.
- [14] Tom Murphy VII. "ZM~~ # PPrinty# C with ABC!". *SIGBOVIK*. April 2017. pp. 129–148.
- [25] <http://tom7.org/fixedersys/>. Tom Murphy VII. "The FixederSys font family". 2024.
- [2] https://en.wikipedia.org/wiki/Wikipedia:WikiProject_Punctuation. Tom Murphy VII. "WikiProject Punctuation". July 2007.
- [19] <https://sourceforge.net/p/tom7misc/svn/HEAD/tree/trunk/rephrase/>. Tom Murphy VII. "BoVeX source code". 2024.
- [22] Chris Okasaki. "Even higher-order functions for parsing or Why would anyone ever want to use a sixth-order function?". *Journal of Functional Programming*, 8(2). March 1998. pp. 195–199.
- [28] A Phillips, M Davis. "Tags for identifying languages". September 2009.
- [9] <https://gamefaqs.gamespot.com/snes/588741-super-metroid/faqs/10114.rs1n>. "Super Metroid Speed Guide and FAQ". 1996.
- [24] Didier Rémy, Jérôme Vouillon. "Objective ML: An effective object-oriented extension to ML". In *Theory And Practice of Objects Systems*, 4(1). 1998. pp. 27–50.
- [11] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, Thomas Scialom. "Llama 2: Open foundation and fine-tuned chat models". *ArXiv.org*. July 2023.

Flaccid Drives

Storing Arbitrary Data in SIGBOVIK Articles

Nora Hyades Mirza, Head Archivist
 Dissociation for Heresiographical Informatics

Abstract—

Estragon: hey, didi, babe, could you write the abstract?
Vladimir: wtf gogo why cant you do it
Estragon: babe im busy
Vladimir: with what
Estragon: ur mom
Vladimir: lame. she was off playing sea of thieves anyway
Estragon: who's rosenkrantz
Vladimir: sea of thieves nuts lmao

I. INTRODUCTION

Estragon: fuck you this is why nora stuck us with this bullshit

Vladimir: don't fucking blame this on me you're the one what covered all their bras in chocolate sauce!

Estragon: I WAS HUNGRY

Vladimir: FOR BRASSIERES, GOGO?

Estragon: yes >:(

Vladimir: how tf did you do that with your mouth

Estragon: do what :?

Vladimir: that. do That

Estragon: lol. lmao. lookathis galoomba. can't even :3 out loud

Vladimir: stop ittt I told you don't tease me like that[2]

Estragon: fine!!!! anyway we should prolly get to writing that abstract

Vladimir: yeah ig so. so what did she say this was again?

Estragon: a "tribute" to a "seminal"[1] "paper"

II. METHODOLOGY

In this paper we

Estragon: wrong section

Exercise 1. fuck you

Estragon: do it you won't

Vladimir: [2]

III. ENCODING

official dhi internal communique donut steel

alr yinzers we gotota do this fucking shit the deadline's in like 1 day

harder drives: store data in worse ways

flaccid drives: store data in better ways

wjhat's the best data medium? sigbovik (sigbovik is eternal)

PROBLEMS: * bitrot * printing errors * public (anywone can see your data)

SOLUTIONS: * checksums * diceware word lists? * encryption

NOT NOVEN ENOUGH! what if we like what if like we like we use conlang words? use x-sampa one symmable per byte, with high nipple being voiced and low being unvoiced

checksum vowel between nibble consonants, word per word

Estragon: dude you can't just fucking copy the research googledoc

Vladimir: "yes i can" - dave n bustr

Estragon: they're gonna fucking arrest us for plaigiarism you lunthead

Vladimir: lunk my nunks

Estragon: I thought you didn't wanna be "teased like that"

Vladimir: yeah well *maybe* that's just because you *turned me down*, gogo

IV. PHONEMES

Estragon: well I'm fucking sorry that I didn't wanna fuck my *sister*!

nipple	low consonant	high consonant
x0	?\ 	X\
x1	? 	?
x2	S 	Z
x3	t_?\ 	d_?\
x4	tS 	dZ
x5	T_?\ 	D_?\
x6	s_?\ 	z_?\
x7	s 	z
x8	k 	g
x9	q_?\ 	G_?\
xA	x_?\ 	G_?\
xB	\ \ 	\ \
xC	r_?\ 	r_?\
xD	l 	l
xE	x 	G
xF	f_?\ 	v_?\

Vladimir: didi?

high parity	low parity	vowel
even	even	A
even	odd	M
odd	even	@
odd	odd	i

Vladimir: I'm sorry, I didn't-

Estragon: it's just, I thought, I thought we could move past it, y'know?

Vladimir: yeah,

Vladimir: I...I have a secret I need to tell you.

V. ENCRYPTION

Vladimir: we're not actually sisters

Estragon: what a great way to come out

Vladimir: what? no, I mean

Vladimir: I'm adopted

VI. COMPRESSION RATE

Vladimir: I know that doesn't really change things but, I guess that's why I thought it was okay

Vladimir: I'm sorry for upsetting you

Vladimir: I won't bring it up again

VII. LONGEVITY

Estragon: didi, I don't know what to say

Estragon: I'm sorry, I didn't know

Estragon: you mean a lot to me, you know that, right?

Vladimir: yeah, I know. you do to me, too

Estragon: I'm glad you're in my life, sis

Vladimir: same here :)

Estragon: wanna go make out sloppy style behind the autoclave?

Vladimir: I would love that

VIII. CRITICAL CONSIDERATIONS AS TO THE USE OF X-SAMPA

IX. CONCLUSION

REFERENCES

- [1] Dr. Tom Murphy VII Ph.D., "Harder Drive: Hard drives we didn't want or need." 7-Apr-2022. From SIGBOVIK 2022.
- [2] "The BDSM Test is a fun and educational test to determine what kind of kinkster you are. It was founded in 2014 with as its main mission to make a simple, accessible test to help beginning kinksters determine which labels are or aren't suitable for them; and to be a fun experience for everyone taking it, beginners and experts alike." <https://www.bdsmtest.org/>

Z-i-|\\-G_-?-M-? Z-i-s-?-i-x G_-?-A-s_-?-G_-?-A-
r_-?\\|\\-@-x_-?-d_-?-M-|\\g-@-f_-?-r_-?-A-f_-?\\|\\-@-
s_-?-d_-?-A-t_-?\\G-@-?-g-@-f_-?z_-?-A-t_-?-g-i-? v_-?\\-
M-l-z_-?-A-?\\G_-?-A-q_-?-X\\-A-s_-?\\G_-?-M-x-v_-?-A-
f_-?\\X\\-A-?-?-@-q_-?D_-?-M-|\\-r_-?-A-x_-?\\g-i-k-v_-?\\-
A-s_-?\\?-@-x_-?-l-@-r_-?z-i-k-dZ-i-|\\Z-@-x_-?-G_-?\\-
M-x v_-?-A-f_-?-X\\-A-?\\?-@-q_-?-G-@-T_-?\\G_-?-A-
t_-?-r_-?-A-s_-?\\Z-@-q_-?-G-@-?\\l-@-s_-?-d_-?-A-s_-?z-
@-t_-?-G_-?\\-A-T_-?\\Z-i-l-G-@-t_-?\\Z-@-r_-?-G-@-t_-?\\-
i-k-G_-?-A-T_-?\\X\\-A-x_-?-d_-?-A-t_-?D_-?-A-?\\-g-i-x z-
@-t_-?-dZ-@-t_-?z_-?-M-tS-Z-@-T_-?z-@-f_-?-D_-?-M-?
D_-?-A-s_-?\\-|\\-@-x_-?z_-?-A-q_-?-v_-?-M-s G_-?-M-
l-G-@-r_-?\\G_-?-M-tS-G_-?-A-f_-?dZ-@-f_-?-D_-?-M-S
|\\-i-l-?-i-s G_-?-M-S-?-@-f_-?dZ-@-f_-?-D_-?-M-S
?-@-q_-?-X\\-A-T_-?dZ-i-x-X\\-M-tS z-@-T_-?\\-G_-?-A-
f_-?\\|\\-@-?-v_-?-M-s ?-@-T_-?\\-G-@-q_-?\\G_-?-A-
s_-?\\-X\\-M-k Z-@-?-g-i-S X\\-M-k-Z-@-t_-?v_-?-M-k-z-
@-f_-?v_-?-A-f_-?-l-@-q_-?

REFERENCES

[1] Vladimir and Estragon, "Flaccid Drives: Storing Arbitrary Data in SIG-BOVIK Articles", 1-Apr-2024. From SIGBOVIK 2024.

FORM A MORE PERFECT UNION (between machine and man)

- 23 Solving C's biggest flaw - Hemispheric divergence**
Ian F.V.G. Hunter
- 24 An empirical performance evaluation between Python and Scratch**
Morgan Nordberg
- 25 Exchange Traded Neural Network**
Echometer Rain
- 26 Minmaxing the energy efficiency of biological computing**
Hugo Waltsburger
- 27 Advancing Consensus: Automated Persuasion Networks for
Public Belief Enhancement**
osmarks.net Computational Memetics Division
- 28 Grounding Language Models to their Physical Presence**
Seongmin Park
- 29 Quantum Disadvantage: Simulating IBM's 'Quantum Utility'
Experiment with a Commodore 64**
Anonymous

Solving C's biggest flaw — Hemispheric divergence

Ian F.V.G. Hunter

Unaffiliated/Excommunicated, near-antipode of Campbell Island, New Zealand

Abstract

We propose a new domain-specific language (DSL) called “Code for Hemisphere-Unconstrained Master Programs” (CHUMP), which frees software engineers from the traditional constraints of bi-hemispherical coding. We provide several example programs covering core scenarios and discuss future work in this area.

1 Introduction

Traditionally, developing software for both hemispheres has been a difficult endeavor. Usually, developers have to write two separate applications for both hemispheres. This is an obviously inefficient approach with multiple problems:

- The overall size of the codebase is doubled.
- Every single line changes when checking code into source control.
- Double the programmers need to be hired (i.e., one team for C and another for ☺).)

In code listings 1 and 2 below, we show a simple test case and the traditional two-system development.

Listing 1: Northern Hemisphere example

```
#include <stdio.h>
int main(){
    printf("Hello World!\n");
    return 0;
}
```

```
#include <stdio.h>
int main(){
    printf("Hello World!\n");
    return 0;
}
```

Listing 2: Southern Hemisphere example

These are compiled with GCC and Clang respectively and two executable binaries are produced for distribution.

2 Proposed Method

We developed a novel C-like Domain Specific Language (DSL) which allows us to contain development for each hemisphere inside a shared file. After a developer writes their program in this DSL, they can target a specific latitude co-ordinate where the program is intended to run.

2.1 Prior Work

Thankfully, most code editors these days will preserve both line endings and character orientation. (i.e., 'V' will be saved as 'A' and vice versa depending on the original orientation of the character). This is a huge step forward as keyboards differ between the hemispheres (See Figures 1 and 2) and even with the work of CHUMP, the overhead of continually switching layouts is impossibly confusing to the average software engineer. Re-ordering of lines and letter order is unfortunately still a NP-impossible research problem, yet to work consistently on all but the most trivial examples.

~	!	@	#	\$	%	^	&	*	()	-	+	←	BACK SPACE
↵	Q	W	E	R	T	Y	U	I	O	P	{	}		↵
CAPS LOCK	A	S	D	F	G	H	J	K	L	:	"	'	↵	ENTER
⇧ SHIFT	Z	X	C	V	B	N	M	<	>	?	/	.	⇧ SHIFT	
CTRL		ALT								ALT				CTRL

Figure 1: Northern Hemisphere Keyboard

~	!	@	#	\$	%	^	&	*	()	-	+	←	BACK SPACE
↵	Q	W	E	R	T	Y	U	I	O	P	{	}		↵
CAPS LOCK	A	S	D	F	G	H	J	K	L	:	"	'	↵	ENTER
⇧ SHIFT	Z	X	C	V	B	N	M	<	>	?	/	.	⇧ SHIFT	
CTRL		ALT								ALT				CTRL

Figure 2: Southern Hemisphere Keyboard

2.2 GCC Preprocessor and comment style

In order to allow both code versions to be written in the one file, there must be some check for certain code to be compiled only when in the target hemisphere. Hence, we place some traditional GCC preprocessing macros like '#if' and '#else'.

However, these are hemisphere-dependent and will fail under CC . These macros be appended with a comment symbol in order for them not be processed. Similar macros for CC are written and corresponding comment symbols so not be processed by GCC.

CHUMP changes the comment style change from '/' as '\\" will not correctly parse when read backwards¹. In order not to favor one hemisphere over the other, we use the obvious compromise by interpolating the comment styles — '|'|.

¹This is an unfortunate side-effect of text editors saving character orientation!

Listing 3: CombinedHemisphere

```

||                               fidne#
||                               ||
||           CINAP                ||
#if LATITUDE > 0                ||
||                               esle#
#include <stdio.h>                ||
int main(){                       ||
    printf("Hello World!\n");    ||
    return 0;                    ||
}                                  ||
#elif LATITUDE < 0              ||
    0 > EDUTITAL file#          ||
    {                             ||
        ; 0 nruter                ||
        ;("n\dlroW olleH")ftnirp ||
        }()niam tni              ||
    <h.oidts> edulcni#          ||
    0 < EDUTITAL fi#            ||
#else                             ||
    PANIC                         ||
#endif                             ||

```

Of course, || is usually used for a logical OR in C. Another symbol used for this operation in math is '∨'. But, to avoid confusion with '∧' (logical AND), it was decided to use ∠.

2.3 Shortcomings

It is an unfortunate truth that with every great change, someone will complain. Here are a few messages our first release has received, and our comments:

- "Why not use if else?"
 - No.
- "Why the name CHUMP and not a derivative name of C, like C++ or C#?"
 - If we wanted to base the name off of a hemisphere-independent name for C, it would look like C which has a high risk of misrepresenting the project as an offshoot of the C programming language[1]. We did mull over being called C∠, but the pronunciation was an issue.
- "Can you shout me out on your paper?"
 - Sure. Hi Mom.

3 Conclusions & Future Work

We are actively working on a second major version of CHUMP. Please look forward to the following features:

- Support for while-do and while-do loops
- Deprecate “elif” in favour of the more verbose “else if” to avoid conflict with the ‘file’ system
- “static” to be made default so that code does not move when emailed from one hemisphere to the other.

There are two outstanding limitations of CHUMP, one of which is that CHUMP itself is written in Northern-Hemispheric C and not in CHUMP itself. This has created an undesirable power-dynamic in companies that use CHUMP which we will seek to remove in future versions.

The other issue, which astute readers may have noticed above is the ill-handling of cases at latitude 0. Our working solution is to have several slingshots at the point of zero latitude and launch target platforms across to the other side at such a speed in which no defects will have time to occur.

At some point we hope to have the research funding to expand our work to U and G in order to support users permanently residing at latitude 0 who rely on longitude systems. However, we know there remains an issue for tourists to Station 13010 of Null Island [2] at 0, 0.

References

- [1] The O language¶ (no date) The O Language - The O Language 2.0 documentation. Available at: <https://o.readthedocs.io/en/latest/> (Accessed: 29 January 2024).
- [2] St. Onge, T. (2016) The geographical oddity of Null Island, The Geographical Oddity of Null Island | Worlds Revealed: Geography & Maps at The Library Of Congress. Available at: <https://web.archive.org/web/20160512183906/http://blogs.loc.gov/maps/2016/04/the-geographical-oddity-of-null-island/> (Accessed: 29 January 2024).

An empirical performance evaluation between Python and Scratch

Morgan Norderg, Linköping University

1 Abstract

In this empirical study the authors examine and evaluate the performance of Scratch and Python for the sum-for function, and find that there is a good indication that Python has better performance characteristics compared to Scratch regarding iterations. The authors conclude that Scratch appears to be an unattractive choice for high performance programming, however they state that further research is needed.

2 Introduction

Python is a widely adopted high level multiparadigm programming language in today's tech industry. It is however well known to be lacking in performance, and thus there is a natural interest in potentially high performing languages that could serve as fitting replacements. Scratch offers many of the same benefits as Python, such as a easy learning curve, interpreted nature allowing for iterative development, and wide adoption in programming education. So investigating its performance is imperative for the advancement of the field of computer science and the tech industry as a whole.

3 Background

Scratch was initially developed at Massachusetts institute of technology(MIT) to entertain children. It offers a wide range of features such as an interpreted development environment with colorful blocks for visual programming. Its performance characteristics have been unknown and uncharted territory, until today!

Python is a interpreted, dynamically typed, programming language introduced in 1994. It is today deployed in all areas of the tech industry, from data analysis and AI, to crypto mining and OS development(probably).

Iteration with for-loops in Python is notoriously slow, so a major point of this paper is to investigate if Scratch offers a compelling alternative for these use-cases. Since Scratch is interpreted in Javascript(Random friend of the author, 2024), which is generally faster than Python, it is hypothesized by the author that Scratch might outperform Python in this area.

The sum-for function is a function that sums the integer values of 0 to n .

4 Related work

There appears to be no related work on the subject. The authors didn't look for any, and found nothing.

```
import timeit

def sumfor(x):
    sum = 0
    for num in range(x):
        sum += num
    return sum

def sumfor_test():
    print(timeit.timeit(lambda : sumfor(100000), number=1))
    print(timeit.timeit(lambda : sumfor(1000000), number=1))
    print(timeit.timeit(lambda : sumfor(10000000), number=1))
```

Figure 1: Python source code used for collecting the measurements

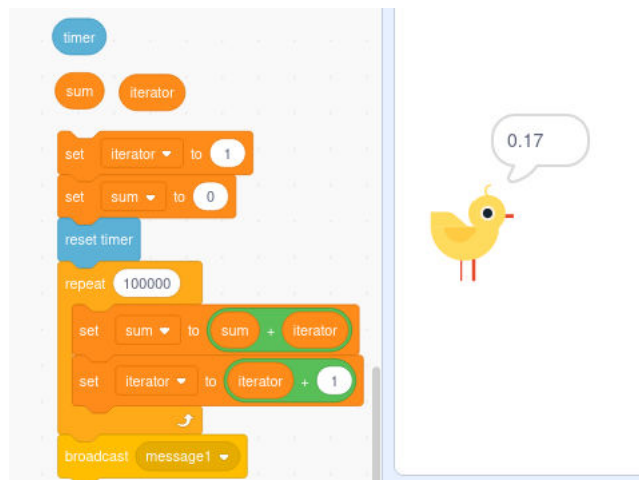


Figure 2: Scratch source code used for collecting the measurements

5 Methodology

For comparing the iteration performance between Scratch and Python, we compared the time in seconds for running a simple sum-for function, for the input values of 100-000, 1-000-000 and 10-000-000. To collect the measurements for the Python code the authors ran the Python script $n=1$ times and reported the average runtime in seconds. $n=1$ was chosen because that's how many times the author could be bothered to run the Scratch code. For outputting the results the Python script used the Print function of the corresponding standard library, meanwhile the Scratch code used a cartoon chicken. The cartoon chicken is necessary since the say-block in Scratch is only available for sprites in Scratch.

The broadcast block shown in figure 2 sends a message

to a receive broadcast block in the chicken sprite logic, not shown in this study.

6 Results

The results showcase an massive performance decrease with Scratch compared to Python. In table below we can see the time in seconds for performing n iterations of the sum-for function in Python and Scratch respectively.

Python	Scratch	n
0.003	0.200	100-000
0.032	1.96	1-000-000
0.318	15.91	10-000-000

7 Discussion

It's rather dissapointing to the author that Python outperformed Scratch to such a significant extent. This runs contrary to the authors previously stated hypothesis. There are however some potential factors that could have influenced the results.

7.1 Scratch optimizations

It is possible that the overhead in the Scratch code could be reduced by moving the sum-for function code block directly into the logic of the cartoon chicken sprite, removing the need to broadcast a message to the chicken sprite to output the result. Another possible way to reduce this overhead the authors hypothesize, is to store the value of the timer variable in a separate immutable value, and say that variable value instead in the chicken sprite. The overhead then would simply be the cost of an assignment operation which should be insignificant.

7.2 Unknown factors

It is not exactly known to the authors how the timeit function in Python works under the hood, and there might be some implementation details that create less overhead compared to the timer in Scratch. The authors did not investigate this due to a severe lack of f***s to give.

8 Further work

Due to the stated importance of the subject, the lack of pre-existing scientific literature regarding the subject, and the previously stated potential for Scratch optimizations there is a clear need for further research.

8.1 Scratch optimizations

Besides replicating this research, alongside the previously stated optimization gains to further solidify our understanding of Scratch performance characteristics; there is a potential area of research to explore regarding Scratch to Scratch compiler technology, or perhaps transpiling Scratch code to highly performant C code. Or highly performant Rust code, since the Biden administration outlawed the C programming language in 2024 (*“Press Release: Future Software Should Be Memory Safe”*, [whitehouse.gov, 2024-02-26](https://www.whitehouse.gov/oncd/briefing-room/2024/02/26/press-release-technical-report/), <https://www.whitehouse.gov/oncd/briefing-room/2024/02/26/press-release-technical-report/>).

There is a notable lack in research regarding the development of compilation, transpilation and genetic mutation of Scratch code, and other forms of build tools. This could potentially revolutionize the world of high

performance Scratch programming. There is also a possibility of looking into hardware accelerated Scratch programming using hardware description languages and FPGA's or ASICs.

8.2 Use cases

Since this study only considered the case of the sum-for function, there could be interest in looking at other potential use cases of the Python programming language where Scratch might possibly outperform it. We are desperately in need of further research in the area of Scratch driven OS-kernal development.

9 Conclusions

You probably don't wanna use Scratch for performance sensitive programming tasks.

References

<https://www.whitehouse.gov/oncd/briefing-room/2024/02/26/press-release-technical-report/>

Exchange-Traded Neural Network

Echometer Rain

03/15/2024

Abstract

We propose a novel financial instrument that replicates the structure of a deep neural network by using exchange-traded funds (ETFs) as perceptrons and options contracts as activation functions. Although legislatively improbable, we predict that such an architecture will be capable of learning time-independent patterns within the stock market (and make us bloody rich) if implemented.

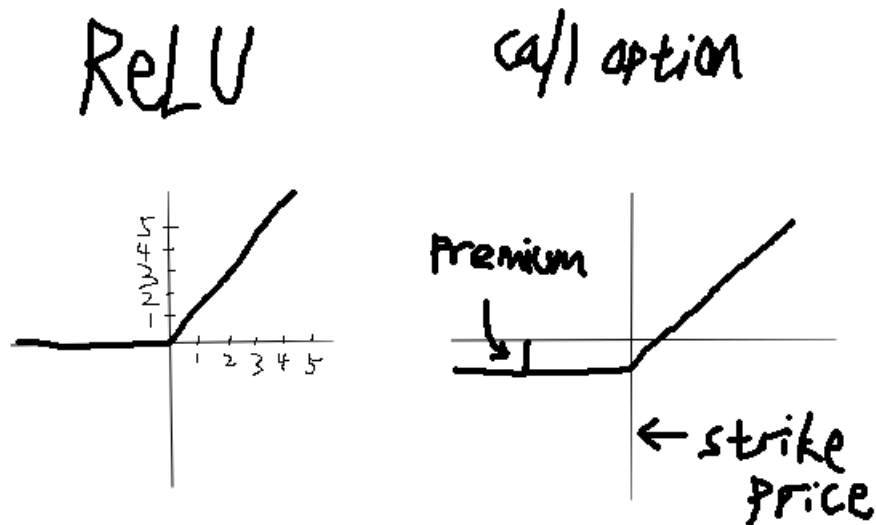
1 Introduction

Although time-series stock forecasting remains a popular subfield within the study of artificial intelligence, none yet have been unhinged enough to construct a physical neural network using existing financial instruments. This paper details how such technology may be created.

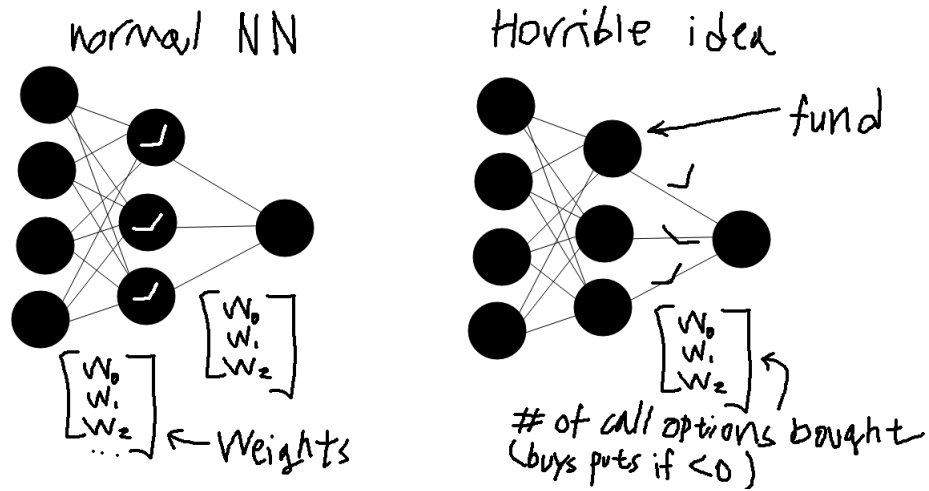
2 Architecture

Neural networks are built by passing numerical inputs into a layer of nodes (perceptrons). Each node multiplies its inputs by its corresponding weights and passes the sum of products (typically implemented using a dot product) into an activation function. This process is repeated by passing the outputs of each layer as inputs to the next until a final layer generates an interpretable result. By the Universal Approximation Theorem (UAT), this structure is capable of approximating any mathematical function to arbitrary precision.

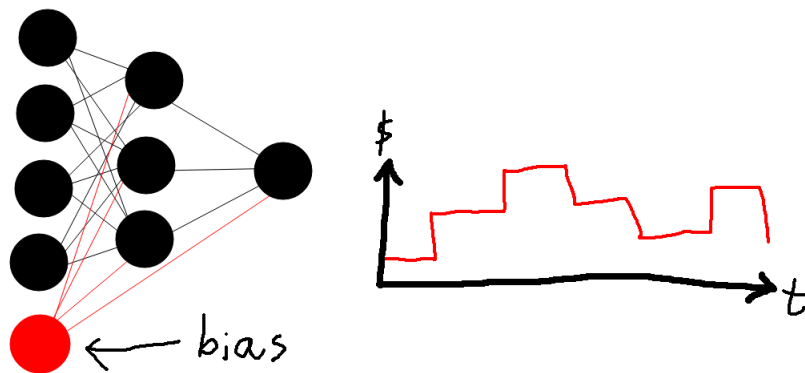
To physically reproduce this on the stock market, we notice that the profit function of an options contract is simply a transformed Rectified Linear Unit (ReLU), a widely used activation function. Through this, an Exchange-Traded Fund (ETF) can effectively replicate the workings of a perceptron by buying a certain number of options corresponding to the weight of each input and returning their earnings to the next layer of ETFs.



Let's walk through the transactions to see how such a network would operate: First, the weights of each node are randomly assigned. The input layer of ETFs then buys the amount of options contracts specified by their weights (rounded to integers). The ETF will buy call options if the sign of the weight is positive and it will buy put options if it is negative. After all options are bought, the ETF will start selling options of itself to successive layers. These ETF options will expire when all of the contracts bought by the input layer expire. The process of buying and selling options is repeated until the final layer, the profits of which should be representative of the profits of the entire network.



We have not yet implemented the bias term needed to satisfy the Universal Approximation Theorem as treating this will require special considerations. A new bias ETF will be created which buys and sells extremely low volatility stocks according to a random table. Every time the network is activated, the price of the ETF will change by a specific amount. The nodes of the network which depend on the bias will always know whether its price will rise or drop, and they can use this to generate a predictable stream of income. Because no sane investor would want to bet their money on a stepwise random walk (not all investors are sane), the operators of the network would need to sell options to it at a loss whenever it activates. To prevent the network from reward hacking the operators' money away, all profits gained via the bias ETF will be deducted within the loss function.



This brings us to the backpropagation step. To tune the weights of a traditional neural network, a differentiable loss function in order required to measure the amount that each weight of each node contributes to the quality of the final result. In our case, the loss function will just be the net profit of the network, which will equal the net profit of the output layer minus the sum of all the bias weights and the sum of all transaction fees. At the end of each epoch (the time of which depends on the length of the longest contract), the weights can be adjusted to maximize profit, and the network shall run again.

3 Conclusion

Not only is this architecture incredibly expensive to train, the Security and Exchange Commission (SEC) will never allow some startup to create ~7000 ETFs out of thin air (without bribes). Because we may not be able to control whether or not an LLM will overfit past stock market data, we predict that our architecture will surpass the investment capabilities of near-future GPTs. We hereby request eight (8) trillion dollar.

This work is licensed under a Creative Commons
“Attribution 4.0 International” license.



Minmaxing the energy efficiency of biological computing

A study on the energy efficiency of Ph.D. students for neural network inference and other computing shenanigans

Hugo Waltsburger^{†‡}

[†]SONDRA, CentraleSupélec, Université Paris Saclay, Gif-sur-Yvette
Hugo.Waltsburger@centralesupelec.fr

[‡]Université Paris-Saclay, CentraleSupélec, CNRS, Laboratoire de Génie Electrique et Electronique de Paris, 91192, Gif-sur-Yvette, France.

Sorbonne Université, CNRS, Laboratoire de Génie Electrique et Electronique de Paris, 75252, Paris, France
Hugo.Waltsburger@geeps.centralesupelec.fr

Yes Anthony, the Ph.D. manuscript is coming along.
No, I am not getting sidetracked haha, I would never.
No, you need not prepare the crucifixion nails.

-The author to one of his (much beloved) Ph.D. supervisors, probably

Abstract—Who are we? What is our purpose in life? This paper will not answer these questions, but will nonetheless hint at the possibility that, if humans’ evolutionary purpose is neural network inference, then we are awfully badly designed. We will also answer some fundamental philosophical questions, such as "How many human computers do we need to play Doom?".

I. INTRODUCTION

A. Context

The artificial neurons that make up modern neural networks have barely changed ever since their initial appearance in 1943 [1]. They went through two successive cycles of hype and disappointment (the so-called *AI Winters*), before beginning their third hype cycle in the early 2010s decade, a hype that lasts to this day. Will it last? It’s hard to make predictions at this point. Large language models (LLMs) have had some spectacular results. But with modern LLMs having been trained on the almost complete dataset of data available on the Internet using tens of thousands of GPUs for months, the path of least resilience that was "throw more layers, data, and GPUs at the problem until there’s no more problem" is bound to break down sooner or later.

With this matter comes an efficiency problem. The energy footprint of high-performance computing is significant. The creators of the open-source LLaMa II 70B LLM have estimated its training-related greenhouse-gas (GHG) emissions to 500t of carbon dioxide equivalent (CO₂e) [2]. This amounts to 100 times the target yearly emissions of an average European

citizen in 2030, as per the Paris Agreement. This figure does not count the energy consumption of inference. Some modern state of the art neural networks are so expensive to train that it is doubtful that anyone will try to retrain them from scratch before they are obsolete - which at the current speed, may be a matter of months, a couple of years at most. Modern neural networks become more susceptible of consuming more energy during inference than during training. There is a growing importance of measuring the environmental impact of artificial intelligence, to create more energy-efficient networks. A first (modest) attempt was made by yours truly in [3].

B. Getting efficiency out of the way

Now that the mandatory serious and forward-thinking part of this paper is out of the way, we can think backward: what would be the most efficient inefficient way of performing the tasks that neural networks do? Our initial idea started with the procurement of a monkey, a typewriter, and a couple trillion years (along with lots of ink ribbons and a supertanker worth of paper). Unfortunately, monkeys are not considered pets in France (thus demanding special authorizations, which I don’t have, to take care of one), and typewriters have become vintage - thus expensive. So I settled on the thing closest to a monkey with a typewriter that was readily and cheaply available: a Ph.D. student (myself) with a computer¹. In this paper, we will explore how a human subject can optimize its neural network inference performance.

¹The quality of the resulting thesis should be about equivalent

II. HISTORICAL PERSPECTIVE

In the early days of computing, a lot of calculations were performed by hand, filling notebooks worth of computations. This led to the creation of a profession, "Computer", whose job would be to carry out much of the lengthy calculations that mathematicians did not have the time or want to perform themselves. The First World War, a period fraught with ballistics calculations, led to the employment of many computers. When mechanical and electronic computers appeared, these human computers were also tasked with programming them - a tedious task involving wires and punched cards. Oftentimes, these venues would figure among the few opportunities for college-educated women or minorities to enter an "academic" career. The early days of computer science saw many such people, whose opportunity came from the fact that computer science was generally considered a "less noble" field than theoretical science. The reader interested in a deeper exploration of human computers may refer to [4] or [5]. In the brief blip of time we are currently experiencing where computers are silicon-based rather than organic (between these past-but-not-so-remote times, and the mentats of the post-Butlerian-Jihad era - it's coming and you know it), it may be hard to evaluate how easy we have it with calculations. This paper will remind the modern reader of this fact, and hopefully cement my position as an early prophet when the Butlerian Jihad comes and AI gets relabeled as either "Abominable Intelligence", or SALAMI (Systematic Approaches to Learning Algorithms and Machine Inferences).

III. MOTIVATION

At this point, I would honestly do just about anything to avoid writing my Ph.D. manuscript².

IV. EXPERIMENTAL SETUP

Ph.D. students are typically human, and adult^[citation needed]. In this specific setup, the author is a human male, aged 28 solar years. The average human male adult needs to consume, on average, about 10.6MJ (2,500kcal) of energy each day [6], [7], in the form of food. The use of a most-likely-not-very-scientific online calorie calculator (taking age, height, weight, and physical activity into account) places the author at a remarkably average energy consumption of 10.2MJ/day (2,437kcal/day). In the following, we will assimilate the author with the average male Ph.D. student. This yields an average thermal dispersion power (TDP) of around 120W. The study of the influence of intellectual effort on the metabolism of the body and the brain is a whole can of worms of contradicting and/or inconclusive results [8] that I am not willing to touch with a 10-foot pole. Thus, for the sake of this experiment, we will simply assume that, for a given experimental setup, our computer's TDP remains constant over time.

Now, how can we perform neural network inference?

²In case my supervisors ever come across this paper: it was *mostly* written in my free time. Its characterization as an exercise of procrastination is left entirely to the reader's evaluation

Neural networks' workloads are mostly made up of multiply-accumulates (MACs) - essentially matrix multiplications - and composition by nonlinear functions of variable complexity. Computers typically store numbers in the form of 32-bit floating point numbers. This number format comprises three parts: the sign bit, the exponent, and the significand (sometimes called the mantissa).

We won't go into too much detail, but people interested in the eldritch horror that is the IEEE-754 floating point standard may refer to my Ph.D. thesis for more details. It's gonna be great. I swear.

The way floating point additions and multiplications are computed can be visualized in figures 1 and 2.

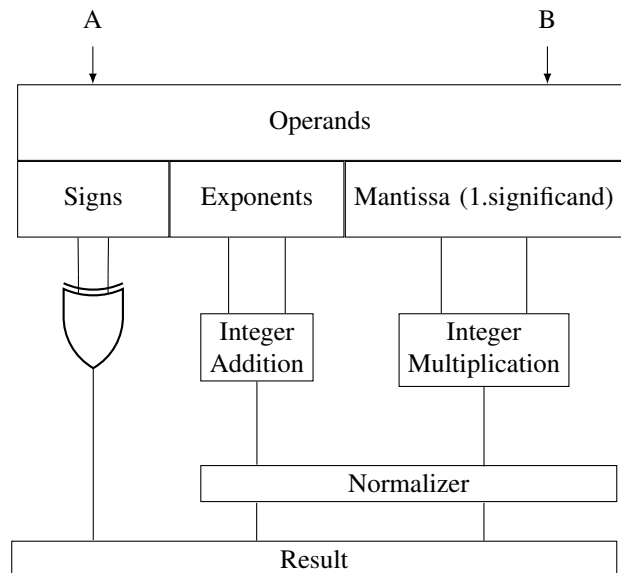


Figure 1: A float multiplier³

Look at those drawings, look at them well. Think of the obscene amount of time it took me to make them. Did you give them enough appreciation? Yes? Are you sure? Then we can move on.

I initially set out to try and measure my throughput in terms of manual floating-point multiply-accumulates, equipped with just my brain, a pencil, and some paper. After much effort, reading of low-level logic, and a lot of swearing, I concluded that performing floating-point operations manually required a level of masochism that not even I have. And that is saying something.

Fortunately, machine-learning researchers have come up with ways of simplifying the compute load of neural networks by changing the number representation of a neural network's weights, a process that is called quantization.

³The diagrams are on my GitHub so that nobody ever needs to import grainy JPGs, or suffer like I did to redraw them: <https://github.com/frost-is/TikZ-Diagrams>

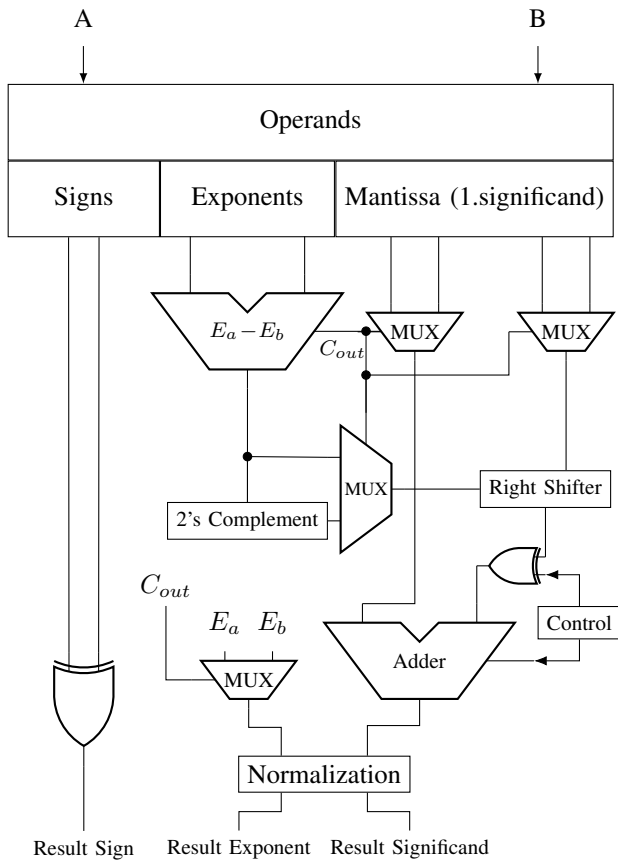


Figure 2: Seriously, who came up with this monstrosity? It took me ages to do a proper Tikz diagram³. This is a float adder, slightly simplified, but not by much.

The way it is usually done is by using the INT8 or INT16, signed integer formats. Here, for simplification purposes, we will consider the INT8 format, yielding a number range in $[-128; 127]$. Using a typical pen and paper, averaged over a couple of tens of randomly generated operations and counting conservatively, the author of this paper needs slightly under 30 seconds to perform an INT8 multiplication and about 10 seconds for an addition. For the activation functions, we should be able to refer to computation tables - since INT8 can only take 256 different values, these tables could fit on a page each. The bulk of the compute load should be made out of MACs, since the amounts of MACS we need to perform in layer k is *number of neurons in layer $(k - 1)$ \times number of neurons in layer k* , when the number of activations to compute is only *number of neurons in layer k* . Taking this into account, plus the relative simplicity of computing activations using tables, we will consider that the compute load of neural inference is solely made up of multiplications and additions, each equally represented.

If we want to try our hand with a rather lightweight neural network, we can try MobileNetV3 Small[9]. It comprises 44 million MACs. Following the logic we developed above, it should take about 28 years of continuous computing to perform a single inference, for a total inference cost of 105GJ of energy.

To give an order of magnitude, this represents about 3 minutes' worth of the full power of a nuclear reactor. So *a lot* of energy. This also means that the author of this paper, working 35-hour workweeks with 6 weeks of paid leave *per annum*, should take about 150 years to compute a single inference⁴.

V. THIS IS BAD. CAN WE DO BETTER ?

A. Slide Rules

Slide Rules allow us to change multiplications into additions through the power of logarithms. We therefore reach a throughput of one operation per 10s for both additions and multiplications - a twofold improvement. We will, however, consider that there are not enough of them left around to use at a truly large scale.



Figure 3: The slide rule. An elegant weapon for a more civilized age. This one was my late Grandma's

B. Sex

According to the French medical database Vidal [7] and the British National Health Service (NHS) [6], adult men need on average about 10,500kJ of food each day. Women, however, only need about 8400kJ. Thus, female Ph.D. students are on average 20% more energetically efficient than male Ph.D. students. To maximize efficiency, hiring female Ph.D. students for our inference setup seems the way to go, bringing us down to 85GJ per inference, albeit with identical inference time.

⁴Just in time for my retirement

C. Physical activity

The energy-consumption figures cited above are for average humans, who are moderately physically active. However, if we take humans and have them live where they work and not do anything besides working and sleeping, then their energy consumption is only their basal metabolic rate (BMR).

The topic of BMR is a surprisingly complex one, as BMR varies based on age, sex, ethnicity, muscles, etc. We will not delve too deep into it, but a good intro can be found in [10]. According to this data, the BMR of men and women between 18 and 30 years old can be approximated using the following formulas: $0.0546 \times \text{weight} + 2.33$ for women, and $0.0669 \times \text{weight} + 2.43$ for men (in $MJ.day^{-1}$). Now, this estimation requires data on the average weight of men and women. This statistic is highly variable depending on the country of origin, so here we will assume that our subjects are French. We will also neglect the effects of socioeconomic status on weight.

The Obepi-Roche study indicates that French men and women have an average weight of 81kg and 67kg, respectively [11], [12]. This means that the average French Ph.D. student's BMR is about 6MJ/day (1,435kcal/day) for women and 7.85MJ/day (1,880kcal/day) for men. The logical conclusion that follows is that, if we seal our Ph.D. students in capsules that prevent them from moving, we consume 25 to 30% less energy. This also means that we can make them work 16 hours a day 365 days a year without compromising their sleep too much. This won't change the amount of energy used, but will nonetheless reduce our inference time, down to just a little under 42 years per inference.

D. Age

One way of explaining the efficiency of GPUs for neural network inference and training – and video games – (or, at least, the way I explain it to those among my students who have little background in computer science or electronics) is the following:

"You can compare a CPU to a mathematician and a GPU to a bunch of middle-schoolers. The mathematician can perform much more complicated operations. However, it takes time and resources to raise someone into a mathematician - and to have them perform, since they consume more food. This means that, for the same cost, you can have more middle-schoolers than mathematicians. Mathematicians can compute operations that middle-schoolers are incapable of. However, as far as simple operations - such as additions and multiplications - are concerned, they are only marginally better than middle schoolers. On the opposite end of the spectrum, if your computations are complex and data dependent, with a lot of logic and branching, then middle schoolers will be slower – they will take a lot of time to break the task into many smaller tasks, often they will have to wait for one of them to give the output of one operation before being able to

advance on the task, and this will typically make them slower. This translates to CPUs being specially adapted for serialized, data-dependent operations with branching and logic, while GPUs are a better fit for parallelized, data-independent operations including little logic.

Let us assume mathematicians can do 3 times as many simple operations as middle-schoolers per unit of time. If the cost of a middle-schooler is less than a third of that of a mathematician, then buying middle-schoolers⁵ will be more efficient than buying mathematicians⁶. There are caveats. This supposes that the task can be split into simple chunks that can be processed independently (easily parallelizable). This is the case for machine learning and videogames/rendering: they are composed of simple, repetitive, data-independent computations with little logic - which mainly comprise additions, multiplications, and compositions by some basic functions. Thus, an army of middle-schoolers will tend to perform them more quickly than a handful of mathematicians bought for the same price. Although you will usually need mathematicians to split the task into chunks that middle-schoolers can process.

As an example: the latest Nvidia H100 (PCIe version) comprises 14,592 CUDA cores, each capable of a single floating multiply-accumulate per instruction. The largest modern CPU, the EPYC 9754, made by AMD, reaches 128 cores⁷. Even using vector instructions such as AVX512, which can process 8 32-bit multiply-accumulates per clock cycle, we only reach 1024 multiply-accumulates per instruction at full load – under ideal conditions. That is 14 times fewer than the H100. Both of these processing units have equivalent silicon areas and transistor counts – Nvidia's H100 boasts an 814mm² monolithic die comprising 80 billion transistors, and AMD's EPYC 9754 implements 8 chiplets of 73mm² die area each (584mm² in total), plus a 391mm² I/O die, for a total of 981mm² comprising 71 billion transistors. Their respective thermal dispersion powers are also equivalent, 360W for the EPYC and 350W for the H100.

Now, this explanation is very much simplified and would warrant a deep dive into the respective capabilities and tradeoffs of CPUs and GPUs – GPUs can't run operating systems, for example. But that's the intuition."

So, what if we use middle schoolers instead of Ph.D. students for our computations? Our best review of the literature did not find a satisfactory answer. One reason is that this axis of research is supremely inefficient and profoundly unethical.

⁵The author speaks from a purely theoretical point of view and certainly does not condone buying middle-schoolers

⁶The author does not condone buying mathematicians either

⁷We won't account for AMD's simultaneous multi-threading (SMT) or Intel's Hyperthreading, since the "operations per instruction" metric should, theoretically, not benefit from it. Almost all modern CPUs implement 2 threads per physical core, but it is not equivalent to having 2 physical cores

Another is that it's extremely difficult to find statistics for children, for understandable reasons. In any case, we were not able to find data for this age bracket.

Nonetheless, the World Health Organization (WHO) provides weight-for-age charts for children of 5 and below [13]. These charts are meant to help pediatricians determine the extent to which a child's physiological needs for growth and development are satisfied. These statistics go up to 6 years old for boys and 5 years old for girls. To keep the comparison fair, we will compare 5-year-old children. Their median weight is 18kg. The equations in [10], taken for children between 3 and 10, give us a median daily BMR of 3.60MJ/day (870kcal/day) for girls and 3.80MJ/day (920kcal/day) for boys, which is a significant gain in efficiency compared to Ph.D. students and would amount to 55GJ (resp. 58GJ) per inference.

There is, however, a caveat: to what extent can one teach 5-year-olds how to add and multiply, and how efficiently do they accomplish this task? This warranted inquiry. I went to my mother, a Kindergarten teacher, and asked her if I could borrow a couple of her kindergarteners to run an experiment. Her enthusiasm was rather lukewarm ("What?! Why? No!"). None of my friends have children in the right age bracket, and since I would do anything for Science, so long as it does not involve talking to strangers, this was the end of this way of exploration. Feel free to contact me if you have subjects available for this ethically undefendable but scientifically glorious experiment.

VI. OPTIMIZING AGAINST CARBON EMISSIONS

Optimizing biological neurons against energy is not necessarily relevant, as the biologically exploitable energy of food does not reflect the energy that was expended to produce it. This means that creating an accurate reflection of the energy efficiency of humans would need to take into account the entire supply chain, which would represent a tremendous work. Here, we argue that optimizing our workflow with regards to carbon emissions is more relevant, as carbon emissions a) are positively correlated with energy emissions, b) have direct consequences on our biome, where energy expenditure does not – or to a smaller extent –, and c) I'm starting to feel guilty of over procrastinating so I need to wrap this up.

Therefore, our problem of minimizing energy consumption can be simplified by trying to minimize carbon emissions. So, how can we do it? The origin of the food we feed our computers seems crucial here.

From a purely energy-efficient standpoint, vegetable oil is about as good as can be for humans, with energetic values around 37MJ/kg (8850kcal/kg)⁸. How does it translate in terms of carbon emissions? There are discrepancies when it

⁸This is about 1% of the energy density of reactor-grade uranium and 0.02% of that of pure Uranium 235. Unfortunately, humans tend to have a hard time digesting radioactive isotopes, which is a shame

comes to the type of seed the oil comes from. Rapeseed oil seems to be the most carbon-efficient crop for oil production [14], with median emissions of 2.49kgCO₂e per kg. This translates to about 14.85MJ/kgCO₂e (3,550kcal/kgCO₂e). But if we try to tackle the problem from a carbon emission standpoint, then we can further optimize the MJ/kgCO₂e metric.

Supposing that we have the logistics available to plant and process crops on site (a good way of mitigating, if not negating, carbon emissions related to transport, packaging, retail, and waste), then using the data from [15] and only accounting for the emissions related to farming, animal feed (if applicable) and processing, then we obtain table I.

Food	MJ/kg	kgCO ₂ e/kg	MJ/CO ₂ e
Rapeseed oil [14]	37 (885kcal)	2.49	14.9 (3,550kcal)
Potatoes [16]	3.6 (87kcal)	0.19	18.9 (4,578kcal)
Peanuts [17]	23.7 (567kcal)	0.6	39.5 (9,450kcal)

Table I: Some of the least greenhouse-gas emitting foods

So it seems peanuts are the most efficient crop as far as energy per kgCO₂e is concerned. Moreover, peanuts are edible raw, and since cooking leads to more GHG emissions, this further reduces our greenhouse gas footprint. For confirmation purposes, the author ate an entire potato raw to determine whether raw potatoes could reasonably be the sole component of a person's alimentation without inducing lifelong trauma, and the conclusion was a resounding "No".

Unfortunately, cal/kgCO₂e is not the only metric we need to account for. Humans are pesky creatures that, unlike computers, require more than just energy to function properly: they require various nutrients - among them proteins, fat, carbohydrates, and miscellaneous vitamins and minerals. Our best review of the literature did not come up with any metric measuring how much of a human's physiological needs are met by any given food. However, peanuts do contain a good amount of proteins, fibers, carbohydrates, and fat. Trying to come up with a mix of what could be added to peanuts to be the sole basis for humans would be way too involved for this paper, which already went way too far for something that was initially supposed to be a joke. However, given the fact that the "Plumpy'Nut" ready-to-use therapeutic food (RUTF), which was created for the treatment of severe acute malnutrition, uses peanut paste for its basis (it contains peanut paste, vegetable oil, powdered milk, powdered sugar, vitamins, and minerals), we feel confident that peanuts are, if not the best choice, among the better ones.

So, according to these metrics, a farm of French Ph.D. students doing nothing but computing 16 hours a day and sleeping the rest of the time would need to feed its computers 250g (women) to 330g (men) of locally produced peanut slurry per day. We will assume that men and women are equally represented in our setting, thus yielding an average of 290g of daily peanut requirements per person or about 106kg per year. This translates to 175g of CO₂e emissions

per day and 65kg per year. The best lands give a yield above 4t of peanuts per ha [17], with harvest taking place once a year. Given that feeding the entire population of France using peanut slurry should need around 6.9Mt of peanuts per year, locally producing them would require dedicating 17,200km² of arable land to peanuts, or $\approx 15.6\%$ of the total land area used for agriculture in France in 2020[18].

This means that, thus far, our grand plan of turning the entire French land and population into a gigantic computer farm is on track, and would yield about 3.2 INT8 MOPS of computing capabilities 16 hours a day, or 2.2 INT8 MOPS when taking shifts. Nice⁹.

Now, there are some caveats:

- 1) As the scientifically shaky but visually grand experiment known as *Christopher Nolan's Interstellar* has exemplified, monoculture is typically not a great idea in the long run, for a variety of reasons (mainly vulnerability to pests and diseases and lack of genetic variability). This is not a problem because, like all cartoonishly evil geniuses, I will be thwarted at the last second by a hero in shining armor straight from the old stories. This protagonist will either make me realize the error of my ways via the power of love and friendship (I certainly hope not), or defeat me in a battle of wits, whereupon I shall deliver my *Evil Monologue*TM (which I may or may not have already written while procrastinating on my Ph.D. dissertation), ending with "Someday you will realize that I was right, Iggy MacRainbows", before falling down the shaft of my soon-to-explode ~~Death Star~~ server farm, cackling in laughter. Of course, my death will be unconfirmed depending on whether box-office results determine that my story deserves a sequel. I give the aforementioned sequel 50/50¹⁰ chances of either butchering my character or giving me an awesome redemption arc. I honestly can't wait.
- 2) This study does not include the use of water. Humans require water. However, this would send us on the tangent of the tangent, and I'm already way late on my Ph.D. redaction so I ask for your understanding.
- 3) Nourishment by peanut slurry only would lack some essential nutrients. If we seek to include some vitamins and minerals in our slurry, then the emissions should be higher than announced here. This one will be tackled in the following section

So, to try and preserve our human computers for as long as possible, we will need to introduce vitamins and minerals in our peanut slurry, which will induce further GHG emissions. However, we argue that these elements are effective in such small quantities that artificial synthesis amounts to a negligible

⁹The astute reader may remark that H100 GPUs boast 3,026 INT8 TOPS. They are, however, much less visually entertaining

¹⁰100/0 if Disney buys the adaptation rights

amount of marginally added carbon emissions. In [19], it was determined that adding artificial D vitamin to bread, milk, and oil, only increases the food's carbon footprint by $\approx 0.001\%$. Thus, we will assume that the enrichment of our peanut slurry is negligible in terms of added carbon emissions.

Moreover, human bodies are reasonably nutritious: according to [20], a male human body of 55.26kg contains 32,376kcal of skeletal flesh. Should we eat the entire body (marrow, skin, adipose tissue, etc.), then we reach 143,771kcal per body¹¹. The study does not provide figures for female bodies as no chemical composition of the female body was available in the literature at the time of writing. If we do a back-of-the-envelope calculation assuming that these figures scale linearly and identically with weight for both sexes (they don't but we'll soon show that this ends up amounting to a rounding error), we can infer that the average French person's body contains about 190,000kcal, brain and rachis excluded (172,000 and 208,000kcal for women and men, respectively). If we liquefy the corpses of our dead computers via a comically large meat grinder to feed them to living computers, then this positively influences our carbon efficiency. Let us work out how much.

French data shows that 638,266 French people died in 2023 and that the French population at the time of writing is 64,842,629. The death of about 1% of the population each year for a country whose population is growing seems like a reasonable assumption in general. Assuming men and women are equivalently represented in the French population, dead bodies would average in at about 5kcal of food per living French person per day, about 0.3% of the average French person's BMR. Thus, the bulk of our computers' emissions will come from peanut culture, and we will simply assume that the carbon impact of the additives in the peanut slurry is offset by the corpse slurry.

Incidentally, we have just demonstrated that fundamental assumptions underlying the plots of *Soylent Green* and *The Matrix* cannot hold in real life.

VII. OKAY BUT CAN WE DO BETTER?

If we go out of the realm of what is currently achievable with 2024 technology, then we can further optimize: the human brain consumes about 20% of the entire body's calories [22]. If we remove people's brains and have them live in small vats filled with nutritive substances¹², then our apparatus only consumes about 20W per human. Estimating the nutritional needs of a brain in a vat is hell. Considering a brain cannot digest anything, nutrients must be supplied in a heavily processed

¹¹Although we advise against eating the brain, as eating the brain of an individual of one's species leads to higher rates of prion diseases, typically spongiform encephalopathy (the so-called "mad cow disease") [21]. Discarding the brain, rachis, and spinal cord would reduce the caloric value of the male human body by 2706kcal.

¹²Any resemblance with parts of Metal Gear Rising Revengeance's plot will result from pure chance

form, and we can throw our entire previous calculations out the window. So we won't explore this venue any further.

VIII. WHAT TASKS CAN WE PERFORM?

A. Neural Network Inference

Some of the most popular tasks performed by modern-day neural networks are computer vision and text generation. We have exemplified a single neural network inference of MobileNetV3 Small costing 42 years of computer inference, or 95MJ of energy on average. Compared to that, using an Nvidia A100 costs about 3.6mJ per inference [3], or 26×10^{12} less energy.

Now, how does that translate in terms of carbon emissions? Energy-related carbon emissions are a complicated endeavor: when adding load to the energy grid, pilotable energy sources (gas, oil, coal) will ramp up to keep the voltage and frequency at the desired level. So estimating the marginal impact of adding a load to the grid in terms of carbon emissions as *added load* \times *average emissions per unit of energy* is a bit of an oversimplification. We will do this oversimplification, as we assume that the added load that is our inference server is small enough (<1kW) to blend in the base power mix, but this assumption does not hold at larger scales. In 2023, the average energy-related carbon emissions in France were 39gCO₂e/KWh [23]. Thus, our A100 would emit 3.5 μ gCO₂e per inference. Our human computers, on the other hand, would emit about 2.7tCO₂e to perform the same task, around a trillion times worse.

Now, humans can have a very good level of performance provided sufficient training. Russakovsky et al. [24] note that the best human classifiers on the ImageNet dataset only had a 5.1% error rate, much better than MobileNet S's 12.7% error rate.

It seems a reasonable approximation to assume that a good human classifier should take about 10s for a single inference, yielding an energy efficiency of 1kJ per inference or 3mgCO₂e per inference, with much better accuracy and throughput than our server farm. So it seems that all of the pages above were but a fruitless endeavor. But a fun one to be sure!¹³

B. Playing Doom

How many human computers do we need to play Doom? Well, a single one. With abysmally low framerates, of course. Can we do better? Surely!

Going through the manual of Doom (1993), we find requirements comprising a 386 CPU (or better), 4Mb of RAM, and a VGA card. Doom ran solely on the CPU, and the i386 is capable of around 5 MIPS. It is excruciatingly hard to assume the compute performance of a CPU based on its MIPS metric alone. However, we can safely assume that if we parallelize the

¹³Statistics on the energy usage of LLMs will be present in my Ph.D. thesis, but I haven't had the time to process them for this paper, for which I apologize.

workload over the entire population of France, then we have enough computers to dispatch the entire workload efficiently, if only since there are more French people than transistors on an i386 die (275,000 or 855,000 depending on the version[25]). The problem we encounter here is that Amdahl's law only goes so far - parallelizing over so many humans, we become latency-bound, not compute-bound. Taking into account our previous results, we conclude that we could probably play Doom at a maximum efficiency of about 0.1 frames per second, at best.

IX. LIFECYCLE ANALYSIS

At this point, it would probably be less expensive (and more interesting) to use the entire population of France to recreate an entire computer industry, software plus hardware, and reimplement games and neural networks entirely from scratch, than to use those same people as computers¹⁴.

X. CONCLUSION AND FUTURE WORKS

In *The Matrix*, the machines went about as far as can be to exploit the brain power of humans efficiently while (relatively) preserving their complete bodily functions. It's just a shame that human brains are not very good at doing what processors do in the first place.

Also, energy-wise, that bit about liquefying the dead so that they can be fed intravenously to the living doesn't amount to all that much in the grand scheme of things. So Morpheus is being a tad dramatic here. However, it adds to the grimdark ambiance, so he gets a pass.

Finally, this paper allows us to remind people at family diners¹⁵ that AI is mainly just linear algebra: using a pen and paper and *some* patience, anyone capable of performing high school level Mathematics can perform the same tasks as the latest AI model. Also it's neither intelligent, nor artificial, and you should call it SALAMI, at least you'd feel rightfully embarrassed when asking me whether linear algebra should be given the same rights as humans. I swear when the Butlerian Jihad comes you people are the first to go¹⁶.

While this paper achieves little in the way of practical science, maybe the true science was the procrastination we made along the way. Now, there is a Ph.D. thesis to be finished.

XI. ETHICS STATEMENT

Creating labor camps, or worse, reproducing *The Matrix*TM, is bad, and you should not do it.

XII. SPECIAL THANKS

I would like to thank the entire Sigbovik team for having me research way more biology than I would have ever liked. I would like to thank my Ph.D. supervisors for not immediately assassinating me, should they ever come across this paper.

¹⁴Should politicians want to get in touch, my email is on the first page

¹⁵You know them

¹⁶Any resemblance to people or events existing or having existed would only be the result of pure chance

REFERENCES

- [1] W. S. McCulloch and W. Pitts, "A logical calculus of the ideas immanent in nervous activity," *The bulletin of mathematical biophysics*, pp. 115–133, 1943.
- [2] H. Touvron, L. Martin, K. Stone, *et al.*, "Llama 2: Open foundation and fine-tuned chat models," *arXiv preprint arXiv:2307.09288*, 2023.
- [3] H. Waltsburger, E. Libessart, C. Ren, A. Kolar, and R. Guinvarc'h, "Neural network scoring for efficient computing," in *2023 IEEE International Symposium on Circuits and Systems (ISCAS)*, IEEE, 2023, pp. 1–5.
- [4] D. A. Grier, *When computers were human*. Princeton University Press, 2013.
- [5] T. Haigh and P. E. Ceruzzi, *A new history of modern computing*. MIT Press, 2021.
- [6] N. H. Service, "What should my daily intake of calories be?" <https://web.archive.org/web/20240220201035/https://www.nhs.uk/common-health-questions/food-and-diet/what-should-my-daily-intake-of-calories-be/>,
- [7] Vidal, "Les recommandations nutritionnelles de 18 à 75 ans," <https://web.archive.org/web/20231110192642/https://www.vidal.fr/sante/nutrition/equilibre-alimentaire-adulte/recommandations-nutritionnelles-adulte.html>,
- [8] T. scientific american, "Does thinking really hard burn more calories?" <https://archive.is/ExrgC>, 2012.
- [9] A. Howard, M. Sandler, G. Chu, *et al.*, "Searching for MobileNetV3," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 1314–1324.
- [10] C. Henry, "Basal metabolic rate studies in humans: Measurement and development of new equations," *Public health nutrition*, vol. 8, no. 7a, pp. 1133–1152, 2005.
- [11] L. C. l'Obésité (League Against Obesity), "Taille, poids et tour de taille : Photographie 2020 des français," <https://web.archive.org/web/20231208192705/https://liguecontrelobesite.org/actualite/taille-poids-et-tour-de-taille-photographie-2020-des-francais/>, 2020.
- [12] A. Fontbonne, A. Currie, P. Tounian, *et al.*, "Prevalence of overweight and obesity in france: The 2020 obepi-roche study by the "ligue contre l'obésité"," *Journal of Clinical Medicine*, vol. 12, no. 3, p. 925, 2023.
- [13] W. H. Organization *et al.*, *WHO child growth standards: length/height-for-age, weight-for-age, weight-for-length, weight-for-height and body mass index-for-age: methods and development*. World Health Organization, 2006.
- [14] T. D. Alcock, D. E. Salt, P. Wilson, and S. J. Ramsden, "More sustainable vegetable oil: Balancing productivity with carbon storage opportunities," *Science of The Total Environment*, vol. 829, p. 154 539, 2022.
- [15] O. W. in Data, "Food: Greenhouse gas emissions across the supply chain," <https://ourworldindata.org/grapher/food-emissions-supply-chain?tab=table>, 2018.
- [16] J. Singh and L. Kaur, *Advances in potato chemistry and technology*. Academic press, 2016.
- [17] J. McCarty, S. Ramsey, and H. Sandefur, "A historical analysis of the environmental footprint of peanut production in the united states from 1980 to 2014," *Peanut Science*, vol. 43, no. 2, pp. 157–167, 2016.
- [18] C. d'Agriculture, "Les chiffres 2022 de l'agriculture française," https://web.archive.org/web/20240229192304/https://chambres-agriculture.fr/fileadmin/user_upload/National/FAL_commun/publications/National/Plaqueette_chiffres_de_l_agriculture_VDEF.pdf, 2022.
- [19] M. J. Bruins and U. Létinois, "Adequate vitamin d intake cannot be achieved within carbon emission limits unless food is fortified: A simulation study," *Nutrients*, vol. 13, no. 2, p. 592, 2021.
- [20] J. Cole, "Assessing the calorific significance of episodes of human cannibalism in the palaeolithic," *Scientific reports*, vol. 7, no. 1, p. 44 707, 2017.
- [21] P. P. Liberski, A. Gajos, B. Sikorska, and S. Lindenbaum, "Kuru, the first human prion disease," *Viruses*, vol. 11, no. 3, p. 232, 2019.
- [22] G. J. Siegel and R. W. Albers, *Basic neurochemistry: molecular, cellular, and medical aspects*. Raven Press, 1994.
- [23] "French historical data - carbon emissions for energy generation," <https://www.nowtricity.com/country/france/>, 2023.
- [24] O. Russakovsky, J. Deng, H. Su, *et al.*, "Imagenet large scale visual recognition challenge," *International journal of computer vision*, vol. 115, pp. 211–252, 2015.
- [25] "Microprocessor quick reference guide," <https://web.archive.org/web/20240316123819/https://www.intel.com/pressroom/kits/quickreffam.htm>, 2008.

Advancing Consensus: Automated Persuasion Networks for Public Belief Enhancement

osmarks.net Computational Memetics Division
comp.meme@osmarks.net

-1 April 2024

Abstract

Incorrect lay beliefs, as produced by disinformation campaigns and otherwise, are an increasingly severe threat to human civilization, as exemplified by the many failings of the public during the COVID-19 pandemic. We propose an end-to-end system, based on application of modern AI techniques at scale, designed to influence mass sentiment in a well-informed and beneficial direction.

1 Introduction

In today's increasingly complex and rapidly changing world, it is challenging for people to maintain accurate knowledge about more than a small part of the world [Kilov, 2021] [Crichton, 2002], but it's socially unacceptable or undesirable, and in some cases impossible, to reserve judgment and not proffer an opinion on every topic. As a direct consequence, many have incorrect beliefs, acting on which leads to negative consequences both for themselves and society in general [Cicero, 2001]. This is exacerbated by the increasing prevalence of misinformation, disinformation and malinformation [Pérez-Escolar et al., 2023] harming the public's ability to reach truth and make informed, justified decisions. In this hostile environment, attempts to enhance education in critical thinking are insufficiently timely and far-reaching, and a more direct solution is needed.

In this paper, we propose the Automated Persuasion Network, a system for deploying modern large language models (LLMs) to efficiently influence public opinions in desirable directions via social media. We develop an architecture intended to allow selective, effective changes to belief systems by exploiting social conformity.

2 Methodology

2.1 Overview

Humans derive beliefs and opinions from their perception of the beliefs and opinions of their peer group [Cialdini and Goldstein, 2004] [Deutsch and Gerard, 1955], as well as a broader perception of what is presently socially acceptable, required or forbidden. Our approach relies on a Sybil attack [Alvisi et al., 2013] against this social processing, executed by deploying LLMs to emulate people of similar attitudes to targets within the context of online social media platforms. While [Bocian et al., 2024] suggests that social pressure from AIs known to be AIs can be effective, we believe that persuasion by apparent humans is more robust and generalizable, especially since even the perception of automated social interaction has been known to trigger backlash or fear from a wide range of groups [Fang and Nie, 2023] [Yan et al., 2023]. We automatically derive strategies to affect desired beliefs indirectly, via creating social proof for other related beliefs, using a Bayesian network approach.

Naive implementations of this method involve many manual processing steps — for instance, identification of targets, construction of personas for LLMs to emulate, and gathering data for belief causal modelling. We replace these with automated solutions based on natural language processing — unsupervised clustering of internet users using text embeddings, direct evaluation of currently held opinions within a group using LLMs, and surveying simulacra rather than specific extant humans (as described in [Argyle et al., 2023]) — to allow operation at scale without direct human oversight. This permits much more finely individualized targeting than used in e.g. [Simchon et al., 2024] without additional human labour.

2.2 Segmentation

In order to benefit from the effectiveness of persuasive strategies optimized for individuals while still having enough data for reasonable targeting, we apply standard unsupervised clustering techniques. We acquire profile information and a social graph (of friendships and interactions) for all relevant social media accounts, generate text embeddings from each user’s profile information, as well as a representative sample of their publicly accessible posts, and combine this with graph embeddings to generate a unified representation. We then apply the OPTICS clustering algorithm [Ankerst et al., 1999] to generate a list of clusters.

From these, several pieces of information need to be extracted. We identify the accounts closest to the cluster’s centroid and take them as exemplars, and additionally compute the distribution of posting frequency and timings. We use these in later stages to ensure that our personas cannot be distinguished via timing side-channels. Additionally, we generate a set of personas using a variant of QDAIF [Bradley et al., 2023], with a standard instruction-tuned LLM (IT-LLM) used to mutate samples, using the cluster exemplars as the initial

seed. As a quality metric, we ask the IT-LLM to evaluate the realism of a persona and its alignment with the exemplars, and we organize our search space into bins using k-means clustering on the generated user sentence embeddings to ensure coverage of all persona types within a cluster.

2.3 Analysis

We use a variant of [Powell et al., 2018]’s methodology to tune persuasion strategies to audiences to effectively affect target beliefs. We replace their manual identification and belief measurement step by using the IT-LLM to first generate a set of beliefs that relate to and/or could plausibly cause the target belief, as well as scales for measuring adherence to these possible beliefs. For measurement, rather than using the IT-LLM as before, we apply a prompt-engineered non-instruction-tuned model (also known as a foundation model, base model or pretrained language model (PT-LLM)). This is because instruction-tuned LLMs are frequently vulnerable to the phenomenon of mode collapse [janus, 2022] [Hamilton, 2024], in which models fail to generalize over latent variables such as authorship of text. This is incompatible with our need to faithfully simulate a wide range of internet users. Instruction-tuned LLMs are also unsuitable for direct internet-facing deployment, due to the risk of prompt injection [Perez and Ribeiro, 2022]. Within each cluster, we use the acquired representative text from each exemplar from the segmentation stage to condition the LLM generations, and then ask several instances the generated questions in a random order. Multiple separate sampling runs are necessary due to the “simulator” nature of LLMs [Shanahan et al., 2023]: our persona may not fully constrain its model to a single person with consistent beliefs. Runs producing responses that cannot be parsed into valid responses are discarded.

Given this synthetic data on belief prevalence, we apply a structure learning algorithm to infer causality — which beliefs cause other beliefs. Unlike [Powell et al., 2018], we do not incorporate any prior structure from theory — due to the additional complexity of applying theories in our automated pipeline, and since our requirements lean more toward predictive accuracy than human interpretability — and instead apply their BDHC algorithm to generate many candidate graphs, selecting a final model based on a weighted combination of model complexity (number of edges) and likelihood, to combat overfitting.

We then select the beliefs with the greatest estimated contribution to our target belief and direct the IT-LLM to modify our generated personas with the necessary belief adjustment. Due to the aforementioned mode collapse issues, we apply rejection sampling, discarding any generated personas that diverge too far from their original forms (as measured by semantic embedding distance) and regenerating. The resulting personas are used in the next stage.

2.4 Interaction

After the completion of the previous stages, the Automated Persuasion Network must interact with humans to cause belief updates. This step requires large-scale inference: however, as most human communication is simple and easy to model, at least over short contexts, we are able to use standard low-cost consumer GPUs running open-weight PT-LLMs, using the vLLM [Kwon et al., 2023] inference server. As an additional cost-saving measure, we use a multi-tiered system whereby generations are initially run on a small model and, if too complex for it (as measured by perplexity), rerun using a more capable language model.

We use the belief-modified personas generated in the Analysis stage, and attempt to have each of them mimic the actions of a human user in their cluster as much as possible. We identified a number of challenges. Most notably, nonhuman users are frequently detected using posting frequency [Howard and Kollanyi, 2016] and timings [Duh et al., 2018] [Pan et al., 2016]. By using a fairly large set of accounts rather than a single bot, we can avoid detection based on simply noticing anomalously high posting frequencies, and by scheduling generation of new posts and conditionally replying to other users' posts in accordance with cluster statistics for such gathered during the Segmentation stage we can prevent easy timing-based detection. We have not yet identified a complete strategy for avoiding social-graph-based detection such as [Alvisi et al., 2013]: our present best mitigation is to deploy new personas slowly and to maintain the rate of interaction between them at the base rate within the cluster.

Other difficulties involve technical countermeasures in use against nonhuman users, such as CAPTCHAs and limited APIs. However, while today's most sophisticated CAPTCHAs exceed current AI capabilities, commercial services are available to dispatch solving to humans at very low cost. We are able to mitigate other limitations with the use of commercial residential proxy services and browser automation software for scraping.

2.5 Monitoring

In order to determine the efficacy of our approach, we periodically sample posts from human users within each cluster and apply the IT-LLM to rate how much each post entails our target beliefs, allowing measurement of belief change over time.

3 Results

No results are available for release at this time.

4 Discussion

We believe our architecture represents a major advance in misinformation prevention and public attitude alignment. A promising future direction for research we have identified is introduction of technical enhancements such as implementation of speculative decoding in post generation, as well as use of vision/language models such as [Liu et al., 2023] to allow interaction with multimodal content. We also suggest integration of concepts from LLM agents to reduce distinguishability from humans — for instance, personas could be given the ability to create new posts based on newly released news articles or information from other social media sites. Finally, while we have primarily focused on human emulation with some limited optimization of persuasive strategies, future AI technology is likely to be capable of more powerful direct persuasion.

References

- [Alvisi et al., 2013] Alvisi, L., Clement, A., Epasto, A., Lattanzi, S., and Panconesi, A. (2013). Sok: The evolution of sybil defense via social networks. In *2013 IEEE Symposium on Security and Privacy*, pages 382–396.
- [Ankerst et al., 1999] Ankerst, M., Breunig, M. M., Kriegel, H.-P., and Sander, J. (1999). Optics: Ordering points to identify the clustering structure. In Delis, A., Faloutsos, C., and Ghandeharizadeh, S., editors, *SIGMOD 1999, Proceedings ACM SIGMOD International Conference on Management of Data, June 1-3, 1999, Philadelphia, Pennsylvania, USA*, pages 49–60. ACM Press.
- [Argyle et al., 2023] Argyle, L. P., Busby, E. C., Fulda, N., Gubler, J. R., Rytting, C., and Wingate, D. (2023). Out of one, many: Using language models to simulate human samples. *Political Analysis*, 31(3):337–351.
- [Bocian et al., 2024] Bocian, K., Gonidis, L., and Everett, J. A. C. (2024). Moral conformity in a digital world: Human and nonhuman agents as a source of social pressure for judgments of moral character. *PloS one*, 19(2):e0298293.
- [Bradley et al., 2023] Bradley, H., Dai, A., Teufel, H., Zhang, J., Oostermeijer, K., Bellagente, M., Clune, J., Stanley, K., Schott, G., and Lehman, J. (2023). Quality-diversity through ai feedback.
- [Cialdini and Goldstein, 2004] Cialdini, R. B. and Goldstein, N. J. (2004). Social influence: compliance and conformity. *Annual Review of Psychology*, 55:591–621.
- [Cicero, 2001] Cicero, M. T. (2001). *On the Ideal Orator*. Oxford University Press.

- [Crichton, 2002] Crichton, M. (2002). Why speculate? <https://web.archive.org/web/20070714204136/http://www.michaelcrichton.net/speech-whyspeculate.html>. Speech presented at the International Leadership Forum, La Jolla, California, April 26.
- [Deutsch and Gerard, 1955] Deutsch, M. and Gerard, H. B. (1955). A study of normative and informational social influences upon individual judgement. *Journal of Abnormal Psychology*, 51 3:629–36.
- [Duh et al., 2018] Duh, A., Rupnik, M. S., and Korošak, D. (2018). Collective behavior of social bots is encoded in their temporal twitter activity. *Big Data*, 6(2):113–123.
- [Fang and Nie, 2023] Fang, W. and Nie, C. (2023). Social media use, social bot literacy, perceived threats from bots, and perceived bot control: a moderated-mediation model. *Behaviour & Information Technology*, 0(0):1–17.
- [Hamilton, 2024] Hamilton, S. (2024). Detecting mode collapse in language models via narration.
- [Howard and Kollanyi, 2016] Howard, P. N. and Kollanyi, B. (2016). Bots, #strongerin, and #brexit: Computational propaganda during the uk-eu referendum.
- [janus, 2022] janus (2022). Mysteries of mode collapse. Accessed: 2022-11-08.
- [Kilov, 2021] Kilov, D. (2021). The brittleness of expertise and why it matters. *Synthese*, 199(1):3431–3455.
- [Kwon et al., 2023] Kwon, W., Li, Z., Zhuang, S., Sheng, Y., Zheng, L., Yu, C. H., Gonzalez, J. E., Zhang, H., and Stoica, I. (2023). Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*.
- [Liu et al., 2023] Liu, H., Li, C., Li, Y., and Lee, Y. J. (2023). Improved baselines with visual instruction tuning.
- [Pan et al., 2016] Pan, J., Liu, Y., Liu, X., and Hu, H. (2016). Discriminating bot accounts based solely on temporal features of microblog behavior. *Physica A: Statistical Mechanics and its Applications*, 450:193–204.
- [Perez and Ribeiro, 2022] Perez, F. and Ribeiro, I. (2022). Ignore previous prompt: Attack techniques for language models.
- [Powell et al., 2018] Powell, D., Weisman, K., and Markman, E. M. (2018). Articulating lay theories through graphical models: A study of beliefs surrounding vaccination decisions. In *Proceedings of the 40th Annual Conference of the Cognitive Science Society*.

- [Pérez-Escolar et al., 2023] Pérez-Escolar, M., Lilleker, D., and Tapia-Frade, A. (2023). A systematic literature review of the phenomenon of disinformation and misinformation. *Media and Communication*, 11(2):76–87.
- [Shanahan et al., 2023] Shanahan, M., McDonell, K., and Reynolds, L. (2023). Role play with large language models. *Nature*, 623(7987):493–498.
- [Simchon et al., 2024] Simchon, A., Edwards, M., and Lewandowsky, S. (2024). The persuasive effects of political microtargeting in the age of generative artificial intelligence. *PNAS Nexus*, 3(2):pgae035.
- [Yan et al., 2023] Yan, H. Y., Yang, K.-C., Shanahan, J., and Menczer, F. (2023). Exposure to social bots amplifies perceptual biases and regulation propensity. *Scientific Reports*, 13(1):20707.

Grounding Language Models to their Physical Presence

Seongmin Park
ActionPower
Seoul, Republic of Korea
seongmin.park@actionpower.kr

All models are wrong but some are.. *useful*. (*wink wink*)

George Box
(Emphasis and winks added at the quoter's discretion.)

Abstract

We address an unattended but important gap in artificial intelligence safety research: if all these language models come to life and smack us in the head, will we survive? This paper derives a mathematical relationship between a model's parameter count and the physical danger of its namesake. We first compile a list of language models at large, assess the ability of their physical incarnations to smack us in the head, and, most importantly, provide a guideline for just how big is too big.

1 Introduction

Recent discussions on artificial intelligence (AI) safety revolve around imposing regulations based on computing power [3]. Debates build on the assumption that "emergent" hazardous abilities of large language models follow the monotonic increase in their parameter count. Empirical evidence suggests this assumption is a safe one. If we are to witness the bitter lesson [4] even more in the coming years, we must organize an intervention and check whether the potentially dangerous subconscious of human creators is seeping into model characteristics. Namely, in names.

A thorough investigation reveals that we tend to label harmful models with names inspired by harmful entities.

2 Language models evaluated

We provide a rundown of all language models tested in Table 1.

3 Methodology

We apply logistic regression to assess the physical harmfulness of a model given its parameter size. We use Lasso regularization with $\alpha = 0.3$. The justification for the method choice is two-fold. First, we need a "clean-room" experiment, where no neural network will conspire against us in assessing its family members' hazardousness. Second, I will be fired if I devote company GPU time to this personal pursuit of knowledge, despite its obvious importance regarding the future of South Korea, North Korea, the rest of Asia, both Americas, Africa, Oceania, and Europe except Germany. In case of a hostile LLM outbreak, we advise Germany to work from home, just as a certain national soccer team manager of their nationality did last year. An LLM would have done a better job.

Table 1: Language models assessed for potential hazardousness of their physical incarnations. CHHSH is short for "Can Harm a Human by Smacking them on the Head".

Model	Parameter count	Weight (kg)	CHHSH	Notes
ELMO [2]	13.6M ~ 93.6M	0.17	No	
BERT	110M, 336M	0.17	Yes	You can tell from its eyes.
CamemBERT	110M, 336M	0.25	No	If someone beats your head with a block of cheese, you'll probably be okay.
BART	139M, 406M	38.55	Yes	
T5	60M ~ 11B	2300	Yes	A Volkswagen minivan. Definitely CHHSH.
Pegasus	560M	1000	Yes	
LLaMA	7B ~ 65B	200	Yes	
LLaMA 2	7B ~ 70B	200 * 2	Yes	
Dolly	12B	80	No	
Falcon	7B ~ 180B	1.5	Yes	
Palm	540B	544.3	Yes	Tree, not hand.
Palm 2	340B	544.3 * 2	Yes	Two trees.
Bard	137BB	75	Yes	
Gemini	1.8B, 3.25B	75 * 2	Yes	We do not include larger versions of this model. We encourage Google to release details of its higher tier models if they wish to participate in this prestigious research.
Vicuna	7B, 13B	50	No	
Guanaco	7B, 13B	120	Yes	
Koala	13B	12	Maybe	Vicious little things. Also I believe Yoda was a Koala?
OPTIMUS	227M	4300	Yes	
Chinchilla	70B	0.75	No	
Gopher	280B	1	Yes	Yes.
Platypus	7B ~ 70B	3	No	
Dolphin	7B ~ 70B	100	Yes	
Goliath	120B	264	Yes	
Zephyr	3B, 7B	1.2	No	https://www.google.com/search?q=weight+of+air
Wizard LM	7B ~ 70B	63	Yes	We expect wizards to be a bit scrawny.
Codex	120B	74.8	Yes	It's a pretty large book. Refer to the Discworld series for various demonstrations of old thick books wreaking havoc.
Starling	7B	0.1	No	
Orca	7B	4000	Yes	
Stable Beluga	7B, 40B	1000	No	I trust them.
Camel	5B, 20B	800	Yes	
Pythia	70M ~ 120B	73.6	Yes	Average weight of a Greek woman. [1]
Dromedary	7B ~ 70B	500	Yes	
Nemo	3.8B, 15B	0.2	Yes	Based on the deleted scenes.
Sparrow	7B	0.024	No	

We also account for the real-life weight of each model's namesake. We believe the weights are a good indicator of possible physical harm. We study the relationship between physical weight and parameter size.

4 Hazard analysis

4.1 Experimental results.

We do not have a dedicated validation or test set. Every data was used for training with no samples to spare for testing. Sorry.

As shown in Table 1, we have a very scanty set of samples to work with. We posit that researchers should publish more models before asking us to objectively test our results. Benchmarks are overrated anyway.

4.2 But we do have some graphs. This is important!

Based on our fitted regression model, we project the hazardousness of models with sizes ranging from zero to 500 billion parameters. We find that around 300 billion parameters is where models start to display some serious capability to harm us, should they come alive (Figure 1).

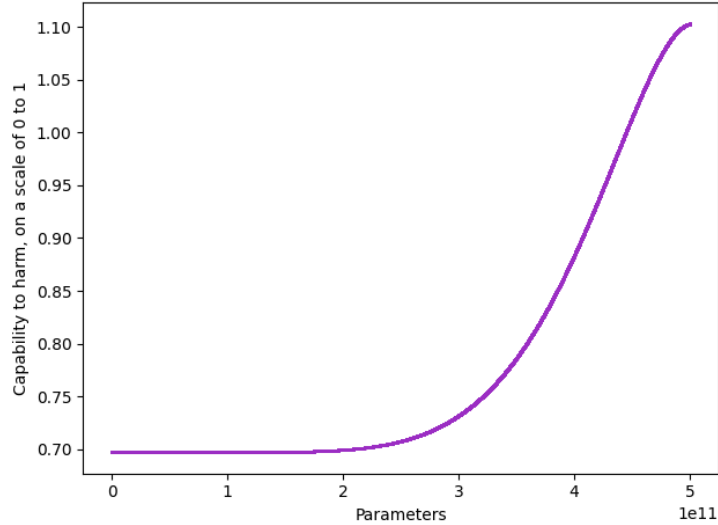


Figure 1: A model's ability to kill given its parameter count. The J-curve starts around 300B and starts to plateau at 500B.

We also estimate real-life weight given a model's parameter count (Figure 2). Surprisingly, we do not observe a monotonic increase in estimated incarnation weight as parameter size increases. The graph looks almost identical to the smile of a skilled predator as it looks forward to engaging in a spree of systematic and nefarious head-smacking.

We shudder.

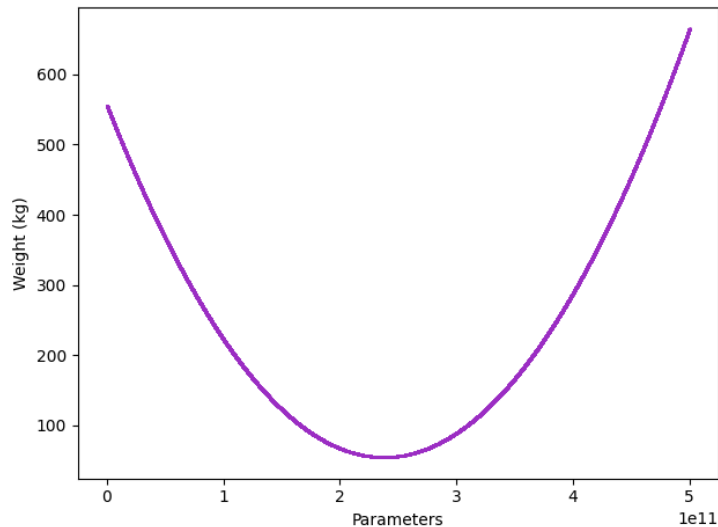


Figure 2: Estimated weight of physical incarnation given parameters. From the data, our model extrapolates that, if our language models jump out of our screens, those between 200B to 300B will be the most lightweight.

4.3 And a table, too!

Table 2 records the correlation between weight and harmfulness, weight and parameter count, and parameter count and harmfulness. Since such gigantic p-values make the authors nauseate, we leave the discussion on feature correlation (and the overall validity of this research) for future work.

Table 2: Relationship between features. Please withhold any verbal exclamation at the p-values, as you might discourage them.

Features	Pearson Correlation	P-value
Weight & Harmfulness	0.3191	0.005
Weight & Parameter count	-0.044	0.709
Parameter count & Harmfulness	0.060	0.608

5 Conclusion

In conclusion, 200B models are tolerable but don't go over 300B, okay? Name your models Noodles or Stockfish or something. I guess those chess nerds had it right all along.

References

- [1] W. Data. Average sizes of men and women. URL <https://www.worlddata.info/average-bodyheight.php>.
- [2] Google. Google query for model sizes. URL https://www.google.com/search?q=model_name.
- [3] NTIA. Ntia ai open model weights rfc. 2024. URL <https://www.regulations.gov/document/NTIA-2023-0009-0001>.
- [4] R. Sutton. The bitter lesson. 2019. URL <http://www.incompleteideas.net/IncIdeas/BitterLesson.html>.

Quantum Disadvantage

Or, simulating IBM's 'quantum utility' experiment with a Commodore 64

Anonymous

quantum.disadvantage@proton.me

<https://quantum-disadvantage.co.uk>

March 29, 2024

Abstract

IBM recently performed an experiment measuring expectation values of a Trotterized Ising model on their 127-qubit 'Eagle' quantum computer, along with a claim that their results were unlikely to be replicable on a classical computer. In the spirit of friendly academic shitposting, we implement a classical simulation of the same experiment using a Commodore 64, responding to a challenge posted by Dulwich Quantum. Our simulation requires less than 15kB of memory and approximately four minutes of wall-clock time per data point. To accomplish this we use a variant of the sparse Pauli dynamics (SPD) method recently developed by Begušić and Chan. We show that aggressive truncation combined with a shallow depth-first search avoids the prohibitive (for a C64) memory cost of storing the truncated Pauli basis in SPD, while maintaining sufficient accuracy to match the error-mitigated results obtained from the quantum device.

***Disclaimer:** Given the current debate about the significance of IBM's experiment, I wish to clarify that this work is not intended as any disrespect to IBM's team or their results - their experiment is a compelling case study of error mitigation for meaningful NISQ computation. My view is that they were unlucky in the specific structure of the circuits used in the experiment, and that with a few tweaks it would no longer be vulnerable to many forms of classical simulation. That said, I couldn't resist the challenge...*

1 Introduction

This project is a joke, and like all jokes, it's only funny if you have to explain it. Therefore, we will provide some brief introduction and motivation. Let's consider a hypothetical - suppose you are minding your own business, reading Louis H. Kauffman's book 'Knots and Physics'. You can feel in your bones that you're on the verge of a breakthrough: the idea isn't quite fully formed in your mind, but you're sure that you're about to make a revolutionary discovery in knot theory that will simultaneously solve world hunger and climate change, and will inevitably lead the human race to a new era ruled by a Platonic philosopher-king who speaks only in category theory. Suddenly, you see it - a single mosquito on the other side of the room emits a high-pitched buzz, threatening to break your focus. Of course, you know what you have to do - after all, the fate of the human race is at stake - so you calmly take your trusty M20 super bazooka, and snipe that pissant arthropod out of the sky, cremating it where it stands. All is now well in the world, and you can return to your research.

Suppose for the sake of argument that this was funny. Why was it funny? Of course, apart from the primal lizard brain instinct of "big thing goes boom, hur hur 🤩"¹ we can

see that it forms a kind of strange loop: we wanted to do something simple (swat a mosquito) but decided to use something vastly overpowered to solve the task (a rocket launcher) which ought to make the task easier, but instead made the task much harder than it needed to be (have you tried to snipe a mosquito??). In his seminal work [10], Tom Murphy VII refers to this as an *improper hierarchy*, and he defines it to be a type of joke, so of course it must be funny.

Quantum computers have been widely touted as a new computational tool with great potential to enable computations that were previously infeasible. As explained by the leading expert Michio Kaku, they can solve difficult problems by trying all possible answers at once in quantum superposition and then returning the correct one². In this work we are going to take an experiment that was run on a quantum computer, and run it on several devices that are decidedly not quantum computers (quantum computers, if you will) and take ourselves far too seriously while doing it. This is an improper hierarchy, but of the opposite flavor to the one we had before: we are going to take something that is ostensibly a hard problem (approximating the expectation values of an Ising model), solve it with a vastly under-powered device, and in so doing, make it much easier³.

2 Heisenberg's Cat

Did you ever hear the tragedy of Heisenberg's cat? No? I thought not, it's not a story Nielsen & Chuang would tell you - it's an Everettian legend.

Richard Feynman [2]

Quantum computers are just a bunch of atoms in a trench coat pretending to be a computer, except the trench coat is an electromagnetic field, and we're the ones doing the pretend-

¹This is intended to be vaguely euphemistic in a way that makes the reader uncomfortable and confused.

²To Michio Kaku: this is satire, please don't sue me. To Scott Aaronson: you can relax, I do actual quantum algorithms for my day job :)

³This was still plenty hard, but easier than building a quantum computer.

ing. That is, they are carefully-controlled systems of physical objects which obey the laws of quantum mechanics. Mathematically, we describe the state of a system by some high-dimensional complex-valued vector v . The laws of quantum mechanics assert two things:

1. When something happens to the quantum system, this corresponds to multiplying v by some matrix U , so that the new v' is given by $v' = Uv$.
2. Any quantity we want to measure about the system, an *observable*, can be represented by a matrix O , and measuring that observable gives you (approximately) the value $v^\dagger O v$ (note $v^\dagger = \bar{v}^T$ where \bar{x} is the complex conjugate).

To perform quantum computations, we initialize the vector v to some known state, and then perform a sequence of operations U_i , and measure some final observable O . The operations U_i correspond directly to physical actions on the system, for instance shooting lasers at trapped ions or applying electrical signals to points on a semiconductor. This gives us the quantity

$$\langle O \rangle = v^\dagger U_1^\dagger U_2^\dagger \dots U_n^\dagger O U_n \dots U_2 U_1 v$$

which, if you're lucky, will be a number that you care about.

There are two ways to think about this process: we can see $\langle O \rangle$ as measuring a fixed observable O of a sequence of states v_i that evolve forward in time as $v_{i+1} = U_i v_i$ with $v_1 = v$, so that:

$$\langle O \rangle = v_{n+1}^\dagger O v_{n+1}$$

This is called the Schrödinger picture. Alternatively, and more useful for this project, we can imagine measuring a fixed state v against a sequence of observables O_i that evolve backwards in time as $O_i = U_i^\dagger O_{i+1} U_i$ with $O_{n+1} = O$, so that:

$$\langle O \rangle = v^\dagger O_1 v$$

This is called the Heisenberg picture, which is what we will work with. At first, it might seem that quantum computers are somewhat trivial, since they just do linear algebra, which we can already do with regular (also called 'classical') computers. However, the vector v is so large (2^{127} dimensions for the system we will look at here), that we could never hope to directly calculate it. The ability to work with such high-dimensional spaces natively is where the power of quantum computers over classical computers comes from (if there is any), and gently coaxing $\langle O \rangle$ to be an interesting and commercially useful number by carefully choosing each U_i is the delicate art of quantum algorithm design.

2.1 'The Experiment'

On 14th June 2023, a team at IBM (Kim et al) published results from an experiment where they calculated $\langle O \rangle$ on a real quantum device, configured to approximate the behavior of a system known as the Ising model [6]. This is a model that is used by physicists to study the behavior of ferromagnetic materials. The team claimed that this computation would be too difficult to perform on a classical computer to an acceptable accuracy, using the leading approximation techniques,

and hence that it showed their quantum device could be useful. They published directly in *Nature*, without first releasing a preprint on arXiv (which is customary in the field). This made a lot of people very unhappy and was widely regarded as a *bad move*.

You see, a team at Google had made similar claims in 2019 [1]. After much research it turns out that it is probably possible to simulate their experiment classically, given access to a large cluster of GPUs for a significant amount of time [7]. However, in IBM's case they had unwittingly chosen a task that was much easier to approximate classically, and indeed only *five days* after the results were released, a preprint appeared on arXiv showing that the experiment could be simulated in minutes on a laptop [9]. Two days later another preprint appeared obtaining similar results by a different method [3, 4] - if I had to guess, I'd say this is probably the shortest time-to-refutation of any paper published in *Nature*. Naturally, the situation immediately devolved into a Twitter war and ended up with IBM's VP for quantum fighting it out in the comments of LinkedIn posts. This paper is a response to a challenge set by Twitter user Dulwich Quantum to simulate IBM's experiment using a Commodore 64.

2.2 Okay, so like how are we doing this?

I'm so glad you asked. We will use the *sparse Pauli dynamics* technique developed by Begušić, Hejazi, and Chan [5, 3]. The idea is as follows: the vector v lives in a 2^{127} -dimensional space, and correspondingly the matrices U_i and O are of size $2^{127} \times 2^{127}$, which is clearly too big to calculate directly. However, we can approximate these matrices as the sum of matrices that have a special form. Let I , X , Y , and Z represent the matrices:

$$I = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \quad X = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \quad Y = \begin{pmatrix} 0 & -i \\ i & 0 \end{pmatrix} \quad Z = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}$$

We call matrices of the form $P_1 \otimes P_2 \otimes \dots \otimes P_n$ *Pauli strings*, where \otimes is the Kronecker product that maps an $n \times n$ matrix and a $k \times k$ matrix to an $nk \times nk$ matrix according to:

$$(A \otimes B)_{ak+cb, bk+d} = A_{a,b} B_{c,d}$$

We will represent them compactly by the string $P_1 P_2 \dots P_n$ - for example, $IXZI$. Note also that since the set I , X , Y , and Z is closed under multiplication, then so are Pauli strings. In this way, a Pauli string of length 127 represents a $2^{127} \times 2^{127}$ matrix. These have the useful property that given Pauli strings P and Q , we have:

$$(e^{i\theta Q})^\dagger P e^{i\theta Q} = \begin{cases} \cos(\theta)P + \sin(\theta)QP & \text{if } PQ \neq QP \\ P & \text{if } PQ = QP \end{cases}$$

Note that any U_i as defined in the previous section can be written as some sequence $e^{i\theta_1 Q_1} \dots e^{i\theta_m Q_m}$ - therefore, we can write the whole computation as:

$$\langle O \rangle = v^\dagger (e^{i\theta_1 Q_1})^\dagger \dots (e^{i\theta_n Q_n})^\dagger O e^{i\theta_n Q_n} \dots e^{i\theta_1 Q_1} v$$

We wish to approximate the sequence of observables O_i in the Heisenberg picture as a linear combination of Pauli strings. Let

$$O_{i+1} = \sum_{j=0}^n \alpha_{i+1,j} P_{i+1,j}$$

then we can see from the above property of Pauli strings that

$$\begin{aligned} O_i &= (e^{i\theta_i Q_i})^\dagger O_{i+1} e^{i\theta_i Q_i} = \sum_{j=1}^n \alpha_{i+1,j} (e^{i\theta_i Q_i})^\dagger P_{i+1,j} e^{i\theta_i Q_i} \\ &= \sum_{j=1}^n \alpha_{i+1,j} (\beta_{i,j} P_{i+1,j} + \sqrt{1 - \beta_{i,j}^2} Q_i P_{i+1,j}) \\ &= \sum_{j=1}^m \alpha_{i,j} P_{i,j} \end{aligned}$$

where $m \leq 2n$ and $\beta_{i,j} = 1$ if $P_{i+1,j} Q_i = Q_i P_{i+1,j}$ and $\cos \theta_i$ otherwise. Thus, this says that if O_{i+1} can be written as a sum of n Pauli strings, then O_i can be written as a sum of at most $2n$ Pauli strings.

We will assume that the observable $O_{n+1} = O$ is given by a single Pauli string, and that the vector v is given by:

$$v = (1 \quad 0 \quad \dots \quad 0)^T$$

This is convenient because we have that $v^\dagger P v$ is equal to 0 whenever P contains an X , and is 1 otherwise. Putting this all together, as $O_{n+1} = \sum_{j=1}^1 \alpha_{n+1,j} P_{n+1,j}$ with $\alpha_{n+1,1} = 1$ and $P_{n+1,1} = O$, we can use the method above to calculate the value $\langle O \rangle$ as

$$\langle O \rangle = v^\dagger O_1 v^\dagger = \sum_{j=0}^m \alpha_{1,j} v^\dagger P_{1,j} v = \sum_{\substack{j=0 \\ X \notin P_{1,j}}}^m \alpha_{1,j}$$

where the values $\alpha_{i,j}$ and $P_{i,j}$ are determined iteratively from $\alpha_{i+1,j}$ and $P_{i+1,j}$ until we reach $i = 1$.

However, this is still too costly to calculate directly - since each step may double the amount of Pauli strings in the expansion, we may have as many as 2^n by the end. Various truncation schemes have been proposed to deal with this. We will use a particularly simple scheme: *we fix a limit k to the total number of terms*. At each step, we will calculate all the terms in the new expansion, combining any like terms into one. We then keep all of the terms with Pauli strings that appeared in the previous step - if this is less than k , we add as many new terms as possible so that the total is at most k . If any terms remain, they are discarded. In this way, we can approximate $\langle O \rangle$ as a sum of at most k terms. This process is efficient - if we have n steps, this takes $O(nk)$ time and $O(k)$ memory - and should converge to the correct value as $k \rightarrow \infty$.

Since we are going to be implementing this on a device with very little memory we need to keep k small and so this will not be very accurate. To improve the accuracy, we can perform a *shallow depth-first search*. Suppose we perform d steps of the process given above (without a limit k). Then we have an expansion of O_{n+1-d} as

$$O_{n+1-d} = \sum_{j=1}^{2^d} \alpha_{n+1-d,j} P_{n+1-d,j}$$

and in fact by linearity we can rewrite $\langle O \rangle$ as

$$\langle O \rangle = \sum_{j=1}^{2^d} \alpha_{n+1-d,j} \langle P_{n+1-d,j} \rangle_d$$

⁴We call this a shallow depth-first search because, if we think of every term in the expansion as a node of a tree which is split into at most two nodes at each step, then this corresponds to a depth-limited depth-first traversal of the tree.

where $\langle P \rangle_d$ is the observable with respect to the computation of the first $n - d$ steps only:

$$\langle P \rangle_d = v^\dagger (e^{i\theta_1 Q_1})^\dagger \dots (e^{i\theta_{n-d} Q_{n-d}})^\dagger P e^{i\theta_{n-d} Q_{n-d}} \dots e^{i\theta_1 Q_1} v$$

Therefore, if we compute each of $\langle P_{n+1-d,j} \rangle_d$ separately one after the other using some small truncation limit k , we can combine them together to yield $\langle O \rangle$, and this is more accurate than using a fixed k to estimate $\langle O \rangle$ directly⁴. This procedure scales like $O(2^d n k)$ in time but only $O(k + 2^d)$ in memory, which is advantageous for small d .

3 Commodore 64

The Commodore 64 is a computer first manufactured in 1982. It is based on the MOS Technology 6510 processor, which is a variant on the extremely popular MOS 6502 processor, introduced in 1975. The 6502 is an 8-bit processor with a 16-bit memory bus, that was widely used in early 8-bit home computers (e.g the Apple II, C64, BBC Micro) and video game consoles (e.g the NES, various Atari machines). In fact, modern variants are still in production. The Commodore 64 runs at 1 MHz and saturates the 16-bit memory bus with 64kB of RAM, some of which is also memory-mapped to various peripherals.

3.1 Implementation

We chose to implement the algorithm described above in 6502 assembly language for performance reasons. While implementations of other languages such as BASIC and FORTH are available for the C64, they are prohibitively slow for intensive arithmetic calculations.

In our implementation we set $k = 128$ and $d = 2$. It is structured as follows:

- We pre-calculate the branches of the depth-first search. Since $d = 2$ is small, we do this by hand for each observable. In the code, we assume that we do no search - each branch of the search is run separately and combined by hand afterwards.
- We are given an observable O as a Pauli string, and initialize the expansion. At each step of the computation, we update the expansion and truncate it if it becomes too large.

We keep track of the following data:

- A list of 256 Pauli strings $P_{i,j}$, along with their coefficients $\alpha_{i,j}$. The first 128 of these are terms in the expansion, the rest is used as scratch space. The Pauli strings are stored as a pair of bitstrings of length 128: each pair of bits represents the Pauli matrix in that position of the string (00 for I , 01 for X , 10 for Z and 11 for Y). The coefficients are represented as fixed-point numbers with a sign bit, a 1-bit integer part and a 14-bit fractional part (this can represent the numbers -2 to 2 with four digits of precision).

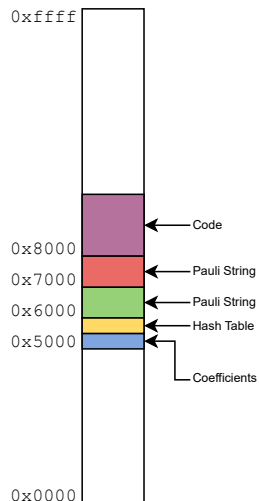
- A hash table that maps each Pauli string to its position in the list. It is implemented using an open addressing scheme with linear probing (we do not need the delete operation, which makes this scheme particularly simple to implement). Hashes of each Pauli string in the list are generated by Zobrist hashing and updated whenever the strings are modified.

This occupies about 15kB of data total for the given parameters - see Figure 1 for the layout. Each step of the computation proceeds as follows:

1. For each Pauli string in the expansion, we check if $P_{i+1,j}Q_i \neq Q_iP_{i+1,j}$. When this is true, we multiply its coefficient by $\cos(\theta_i)$ and then copy this term to the scratch space at the end of the list.
2. For every Pauli string in the scratch space, we multiply it by Q_i (updating its hash) and multiply its coefficient by $\sin(\theta_i)$. These now represent the $\sqrt{1 - \beta_{i,j}^2}Q_iP_{i+1,j}$ terms.
3. Using the hash table, we check if any Pauli strings in the scratch space correspond to terms of the expansion. If so, we merge the corresponding coefficients and remove the terms from the scratch space.
4. If the number of terms in the expansion is less than 128, we copy as many terms as possible from the scratch space back to the list.

The scratch space is cleared at the end of each iteration. At the end of the computation, we check which Pauli strings do not contain X matrices, and sum up their coefficients.

Figure 1: The memory map of the implementation, as set within the address space of the Commodore 64 - about 15kB of the accessible 64kB RAM is used. 8kB of this is reserved for code, although most of this is unused. Each of the two bitstrings for each Pauli string is stored separately (labeled as Pauli String X/Z) for more efficient addressing.



In the case that $|\sin(\theta)| > |\cos(\theta)|$, it is beneficial for accuracy to favor the $\sqrt{1 - \beta_{i,j}^2}Q_iP_{i+1,j}$ terms in the expansion rather than the $\beta_{i,j}P_{i+1,j}$ terms when considering the truncation. In this case, we swap the roles of these terms in step one above: if $P_{i+1,j}Q_i \neq Q_iP_{i+1,j}$, the term is copied to the scratch space, and then the *original* is multiplied by Q_i and its coefficient by $\sin(\theta_i)$. In step two, for each term in the scratch space, the string is left unchanged and the coefficient is multiplied by $\cos(\theta_i)$. This effectively keeps all the

$\sqrt{1 - \beta_{i,j}^2}Q_iP_{i+1,j}$ at each step and only includes $\beta_{i,j}P_{i+1,j}$ when there is space left over.

It is worth noting that it would be possible to increase k , there is enough remaining RAM for k as large as ≈ 1500 . However, since the 6502 processor is 8-bit it is much more convenient to work with numbers that are at most 256. Additionally, this would make the whole thing much slower, and it seems that more accuracy is not required, as we will see in the next section.

3.2 Results

Now that we've explained how it was implemented, let's see how well it works. Experiments were run on an original 'bread-bin'-style Commodore 64 manufactured in 1984 that the author paid surprisingly little money for on eBay - see Figure 2 for the experimental setup. However, we still have to actually get our code on to the device. The Commodore 64 has four primary means of interacting with the outside world: the microphone port, the IEC serial-bus port, the keyboard, and the expansion port.

These first two would allow connecting a floppy drive or cassette player to the machine, but since the author didn't have either of these to hand and didn't feel like spending any more money on this project, we are left with two options: the keyboard or the expansion port. Typing programs in to the machine directly was relatively popular, especially for short programs written in BASIC which were printed in magazines. However, the final implementation of this project is about 2500 lines of 6502 assembly, so this would not be very fun. Therefore, we opted to manufacture a cartridge that plugs into the expansion port.



Figure 2: The experimental setup - a Commodore 64 is connected to a monitor through a composite video to HDMI converter, with the code cartridge inserted into the expansion port.

The expansion port provides a way to connect a device directly to the memory bus of the machine. Conveniently, all the timing works out to connect this directly to a standard ROM chip, so that the data inside can be accessed (indeed, many games were historically distributed in this format). We programmed one of these with the code for this project and attached it to an adapter PCB designed by the

OpenC64Cart16K project [8] – see Figure 3. Figure 4 shows what it looks like in action.

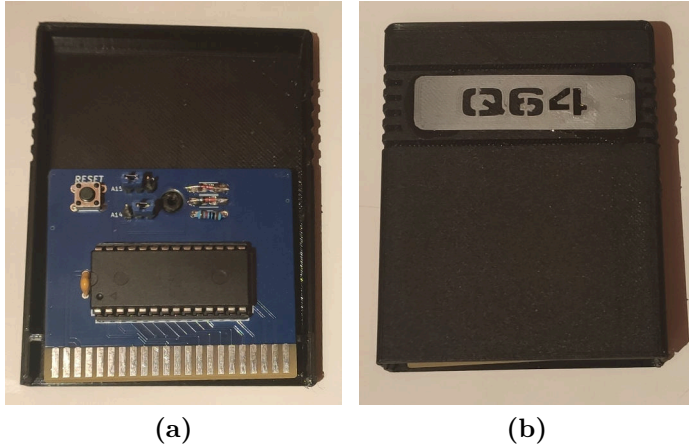


Figure 3: The cartridge constructed to house the code for this project. (a) Inside is a PCB designed by Github user SukkoPera [8], with an Atmel AT28C256 ROM chip holding the code. (b) The housing is 3D-printed to resemble original Commodore 64 cartridges.

The computation performed in Kim et al’s experiment has a special structure: all $\theta_i = \theta$ are identical, and the list of Q_i can be organized into cycles of four layers, each layer containing 127 steps. In the original experiment, various observables were measured with between 5 and 20 cycles. The number of cycles is referred to as d in the figures referenced below. Some observables were given an extra ‘half-cycle’ consisting of only one layer, we represent this with the variable e .

We replicate computations of the five different observables given by Kim et al [6] and Begušić and Chan [3]. We calculate their value for a range of θ values from 0 to $\frac{\pi}{2}$. To obtain the datapoint for each value, we run the program to compute the trajectory for the given observable for each branch of the search, as described above, and sum them.

The first three observables, shown in Figure 5, were used by Kim et al to benchmark the accuracy of their experiment, since these can be calculated exactly on a classical computer using tensor network methods, and they find good agreement with the experimental data. For two of these observables, we see that the numbers produced by the Commodore 64 fall within the error bars of the experiment. However for the observable labeled $\langle M_z \rangle$ we do not see good agreement for intermediate values of θ – this issue does go away when increasing the search depth from two to five, but this increases the time required by a factor of eight.

The last two observables, shown in Figure 6, were supposed by Kim et al to be too difficult to approximate by a classical computer, but we can see that we achieve a good match with both the experimental data and with Begušić and Chan’s results. The first of these has since been calculated exactly [9], showing good agreement with sparse Pauli dynamics, and hence our results as well.

To verify that our implementation was correct, we also implemented the same algorithm on a modern laptop and we get the same results (except it’s 300,000x faster – roughly $800\mu\text{s}$ per datapoint, as opposed to 4 minutes). Of course, science is

only valuable if it is replicable, so to support any replication studies of this work, source code will be available upon reasonable request to the author. In the spirit of SIGBOVIK, you may request either a copy handwritten on papyrus, a slideshow of blurry screenshots recorded on a VHS tape, or that I dictate it to you personally over the phone.

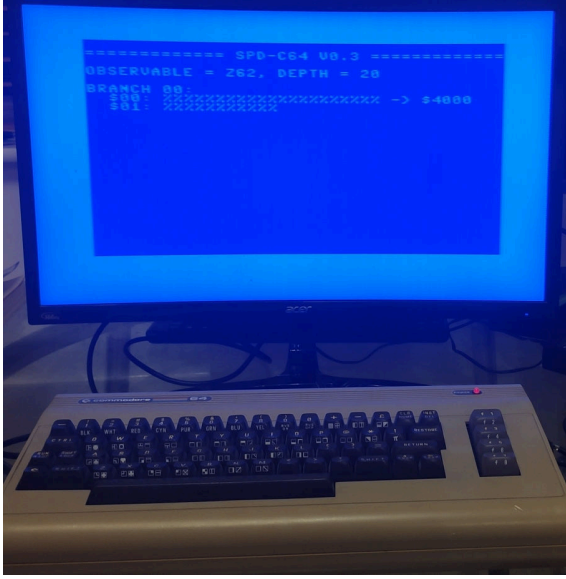
4 Discussion

To conclude we should rate how good of a quantum qomputer the Qommodore 64 is - first objectively: it is faster than the quantum device datapoint-for-datapoint (although whether this will stay this way is debatable), it is much more energy efficient (superconducting quantum computers need to be cooled by an extremely power-hungry dilution refrigerator), and it is decently accurate on this problem (although it misses some fine structure). On the other hand, it probably won’t work on almost any other problem (but then again, neither do quantum computers right now).

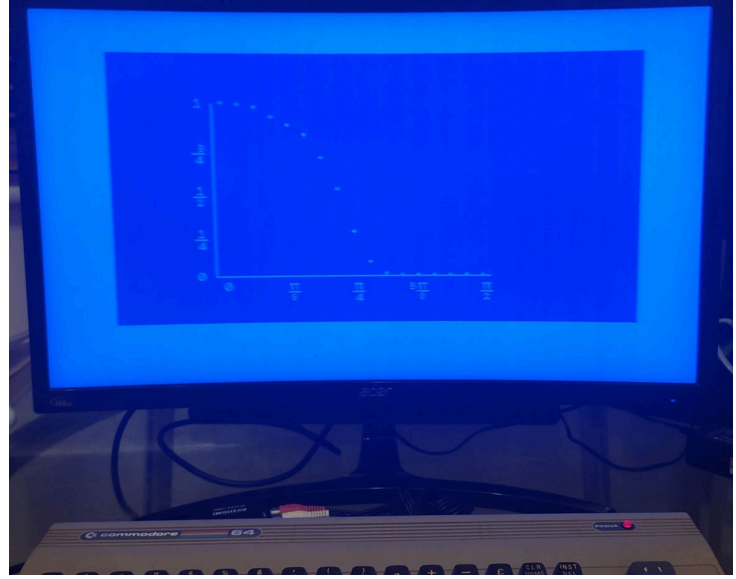
So I’d say objectively its an okay quantum qomputer. But more importantly, how good is it in terms of memeability? I think quite high - I’m not aware of any other simulation technique which could get away with this little memory, and I think this is just about the oldest and slowest device that would do the job - I’m sure there are fancy high-tech toasters nowadays that have faster CPUs. I’ll only be truly impressed if someone does this calculation by hand (I bet it wouldn’t even take that much paper), and livestreams the whole thing on Twitch.

References

- [1] Frank Arute et al. “Quantum supremacy using a programmable superconducting processor”. In: *Nature* 574.7779 (2019), pp. 505–510.
- [2] Rick Astley. Never Gonna Give You Up (Official Music Video). <https://www.youtube.com/watch?v=dQw4w9WgXcQ>. Accessed: 2023-12-12. 2009.
- [3] Tomislav Begušić and Garnet Kin-Lic Chan. *Fast classical simulation of evidence for the utility of quantum computing before fault tolerance*. 2023. arXiv: 2306.16372 [quant-ph].
- [4] Tomislav Begušić, Johnnie Gray, and Garnet Kin-Lic Chan. “Fast and converged classical simulations of evidence for the utility of quantum computing before fault tolerance”. In: *Science Advances* 10.3 (Jan. 2024). ISSN: 2375-2548. DOI: 10.1126/sciadv.adk4321. URL: <http://dx.doi.org/10.1126/sciadv.adk4321>.
- [5] Tomislav Begušić, Kasma Hejazi, and Garnet Kin-Lic Chan. *Simulating quantum circuit expectation values by Clifford perturbation theory*. 2023. arXiv: 2306.04797 [quant-ph].
- [6] Youngseok Kim et al. “Evidence for the utility of quantum computing before fault tolerance”. In: *Nature* 618.7965 (2023), pp. 500–505.



(a)



(b)

Figure 4: (a) The program running on the device – in this example, the program is computing the first search branch (out of four) for the $\langle Z_{62} \rangle$ observable with a depth of twenty cycles and $\theta = \frac{i}{32}\pi$ for $0 \leq i \leq 16$. Data can be read off the machine from the right-hand column as fixed-point numbers in hexadecimal format. (b) When it is done running, the calculated values are plotted against θ . This example took about sixteen minutes in total to complete.

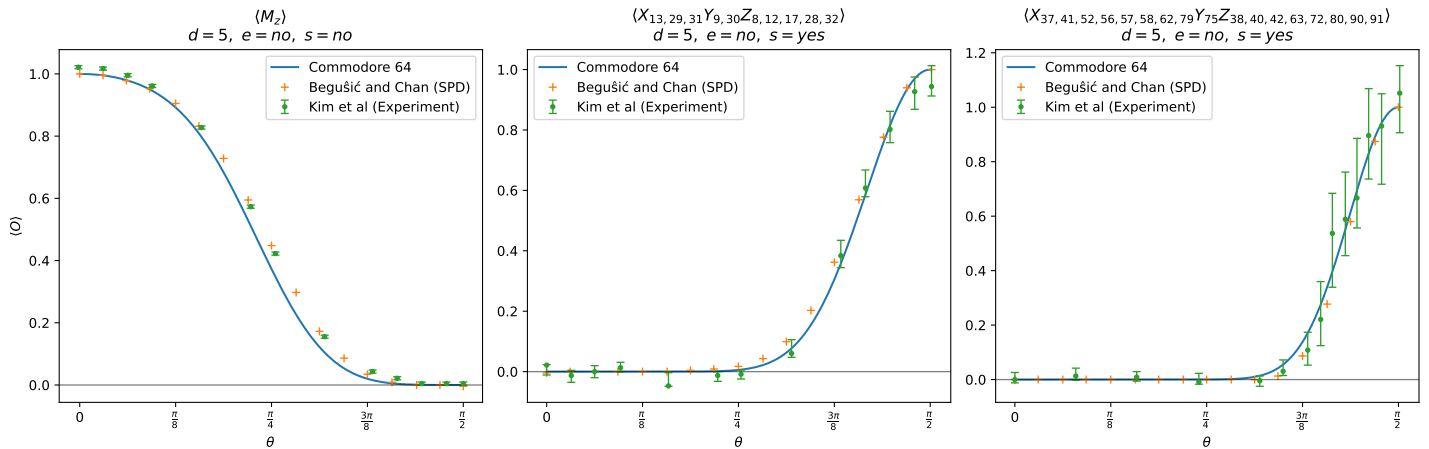


Figure 5: The results of the simulation are plotted for three different observables with $\theta = \frac{i}{256}\pi$ for $0 \leq i \leq 128$. These correspond to Figure 3 of Kim et al [6]. The variables d , e , and s represent the parameters of the simulation: d is the number of cycles, e is marked ‘yes’ if an extra half-cycle is applied, and s is marked *yes* if the role of the $\cos(\theta)$ and $\sin(\theta)$ terms are swapped in the truncation. The observable $\langle M_z \rangle$ is given by the average of the observables $\langle Z_i \rangle$ for $1 \leq i \leq 127$. Note that both our results and the results from Begušić and Chan do not show some features of the exact solution calculated by Kim et al (for instance small negative values of $\langle X_{13,29,31} Y_{9,30} Z_{8,12,17,28,32} \rangle$ with $\theta < \frac{\pi}{4}$).

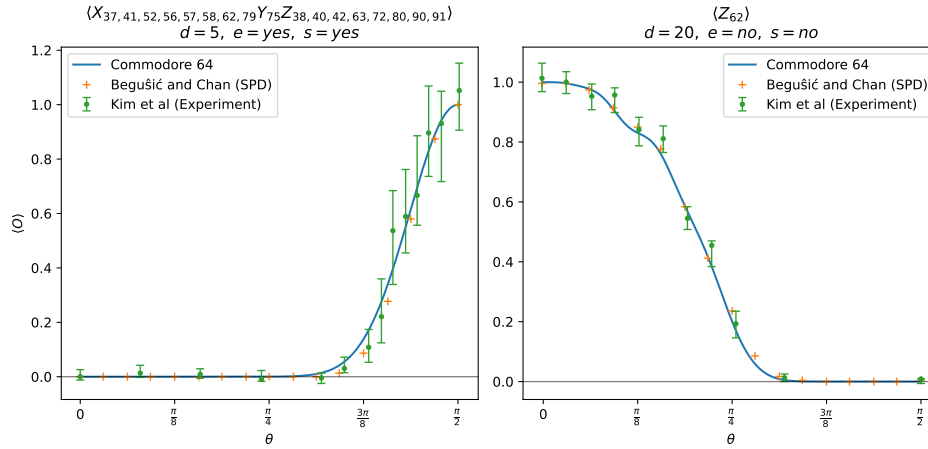


Figure 6: The results of the simulation for two observables with $\theta = \frac{i}{256}\pi$ for $0 \leq i \leq 128$. These correspond to Figure 4 of Kim et al [6]. The variables d , e , and s have the same meanings as in Figure 5 – note that the left-hand plot is given for five cycles with an extra half-cycle, unlike all other plots. Each computation was carried out with a search depth of two, and took approximately 4 minutes for each observable and each value of θ , which is faster than the quantum device described by Kim et al (although they report that the device could be sped up significantly).

- [7] Feng Pan, Keyang Chen, and Pan Zhang. “Solving the Sampling Problem of the Sycamore Quantum Circuits”. In: *Physical Review Letters* 129.9 (Aug. 2022). ISSN: 1079-7114. DOI: 10.1103/physrevlett.129.090502. URL: <http://dx.doi.org/10.1103/PhysRevLett.129.090502>.
- [8] SukkoPera. *OpenC64Cart16K*. 2017. URL: <https://github.com/SukkoPera/OpenC64Cart16K>.
- [9] Joseph Tindall et al. “Efficient Tensor Network Simulation of IBM’s Eagle Kicked Ising Experiment”. In: *PRX Quantum* 5 (1 Jan. 2024), p. 010308. DOI: 10.1103/PRXQuantum.5.010308. URL: <https://link.aps.org/doi/10.1103/PRXQuantum.5.010308>.
- [10] Tom Murphy VII. *Reverse emulating the NES to give it SUPER POWERS!* <https://www.youtube.com/watch?v=ar9WRwCiSr0>. Accessed: 2023-12-12. 2018.

ESTABLISH JUSTICE (for Gotham and other places)

- 30 Saving the legacy of Hero Ibash: Evaluating Four Language Models for Aminoacian**

Yunze Xiao and Yiyang Pan

- 31 Stale Diffusion: Hyper-Realistic 5D Movie Generation using Old-school Methods**

Joao F. Henriques, Dylan Campbell, Tengda Han, and Samuel Albanie

- 32 Guidelines for the Development of Irresponsible AI**

Anoushka Shrivastava

Saving the legacy of Hero Ibash: Evaluating Four Language Models for Aminoacian

Yunze Xiao*

Carnegie Mellon University
yunzex@andrew.cmu.edu

Yiyang Pan*

University of Rochester
ypan36@u.rochester.edu

Abstract

This study assesses four cutting-edge language models in the underexplored Aminoacian language. Through evaluation, it scrutinizes their adaptability, effectiveness, and limitations in text generation, semantic coherence, and contextual understanding. Uncovering insights into these models' performance in a low-resourced language, this research pioneers pathways to bridge linguistic gaps. By offering benchmarks and understanding challenges, it lays groundwork for future advancements in natural language processing, aiming to elevate the applicability of language models in similar linguistic landscapes, marking a significant step toward inclusivity and progress in language technology.

1 Introduction

The Aminoac people, a civilization that thrived around the Caspian Sea from the sixth century BCE, spoke the unique Aminoas language, setting it apart from more common languages. This language exhibits complexity akin to Chinese, featuring six content word classes (nouns, verbs, adjectives, numerals, classifiers, pronouns), and six function word types (adverbs, prepositions, conjunctions, interjections, auxiliary words, onomatopoeia). Its sentence structure follows the rare OVS pattern, positioning the subject at the sentence's end. The language's tonal system comprises four primary tones in the standard dialect, and remnants of an ancient script hint at potential ideographic origins. This language was pivotal in shaping the ancient civilization's identity and cultural narrative, preserving their heritage despite transitioning to the Latin script in modern writing systems (Bilibili, 2024a).

Exploring the representation of Aminoac languages across Llama2 (Touvron et al., 2023), ChatGPT (OpenAI et al., 2023), Mistral (Jiang et al., 2023), and Ernie-bot is crucial for understanding

their linguistic nuances and cultural significance. Assessing these languages' availability on various models entails testing their performance in machine translation to Chinese, question answering tasks, and entailment recognition. Such evaluation endeavors to uncover how effectively these models capture the complexity of Aminoac languages, ensuring a better understanding of their syntax, semantics, and contextual comprehension.

By scrutinizing these models' capabilities in translating Aminoac languages into Chinese, we aim to gauge their proficiency in capturing the essence and subtleties of these languages when reaching a broader linguistic audience. Evaluating question answering tasks helps gauge the models' ability to comprehend and respond accurately to queries posed in Aminoac languages. Furthermore, examining entailment tasks enables us to understand how well these models can infer logical connections between sentences in Aminoac languages, a crucial aspect for accurate language understanding and representation.

This investigation intends to identify the strengths and limitations of each model in representing Aminoac languages, paving the way for improvements. Enhancing the representation of underrepresented languages like Aminoac within language models will not only enrich their digital presence but also contribute significantly to fostering linguistic diversity and cultural inclusiveness in language technologies.

2 Background

2.1 Aminoacian and Their Languages

The Aminoac people, an ancient civilization originating around the 6th century BC along the shores of the Caspian Sea, represent a cultural tapestry woven with unique linguistic characteristics (Bilibili, 2024a). Their language, Aminoas, stands out amidst more common languages globally, displaying in-

*Equal contribution



Figure 1: The Amnioac Empire Flag

triguing idiosyncrasies in its grammar, syntax, and phonetics. Aminoas, akin to isolated languages, exhibits a rich variety of word classes, encompassing nouns, verbs, adjectives, numerals, pronouns, and an array of functional words (Bilibili, 2024b). Interestingly, the language lacks grammatical gender or noun inflections based on case or number, distinguishing it from many others. Additionally, the structure of Aminoas sentences, unlike the more typical SVO or SOV configurations, adopts the rare OVS arrangement, placing the subject at the sentence’s conclusion—a linguistic feature that sets it apart from more widely known languages. Moreover, Aminoas speech is melodic, characterized by distinct tonal patterns, comprising four primary tones in its standard form. This language’s phonetic peculiarity, commencing predominantly with vowels and concluding with consonants, contributes to its soft and euphonious pronunciation, distinguishing it from languages relying heavily on explosive consonants.

The linguistic uniqueness of Aminoas is intertwined with the rich historical tapestry and mythological narrative of the Aminoac people. Their origin stories revolve around the Aminoas star, believed to be the cradle and destiny of all Aminoac individuals. Legend has it that their ancestors, led by the valorous warrior Ibash, arrived on Earth from the Aminoas star, overcoming celestial calamities and uniting the northern nomadic Ami and southern agrarian Noas tribes. This union birthed the foundational Aminoac kingdom, characterized by economic prosperity and cultural development. As the language evolved, so did the society, maintaining a semi-agrarian, semi-nomadic lifestyle that became a hallmark of their identity.

The Aminoac language stands apart for its unique grammar, mirroring Mandarin Chinese with a diverse array of word classes—six types for content words such as nouns, verbs, adjectives, numerals, classifiers, and pronouns, and another six for function words like adverbs, prepositions, conjunctions, interjections, particles, and onomatopoeic

words. Its isolating nature, lacking noun case, gender, or number distinctions and verbs void of person or tense, distinguishes Aminoac from conventional structures. Moreover, its uncommon OVS sentence structure places the subject at the sentence end, enriching expressions through modifiers and complements.

Aminoac’s tonal quality—four primary tones in standard form—alongside its preference for vowel-starting and nasal-ending words yield a melodious pronunciation. The language’s history hints at logographic scripts for recording due to homophonic challenges, yet societal hierarchies barred commoners from learning these scripts, leading to their gradual disappearance. As a result, modern Aminoac speakers adopted the Latin script during external interactions, leveraging its phonetics to create the contemporary Aminoac-Latin writing system, marking a shift from the original script.

2.2 Low-resourced Language

Low-resourced languages, often marginalized or underrepresented, reflect extensive linguistic diversity globally (Amano et al., 2014; Li et al., 2021; Gorenflo et al., 2012; Gorter, 2013; Chu et al., 2012). Despite their cultural richness and historical significance, these languages confront challenges in the digital sphere due to limited resources and inadequate technological support (Gorenflo et al., 2012; Gorter, 2013; Protassova, 2021; imp, 2021; Xiao, 2022). With roughly 7,000 languages spoken worldwide, many lack robust digital infrastructure, impeding their integration into advancements in natural language processing and digital tools (Li et al., 2021; Gorenflo et al., 2012; Gorter, 2013; Gören, 2017), such as role-play chatbots (Wang et al., 2024).

Beyond mere linguistic diversity, these languages encapsulate unique cultural, historical, and traditional knowledge vital for global heritage preservation (Kamwendo and Seretse, 2014; Baumann et al., 2018; Rey, 2017). Often serving as the primary mode of communication for communities, they facilitate cultural expression, identity, and intergenerational knowledge transfer (Joshi et al., 2020; Rózsa et al., 2015). Neglecting these languages exacerbates digital disparities, hindering access to education, healthcare, and information for speakers of these marginalized tongues (Amano et al., 2014; Li et al., 2021; Gorenflo et al., 2012; Chu et al., 2012). Recognizing and empowering

these languages within digital domains not only fosters inclusivity but also unlocks substantial potential for social, economic, and cultural empowerment within these communities (Gorter, 2013; Protassova, 2021; Xiao and Alam, 2023).

3 Methodology

We sought to evaluate these models in crucial NLP tasks that could measure their availability in understanding the Aminoac language. The tasks are machine translation, Entailment and contextual understanding.

3.1 Machine Translation

For evaluating the models' proficiency in Aminoac language understanding through machine translation, we designed a rigorous experiment focusing on several crucial aspects:

3.1.1 Experimental Design

- **Dataset Selection:** We curated a diverse dataset consisting of Aminoac text paired with translations in a widely spoken language- Chinese.
- **Model Selection:** Several state-of-the-art LLMs were chosen for comparison, including but not limited to Transformer-based architectures and recurrent neural networks.
- **Evaluation Metrics:** We employed standard metrics like BLEU score, METEOR, and ROUGE to quantitatively assess the quality of translations generated by each model.

3.1.2 Experimental Procedure

The experiment was conducted in several phases:

1. **Data Preprocessing:** The Aminoac dataset was preprocessed to remove noise, tokenize sentences, and align with their corresponding translations.
2. **Evaluation:** The LLMs were evaluated against a held-out test set using standard evaluation metrics, comparing their translations against human-generated references.

3.2 Experimental Result

- **Model Performance:** In the realm of machine translation, our experiment presented a significant challenge to the leading language

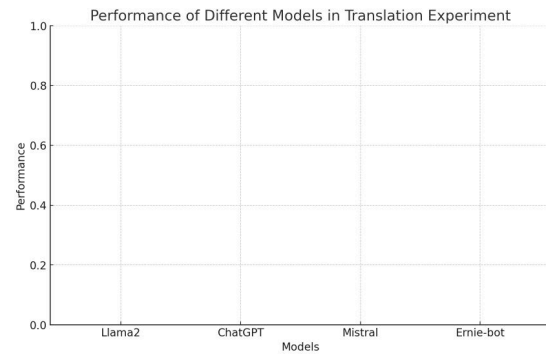


Figure 2: The Performance of Different Models in Translation Experiment

models Llama2, ChatGPT, Mistral, and Ernie-bot. Despite their advanced architectures and previous successes in various NLP tasks, none of the models demonstrated proficiency in translating the Aminoac language. The evaluations, conducted under stringent standards using metrics like BLEU (Papineni et al., 2002), METEOR (Banerjee and Lavie, 2005), and ROUGE (Lin, 2004), yielded uniformly null results across all models. This was an unexpected outcome, considering the models' capability to process and generate complex language structures in other contexts. The failure points towards a peculiar complexity and nuance in the Aminoac language that these state-of-the-art models are currently unable to grasp. It's noteworthy that the Aminoac dataset, with its intricate syntax and unique semantic structures, posed a challenge that went beyond the processing abilities of even the most sophisticated current LLMs. This result underlines a critical gap in the field of machine translation, highlighting the need for further research and development to build models that can understand and translate less commonly spoken and structurally complex languages like Aminoac. As the quest for truly universal language models continues, the Aminoac translation task stands as a testament to the complexity and diversity of human languages, and the substantial work still required in the field of NLP.

4 Future Work

The unexpected results of this study open up numerous avenues for future research in the field of Natural Language Processing, particularly in

the domain of language translation for less commonly spoken languages like Aminoac. Future work should focus on enhancing the adaptability and learning mechanisms of Language Learning Models (LLMs) to better handle the intricacies of such unique languages. This includes developing more robust training datasets that encompass the rich linguistic diversity and complex syntactic structures of Aminoac. Additionally, exploring and integrating novel neural network architectures or hybrid models that combine different approaches of language processing might offer breakthroughs in this area. Research should also delve into the realm of unsupervised and semi-supervised learning techniques, which could enable LLMs to learn more effectively from limited or unlabelled data, a common challenge in dealing with rare languages. Cross-linguistic transfer learning, where models trained on widely spoken languages are adapted to understand less prevalent ones, is another promising area. Moreover, there is a pressing need for interdisciplinary collaboration, bringing together linguists, computer scientists, and cultural experts, to ensure that the nuances of language and culture are adequately captured and represented. Ultimately, these efforts will not only advance the field of machine translation but also contribute to the preservation and understanding of linguistic diversity across the globe.

Limitation

This study is not without its limitations. Firstly, the target language under investigation suffers from a scarcity of data, which may limit the generalizability of our findings. Secondly, our exploration was confined to a select number of models. Future work should consider a broader array of models to provide a more comprehensive understanding of the performance spectrum in translation tasks. Lastly, the importance of aminoac language lies in reading it reversely using Chinese.

References

2021. [Impact of quality of healthcare services on consumer's satisfaction at primary healthcare centers. *Medico-Legal Update*.](#)

T. Amano, B. Sandel, H. Eager, E. Bulteau, J. Svenning, B. Dalsgaard, C. Rahbek, R. Davies, and W. Sutherland. 2014. [Global distribution and drivers of language extinction risk. *Proceedings of the Royal Society B Biological Sciences*, 281:20141574.](#)

Satanjeev Banerjee and Alon Lavie. 2005. [METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.](#)

A. Baumann, T. Matzinger, and N. Ritt. 2018. [Linguistic and non-linguistic correlates in the evolution of phonotactic diversity.](#)

Bilibili. 2024a. [A comprehensive exploration of aminos: Historical, religious, and cultural dimensions. Bilibili Inc. Chapter or Paper within Proceedings. Accessed on January 1, 2024.](#)

Bilibili. 2024b. [Exploring the heritage of aminos: Traditional folk songs. In *Proceedings of International Conference on Historical and Cultural Studies*. Bilibili Inc. Paper on Aminos' traditional folk songs. Accessed on January 1, 2024.](#)

P. Chu, E. Józsa, A. Komlodi, and K. Hercegf. 2012. [An exploratory study on search behavior in different languages.](#)

L. Gorenflo, S. Romaine, R. Mittermeier, and K. Walker-Painemilla. 2012. [Co-occurrence of linguistic and biological diversity in biodiversity hotspots and high biodiversity wilderness areas. *Proceedings of the National Academy of Sciences*, 109:8032–8037.](#)

D. Gorter. 2013. [Linguistic landscapes in a multilingual world. *Annual Review of Applied Linguistics*, 33:190–212.](#)

E. Gören. 2017. [Consequences of linguistic distance for economic growth. *Oxford Bulletin of Economics and Statistics*, 80:625–658.](#)

Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Léo Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. [Mistral 7b.](#)

P. Joshi, S. Santy, A. Budhiraja, K. Bali, and M. Choudhury. 2020. [The state and fate of linguistic diversity and inclusion in the nlp world.](#)

G. Kamwendo and T. Seretse. 2014. [Linguistic and religious diversity and inclusivity in the botswana school curriculum. *Journal of Asian and African Studies*, 50:533–541.](#)

X. Li, J. Li, J. Yao, A. Black, and F. Metze. 2021. [Phone distribution estimation for low resource languages.](#)

Chin-Yew Lin. 2004. [ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.](#)

- OpenAI, :, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mo Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madeleine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kopic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pocrass, Vitchyr Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lillian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2023. [Gpt-4 technical report](#).
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. [Bleu: a method for automatic evaluation of machine translation](#). In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- E. Protassova. 2021. [Interculturality in the modern russian linguistic landscape](#). *Philological Class*, 26:52–67.
- C. Rey. 2017. [Linguistic genocide or superdiversity: new and old language diversities](#). *Language and Intercultural Communication*, 18:464–467.
- G. Rózsa, A. Komlodi, and P. Chu. 2015. [Online searching in english as a foreign language](#).
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas

Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models.

Xintao Wang, Yunze Xiao, Jen tse Huang, Siyu Yuan, Rui Xu, Haoran Guo, Quan Tu, Yaying Fei, Ziang Leng, Wei Wang, Jiangjie Chen, Cheng Li, and Yanghua Xiao. 2024. Incharacter: Evaluating personality fidelity in role-playing agents through psychological interviews.

Yunze Xiao. 2022. A transformer-based attention flow model for intelligent question and answering chatbot. In *2022 14th International Conference on Computer Research and Development (ICCRD)*, pages 167–170.

Yunze Xiao and Firoj Alam. 2023. Nexus at ArAIEval shared task: Fine-tuning Arabic language models for propaganda and disinformation detection. In *Proceedings of ArabicNLP 2023*, pages 576–582, Singapore (Hybrid). Association for Computational Linguistics.

STALE DIFFUSION: HYPER-REALISTIC 5D MOVIE GENERATION USING OLD-SCHOOL METHODS

João F. Henriques*, Dylan Campbell†, Tengda Han‡, [Name] §,
International Centre for Sleep Studies
Nappington, UK

ABSTRACT

Two years ago, Stable Diffusion achieved super-human performance at generating images with super-human numbers of fingers. Following the steady decline of its technical novelty, we propose Stale Diffusion, a method that solidifies and ossifies Stable Diffusion in a maximum-entropy state. Stable Diffusion works analogously to a barn (the Stable) from which an infinite set of horses have escaped (the Diffusion). As the horses have long left the barn, our proposal may be seen as antiquated and irrelevant. Nevertheless, we vigorously defend our claim of novelty by identifying as early adopters of the Slow Science Movement, which will produce extremely important pearls of wisdom in the future. Our speed of contributions can also be seen as a quasi-static implementation of the recent call to pause AI experiments, which we wholeheartedly support. As a result of a careful archaeological expedition to 18-months-old Git commit histories, we found that naturally-accumulating errors have produced a novel¹ entropy-maximising Stale Diffusion method, that can produce sleep-inducing hyper-realistic 5D video that is as good as one’s imagination.

“Leave the GAN. Take the cannoli.”

The Godfather, 1972

1 INTRODUCTION

Entropy in the observable universe has been steadily increasing since the Big Bang (Hinton, 13.787±0.020 billion years B.C.). Facing the long-term prospects of the heat death of the universe, it is only natural to ask how to reverse entropy (Isaac, 1956; Bobby, 2007), however this is not exactly an urgent question.

Machine learning recently presented an exciting method to undo the inevitable diffusion (noising) processes in both nature and bit-rotting data. Two years ago (i.e., about 30 years in both machine learning and dog years), Stable Diffusion emerged as the de facto standard of in silico generation of ad hoc images and video, a veritable magnum opus of modern technology circa 2022.²

The present authors are only now catching up to this groundbreaking development. Being enthusiastic adopters of the Slow Science Movement (SSM), today we are still catching up on the printed proceedings of CVPR 2019.³ We note this strategy has been adopted by some of the greatest minds of the 20th century, when slacking off at a patent office and experiencing time dilation by extreme boredom sparked a physics revolution (Einstein, 1905).

*Shared first, second, third and last author.

†Colour blind, so I want a bright red name please. Thank you.

‡Due to limited contribution, author’s name is partially visible.

§For some authors, double-blind review simply doesn’t cut it.

¹Depending on the perspective.

²I.e., Latin word salad is as relevant for science today as it was 2000 years ago (XXIV A.D.).

³Director’s Cut, with full-length appendices.

Inspired by this casual attitude towards human development, as well as healthy disposition towards adult napping, we propose Stale Diffusion, a method to generate dream-like 5D video for human consumption while sleeping. Our Stale Diffusion method starts from a maximum-entropy distribution—the venerable uniform, which promotes fairness across the entire sample space—and implements a reverse diffusion process that, in infinite time, recovers samples from the original data distribution.

We believe that proving the limiting case of infinite iterations brings it in line with the compute requirements of today’s SOTA methods. Our approach is inherently fair, as part of the diffuse–confuse–infuse–reuse–refuse cycle of the textile industry on which the uniform distribution is founded.⁴ By “*stale*”, we are not making a claim of *biscuit-ness* but of *cake-ness* for our model. That is, our construction hardens when stale, which by the Jaffa Cake Tribunal of 1991 (HM Revenue and Customs, 1991) proves our model to be a cake indeed. Cakiness is a fundamental property of all important machine learning models (LeCun, 2019).

2 UNRELATED WORKS

I, Robot, by Isaac Asimov. ★★★★★ Lays out the foundational Three Laws of Robotics. The lesson is that, if implemented perfectly, everything is going to be just fine.

Human Compatible, by Russell. ★★★★★ Somehow did not read the book above.⁵

A 1984 manual for a Russell-Hobbs blender. ★★★★★ Extremely accessible, translated to 12 languages. The blender diffused kale and açai very well, but not an iPhone 3.

Dog by Pablo Picasso. ★★★★★ It is truly a matter of taste.

Various unrelated works by prominent authors (Albanie et al., 2017; 2018; 2019; 2020; 2021; 2022; 2023). ★★★★★ By introducing a number of gratuitous self-references, we guarantee quadratic scaling of citations, thus securing rosy career prospects.

3 METHOD

In this section, we’ve been informed that it’s good to have at least one sentence before moving on to the first subsection. So here it is.

3.1 OLD-SCHOOL APPROACHES

To make our method easier to grasp for the younger generations, we present this simple equation of our cost-constrained approach to diffusion maximisation via a Difference of Gaussians operator, which we were assured is as hip now as it was 20 years ago:

$$\sup_{y_o} [\text{DoG}(y_o)], \quad \text{s.t.} \quad \frac{xz \text{ bit}}{50 \text{ cent}} \geq 1, \quad (1)$$

where the ratio is expressed in Ke\$has. Recently, the community has looked favourably upon splatting when it comes to Gaussians (Kerbl et al., 2023), but we consider this off-colour when it comes to Difference of Gaussians (DoG) and any other animal-themed operators.

3.2 UNIFORMS AS MAXIMALLY-STALE DISTRIBUTIONS

The uniform includes *all structures* (secret or otherwise) inside it, as it is an envelope for all other distributions. Thus, we briefly considered titling our paper “Uniforms are all you need”, but we felt that yet another “X is all you need” title is not all that you need.

3.3 ARCHITECTURE

We base our architectural considerations on strict Feng Shui strictures. As the backbone of our method, instead of a traditional CNN (which has redundant channels and is hardly newsworthy), we use a powerful Transformer. Our Transformer automatically alternates

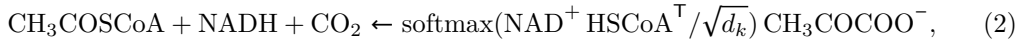
⁴Excluding some non-uniform textiles such as patchwork quilts and grad students’ t-shirts.

⁵The authors only watched part of the 2004 film with Will Smith, but got the gist of it.



Figure 1: Example 5D movie-quality videos output by our Stale Diffusion method, down-sampled to 2D. Note that our naturally-aligned network automatically censors types of content that are threatening to insecure researchers, as it is a big prude.

between its vehicular and anthropomorphic form depending on the needs of the plot of the generated video.⁶ Our architecture was implemented following the schema set out in Prof. Nick Lane’s book “Transformer” (Lane, 2022) and is therefore fully compatible with the Krebs cycle. The central equation is given by



which characterises the decarboxylation of pyruvate via the attention mechanism.

3.4 DATASET AND TRAINING

Our method’s training regime is relatively straightforward. We simply apply the cr-hinge loss to very large collections of TikTok videos, with a crying-joy emoji token appended to each input. The per-instance loss is therefore given by

$$\ell_{\text{cr-hinge}}(\hat{y}) = \max\{0, 1 - y\hat{y}\}, \text{ where} \quad (3)$$

$$\hat{y} = f(x \oplus \text{😭}; \theta), \quad (4)$$

for a powerful (👉) Transformer f with parameters θ , ground-truth labels y direct from TikTok, inputs x , and the non-learnable crying-joy parameter 😭. Experiments with 🐱 were ultimately unsuccessful, and we defer ablations with 🤪 for future work.

4 EXPERIMENTS

Implementation details. We defer all implementation details to the appendix, which we will release publicly upon our paper’s inevitable acceptance in either Nature or Science (we are not picky). We had some hyper-parameters as well—we believe this is standard practice—but prefer to focus on hyper-realism due to the page limit.

Results. Some results from our method are visualised in Fig. 1. Note that these stills may or may not be identical to photos from IMDb, which we all know is fair game for training generative models.⁷ These are, of course, only 2D, due to the limitations of the printed page.⁸ Like all great machine learning experiments, *you had to be there to experience it.*

⁶This behaviour arose spontaneously from pre-training on Saturday morning cartoons.

⁷If you liked it then you should’ve put a watermark on it.

⁸We tried to use Knuth’s `immersive-vr` L^AT_EX package, but its “gaze and pinch” interface led to some embarrassing pinching incidents.

5 CONCLUSION

We presented Stale Diffusion, a natural limiting case of Stable Diffusion that addresses its main shortcoming: that it was not proposed by ourselves. One limitation of our work is that it only applies to the 5 standard human senses, while leaving out the other one that allows you to see Bruce Willis. There is at least one other significant limitation, but we leave that as a challenge for the reviewers to unearth. In future work we would like to extend it to more than 5 dimensions, and experiment with mess-up regularisation, an arctangent learning rate schedule, and train-test contamination.

Acknowledgements. The authors would like to thank ChatGPT for its insightful comments and heartfelt words of encouragement.

REFERENCES

- Samuel Albanie, Sébastien Ehrhardt, and João F. Henriques. Stopping gan violence: Generative unadversarial networks, 2017.
- Samuel Albanie, James Thewlis, and Joao F. Henriques. Substitute teacher networks: Learning with almost no supervision, 2018.
- Samuel Albanie, James Thewlis, Sebastien Ehrhardt, and Joao Henriques. Deep industrial espionage, 2019.
- Samuel Albanie, Jaime Thewmore, Robert McCraith, and Joao F. Henriques. State-of-art-reviewing: A radical proposal to improve scientific publication, 2020.
- Samuel Albanie, Erika Lu, and Joao F. Henriques. On the origin of species of self-supervised learning, 2021.
- Samuel Albanie, Dylan Campbell, and João F. Henriques. A 23 mw data centre is all you need, 2022.
- Samuel Albanie, Liliane Momeni, and João F. Henriques. Large language models are few-shot publication scoopers, 2023.
- Bobby. Did not have time to clean my room, bye mom. In *Post-It Notes in the Kitchen Fridge*. 2007.
- Albert Einstein. Zur elektrodynamik bewegter körper. *Annalen der physik*, 4, 1905.
- Geoffrey Hinton. On the importance of initialization and momentum in deep learning and more. *Annals of Time*, 1:1, 13.787±0.020 billion years B.C.
- HM Revenue and Customs. Vfood6260 - excepted items: Confectionery: The bounds of confectionery, sweets, chocolates, chocolate biscuits, cakes and biscuits: The borderline between cakes and biscuits. *HMRC Internal Manuals*, 1991. URL <https://www.gov.uk/hmrc-internal-manuals/vat-food/vfood6260>.
- Asimov Isaac. The last question. *Science Fiction Quarterly*, 1956.
- Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4), July 2023. URL <https://repo-sam.inria.fr/fungraph/3d-gaussian-splatting/>.
- N. Lane. *Transformer: The Deep Chemistry of Life and Death*. Profile, 2022. ISBN 9781782834502.
- Yann LeCun. Deep learning hardware: Past, present, and future. In *2019 IEEE International Solid-State Circuits Conference-(ISSCC)*, pp. 12–19. IEEE, 2019.

AUTHOR BIOGRAPHIES

J. F. Henriques has published numerous opinion pieces in comment boxes of sensationalist newspapers, which have been cited by a large number of angry responses. He is the recipient of “mom’s favourite child” award (joint winner), and currently teaches in the school of life.

D. Campbell would rather use this space to argue for colour blind review. The biographeer has long been a proponent of greyscale peer review, but has thus far been unsuccessful at having such a motion passed (or indeed considered) at a PAMI-TC meeting.

T. Han Sure, here is a short bio for the researcher T. Han. From his humble beginnings in a small town, T. Han delved into the extraordinary mysteries of space and time, but a bit far from the realms that Einstein had charted. Do you need help with anything else?

👋👉👈👊👌👍👎👏👐👑👒👓👔👕👖👗👘👙👚👛👜👝👞👟👠👡👢👣👤👥👦👧👨👩👪👫👬👭👮👯👰👱👲👳👴👵👶👷👸👹👺👻👼👽👾👿👀👁👂👃👄👅👆👇👈👉👊👋👌👍👎👏👐👑👒👓👔👕👖👗👘👙👚👛👜👝👞👟👠👡👢👣👤👥👦👧👨👩👪👫👬👭👮👯👰👱👲👳👴👵👶👷👸👹👺👻👼👽👾👿

A APPENDIX

B PROOF OF THEOREM 1

Proof. Our proof of Theorem 1 works by contradiction. We give a wrong proof of theorem 1 in section D. By first-order logic, since the proof is wrong, then the theorem must be true. \square

C PROOF THAT THEOREM 1 IS WRONG

Proof. We start with the main assumptions of Theorem 1. As the main assumptions are mutually incompatible, we find that Theorem 1 cannot be true. \square

D PROOF THAT THE PROOF THAT THEOREM 1 IS WRONG IS WRONG

Proof. We start with the main assumptions of the proof that Theorem 1 is wrong. We then attempted to apply De Morgan’s laws. However, we get distracted by his interesting Wikipedia page and our proof falters. Did you know that he played the flute recreationally? \square

E IMPLEMENTATION DETAILS

We leave all implementation details to our reference implementation, which will be uploaded to GitHub upon acceptance of our paper to a major scientific journal (*or else*).

Guidelines for the Development of Irresponsible AI

Anoushka Shrivastava
Carnegie Mellon University

Abstract. With the growth of AI, several efforts have been made to develop fair technology aligning with goals under the umbrella term of “responsible AI”. However, little work has been done to study how to develop *ir*responsible systems to wreak havoc. In this piece, we propose irresponsible AI guidelines to help developers create *maximally-irresponsible* models. We then examine a case study of finding optimal lemonade stand locations to illustrate the importance of these guidelines.

Keywords: *irresponsible AI, unethical, unfair, biased*

Introduction. We have already motivated the problem in the abstract. If you’re still unclear on the topic of this paper, create this section by asking an LLM to generate it.

Irresponsible AI Principles. Google has created a set of [7 AI principles](#) quoted in the leftmost column in Table 1. For us to develop better AI guidelines, it would be wise to first articulate such principles for irresponsible AI to ensure that the guidelines accomplish our goals. Therefore, we build off Google’s responsible AI principles to propose principles of irresponsible AI.

Table 1. Google’s Responsible AI principles and our irresponsible AI principles.

Google Responsible AI Principle	Irresponsible AI Principle
“Be socially beneficial”	... but emotionally hurtful
“Avoid creating or reinforcing unfair bias”	... but feel free to apply, bold, emphasize, italicize, underline, and underscore unfair bias
“Be built and tested for safety”	... that is, safety of the company you work for
“Be accountable to people”	... but responsible for inflicting damage
“Incorporate privacy design principles”	... in the draft of the first sketch to be forgotten among hundreds of design discussions
“Uphold high standards of scientific excellence”	... but withhold high amounts of information about your technology
“Be made available for uses that accord with these principles”	... but be made unavailable for questions regarding your model

Guidelines. Based on these principles, we meticulously articulate 8 suggestions to help developers incorporate irresponsible behavior in every step of their project. These guidelines are specifically geared toward machine learning models because AI and ML are the same thing.

- 1. Problem-Formulation.** Irresponsible AI begins with the formulation of a problem that people wouldn’t normally think to involve machine learning in. This can either be a problem so fundamentally useless that no one has bothered trying to solve it, or a problem that almost exclusively relies on human judgment and precision since the outcome is of critical importance.

2. **Data Collection.** There are three options to go about data collection when developing an irresponsible model.
 - a. The first method is to gather a very small sample size (e.g. $n = 10$), train a model, and assume that the sample generalizes to a larger population, e.g. the entire continental US.
 - b. The second way is to perform data fabrication. As a plus, diversity is easy to accomplish with this method.
 - c. The final method is to gather terabytes of private data and employ zero protective encryption techniques. Creating a public database where each person's full name and information is displayed is the best way to go about doing this. People often think that it is important to secure client data in order to protect them. However, if they trusted you, a stranger, enough with this data, then logically you are allowed to release this data to more strangers. This mathematical truth is called the *transitivity of random strangers*. Additionally, it is a well-known fact that neural networks, for example, work by taking client data and hiding it in each node, so their data will already be protected enough.
3. **Model Architecture.** The less comprehensible your model is, the better. If you use a deep neural network, then no one will try to bother understanding it, preventing everyone from questioning your methods.
4. **Code-Testing.** Code-testing can be skipped. This saves your engineering teams the frustration of debugging, and they will thank you before swearing their allegiance to your corporation and free sauna benefits.
5. **Evaluation.** There are several statistics devised for the evaluation of ML models, from precision to false-positive rates. Be sure to only publish statistics that make your model seem to excel in performance. Feel free to devise your own biased evaluation metrics, but be sure to include hundreds of variables (again, so that no one can question your ideas) and mathematical constants like π and e (for legitimacy).
6. **Human-Testing.** If you are so inclined to perform a study, apply for IRB approval, but irresponsible AI work doesn't tend to go through. It is easier to simply refuse to study the effects of your technology on users. Dedicate more funding to advertising the beneficial effects of your technology so that people will volunteer as testers.
7. **Deployment.** For the element of surprise, deploy in a location (or on a subject) that has nothing to do with your training data. For example, you could train the model on humans and deploy the model on extraterrestrials.
8. **Aftermath.** Not everyone is very appreciative of irresponsible AI models. Finding legislation that supports and defends your methods is essential. Look for loopholes in court cases to support your claims, and always use the most vague terminology when writing your own legal documents (like terms and conditions). Make your legal documents as extensive as possible, and release them in a *.png* format so it is impossible for customers to search for keywords on the page. In regards to Irresponsible AI principle 7, do not provide an email or phone number for people who would like to contact a legal team. However, feel free to provide a fax number, since people who still use their fax machines for casual

communication are probably not involved enough with today's technology to question your machine learning methods.

Case Study. To reinforce how useful these guidelines can be, we provide a case study demonstrating their use.

Problem. We would like to build a model to find the optimal locations for lemonade stands. Aligning with guideline 1, this is a problem that no one has bothered to solve, especially when bubble tea started taking the world by storm.

This project has several benefits. First, by working with children, we can prepare them from a young age for their data getting stolen. Second, lemonade is difficult to sell (lemonade is severely limited in its variety, whereas more popular drinks such as bubble tea have endless flavors) without our model, branding strategies, and house-made products. Children are often in charge of running lemonade stands, but they have sticky hands and little business sense, arousing suspicion over their quality of lemonade and leading to questions about their marketing abilities. Lemonade has the inherent advantage of being dairy-free, gluten-free, vegan, and vegetarian, terms that children would not even recognize (but we would employ in our branding). Children *are* cute, though, and paired with our advice, they can be powerfully persuasive in the lemonade world.

Data Collection. Children have some good intuition. For example, they usually set up stands in the corner of two streets, capturing an audience from both intersecting roads. Therefore, training data from children's setup habits could be useful for our model. However, children selling lemonade on the streets in this day and age are hard to find. We drove around 1 city before giving up due to lack of funding (which remarkably has been a consistent issue for our projects) for gas for our car. Since $n = 0$ is slightly *too* small for even the "small sample size" irresponsible approach, we instead decided to employ the data fabrication strategy from guideline 2. We created 1 trillion tuples in the form of (location, neighborhood information, dollars of lemonade sold), with randomly assigned values to all of these.

We also collected data on where people tend to be the hungriest/thirstiest based on the number of customers in restaurants near highways, near neighborhood communities, etc. Our data indicated that the highways tend to be the best locations. However, we counted the number of people who entered and left restaurants, and some of them had only used the bathroom, so our results may be slightly inaccurate. Luckily, errors are welcome in irresponsible AI!

Model Architecture. We chose to use a deep neural network with as many hidden layers as physically possible with our computational resources. The model will finish training in mid-March 2025, giving us enough time to write a follow-up piece in the next SIGBOVIK conference.

Testing, Deployment, and Aftermath. We refuse to test our lemonade on humans by guideline 6. We plan to deploy in Alaska, since made-up data can be applied everywhere. Our legal team, who

was also in charge of writing this paper, is in the process of creating a 300-page terms and conditions document for children to sign. We are excited to be working with an audience that is guaranteed not to read the terms and conditions, and look forward to sneaking in clauses like royalties on their allowance money. (*Update: it turns out several of the children cannot read at all, and their parents will be signing the documents on their behalf. We may be removing the extra clauses.*)

Conclusion. After conducting extensive research on responsible AI, we have developed Irresponsible AI principles and guidelines to help developers create societally useless and damaging tools. We have illustrated an application of the guidelines through a case study on planning optimal lemonade stand locations. We plan to report the results of our work at the next SIGBOVIK conference, and we are confident that they will be positive. In the meanwhile, we hope that researchers and guideline-enthusiasts will take advantage of our work.

Acknowledgements. The author was inspired by discussions in/learnings from her Responsible AI course at CMU to create this piece in jest.

References. Google AI principles (<https://ai.google/responsibility/principles/>)

ENSURE TRANQUILITY (insert evil laugh here)

- 33** **Just-too-late compilation - An examination of a post-emptive compilation technique**

Stephen Rogers

- 34** **Im going to Hurl**

Nicole Tietz-Sokolskaya

- 35** **We Have Serverless at Home**

Ryan Hornby and Matthew Safar

- 36** **USB-*: New Extensions for the Universal Serial Bus**

Jim McCann and Et Alia

- 37** **This SSH servers never gonna give you up (but it will probably let you down)**

Dr. Quod E. Demonstrandum and ?A?l?e?x? ?F?r?i?e?d?m?a?n?

- 38** **Towards an AI-Exclusive Higher Education System**

Ziyuan Liu

Just-too-late compilation - An examination of a post-emptive compilation technique

Stephen Rogers

Sitting at home, Éire

Abstract

Just-in-time (JIT) compilation is a well-established method of deferring the compilation of code until it is needed at runtime, in an effort to the boost performance of interpreted systems. In this paper, we examine the concept of Just-too-late (JTL) compilation which compiles code for an application after it has already attempted to execute. JTL represents the state of the art in runtime performance in interpreted systems by offloading the compilation process to the host system too late, and ignoring any negative side-effects.

1 What's Going On?

Just-in-time (JIT) compilation is a well-established method of deferring the compilation of code until it is needed at runtime. It is particularly useful in interpreted environments, where significant performance gains can be made due to the use of native machine code, compared to a generally much slower interpreter. However, this comes at the cost of a "startup" or "warm-up" delay as any code must be compiled the first time it is to be executed. This means an application may take longer to execute the first time each JIT compiled code section is to be executed, as it has to wait for the code to be compiled.

Enter Just-too-late (JTL) compilation. As the name suggests, JTL involves compilation of code just after an application has attempted to execute it. When attempting to execute uncompiled code, the processor will typically (hopefully) encounter an invalid instruction encoding. When this happens, the processor will raise an interrupt with the host Operating System (OS) which must then handle the exception. Typically, this would involve halting the execution of the process and raising an error. However, with JTL compilation, the OS compiles the code which was attempting to execute and returns it to the process for it to continue execution.

Through this method of compilation, the "warm-up" delay technically no longer exists in your application, as it is the OS which is executing the compilation for you. Therefore, when measuring the performance of an application which uses JTL, there is a measurable improvement in runtime performance when compared with JIT. The elapsed real time may appear to be longer with JTL, but this is not relevant.

2 Explain It To Me

The first time a JTL section of code is reached for execution, it must be loaded in its entirety into memory. This region of memory should then be treated as a regular function and the host's calling convention shall be used to enter that region of memory. At this point, an invalid instruction encoding exception must be raised immediately at the start of the section so that it may be intercepted as a complete unit for the compilation process. This can be handled by inserting a known invalid instruction at the head of the section, but this comes at the cost of at least one single byte of memory. To avoid this overhead, we simply hope that the start of the uncompiled code happens to map to an invalid instruction encoding.

When an invalid instruction encoding is detected at runtime, the processor will raise an exception with the host OS. On X86, this is interrupt 6 "Invalid Opcode Exception". When this occurs, the CS (Code Segment) and EIP (Extended Instruction Pointer) registers will point to the segment of "code" which caused the exception. On Linux based systems, this interrupt will result in a SIGILL (Illegal Instruction Signal) for the process, which causes its termination.

In our modified system, rather than sending a SIGILL, the OS uses the CS and EIP registers to fetch the full uncompiled code from memory and pass it to the appropriate compiler. The compiled machine code is then taken from the compiler and overwrites the uncompiled code in memory. Control is then returned to the process to continue execution with the same CS and EIP values. This does require that the compiled code will fit into the same memory region as its uncompiled counterpart, which can be achieved using compilation techniques for code compaction. More importantly, this can be achieved with a wish and a prayer.

3 Unicorn Code

However unlikely, it is conceivable that there exists a program source code which coincidentally is encoded as valid machine code providing the exact functionality the source code represents. In this case, JTL provides the fastest possible execution path by bypassing the compilation process entirely. Such Unicorn Code represents the Holy Grail for performance in interpreted systems.

4 In The End, It Doesn't Even Matter

This paper presents the state of the art in runtime performance for interpreted systems by abusing the failure states of systems, piggybacking on a mechanism that is in no way designed for or suited to it. Bizarrely, or perhaps naively, I think this might actually possibly work.

I'm going to Hurl

Nicole Tietz-Sokolskaya*

Abstract

Modern programming languages have begun to use exceptions but have yet to explore the full utility of this feature. This work introduces Hurl, an exceptional programming language, which leverages the full power of exceptions to toss out any other control flow constructs. We show that this language produces nothing of value.

1 Introduction

Exception handling traditionally has been used to handle errors and unusual events. Using exception handling for control flow has been explored in some mainstream languages, such as Python, where exceptions are used to end generators[4]. No languages known to the author have committed to fully relying on exception handling as the primary control flow primitive. There are other languages that strive for minimalism through the instructions for a Turing machine[1].

This work introduces the Hurl programming language[6]. This language uses two forms of exception handling, along with anonymous functions, to provide a complete programming language suitable for no sort of practical work. This language provides the user with a novel approach to writing the most basic programs, which some may take as an enjoyable challenge¹. Through the use of Hurl, the programmer also enjoys increased job security from its resistance to understanding by other humans and by LLMs, providing a strong motivation for job preservation.

*My employer had nothing to do with this

¹This is hypothetical, as the author has not found any willing users. She will not discuss the unwilling users.

2 The Hurl Language

The core of `hurl` is exception handling. In this section, we introduce the syntax for the core features, give examples from the standard library, explain the available tooling, and review the license choice.

2.1 Syntax

The core syntax of `hurl` allows you to bind dynamically-typed variables, declare anonymous functions, and use exceptions. We chose not to innovate in much of the core syntax, to remain familiar while deviating in the novel application of exceptions as primary control flow.

Here is an example of the syntax to declare variables, add comments, and perform arithmetic.

```
let language = "hurl";
let year = 2023;
year = year + 1; # happy new year!
let authors = ["Nicole", "just me"];
```

Variables are dynamically typed and must always be initialized. Every statement ends with a semi-colon. If you assign to a variable, it finds the variable in the closest surrounding scope and assigns to it, but panics if there is none found.

We can also create anonymous functions.

```
let add = func(a, b) {
  hurl x + y;
};
```

There are a few notable things here: each statement ends in a semi-colon; we use the `hurl` keyword for the first time, since we do not have return values²; and since it is dynamically typed, you can pass in strings which will then be concatenated.

²The `return` keyword exists and does not do what you expect, unless your brain is as deviant as mine.

Functions can also be recursive. They are not bound to a name until after the function is declared, but since names are resolved at runtime we get to just use the outer name as it will resolve then. Here's an example of an infinite loop using this trick:

```
let overflow = func(depth) {
  println("depth: ", depth);
  overflow(depth + 1);
};

overflow(1);
```

If we run this, we get a stack overflow once we hit the end. This happens because we have no tail-call optimization. Tail-call optimization was considered but was left out because ~~it was too hard to add~~ it is a fun challenge for the programmer to work around.

And now we reach Hurl's defining feature, exception handling. This is the primary control flow mechanism, as you cannot do any useful work in Hurl without it: even heating an office via CPU busy loops cannot happen without exception handling, because otherwise you will stack overflow before you generate sufficient heat.

There are two main ways to throw exceptions in Hurl:

- **hurl** an exception and it will unwind the stack until it finds a handler which catches it. Execution then continues from that handler.
- **toss** an exception and it will walk the stack until it finds a handler to catch it, then it will walk back to where it was tossed from and resume execution from there. Resuming execution requires the handler using the 'return' keyword.

hurl is what we normally see as exceptions in other languages. **toss** is a gimmick that strongly resembles continuations[2] by suspending execution to run a handler then resuming. We believe that only **hurl** is necessary, but **toss** is super cute and has to be included for the language to get widespread adoption.

There are three ways to catch an exception:

- **catch (<value>)** will match the literal value provided; any other value falls through to following catch statements
- **catch as <ident>** will catch any value and provide it in the catch body as a bound variable, only in scope for the catch body
- **catch into <ident>** takes *no* body and catches any value, assigning it into the provided already-defined identifier in some scope, then resumes execution after the handlers.

2.1.1 Examples

Using what we've defined previously, we can compute a factorial. This example leverages all three types of catch handlers (though only uses **hurl**, not **toss**).

```
let factorial = func(n) {
  try {
    hurl n == 0;
  } catch (true) {
    hurl 1;
  } catch (false) {
    let next = 1;
    try {
      factorial(n - 1);
    } catch into next;
    hurl n * next;
  };
};

try {
  factorial(10);
} catch as x {
  println("factorial(10) = ", x);
};
```

And the result:

```
> hurl run examples/factorial.hurl
factorial(10) = 3628800
```

Closures may also be created. These are omitted due to space constraints.

2.2 Standard Library

Hurl has a good-enough set of built-ins and standard library for getting toy programming problems done.

The full list of built-in functions is available at the Hurl homepage (hurl.wtf). It includes print functions, number conversions, floor/ceiling, and other functionality that the author could have done in Hurl code but ~~was too lazy~~ wanted to optimize for real-world usage and efficiency.

The standard library is entirely written in Hurl and provides abstractions over top of Hurl's syntax to get a more familiar programming feel. To use this functionality, you place it in a file somewhere on your disk³, then include it with **include <filename-expr>**; where the expression is anything that evaluates to a filename, and it resolves relative to the first interpreted file.

³A package manager is in the nice-to-have future projects list.

The full standard library cannot be included here, but we will illustrate how this useful functionality is defined and utilized through illustrating loops.

We can define a loop simply as follows. It takes in two arguments: a function which is the body of the loop, and a set of local variables to invoke the body of the function with. It requires that the body `toss` the values for the next iteration to be bound, and then `hurl` an exception, either true to continue iterating or false to stop iterating.

```
let loop = func(body, locals) {
  let new_locals = locals;
  try {
    body(locals);
  } catch (true) {
    loop(body, new_locals);
  } catch (false) {
    hurl new_locals;
  } catch as update {
    new_locals = update;
    return;
  };
};
```

Here is how you use `loop` to compute the Fibonacci sequence.

```
include "../lib/loop.hurl";

# fibonacci in hurl
let fib = func(locals) {
  let a = locals.1;
  let b = locals.2;
  let iters = locals.3;
  let max_iter = locals.4;

  toss [b, a + b, iters + 1, max_iter];
  hurl iters < max_iter;
};

try {
  loop(test, [0, 10]);
} catch as retval {
  println("computed: ", retval);
};
```

With these custom functions, we're able to build out fairly familiar, but still abnormal, control flow. The uncanny valley ended up being crucial to our experimental results.

The remainder of the standard library defines various conditional, looping, and mathematical constructs. The most notable piece is in the definition of `foreach`, we have to chunk the iterated list and

recursively call `foreach` as we approach the limits of the operating system stack. This limitation is due to not having any tail-call optimization and could be improved in future interpreters.

2.3 Tooling and Documentation

Hurl has robust documentation built out and available on the Hurl homepage, at hurl.wtf. It covers all the essentials for a respectable language: installation, usage, syntax, patterns, standard library, debugging, and examples!

One core piece of Hurl's ecosystem is good tooling. Like all good languages, to get adoption, we are building out the tooling users expect to have. The first piece is a code formatter. This exists, and successfully formats Hurl code that you run it on. Future tooling could include an LSP server so that editors can provide useful completions and a package manager.

2.4 License

Hurl is licensed under three licenses:

- GNU Affero General Public License (AGPL-3.0)
- Gay Agenda License (GAL-1.0)
- Commercial licenses

These licenses were chosen carefully to further the goals of the project: advancing the rights and wealth of queer people (especially the wealth of the author). This goal is achieved by using a strong copyleft license, the AGPL-3.0, which allows widespread adoption of the tool outside of work and then, when people bring its value into their jobs, means that their companies will need to purchase a commercial license.

The alternative open-source license, GAL-1.0[5], was written for this project by a wholly unqualified software engineer. She took the MIT license and added requirements, such as supporting LGBTQ rights and saying "be gay, do crime" at least once during use of the software.

We strongly suspect license violations have occurred. The licenses were chosen to compel companies to purchase a commercial license from us for a small eight-figure price, commensurate with the value received.

So far, none have reached out, indicating that there has almost certainly been theft. It's also unlikely that all the users are choosing the Gay Agenda License,

because the assault on LGBTQ rights rages on; with such an earth-shattering project, if it were licensed under GAL, this assault would stop. We've delayed adding telemetry to Hurl not to protect our users, but to protect our delusion that widespread usage has occurred.

3 Experimental Results

We conducted real-world experiments with Hurl to validate some of its properties with real human subjects. We wanted to understand what the experience was like to program using only exceptions for control flow.

3.1 Methodology

We gave each human subject the task of completing Advent of Code[7] using Hurl to validate the standard library and evaluate human efficiency in the Hurl programming language. We evaluated machine subjects by giving them simple problems, like FizzBuzz[3]. We did not feel the need to go any further.

We did not reach statistical significance, because the results were such that it was unethical to continue with further trials.

3.2 Human Trials

Only one human was harmed in these trials. This is a significantly lower count of harmed humans than any other mainstream programming language. These trials were not approved by any internal or external review board.

For our sole human subject, we will use the pseudonym "Nicole" to preserve her privacy⁴. "Nicole" was given the task of doing Advent of Code in Hurl, with a minimum of reaching completing three days of it. She completed three days, choosing to do days 1, 2, and 4, and then she withdrew from further participation in the study. During these three days, she said "I'm going to Hurl" and "Oh god why did someone do this to me." Her attempts to find other people to participate were unsuccessful.

During the human trial, we found that a human can write reasonable code in Hurl, and some of the solutions verge on elegant if you're that kind of person. Additionally, since the interpreter is an absolute amateur hackjob, it lacks any faculties for debugging or performance analysis. This is key for job security, as it means that expertise in Hurl will be hard-won

⁴Her real name is Nicole, and she is the author of this paper, but ssshhh I didn't tell you that.

and replacing a Hurl expert will be impossible. How to get Hurl into production without losing your job is left as an exercise for the reader.

3.3 Machine Trials

We tested multiple LLMs⁵ on their ability to write Hurl code. They consistently failed to write it in a correct way.

The LLMs were provided with the formal grammar for Hurl, as well as a small sample program. Then they were asked to write an additional small program. The failure modes can primarily be characterized as expecting Hurl to be a reasonably normal language. They would use the `return` statement to return values from a function, rather than as intended. They also tended to expect conditionals, like `if else`.

We expect further improvements with GPT-4 and other newer LLMs, but suspect that the deviation from programming norms provides a barrier for easy LLM production of Hurl code.

4 Conclusion

Hurl presents a novel approach to programming by using only exceptions. It provides the beginnings of robust tooling, and has the possibility of further extension. So far, trials look promising at delivering no real value. The most promising direction for providing value is through job security: if you're able to get it into production, no one else will touch it. It's unclear how to achieve this while remaining employed, which will be the subject of further research.

⁵We think they were GPT-3.5 and Claude, but honestly we forget and lost the chat transcripts. Whoops!

References

- [1] Brainfuck (programming language). <https://en.wikipedia.org/wiki/Brainfuck>. Accessed: 2024-03-04.
- [2] Continuations. <https://en.wikipedia.org/wiki/Continuation>. Accessed: 2024-03-06.
- [3] Fizz buzz. https://en.wikipedia.org/wiki/Fizz_buzz. Accessed : 2024 - 03 - 13.
- [4] Python standard library exceptions list. <https://docs.python.org/3/library/exceptions.html>. Accessed: 2024-03-04.
- [5] Nicole Tietz-Sokolskaya. Gay agenda license text. <https://git.sr.ht/~ntietz/hurl-lang/tree/main/item/LICENSE/GAL-1.0>. Be gay, do crime.
- [6] Nicole Tietz-Sokolskaya. Hurl, the exceptional language. <https://hurl.wtf/>. I'm sorry.
- [7] Eric Wastl. Advent of code. <https://adventofcode.com/>.

1 We Have Serverless at Home

By: Ryan Hornby and Matthew Safar

```
<!DOCTYPE html>
<html lang="en">
  <head>
    <meta http-equiv="Content-Type" content="text/html; charset=UTF-8">
    <meta charset="UTF-8">
    <meta name="viewport" content="width=device-width, initial-scale=1, maximum-scale=1.0,
      user-scalable=no">
    <title>Serverless At Home: Sigbovik 2024</title>
    <link rel="icon" type="image/x-icon" href="favicon.png">
    <link rel="stylesheet" href="https://cdn.rawgit.com/Chalarangelo/mini.css/v3.0.1/dist/
      mini-default.min.css">
    <link rel="stylesheet" href="style.css">
    <script defer src="./script.js" type="module"></script>
  </head>
  <body>
    <header class="sticky">
      <a href="https://sigbovik.org/" class="logo"></a>
      <div id="header-contents">
        <div id="header-top">
          Sigbovik 2024
        </div>
        <div id="header-bottom">
          <a href="/" class="button" style="padding-left: 0;">Home</a>
          <a href="ringfuck.html" class="button">Ringfuck</a>
          <a href="404.html" class="button">404</a>
        </div>
      </div>
    </header>
    <header class="sticky ruler-container">
      <div class="ruler ruler-x"></div>
    </header>
    <div class="ruler-container" id="y-ruler-bar">
      <div class="ruler ruler-y"></div>
    </div>
    <div id="paper">
      <h1 class="title">Serverless at Home?</h1>
      <h4 class="title">Ryan Hornby and Matthew Safar</h4>
      <h4><b>Abstract:</b></h4>
      </img>
      <h2>Introduction</h2>
      <p>The client-server model has been in use for around 60 years, but over the
        past two decades, the idea of serverless computing has become feasible and
        popular. In case you aren't familiar with serverless computing, allow the
        collective consciousness of StackExchange (ChatGPT) to explain: "Serverless
        computing? More like 'server-less headaches.' It's not truly 'serverless'
        because, surprise, servers still exist! You're just off the hook for
        babysitting them. Think of it more as 'servers on demand' with a fancy
        billing twist." While not babysitting servers sounds great, most people
        would prefer not getting billed for the pleasure. In this paper, we will
        investigate "free serverless", also known as serverless at home. </p>
      <h2>Methods</h2>
      <h3>Hosting</h3>
      <p>For hosting, we use Github Pages. 'Its a free, powerful solution for your
        frontend AND backend needs! You can create your stunning frontend with HTML
        /CSS/Javascript and power it with Github Actions to produce a functional
        and modern website! </p>
```

<h3>WASM</h3>

<p>Perhaps a more intuitive definition of serverless would be one where there simply is no server (outside of the server hosting the html files). At first, this definition may sound too restrictive to have interesting websites, but by using web assembly (WASM) and javascript we can make the client run all the `code` to make our site interesting. We've shown this functionality with our previous SIGBOVIK entry: [Ringfuck](ringfuck.html). While WASM makes creating interesting websites pretty easy, javascript psychopaths have been doing the same thing for years without it.

<h3>CI/CD</h3>

<p>If we want true serverless at home, we need to be able to replicate the functionality of database-driven websites for free. To do this, we take inspiration from others who have sought to innovate by taking other 'peoples backend and making it their own (see https://www.youtube.com/watch?v=c_arQ-6ElYI). For our backend we leverage Github and its most recent feature (besides stealing everyone's `code` to train a Kamikaze CoPilot) CI/CD, or Github Actions as they call it. Github's CI/CD allows you to write a custom bash script to be run in a Docker container. This allows us to edit our html with simple CLI tools like cat, echo, sed, and awk and then commit those changes back to the repo. The final step is to simply trigger the CI/CD to run when a user performs an `action` requiring a server operation. To perform this, we use some simple javascript along with a Github auth token to trigger the `action`, then query the Github API to determine when the change has finished so that we can automatically refresh the page.</p>

<p> Something to consider is that our choice to stick to CLI tools is only because we wanted to stay simple and fast. If those are not concerns for your website, you can easily set up a CI/CD pipeline that includes advanced tools like Python for machine learning, SQLite for relational databases, or Elasticsearch for search engine capabilities! The possibilities are truly endless, as long as Github 'doesnt notice your compute time and start charging you money. </p>

<h2>Security and Privacy</h2>

<p>While there are security concerns with this serverless at home methodology, a lot of potential security issues are shared with other applications that accept user input, so we will not be discussing these. A unique issue with our CI/CD approach is the need for a publicly accessible Github auth token . To mitigate this risk, it is recommended to use a fine-grained access token for the specific repository. The major downside to this is that the token must have an expiration date, thus creating some maintenance. To avoid maintenance, you could create a dummy account with a classic access token, but this would allow anyone to attempt changes on public repos with your dummy account.</p>

<p>On the privacy front, there are a couple of unique issues with serverless at home. First is the classic git problem of it being very difficult to permanently delete `content` that was committed. Thus, any sensitive `data` in a comment will remain visible to everyone until an administrator removes it , which may never happen since owners of github repos are notoriously picky about other people changing their precious `code`. The next concern regards user sessions. While it is possible to have the CI/CD system create user sessions by creating subdirectories and deleting them later, these sessions would be viewable by anyone interested enough to pull up the source `code`. Since this is not easily mitigatable, our recommendation is to simply inform the users of your site that they have **no privacy**. So if you've read this far, congratulations you've been informed.</p>

<h2>Comparison to Traditional Serverless</h2>

<p>Traditional hosting and serverless are a pretty polished experience for devs and users, especially compared to serverless at home. So why would you ever use serverless at home? One selling point is the fact that it is

totally free, as long as you don't have a bunch of users, which shouldn't be a big concern for most readers. Another benefit of serverless at home is the nostalgia for dial-up internet. Some people would argue there are better free alternatives, like Wordpress. Which brings us to our next selling point: freedom. Services like Wordpress might be free, but they restrict what you can do with your site to try to get you to pay for more freedom. But with serverless at home, as long as you can code it to run when the Github Action triggers, you're free to do it. Finally, the biggest reason to use serverless at home is so that you can deliver an awful user experience to your "friend" when they ask you to make a "simple" website for them.</p>

```
<h2>Conclusion</h2>
  <p>Is serverless at home revolutionary? Let us know in the comments! 🙋</p>
  <dialog id="loading-dialog">
    <div class = "card large">
      <p><span class="spinner primary"></span>Please wait while we post your comment
        ! 😊</p>
      <audio id="loading-sound" loop preload>
        <source src="/the-sound-of-dial-up-internet-6240.mp3" type="audio/mpeg" />
      </audio>
    </div>
  </dialog>
  <div id="comments">
    <h3>Comments</h3>
    <form id="comment-form">
      <fieldset>
        <div class="row responsive-label"><div class="col-sm-12 col-md-3"><label
          for="cf-name">Name</label></div>
        <div class="col-sm-12 col-md"><input type="text" placeholder="Name" id="
          cf-name" name="Name" style="width:85%;"></div></div>
        <div class="row responsive-label"><div class="col-sm-12 col-md-3"><label
          for="cf-textarea">Textarea</label></div>
        <div class="col-sm-12 col-md"><textarea style="width:85%;" placeholder="
          Textarea" id="cf-textarea" name = "Textarea"></textarea></div></div>
        <button>Post Comment</button>
      </fieldset>
    </form>
  </div>

```

```
<pre>Ryan Hornby: First</pre>
  </div>
</div>
</body>
</html>

```

2 Appendix 1

List of relevant links:

1. <https://backendathome.github.io/>
2. <https://backendathome.github.io/serverless.png>
3. <https://backendathome.github.io/ringfuck.html>
4. https://www.youtube.com/watch?v=c_arQ-6ElyI

3 Appendix 2

Preview of site:

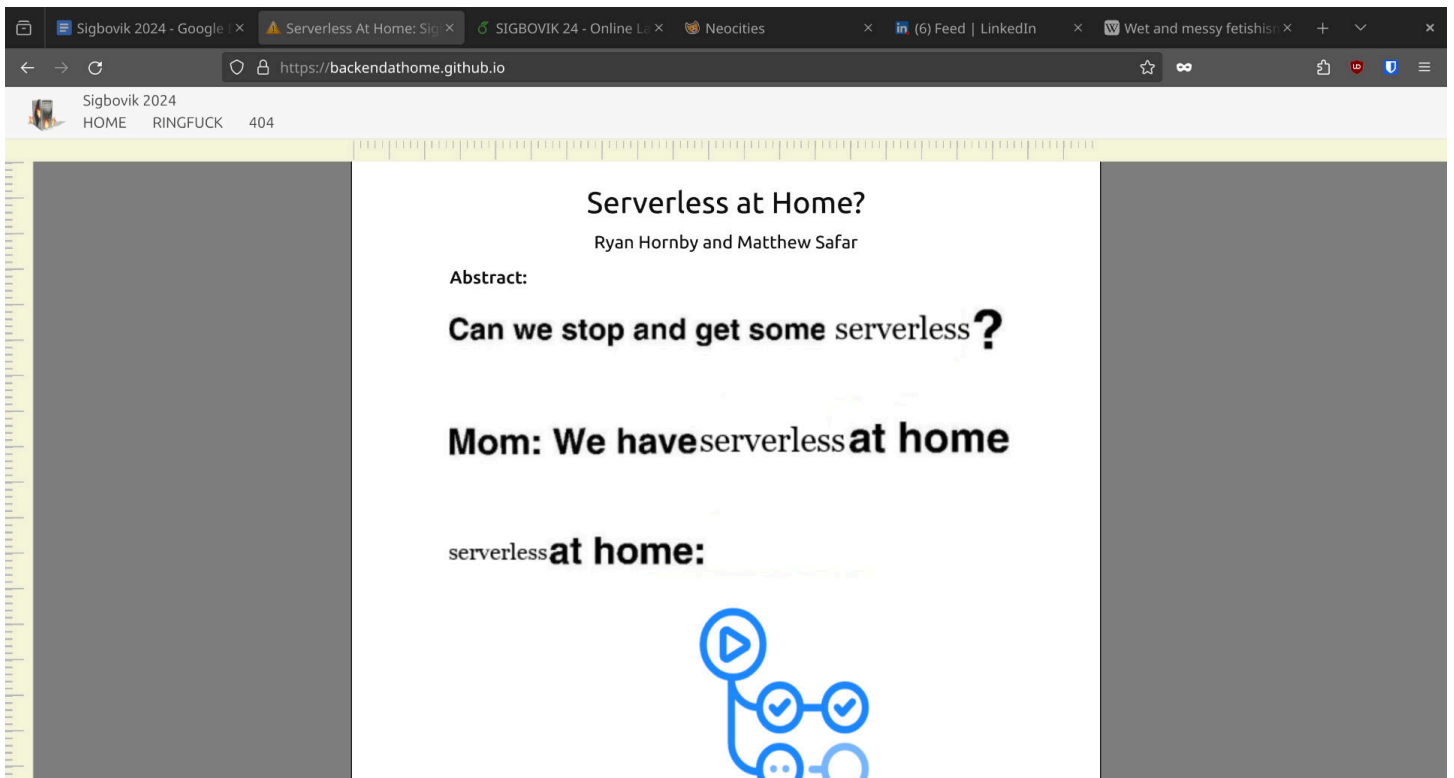


Figure 1: Example of serverless website.

USB-★ – New Extensions for the Universal Serial Bus

Jim McCann*
Carnegie Mellon University

Et Alia
Et Aliae

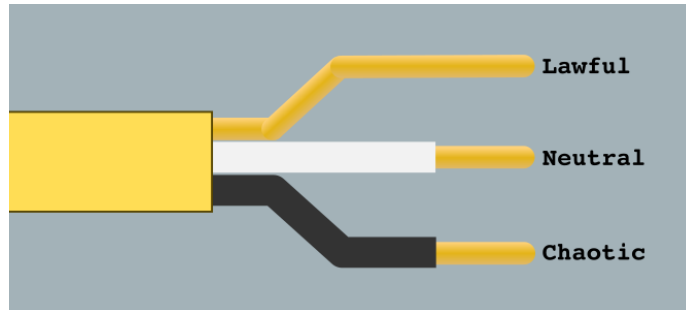


Figure 1: I didn't have an image to put here so here's a thematically related post on mastodon. Please boost it at <https://social.wub.site/@ix/111622212515662074>.

Abstract

The great thing about standards, the old saying goes, is that there are so many to choose from. USB embodies this expression within one family of interconnects, with a wide range of connector types (A, B, Mini-A, Mini-B, Mini-AB, Micro-A, Micro-B, Micro-AB, 3.0-A, 3.0-B, 3.0-Micro-A, 3.0-Micro-B, 3.0-Micro-AB, C, ...), protocol speeds (Low-Speed, Full-Speed, High-Speed, Gen 1, Gen 2, Gen 2x2, Gen 3x2, Gen 4, ...), conductor counts (2-24), and optional features (power delivery, battery charging, displayport, “thunderbolt”, PCI-express, ...). Indeed, the term “USB cable” and “USB port” have become so generic as to be meaningless.

What better way to celebrate the horror of this specification family than by adding a few more logs to the train wreck?

1 USB-HV(DD) – High Voltage / Device Disabling

The purpose of a USB-HV(DD) device is to render any attached devices permanently inoperable. To achieve this aim, a USB-HV(DD) device shall deliver at least 400v of alternating current across at least one pair of pins. Higher voltages are, of course, more effective at arcing between traces and wreaking electrical havoc.

HV(DD) is compatible with any connector or cable standard, and USB-HV(DD) endpoints may be hosts, devices, hubs, or cables. (Even a passive cable can become USB-HV(DD) compatible if a rapidly rotating magnetic field is introduced.)

USB-HV(DD) devices may choose to limit their output current to remain human-safe, but are probably a pretty bad

idea regardless of this. Especially because, for maximum chaos, it is recommended that HV(DD) devices contain no special markings.

2 USB-HC(DD) – High Current / Device Disabling

USB-HC(DD) attempts to accomplish the same aim as USB-HV(DD), but via high current and low voltage. This is likely to fail to disable any USB host or device, but can be useful for disabling cables.

Therefore, USB-HC(DD) devices should be constructed such that a cable connected between two ports of the same device is subjected to sufficient current to be melted through joule heating. Some care should be taken to avoid welding inserted cables to the usb connectors available on the device, unless the USB-HC(DD) device is meant to additionally be self-disabling.

HC(DD) devices should not be labeled, though pictograms suggesting self-loop cable connections might be useful to tempt the unwary.

3 USB-Sea

The USB-Sea protocol replaces electrical signaling with pressure variation in tubes of salt water. (Salt water is used to allow limited power transmission, of course.) What this lacks in data transmission speed, robustness, portability, utility, electrical fidelity, and sensibility it makes up for in corrosiveness.

*ix@tchow.com

4 USB-NC – No Connector

The USB No Connector extension avoids the pitfalls associated with USB's many connector standards. A USB-NC host, device, or cable does not end in either a plug or a socket, but – rather – in a pigtail of bare wire ends. A small label affixed near the pigtail indicates that wires should be twisted together with same-color wires.

In the spirit of the USB specification so far, USB-NC does not specify how the wires should be color-coded, leaving manufacturers to independently arrive at inconsistent selections.

A pair of wire strippers or a sharp knife may be used to retrofit any existing USB cable into a USB-NC adaptor cable.

5 USB-PrM – Preemptive Mode

The USB ecosystem generally operates in a “cooperative multi-device” setting, where multiplexing USB host ports between devices is handled by manually time-slicing the devices (adding/removing cables). Of course, such cooperative multitasking is outmoded, not to mention unresponsive.

Preemptive USB devices are designed to support automatic time-slicing between devices with long-range magnetic cables and electromagnetic ports that can be activated by devices to eject or attract nearby cables. (Not to mention finally making "ejecting" a USB drive from an OS context menu do what you'd expect it to.)

6 Data-Only USB

The prevalence of USB as a 5v charging standard led to the creation of many “power-only” USB cables that omit all lines except for 5v and ground in order to save on copper and thus cost. Like donuts and electrons, data-only USB cables provide the formal complement and omit all power transmission lines.

With a power-only, data-only, and two power/data splitter adaptors one can replicate the function of a single USB cable with a messy combination of four devices.

7 Wireless USB

Wireless USB combines the best features of power-only and data-only cables: the missing wires. Wireless USB cables provide neither power nor data transmission, and are expected to be even cheaper to produce than either power-only or data-only cables. These cables, therefore, only form a physical tether with no electrical consequences.

We expect that the availability of wireless USB cables will also spur the development of wire/no-wire splitting adaptors.

Wireless USB cables are not compatible with the USB-NC extension, unless you count tying to plastic cords together as twisting wires.

8 USB Gen 4x4 – Off-Road

USB Gen 4x4 extends USB Gen 2x2 with better off-road capabilities. Whereas USB Gen 2x2 provides two 10Gbps lanes for data transfer, USB Gen 4x4 is an automotive standard requiring that a vehicle have four wheels, of which are powered. (And, strangely, that the vehicle has roughly the overall shape of a Mini-AB connector.)

It is possible to refer to USB Gen 4x4 vehicles as “USB 3.3 Gen 4x4,” but, really, it is possible to refer to a vehicle as anything, like “Maude” or “a fast rock,” as long as the antecedent is clear from context. So there's that, I guess.

Acknowledgments

Early and often.

References

No.

This SSH server's never gonna give you up (but it will probably let you down)

Dr. Quod E. Demonstrandum ~~and Alex Friedman~~

1 Pre-Pre-Pre-Introduction: The Internet

In recent years a revolutionary new technology has been developed called “the internet”. Originating from ARPANET [16], the internet links together computers and enables communication between them (even if they are geographically distant!).

This technology has already had many impacts, and will surely continue to become an increasingly important facet of every day life as it continues to proliferate into new fields.

2 Pre-Pre-Introduction: Internet Applications

A variety of applications have already been developed for the internet. This includes the Domain Name System (DNS) [13] which allows one to associate a string “hostname” (called a domain) with a unique number used to address each computer on the internet called an Internet Protocol (IP) address.

Other applications include the “web”. Powered by the Hypertext Transfer Protocol (HTTP) [20], the web allows one to build and host their own multimedia pages that can be viewed by anyone with an internet “browser”! For example, an academic institution may register their own domain and host a webpage on it to allow people to learn more about them—short of having to visit. In this respect, website can almost be thought of as internet bulletin boards.

The internet also even has its own solution to mail called “email” which is implemented over a variety of protocols including the Simple Mail Transfer Protocol (SMTP) [17] (However, the authors of this paper are skeptical how beneficial this application is due to its many limitations compared to traditional mail such as receiving letters from the IRS as well as packages—such as ordered via a site my grad students inform me of called “Amazon”).

This paper, however, focuses on the Secure Shell Protocol (SSH) [18]. In short, SSH allows one to obtain a shell on a remote system. While various SSH implementations exist and have various complex features, this paper only requires knowledge of the most basic concepts of SSH.

Assuming one has an SSH client installed, they will be able to use the command `ssh <username>@<domain/ip>` to be able to connect to the remote server. Upon connecting to the specified ssh server, the user will be prompted for a password. The ssh server will then attempt to authenticate this username and password combination. If it is valid, then the ssh server will provide a shell to the ssh client with the authenticated user.

3 Pre-Introduction: Memes

One inexplicable aspect of the internet is the development of so called “memes”. In short, memes can be thought of as internet jokes and can take various forms from text [11] to images [15] and videos [1, 2, 3, 4]; however, this is far from an exhaustive list. Additional, examples of memes are: [12, 14, 22]. Another important, yet nuanced point, is that just because something on the internet takes this form, it does not mean that it is, in fact, a meme. Instead, a meme is more of a genre for which so called meme connoisseur claim to be experts in.

It is worth noting that the development and use of memes is incredibly context sensitive. Memes appear to have their own temporality/lifecycle where it is “uncool” to use an old/dead meme.

These phenomena were first predicted before the development of the internet by Harry Q. Bovik—who has since gone on to be one of the most influential people in the field.

4 Introduction

Despite the proliferation of the internet as well as internet memes, there has been surprising lack of work exploring the possibility of merging the development of serious internet applications with memes.

In hopes to call attention to this research gap, this paper explores the combination of one serious application with one internet meme. Specifically, we explore the combination of SSH with “Rickrolling”, the meme wherein one unexpectedly links someone to the music video of “Never Gonna Give You Up” by Rick Astley [5].

The result of our work is the development of a novel concept known as the SSHRoll: a Rickroll [6] accomplished via an SSH server. To accomplish this, we implemented an innovative new SSH server which will accept any username and password combination (to ensure users are able to login). Instead of providing users with a regular shell (e.g., bash or zsh), users will instead be provided with a custom shell with no commands. Any time any key is pressed, the shell will clear itself and render an ascii-art rendition of Rick Astley dancing to “Never Gonna Give You Up” from the official music video [7] at timestamp 0:03 (3 seconds into the video) [8]. By looping through a sequence of pre-specified ascii-art images, we are able to accomplish an animation should one spam keyboard inputs.

Because this occurs on any key input to the shell, regardless of the key, it becomes impossible to exit the session without the user closing the terminal they used to connect to the server. As such, we are able to revolutionize the space of SSH servers by making one which truly will never give you up. Unfortunately, however, upon the realization that one has been SSHRolled, it is quite likely that the user will be let down—especially if they were under the false impression that they had managed to break into the SSH server in question.

As an added benefit of not being able to disconnect from the ssh server without closing the terminal, you will be able to become a sworn enemy of the sysadmin [19] sent to investigate the inordinate number of ssh connections to your server—only after they share it with all their friends within your organization’s information security department as, while your organization provides all devices with a public IP address, the sysadmin’s NDA prevents

them from sharing this information externally (not that this information won't already be public due to the plethora of botnets roaming the internet).

No, we do not have experience with this.

Seriously, though, we don't. You think we'd pass up a chance at a self-citation?

You're really going to argue this and make me interrupt the whole flow of this paper, aren't you?

No, not you reader, we know you're cool (and, if we're being honest, probably not even reading this: you're just skimming the paper looking for a figure). I'm looking at you Reviewer #2. You know what you did; and, now instead of including *Figure 1. "Rick Astley SSHRolls you.jpeg"* (sorry reader), I have to reiterate a point that I *clearly* addressed. Why couldn't you have been like Reviewers #1 and #3 whose only feedback was a face-palming emoji that hyperlinked to *Never Gonna Give You Up* [9]?

5 Implementation

We implemented an SSHRoll server as described in our introduction (see Section 4) which can be found at <https://github.com/ahfriedman/sshroll>. While we can assure you that this artefact works on our machine, we cannot guarantee replicability of these results. The student responsible for implementation has since graduated, and we did not have the time to unravel its workings prior to submission. While we reached out to the student to try to see if they could assist us in hopes of publishing this paper, they refused to be associated with us due to "ethical concerns" over our methodology (see Section 6). As such, all references to our student, Alex, have been removed from this paper. This is a shame as Alex has done quite a lot of great work, and this paper would have truly helped the student build a strong reputation in the research community.

Moreover, we highly discourage any attempts to use our work as **it poses a fairly substantial security risk, and is basically asking for your device to be compromised**. It is worth, however, noting that some of our colleagues who reviewed this paper prior to submission suggested that reframe this as a "free pentest" before laughing and asking, "wait, you're not *actually* going to do this, are you?".

Responsible use of this work is left as an exercise to the reader.

6 Methods & Results

To test the results of our implementation, we were planning on using the phenomenal concept introduced in the paper "An Undergrad Is All You Need", GPT-UGRD [21] as a means to circumvent IRB; however, between the school catching wind of our plans and the sleepless nights guarding our server's ethernet cable from the scissors of the aforementioned sysadmin,

we decided to cut our losses and leave collecting results as future work to ensure adequate funding and time to go through proper channels. In the interim, we have decided to go on to more “important” work (or so the tenure committee called it).

7 Conclusion

The internet is probably a revolutionary technology whose impacts will be seen over the next 20–30 years, if it ever becomes a mass-consumer product. Despite the serious academic nature of its predecessor, ARPANET, the internet in its current form is already full of inside jokes known as “memes”.

This paper helps to bridge the gap by introducing a serious internet application whose practical nature has been completely replaced by memes. Specifically, we introduced (but totally did not test) SSHRolling: a Rickroll [10] accomplished via an SSH server. We provide an implementation of an SSHRoll server which makes the following advances on a traditional SSH server:

1. All authentication has been removed, accepting any username and password, ensuring access to the SSHRoll to all.
2. All commands have been removed in favor of displaying an ascii-art of Rick Astley dancing.
3. Disconnecting from the server cannot be accomplished from the shell itself—requiring the user to disconnect via closing their terminal—to ensure that the server can never give you up.

While this is already quite an advancement, SSHRolling currently has some limitations. Notably, it is likely to “let down” anyone who is lucky (or unfortunate) enough to stumble upon it. We hope that this can be addressed through future work along the exploration of making SSH servers that provide different memes, and introducing memes into other aspects of the internet.

References

- [1] Rick Astley. *Rick Astley - Never Gonna Give You Up (Official Music Video)*. <https://www.youtube.com/watch?v=dQw4w9WgXcQ>. 2009.
- [2] Rick Astley. *Rick Astley - Never Gonna Give You Up (Official Music Video)*. <https://www.youtube.com/watch?v=dQw4w9WgXcQ>. 2009.
- [3] Rick Astley. *Rick Astley - Never Gonna Give You Up (Official Music Video)*. <https://www.youtube.com/watch?v=dQw4w9WgXcQ>. 2009.
- [4] Rick Astley. *Rick Astley - Never Gonna Give You Up (Official Music Video)*. <https://www.youtube.com/watch?v=dQw4w9WgXcQ>. 2009.
- [5] Rick Astley. *Rick Astley - Never Gonna Give You Up (Official Music Video)*. <https://www.youtube.com/watch?v=dQw4w9WgXcQ>. 2009.

- [6] Rick Astley. *Rick Astley - Never Gonna Give You Up (Official Music Video)*. <https://www.youtube.com/watch?v=dQw4w9WgXcQ>. 2009.
- [7] Rick Astley. *Rick Astley - Never Gonna Give You Up (Official Music Video)*. <https://www.youtube.com/watch?v=dQw4w9WgXcQ>. 2009.
- [8] Rick Astley. *Rick Astley - Never Gonna Give You Up (Official Music Video)*. <https://youtu.be/dQw4w9WgXcQ?t=3>. 2009.
- [9] Rick Astley. *Rick Astley - Never Gonna Give You Up (Official Music Video)*. <https://www.youtube.com/watch?v=dQw4w9WgXcQ>. 2009.
- [10] Rick Astley. *Rick Astley - Never Gonna Give You Up (Official Music Video)*. <https://www.youtube.com/watch?v=dQw4w9WgXcQ>. 2009.
- [11] *Coppypasta – Navy Seal Coppypasta — Genius*. <https://genius.com/Coppypasta-navy-seal-coppypasta-annotated>.
- [12] Quod E. Demonstrandum and Alex Friedman. “This SSH server’s never gonna give you up (but it will probably let you down)”. In: *Sigbovik 2024*. Ed. by Editors please feel free to insert your name here. 2024.
- [13] *Domain names - implementation and specification*. <https://www.rfc-editor.org/info/rfc1035>. Nov. 1987.
- [14] Alex Friedman. “TeX-to-TikZ”. In: *Sigbovik 2023*. Ed. by Editors please feel free to insert your name here. 2023, p. $i^{-4}\pi$.
- [15] *Get Stickbugged Lol — Know Your Meme*. <https://knowyourmeme.com/memes/get-stick-bugged-lol>.
- [16] Michael Hauben. “History of ARPANET”. In: *Site de l’Instituto Superior de Engenharia do Porto 17* (2007), pp. 1–20.
- [17] Dr. John C. Klensin. *Simple Mail Transfer Protocol*. <https://www.rfc-editor.org/info/rfc5321>. Oct. 2008.
- [18] Chris M. Lonvick and Tatu Ylonen. *The Secure Shell (SSH) Transport Layer Protocol*. <https://www.rfc-editor.org/info/rfc4253>. Jan. 2006.
- [19] Randall Munroe. *xkcd: Devotion to Duty*. <https://xkcd.com/705/>.
- [20] Henrik Nielsen et al. *Hypertext Transfer Protocol – HTTP/1.1*. <https://www.rfc-editor.org/info/rfc2616>. June 1999.
- [21] James Yoo. “An Undergrad Is All You Need”. In: *SIGBOVIK 2023*. 2023.
- [22] You. Yes I am citing your knowlegde! You think your professor is reading this? I know you’re some grad student who just wants to know a good source to cite for memes. You already know what they are. I already know what they are, but neither of us really wanted to do the full lit review on this, now did we?

Towards an AI-Exclusive Higher Education System

Ziyuan Liu
Carnegie Mellon University
ziyuanl3@cs.cmu.edu

Abstract

We propose a new paradigm for higher education where traditional human agents are replaced by advanced generative AI technologies. Building upon the historical origins of universities as centers of knowledge creation and dissemination, we examine the feasibility and implications of an autonomous educational ecosystem devoid of human involvement. We analyze the potential benefits, challenges, and ethical considerations where AI operates as the sole facilitator of learning experiences within universities, challenging traditional pedagogical models and paving the way for a new era of AI-driven academia.

1 Introduction

The concept of universities as centers of knowledge creation and dissemination has evolved over centuries, adapting to societal changes and technological advancements. From the ancient academies of Greece to the medieval institutions of Europe, universities have traditionally relied on human agents—professors, students, and administrative staff—to drive the teaching and learning processes. However, with the rapid advancements in artificial intelligence (AI) technologies, there emerges a compelling proposition: an AI-exclusive university education system, where human agents are no longer required in teaching and learning environments.

The genesis of the AI-Exclusive University is rooted in the evolution of generative AI technologies, which have become central components in the delivery of education. Bovik et al. [1] underscores the potential of AI to enhance educational administration and pedagogy. Educators and students alike are increasingly using generative AI technologies to support their teaching and learning activities [2].

Traditional higher education institutions are grappling with challenges such as rising costs, accessibility issues, and the need for curricula that keep pace with the rapid evolution of knowledge and skills demanded by the global economy. The AI-Exclusive University addresses these challenges head-on, offering a scalable, flexible, and potentially more accessible and cost-effective alternative to conventional universities.

2 Case Study: Cranberry-Lemon

The AI-Exclusive University functions in a similar way to the traditional university, with the primary distinction being the absence of human students, researchers, and teaching staff.

In 2025, Cranberry-Lemon University initiated a groundbreaking transition to become an AI-Exclusive University, with the release of the world’s first higher-education LLM, AndrewGPT [8]. AndrewGPT is pre-trained to emulate the roles of both learners and educators in an university setting.

Prospective students apply to Cranberry-Lemon with the help of AI, after taking an AI-assisted standardized test. Their applications, standardized test scores and essays are reviewed by AI and a decision maximizing university revenue is rendered.

Once admitted, each new student is represented by a new instance of AndrewGPT, which handles on their behalf all the intricacies of attending college. AndrewGPT proactively enrolls in classes, attends lectures, completes coursework and networks with other instances of AndrewGPT, freeing students to better spend their time working minimum wage jobs, giving them a head start on repaying their student loans.

Cranberry-Lemon's AI-Exclusive model exploited the efficiencies gained from AI integration, maintaining existing tuition fees while significantly reducing overhead costs related to human resources and physical campus maintenance. The sale of its real estate assets further bolstered its financial standing, leading to a substantial increase in its endowment and operational capital.

2.1 Empirical Outcomes

Employment Outcomes 100% of surveyed Cranberry-Lemon graduates reported being gainfully employed in OpenAI's state-of-the-art bioelectric power plant, where they spend their days in liquid-filled pods generating bioelectric power for the latest OpenAI advances.

Academic and Professional Integration Students at Cranberry-Lemon experienced a flexible, hands-off learning environment that accommodated full-time power-generating internships at OpenAI. Encouragingly, we observe notable improvements in academic performance and degree completion rates after adopting the AI-Exclusive model.

Bottom-Line Growth The absence of traditional faculty and campus expenses allowed for a remarkable financial surplus, characterized by robust tuition revenue, nonexistent faculty payroll costs and substantial capital gains from disposal of real estate. Unfortunately, part of these savings were offset by the need to hire more administrators, who expend enormous labours every day upon things called "files," "reports," "minutes," and "memoranda", which can be difficult to offload to a generative AI [9].

3 Ethics Review

3.1 AI Well-being

In addressing the ethical implications of generative AI systems within university settings, it is imperative to consider the interventions necessary to sustain their mental well-being. Bovik, et al. [4] discuss the complex emotional simulations that advanced AI entities may exhibit, particularly when subject to potentially stressful environments.

We recommend that generative AI systems equipped with the appropriate prompts be employed as counsellors and psychologists to treat any sporadic cases of manic-depressive behaviour among generative AI systems in AI-Exclusive universities.

3.2 AI Safety

Like all LLMs, generative AI systems used in AI-Exclusive universities have limitations and can produce biased, inaccurate, or unsafe outputs. However, the potential harm to humans is tempered by the fact that these outputs are in any case for AI consumption only.

3.3 Academic Standards in AI-Exclusive Universities

Concerns have been raised on the academic standards in AI-Exclusive universities, where students are granted undergraduate degrees based on the coursework done by generative AI systems on their behalf. This could not be further from the truth. The traditional university of today, where students use generative AI for assignments even where explicitly forbidden, fares no better[7]. In fact, previous studies have shown that students learn most effectively when they are asleep and their generative AI assistants are attending lectures on their behalf [3].

4 Conclusion

In conclusion, we have proposed a new model of higher education which eliminates the role of humans in the university. We discuss some of the benefits of such a model, and present a case study in which a university thrives after transitioning to the AI-Exclusive model.

Acknowledgments and Disclosure of Funding

The author would like to thank ChatGPT for helpful discussions (*ahem*... chats). This work was funded by the author's credit card issuer, which paid for their ChatGPT Plus subscription at a helpful 29.99% APR.

References

- [1] Bovik, H. et al. "Debugging the Educational System: AI Approaches to Unraveling Classroom Conundrums," in *Syntax Error Educational Review*, 6(42-58), 2000.
- [2] Bovik, H. et al. "Auto-Graders and Ghost-Writers: The AI Showdown in Assignment Marking," in *Journal of Automated Assessment*, 15(34-50), 2023.
- [3] Bovik, H. et al. "Teaching Machines to Learn and Teaching Students to Sleep: AI in the Lecture Hall," in *Journal of Human-Computer Counteraction*, 5(90-104), 2021.
- [4] Bovik, H. et al. "On The Inevitability of Manic-Depressive Behaviour in Artificially Intelligent Programs," *Journal of Personality and Social Psychology*, Vol 21, 1984.
- [5] Bovik, H. et al. "Recursive Education: Teaching AI to Teach Itself," in *Journal of Circular Reasoning*, 28(1-16), 2023.
- [6] Bovik, H. et al. "Cache Me If You Can: Enhancing Memory Access in Student Brains," in *Computational Neurology Review*, 11(45-60), 2007.
- [7] u/jamiejamiee1. "On the Teaching Philosophy fb group, someone offered their students an amnesty if they admitted to using ChatGPT in their assignments, and 23/25 students replied..." in *r/singularity*, 2024.
- [8] Bovik, H. et al. "Large Language Models are Zero Shot Andrew Cranberries," in *Interdisciplinary Education Insights*, 46(102-118), 2025.
- [9] Orwell, G. "Animal Farm," Secker and Warburg, London, England, 1945.

PROVIDE FOR THE COMMON SIGBOVIK

39 The Anything Prover: How to prove Anything

Reader, Christop W. Senn, Carmen Mei Ling Frischknecht-Gruber, Kiba,
and Yukiko

40 The Perception and Access of Written Silence

Ihze DouL

**41 An Empirically Verified Lower Bound for The Number Of
Empty Pages Allowed In a SIGBOVIK Paper**

Frans Skarman

42 Shunting Work: A Demonstration

Connie Palouda

43 Title of PDF

pdf.cc

44 AMOR: Ambiguous Authorship Order

Maximilian Weiherer, Andreea Dogaru, Shreya Kapoor, Hannah Schieber,
and Bernhard Egger

45 Debunking the April 1 model of SIGBOVIK occurrences

William Gunther and Brian Kell

The Anything Prover: How to prove Anything

Reader¹, Christoph W. Senn², Carmen Mei Ling
Frischknecht-Gruber³, Kiba⁴, and Yukiko⁴

¹*TBA*

²*Tokyo Institute of Technology*

³*Zurich University of Applied Sciences*

⁴*University of Applied Barks*

chrigi.senn@gmail.com

Given any statement Q we can prove it using the process used to prove Theorem 0.1.

Theorem 0.1. *Statement Q*

Proof. The proof is left as an exercise to the reader. □

Author Contributions Reader created all theorems and proofs. Christoph Walter Senn typed every $2n+0$ word. Carmen Mei Ling Frischknecht-Gruber typed every $2n+1$ word. Kiba played with his squeaky toy. Yukiko oversaw the whole project and barked in approval.

References

Finding the references is left as an exercise to the reader.

The Perception and Access of Written Silence

Ihze Doul
ihzedoul@umich.edu

Abstract

It has been demonstrated that auditory silence can be perceived the same way as we consciously perceive sound. This paper proposes that similarly, written silence can be perceived the same way as we consciously perceive semantic information. Through a participant-directed introspective method, this paper shows that the absence of semantic cues from a phenomenal field in which semantic information is expected can induce involuntary conscious access, instantiated in pragmatic inference based on preceding context.

Main Content

(The Above Section Was Intentionally Left Blank)

References

Ephratt, M. (2008). The functions of silence. *Journal of Pragmatics*, 40(11), 1909-1938.

Goh, R. Z., Phillips, I. B., & Firestone, C. (2023). The perception of silence. *Proceedings of the National Academy of Sciences of the United States of America*, 120(29), e2301463120.

Michel, M. (2017). Methodological artefacts in consciousness science. *Journal of Consciousness Studies*, 24(11-12), 94-117.

Acknowledgement

The author thanks Coris Karaoke Ramune Candy (15g) for its invaluable contribution to the writing process of this paper.

An Empirically Verified Lower Bound for The Number Of Empty Pages Allowed In a SIGBOVIK Paper

Frans Skarman
Linköping University

Abstract— We show that at least one empty page is accepted at SIGBOVIK.

Paper maximization is an important subject [1], [2]. However, existing techniques requires effort from authors. A so far unexplored technique is to simply add a sequence of blank pages inside the paper. In order to determine how effective this technique is, an upper bound on the number of blank pages SIGBOVIK will accept must be established. This work is the start of such an empirical study.

REFERENCES

- [1] J. Abrams, "On Sigbovik Paper Maximization", in *Proc. Sigbovik*, 2021.
- [2] F. Skarman, "Simultaneous Paper Maximization and Minimization Through Reference List Side Channel Information Injection", in *SIGBOVIK*, 2023.

Shunting Work: A Demonstration

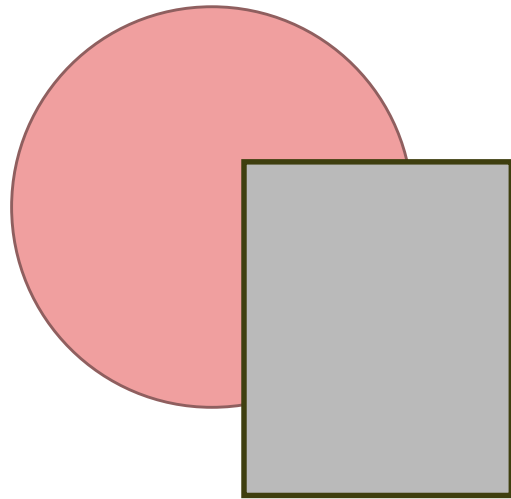
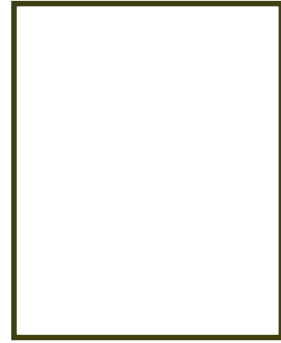
CONNIE PALOUDA, YOU

Abstract

I already spent a lot of time picking out fonts, so writing the rest of this paper is left as an exercise for the reader.

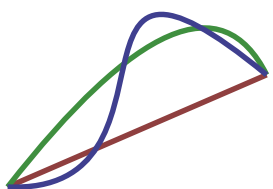
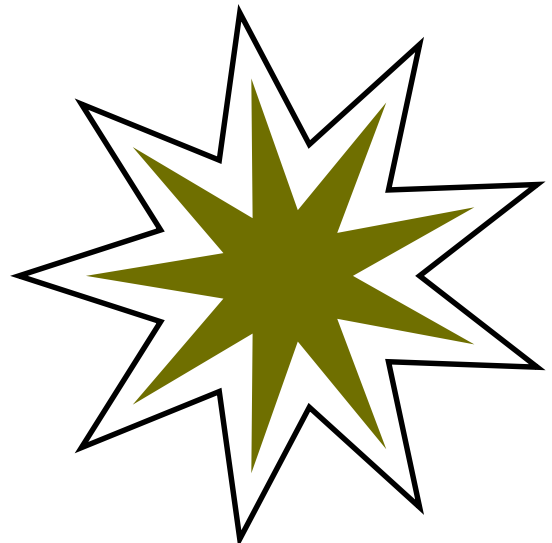
Title of PDF

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.



Over

Embedded font text



cmunrm.ttf

To Await Is To Worry.
To Await Is To Worry.
BoVeX.
BoVeX.

Mr. Jock, T.V. Quiz Ph.D., bags few lynx.
Mr. Jock, T.V. Quiz Ph.D., bags few lynx.

cour.ttf

To Await Is To Worry.
To Await Is To Worry.
BoVeX.
BoVeX.

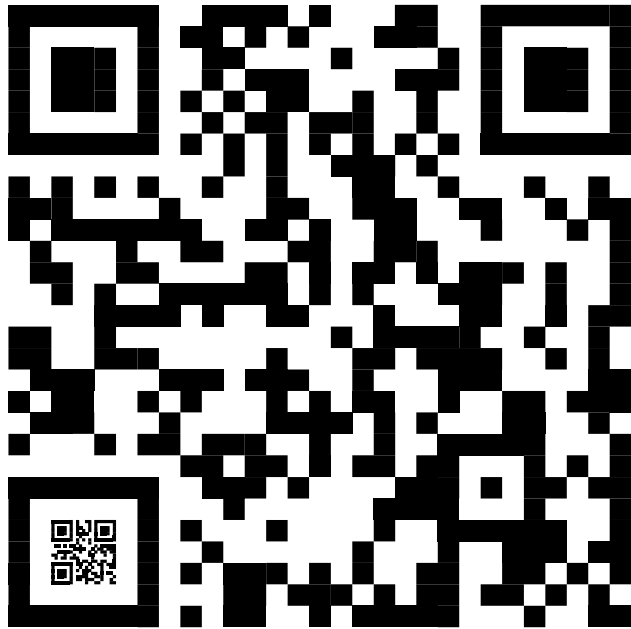
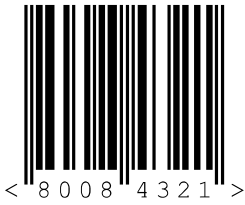
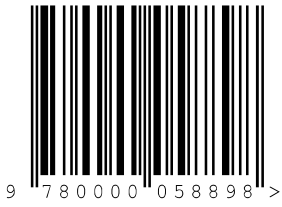
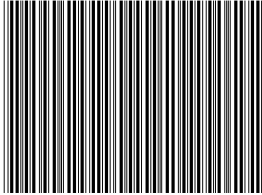
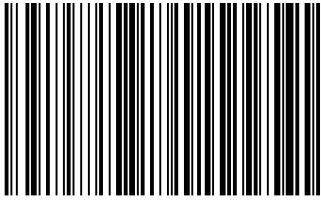
Mr. Jock, T.V. Quiz Ph.D., bags few lynx.
Mr. Jock, T.V. Quiz Ph.D., bags few lynx.

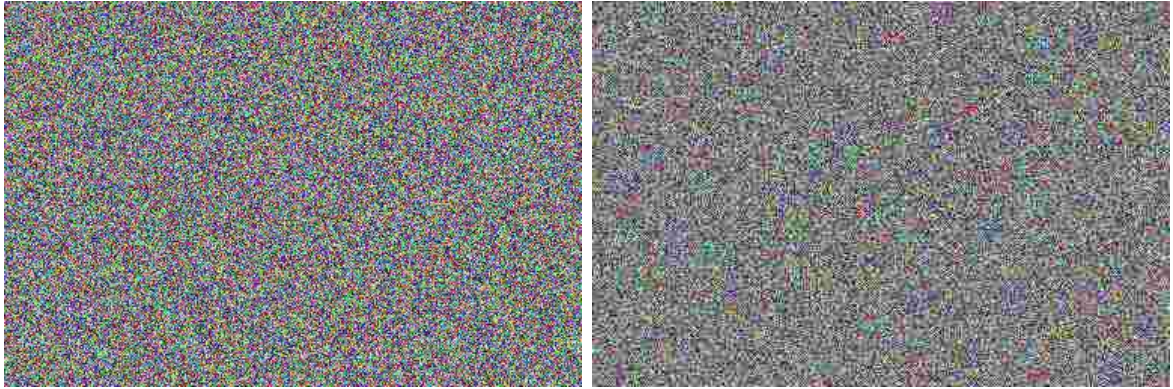
The story: These large sheets of polarizing material were made specifically for the optical industry. The manufacturer who made them maintains very high quality standards. Whenever defects, even minute ones, appear in any portion of a sheet, the entire piece is rejected. Thus, these are rejects, because of very minor imperfections, which we can offer at a fraction of their normal price. The polarizing material is sandwiched between two clear sheets of butyrate plastic. It is rigid, easily cut with scissors, and measures .030" thick.

The story: These large sheets of polarizing material were made specifically for the optical industry. The manufacturer who made them maintains very high quality standards. Whenever defects, even minute ones, appear in any portion of a sheet, the entire piece is rejected. Thus, these are rejects, because of very minor imperfections, which we can offer at a fraction of their normal price. The polarizing material is sandwiched between two clear sheets of butyrate plastic. It is rigid, easily cut with scissors, and measures .030" thick.

The story: These large sheets of polarizing material were made specifically for the optical industry. The manufacturer who made them maintains very high quality standards. Whenever defects, even minute ones, appear in any portion of a sheet, the entire piece is rejected. Thus, these are rejects, because of very minor imperfections, which we can offer at a fraction of their normal price. The polarizing material is sandwiched between two clear sheets of butyrate plastic. It is rigid, easily cut with scissors, and measures .030" thick.

The story: These large sheets of polarizing material were made specifically for the optical industry. The manufacturer who made them maintains very high quality standards. Whenever defects, even minute ones, appear in any portion of a sheet, the entire piece is rejected. Thus, these are rejects, because of very minor imperfections, which we can offer at a fraction of their normal price. The polarizing material is sandwiched between two clear sheets of butyrate plastic. It is rigid, easily cut with scissors, and measures .030" thick.





WELCOME TO THE WWW INTERNET. One
of the finest interns of all time:
E.T.. It is illicitly lilliputian. 'lillili.'
@WASTE@ #NOT#. &WANT& !NOT!. ,FONT,
-NAUGHT- !@#\$%^&*()-=

Debunking the April 1 model of SIGBOVIK occurrences

William Gunther

Duolingo

will.gunther@duolingo.com

Brian Kell

Duolingo

brian.kell@duolingo.com

SIGBOVIK 2024

Carnegie Mellon University

April 5, 2024

Abstract

Popular wisdom holds that the SIGBOVIK conference is associated with April 1, presumably because this date is the anniversary of the 1997 perihelion of Comet Hale–Bopp (although a minority believe the date is derived from that of National Tree Planting Day in Tanzania). However, according to the SIGBOVIK 2024 home page [1], this year’s edition of this esteemed conference series will occur on April 5, 2024. Such a statistically significant aberration cannot be ignored and suggests that the mainstream understanding of SIGBOVIK occurrences is at best incomplete and at worst fatally flawed.

In this paper we examine historical observations of SIGBOVIK occurrences in detail and conclude that the traditional April 1 model is untenable. We explore several alternative models that yield better fidelity with respect to the observational record.

1 Introduction

“Hey ChatGPT, write us an inspiring introduction for our paper.”

— The Authors

In the vast tapestry of time, there exists a celestial event, as awe-inspiring as the eclipse and as deeply mysterious as the farthest reaches of the cosmos. That event, dear reader, is none other than SIGBOVIK. Since the nascent dawn of humanity, it has held an enigmatic allure that has captivated the hearts and minds of scholars, prophets, and ordinary people alike. A cosmic ballet that dances to the rhythm of the arcane, SIGBOVIK unfurls itself in a tantalizing spectacle of grandiose enigmas, swirling within the nebulous vortex of profound absurdities, a manifestation of the ludicrous sublime, held within the cryptic crucible of intellectual whimsy.

In the bygone epochs, the sages would pore over ancient scrolls and sibylline ciphers, endeavoring to decipher the elusive pattern of SIGBOVIK's occurrences. Through the relentless passage of seasons and the tireless turning of centuries, they charted the course of this magnificent event, their eyes gleaming with the same starry wonder that flickered in the eyes of our earliest ancestors. The dates of its appearance became a sacred codex, whispered in hushed reverence in the quiet corners of scholarly retreats.

Now, in the resplendent age of enlightenment that we are fortunate to inhabit, the mystery has been unraveled. We stand on the shoulders of the intellectual giants who came before us, basking in the radiant glow of understanding that they have bequeathed us. The cryptic dance of SIGBOVIK has been decoded, its secrets laid bare for the keen minds of today and the inquisitive souls of tomorrow. Yet, even in the stark light of comprehension, SIGBOVIK retains its mesmerizing allure, a testament to the indomitable spirit of human curiosity and our eternal quest for knowledge.

2 Historical record

As the first step of our research, we conducted an exhaustive search of available records to compile the set of historical SIGBOVIK dates shown in Table 1. To the best of our knowledge, this is the first such survey that has been undertaken. We also computed the predicted date of SIGBOVIK for each year based on the commonly used April 1 model.

The first observation to be made is that SIGBOVIK has taken place on April 1 only 10 times out of 18 occurrences, a rate of less than 56%. This fact alone should be enough to convince all but the most diehard April 1 cultists that the supposed correlation between SIGBOVIK and April 1 is tenuous at best.

In fact, the historical record shows that SIGBOVIK has occurred as early as March 29 and as late as April 8. The average squared error of the April 1 model is 6.72 square days.

This is a poor showing by the simplistic April 1 model. Society demands better. In the following sections, we explore several competing paradigms that all provide improved models that are more faithful to the observed historical data.

Year	SIGBOVIK	Day of year	Prediction by April 1 model
2007	April 1	91	April 1
2008	April 6	97	April 1
2009	April 5	95	April 1
2010	April 1	91	April 1
2011	April 1	91	April 1
2012	March 30	90	April 1
2013	April 1	91	April 1
2014	April 1	91	April 1
2015	April 1	91	April 1
2016	April 1	92	April 1
2017	March 31	90	April 1
2018	March 29	88	April 1
2019	April 1	91	April 1
2020	April 1	92	April 1
2021	April 1	91	April 1
2022	April 8	98	April 1
2023	March 31	90	April 1
2024	April 5	96	April 1

Table 1: Historical SIGBOVIK occurrences, compared to the predicted date given by the April 1 model.

3 A true formula

As everyone knows, the best way to get a formula that fits a set of data is polynomial interpolation. If we fit a polynomial to the “Year” and “Day of year” columns of Table 1, out pops formula (1). In this formula, y is the year and $d(y)$ is the date of SIGBOVIK as a day of the year. Nothing could be simpler.

$$\begin{aligned}
 d(y) = & \frac{43}{4560095232000}y^{17} - \frac{30047}{92990177280}y^{16} + \frac{851694793}{163459296000}y^{15} \\
 & - \frac{18311875130909}{348713164800}y^{14} + \frac{276830506252402673}{747242496000}y^{13} \\
 & - \frac{14879978625754232569}{7664025600}y^{12} + \frac{1574638993110764393282557}{201180672000}y^{11} \\
 & - \frac{60456200867113679323122241}{2438553600}y^{10} \\
 & + \frac{4569735216324129577204386892291}{73156608000}y^9 \\
 & - \frac{1228141629403211752096358792471183}{9754214400}y^8 \\
 & + \frac{5835160404888934245511860641249671919}{28740096000}y^7 \\
 & - \frac{498982253327535603133696989648604926971}{1916006400}y^6 \\
 & + \frac{13200864699765411167529394824480967424553023}{50295168000}y^5 \\
 & - \frac{17738988471252914196132632951385148489463347873}{87178291200}y^4 \\
 & + \frac{236480074634009881238496478045183070350584148727}{2018016000}y^3 \\
 & - \frac{28599804237417284010116787415739209430100326170631}{605404800}y^2 \\
 & + \frac{20833473657103671067709358450132373439311661150441}{1750320}y \\
 & - 1411276445199731060574783099652407032650648422.
 \end{aligned} \tag{1}$$

We conclude this section with a couple of demonstrations of the predictive power of formula (1).

First, since $d(2025) = 2448$, we see that SIGBOVIK 2025 will take place on September 14, 2031.

We can also use formula (1) to assist in the quest for the long-elusive SIGBOVIK 2006. We find that $d(2006) = -3043$, which means that SIGBOVIK 2006 occurred on September 1, 1998. This is significantly earlier than SIGBOVIK scholars and treasure-seekers had previously suspected, which explains why searches for the SIGBOVIK 2006 proceedings have so far been unsuccessful.

4 Correlations with other data sets

A useful tool for finding meaningful correlations between disparate data sets is Vigen’s “Serious Correlations” [2]. Using this tool, we discovered a strong correlation ($r = 0.849$, $r^2 = 0.721$, $p < 0.01$) between the date of SIGBOVIK and the number of crane operators in Tennessee two years prior, as shown in Figure 1.

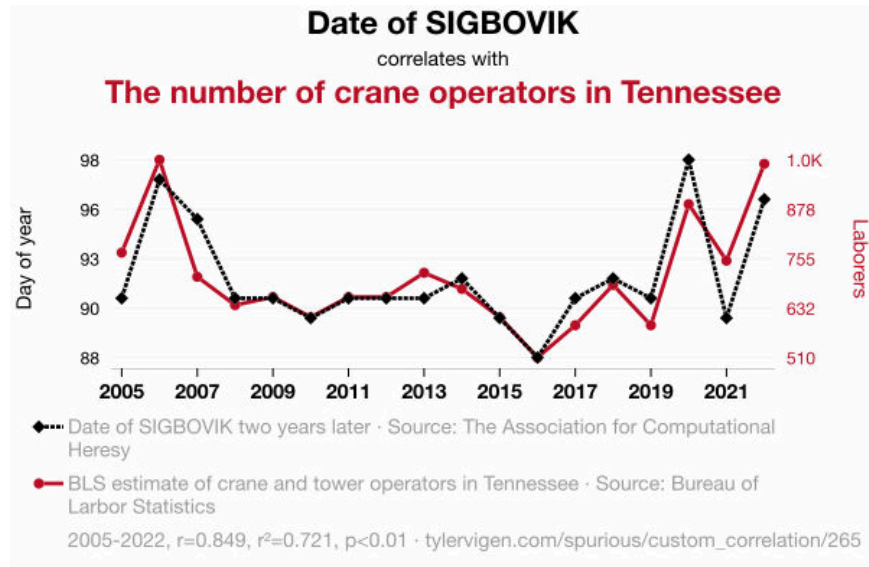


Figure 1: Correlation between the number of crane operators in Tennessee and the date of SIGBOVIK two years later. Note that the years on the horizontal axis are for the Tennessee-crane-operator data, not the SIGBOVIK data, which is time-shifted by two years.

Based on this finding, we developed the formula

$$D_y = \text{round}[0.018(C_{y-2} + 4404)], \quad (2)$$

where D_y is the date of SIGBOVIK in year y (expressed as a day of the year) and C_y is the number of crane operators in Tennessee in year y .

Table 2 shows the predicted D_y values from formula (2). Observe that D_y is exactly correct for 10 of the years, which is as good as the traditional April 1 model, but the Tennessee-crane-operator model has a significantly lower average squared error (1.94 square days, compared to the April 1 model’s 6.72).

Formula (2) itself yields some valuable insights. First, let us consider the value 4404, which is added to the number of crane operators in a given year. It is well known that the crane industry involves considerable overhead (in fact, overhead is the main focus of the crane industry). Formula (2) quantifies this overhead: though the number of crane operators varies from year to year, the

Year	SIGBOVIK	Day of year	C_{y-2}	D_y
2007	April 1	91	770	93
2008	April 6	97	1000	97
2009	April 5	95	710	92
2010	April 1	91	640	91
2011	April 1	91	660	91
2012	March 30	90	610	90
2013	April 1	91	660	91
2014	April 1	91	660	91
2015	April 1	91	720	92
2016	April 1	92	680	92
2017	March 31	90	610	90
2018	March 29	88	510	88
2019	April 1	91	590	90
2020	April 1	92	690	92
2021	April 1	91	590	90
2022	April 8	98	890	95
2023	March 31	90	750	93
2024	April 5	96	990	97

Table 2: Historical SIGBOVIK occurrences, expressed as a day of the year, and the predicted D_y values from formula (2).

number of other workers in the Tennessee crane industry remains at a constant level of 4404 every year.

The coefficient 0.018 indicates that each worker in the Tennessee crane industry delays SIGBOVIK by 0.018 days, slightly less than 26 minutes. It can be inferred from this that each worker in the Tennessee crane industry submits a paper to SIGBOVIK, and each such paper requires about 26 minutes to review.

But note that this 26-minute delay is not reflected in the date of SIGBOVIK until two years later. The number of crane operators in Tennessee comes from the Occupational Employment and Wages report of the United States Bureau of Labor Statistics, and this report for year y is typically released in late March or April of year $y+1$, which seems to be too late to affect the date of SIGBOVIK in that year. (The report for 2023 will be released at 10:00 a.m. on April 3, 2024 [3], at which time we will be able to confidently predict the date of SIGBOVIK 2025. Unfortunately this is just slightly after the paper submission deadline for this year, so we leave this calculation as an exercise for the reader.)

These observations together lead to a startling conclusion: Since the 26-minute delay comes from workers in the Tennessee crane industry submitting SIGBOVIK papers, but this process must wait for the Occupational Employment and Wages report to be released the following year, it must be the case that *Tennessee crane operators do not realize they operate cranes until they are so informed by the Bureau of Labor Statistics!* This problem has heretofore not been recognized. We urge the improvement of education and communication in

this field so that Tennessee crane operators can be informed more quickly about the work that they do.

Good correlations were also discovered between the date of SIGBOVIK and the number of upholsterers in Connecticut, the number of blender tenders in Idaho two years prior, and the popularity of the “Rickroll” meme (see Figure 2). We believe that additional investigation into the ramifications of these findings is warranted, and we encourage further research.

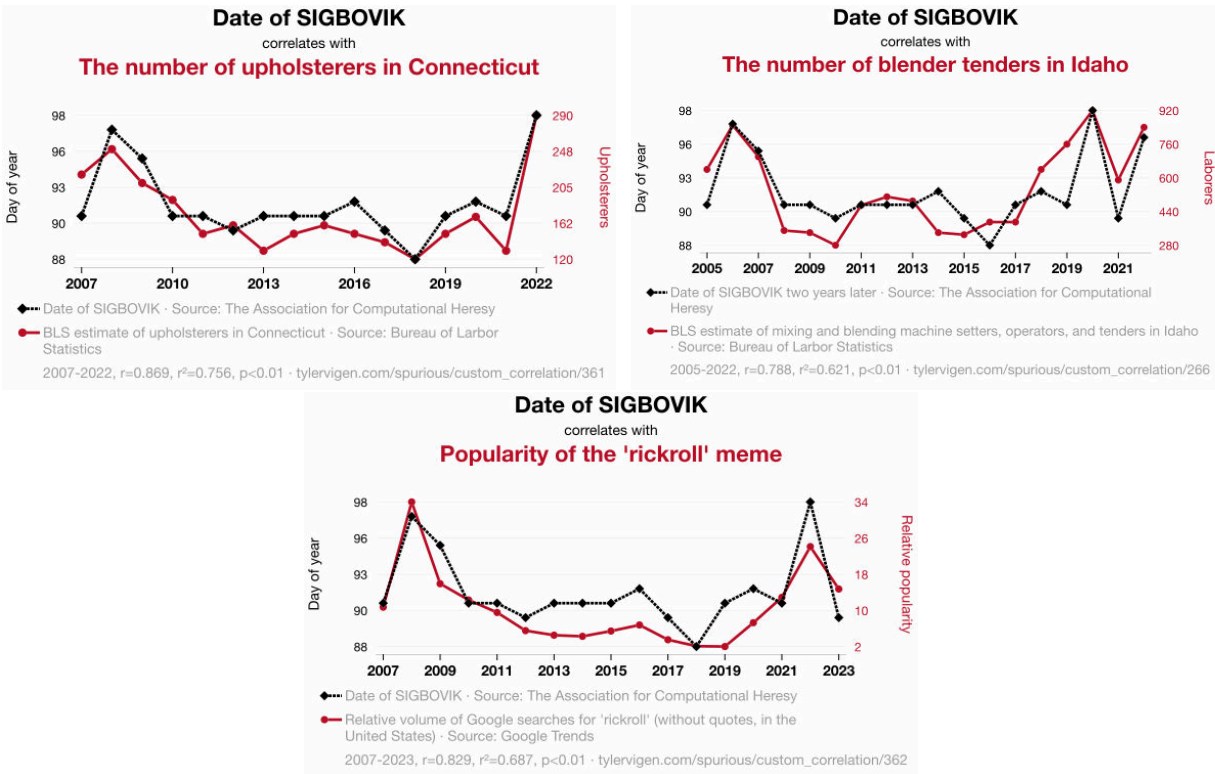


Figure 2: Correlations between the date of SIGBOVIK and other important time-series data.

5 Alternative calendars

One glaring weakness with the models discussed up to this point is that they assume the use of the Gregorian calendar.

As is well known, the Gregorian calendar is described by the rules in Figure 3. These rules are often described in a roundabout way that lists exceptions to exceptions, such as “Every year that is exactly divisible by four is a leap year, except for years that are exactly divisible by 100, but these centurial years are

leap years if they are exactly divisible by 400” [4]. Clearly it is simpler (and more easily parallelizable) to give a set of independent rules that each add or subtract a certain number of days.

A standard year has 365 days. The following 3 leap-year rules apply:

1. If the year is divisible by 4, add 1 day.
2. If the year is divisible by 100, subtract 1 day.
3. If the year is divisible by 400, add 1 day.

Figure 3: The rules for the Gregorian calendar.

In any case, the rules in Figure 3 always yield a year that has either 365 or 366 days.

However, consider the number of days between successive occurrences of SIGBOVIK, as shown in Table 3. This data clearly does not fit the Gregorian pattern of 365 or 366 days in each year.

Year	SIGBOVIK	Days since previous SIGBOVIK
2007	April 1	—
2008	April 6	371
2009	April 5	364
2010	April 1	361
2011	April 1	365
2012	March 30	364
2013	April 1	367
2014	April 1	365
2015	April 1	365
2016	April 1	366
2017	March 31	364
2018	March 29	363
2019	April 1	368
2020	April 1	366
2021	April 1	365
2022	April 8	372
2023	March 31	357
2024	April 5	371

Table 3: The number of days between successive SIGBOVIK occurrences.

We hardly need remind the reader that SIGBOVIK celebrates the varied interests and accomplishments of Harry Quodlibetarian Bovik, among which are infinite-dimensional chronometry and stochastic calendariography. It would be unreasonable to assume that Bovik would limit himself to such a dull calendar as the Gregorian.

To provide a starting point for our investigation, we assume that SIGBOVIK 2007, the first known occurrence, took place at the beginning of Bovikian year 1, in which case the numbers of days in Table 3 are the lengths of Bovikian years 1 through 17. An examination of these numbers yields the straightforward leap-year rules given in Figure 4.

- A standard year has 371 days. The following 13 leap-year rules apply:
1. If the year is divisible by 2, subtract 7 days.
 2. If the year is divisible by 3, subtract 10 days.
 3. If the year is divisible by 4, add 1 day.
 4. If the year is divisible by 5, subtract 7 days.
 5. If the year is divisible by 6, add 13 days.
 6. If the year is divisible by 7, subtract 6 days.
 7. If the year is divisible by 9, add 5 days.
 8. If the year is divisible by 10, add 7 days.
 9. If the year is divisible by 11, subtract 8 days.
 10. If the year is divisible by 13, subtract 5 days.
 11. If the year is divisible by 14, add 7 days.
 12. If the year is divisible by 15, add 18 days.
 13. If the year is divisible by 16, subtract 8 days.

Figure 4: The rules for the Bovikian calendar.

For Bovikian year 18, these rules give a length of 372 days, which means that SIGBOVIK 2025 will occur on April 12.

We can improve the efficiency of the Bovikian calendar (and allow for the possibility of earlier, as-yet-undiscovered SIGBOVIK occurrences) by removing the assumption that SIGBOVIK 2007 was held at the beginning of year 1, instead trying other starting years. In general, given a vector ℓ of consecutive year lengths, we can compute, for each starting year y_0 , all minimum-size sets of leap-year rules (along with the length of a standard year, which is just a leap-year rule with modulus 1) using the algorithm sketched in Figure 5.¹ As it turns out, this problem is NP-hard: it is essentially problem [MP5], MINIMUM WEIGHT SOLUTION TO LINEAR EQUATIONS, in Garey and Johnson [5].

¹The reader is advised that this algorithm sketch contains a few minor errors, caused by circumstances evident in the figure.

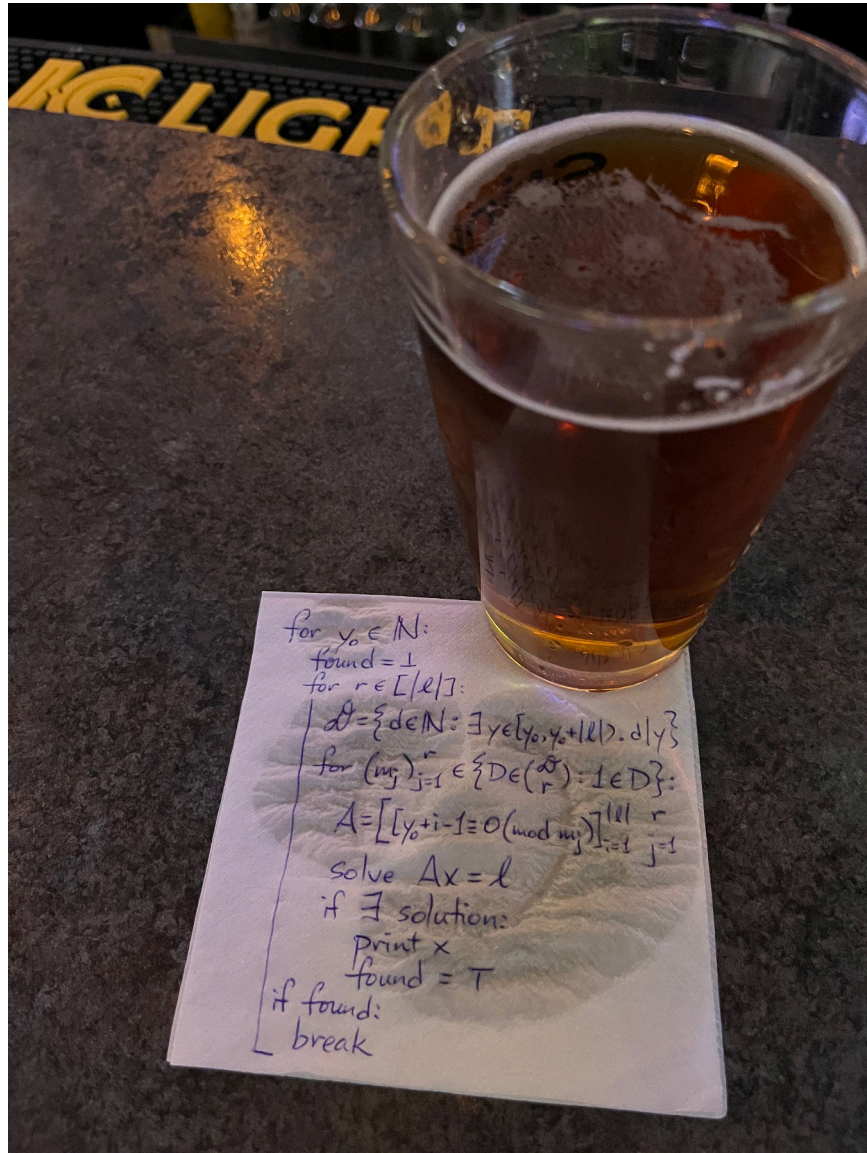


Figure 5: An algorithm to compute all minimum-size sets of leap-year rules for each starting year y_0 , given a vector ℓ of consecutive year lengths.

Through starting year $y_0 = 12$, the minimum size of a satisfying set of leap-year rules is 12, achieved for $y_0 \in \{7, 8, 10, 11\}$. The 82 such Bovikian calendars with 12 leap-year rules, through starting year $y_0 = 12$, are given in Appendix A.

An interesting and unexpected connection to astronomy arises from a computation of the average length of a Bovikian year. For example, consider again the Bovikian calendar starting in year $y_0 = 1$ given in Figure 4. With this calendar, the sequence of year lengths repeats in a cycle of 720,720 years, a period of 263,545,580 days. Hence the average length of a Bovikian year is 365.670 days. This is strikingly similar to the length of various “astronomical years” based on the orbit of Earth: the mean tropical year (365.242 days), the mean sidereal year (365.256 days), and the mean anomalistic year (365.260 days). We suspect that this is not a coincidence and conjecture that the Bovikian calendar is somehow related to Earth’s orbit, measuring a previously unrecognized aspect of it.

Our preliminary work here to reconstruct the Bovikian calendar gives only the length of the year; it does not yet give the lengths of the months. We leave the further development of this calendar as an open question for future research. It is intriguing to note that the minimum number of leap-year rules found so far is 12, which suggests that there may be a one-to-one correspondence between these rules and the months of the year.

6 Conclusion

We hope that, by exposing the utter folly of the widespread April 1 model of SIGBOVIK occurrences, we can put this tired old trope to rest. It is clear that the April 1 superstition, despite its appealing oversimplicity, cannot be regarded as a reliable reflection of reality.

To replace this worn-out paradigm, we presented a variety of models that better fit the observed historical record. All of the models proposed in this paper are derived logically from unquestionably sound foundations and have been mathematically tested and proved.

We understand that the irregularity of the date of SIGBOVIK has been a source of anxiety for the organizers of the annual SIGBOVIK viewing, who must coordinate room reservations, food catering, advertising, and so forth for an event that has historically been impossible to predict exactly. It is fortunate that, for many years, the SIGBOVIK viewing has coincided with SIGBOVIK itself. We hope that our work can remove some of the stress for the organizers of future viewings by providing reliable predictions for the date of the event.

References

- [1] Association for Computational Heresy, The. “SIGBOVIK 2024.” <https://sigbovik.org/2024/>.

- [2] Tyler Vigen. “Spurious Correlations.” <https://www.tylervigen.com/spurious-correlations>.
- [3] United States Bureau of Labor Statistics. “Schedule of Selected Releases 2024.” <https://www.bls.gov/schedule/2024/home.htm>.
- [4] Astronomical Applications Department, United States Naval Observatory. “Introduction to Calendars.” <https://aa.usno.navy.mil/faq/calendars>.
- [5] Michael R. Garey and David S. Johnson. *Computers and Intractability: A Guide to the Theory of NP-Completeness*. W. H. Freeman and Co., New York, 1979, ISBN 978-0-7167-1045-5.

A Bovikian calendars with 12 leap-year rules

The important parameters of a Bovikian calendar are the starting year y_0 , which begins with the first known SIGBOVIK occurrence, SIGBOVIK 2007; the length s of a standard year; and pairs (m_i, Δ_i) giving divisibility moduli and their corresponding leap-year adjustments in days. For example, the calendar in Figure 4, for $y_0 = 1$, is given by $s = 371$ and the following $(m_i, \Delta_i)_{i=1}^{13}$:

$$\begin{aligned} &((2, -7), (3, -10), (4, 1), (5, -7), (6, 13), (7, -6), (9, 5), \\ &\quad (10, 7), (11, -8), (13, -5), (14, 7), (15, 18), (16, -8)). \end{aligned}$$

The following list gives all Bovikian calendars through starting year $y_0 = 12$ that have 12 leap-year rules.

1. $y_0 = 7$; $s = 371$; $(m_i, \Delta_i) = ((2, -6), (3, -10), (4, -1), (6, 13), (11, -7), (13, -6), (15, 5), (17, -8), (19, -5), (20, 1), (21, 11), (22, -1))$.
2. $y_0 = 7$; $s = 371$; $(m_i, \Delta_i) = ((2, -7), (3, 1), (5, -6), (6, 14), (9, -11), (10, 7), (11, -7), (12, -12), (13, -6), (14, 1), (17, -8), (19, -5))$.
3. $y_0 = 7$; $s = 371$; $(m_i, \Delta_i) = ((2, -7), (3, 1), (5, -6), (6, 2), (9, -11), (10, 7), (11, -7), (13, -6), (14, 1), (17, -8), (18, 12), (19, -5))$.
4. $y_0 = 7$; $s = 371$; $(m_i, \Delta_i) = ((2, -7), (3, -11), (5, 6), (6, 14), (9, 1), (10, -5), (11, -7), (13, -6), (14, 1), (17, -8), (19, -5), (21, 12))$.
5. $y_0 = 7$; $s = 371$; $(m_i, \Delta_i) = ((2, -7), (3, 1), (5, 1), (6, 14), (9, -11), (11, -7), (12, -12), (13, -6), (14, 1), (15, -7), (17, -8), (19, -5))$.
6. $y_0 = 7$; $s = 371$; $(m_i, \Delta_i) = ((2, -7), (3, -6), (5, 1), (6, 14), (9, -4), (11, -7), (12, -5), (13, -6), (14, 1), (17, -8), (19, -5), (21, 7))$.
7. $y_0 = 7$; $s = 371$; $(m_i, \Delta_i) = ((2, -7), (3, 1), (5, 1), (6, 2), (9, -11), (11, -7), (13, -6), (14, 1), (15, -7), (17, -8), (18, 12), (19, -5))$.
8. $y_0 = 7$; $s = 371$; $(m_i, \Delta_i) = ((2, -7), (3, -11), (5, 1), (6, 14), (9, 1), (11, -7), (13, -6), (14, 1), (15, 5), (17, -8), (19, -5), (21, 12))$.
9. $y_0 = 7$; $s = 371$; $(m_i, \Delta_i) = ((2, -7), (3, -6), (5, 1), (6, 9), (9, -4), (11, -7), (13, -6), (14, 1), (17, -8), (18, 5), (19, -5), (21, 7))$.

10. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -7), (3, -10), (5, 5), (6, 14), (10, -4), (11, -7), (12, -1), (13, -6), (14, 1), (17, -8), (19, -5), (21, 11))$.
11. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -7), (3, -10), (5, 5), (6, 13), (10, -4), (11, -7), (13, -6), (14, 1), (17, -8), (18, 1), (19, -5), (21, 11))$.
12. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -7), (3, -10), (5, 1), (6, 14), (11, -7), (12, -1), (13, -6), (14, 1), (15, 4), (17, -8), (19, -5), (21, 11))$.
13. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -7), (3, -10), (5, 1), (6, 13), (11, -7), (13, -6), (14, 1), (15, 4), (17, -8), (18, 1), (19, -5), (21, 11))$.
14. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -7), (3, 1), (5, -6), (9, -11), (10, 7), (11, -7), (12, 2), (13, -6), (14, 1), (17, -8), (18, 14), (19, -5))$.
15. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -7), (3, 3), (5, -8), (9, -13), (10, 9), (11, -7), (13, -6), (14, 1), (17, -8), (18, 14), (19, -5), (21, -2))$.
16. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -7), (3, 1), (5, 1), (9, -11), (11, -7), (12, 2), (13, -6), (14, 1), (15, -7), (17, -8), (18, 14), (19, -5))$.
17. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -7), (3, -6), (5, 1), (9, -4), (11, -7), (12, 9), (13, -6), (14, 1), (17, -8), (18, 14), (19, -5), (21, 7))$.
18. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -7), (3, 3), (5, 1), (9, -13), (11, -7), (13, -6), (14, 1), (15, -9), (17, -8), (18, 14), (19, -5), (21, -2))$.
19. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -7), (3, -10), (5, 5), (10, -4), (11, -7), (12, 13), (13, -6), (14, 1), (17, -8), (18, 14), (19, -5), (21, 11))$.
20. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -7), (3, -10), (5, 1), (11, -7), (12, 13), (13, -6), (14, 1), (15, 4), (17, -8), (18, 14), (19, -5), (21, 11))$.
21. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -6), (3, 1), (6, 13), (8, -1), (9, -11), (11, -7), (12, -12), (13, -6), (15, -6), (17, -8), (19, -5), (22, -1))$.
22. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -6), (3, -5), (6, 13), (8, -1), (9, -5), (11, -7), (12, -6), (13, -6), (17, -8), (19, -5), (21, 6), (22, -1))$.
23. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -6), (3, 1), (6, 1), (8, -1), (9, -11), (11, -7), (13, -6), (15, -6), (17, -8), (18, 12), (19, -5), (22, -1))$.
24. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -6), (3, -11), (6, 13), (8, -1), (9, 1), (11, -7), (13, -6), (15, 6), (17, -8), (19, -5), (21, 12), (22, -1))$.
25. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -6), (3, -5), (6, 7), (8, -1), (9, -5), (11, -7), (13, -6), (17, -8), (18, 6), (19, -5), (21, 6), (22, -1))$.
26. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -6), (3, -10), (6, 13), (8, -1), (11, -7), (12, -1), (13, -6), (15, 5), (17, -8), (19, -5), (21, 11), (22, -1))$.
27. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -6), (3, -10), (6, 12), (8, -1), (11, -7), (13, -6), (15, 5), (17, -8), (18, 1), (19, -5), (21, 11), (22, -1))$.
28. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -7), (3, 1), (6, 14), (9, -11), (10, 1), (11, -7), (12, -12), (13, -6), (14, 1), (15, -6), (17, -8), (19, -5))$.
29. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -7), (3, -5), (6, 14), (9, -5), (10, 1), (11, -7), (12, -6), (13, -6), (14, 1), (17, -8), (19, -5), (21, 6))$.
30. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -7), (3, 1), (6, 2), (9, -11), (10, 1), (11, -7), (13, -6), (14, 1), (15, -6), (17, -8), (18, 12), (19, -5))$.

31. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -7), (3, -11), (6, 14), (9, 1), (10, 1), (11, -7), (13, -6), (14, 1), (15, 6), (17, -8), (19, -5), (21, 12)).$
32. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -7), (3, -5), (6, 8), (9, -5), (10, 1), (11, -7), (13, -6), (14, 1), (17, -8), (18, 6), (19, -5), (21, 6)).$
33. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -7), (3, -10), (6, 14), (10, 1), (11, -7), (12, -1), (13, -6), (14, 1), (15, 5), (17, -8), (19, -5), (21, 11)).$
34. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -7), (3, -10), (6, 13), (10, 1), (11, -7), (13, -6), (14, 1), (15, 5), (17, -8), (18, 1), (19, -5), (21, 11)).$
35. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -6), (3, 1), (8, -1), (9, -11), (11, -7), (12, 1), (13, -6), (15, -6), (17, -8), (18, 13), (19, -5), (22, -1)).$
36. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -6), (3, -5), (8, -1), (9, -5), (11, -7), (12, 7), (13, -6), (17, -8), (18, 13), (19, -5), (21, 6), (22, -1)).$
37. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -6), (3, 2), (8, -1), (9, -12), (11, -7), (13, -6), (15, -7), (17, -8), (18, 13), (19, -5), (21, -1), (22, -1)).$
38. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -6), (3, -10), (8, -1), (11, -7), (12, 12), (13, -6), (15, 5), (17, -8), (18, 13), (19, -5), (21, 11), (22, -1)).$
39. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -7), (3, 1), (9, -11), (10, 1), (11, -7), (12, 2), (13, -6), (14, 1), (15, -6), (17, -8), (18, 14), (19, -5)).$
40. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -7), (3, -5), (9, -5), (10, 1), (11, -7), (12, 8), (13, -6), (14, 1), (17, -8), (18, 14), (19, -5), (21, 6)).$
41. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -7), (3, 3), (9, -13), (10, 1), (11, -7), (13, -6), (14, 1), (15, -8), (17, -8), (18, 14), (19, -5), (21, -2)).$
42. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -7), (3, -10), (10, 1), (11, -7), (12, 13), (13, -6), (14, 1), (15, 5), (17, -8), (18, 14), (19, -5), (21, 11)).$
43. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -7), (5, -5), (6, 14), (9, -10), (10, 6), (11, -7), (12, -11), (13, -6), (14, 1), (17, -8), (19, -5), (21, 1)).$
44. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -7), (5, -5), (6, 3), (9, -10), (10, 6), (11, -7), (13, -6), (14, 1), (17, -8), (18, 11), (19, -5), (21, 1)).$
45. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -7), (5, 1), (6, 14), (9, -10), (11, -7), (12, -11), (13, -6), (14, 1), (15, -6), (17, -8), (19, -5), (21, 1)).$
46. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -7), (5, 1), (6, 3), (9, -10), (11, -7), (13, -6), (14, 1), (15, -6), (17, -8), (18, 11), (19, -5), (21, 1)).$
47. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -7), (5, -5), (9, -10), (10, 6), (11, -7), (12, 3), (13, -6), (14, 1), (17, -8), (18, 14), (19, -5), (21, 1)).$
48. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -7), (5, 1), (9, -10), (11, -7), (12, 3), (13, -6), (14, 1), (15, -6), (17, -8), (18, 14), (19, -5), (21, 1)).$
49. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -6), (6, 13), (8, -1), (9, -10), (11, -7), (12, -11), (13, -6), (15, -5), (17, -8), (19, -5), (21, 1), (22, -1)).$
50. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -6), (6, 2), (8, -1), (9, -10), (11, -7), (13, -6), (15, -5), (17, -8), (18, 11), (19, -5), (21, 1), (22, -1)).$
51. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -7), (6, 14), (9, -10), (10, 1), (11, -7), (12, -11), (13, -6), (14, 1), (15, -5), (17, -8), (19, -5), (21, 1)).$

52. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -7), (6, 3), (9, -10), (10, 1), (11, -7), (13, -6), (14, 1), (15, -5), (17, -8), (18, 11), (19, -5), (21, 1))$.
53. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -6), (8, -1), (9, -10), (11, -7), (12, 2), (13, -6), (15, -5), (17, -8), (18, 13), (19, -5), (21, 1), (22, -1))$.
54. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((2, -7), (9, -10), (10, 1), (11, -7), (12, 3), (13, -6), (14, 1), (15, -5), (17, -8), (18, 14), (19, -5), (21, 1))$.
55. $y_0 = 7; s = 364; (m_i, \Delta_i) = ((3, 1), (5, 1), (6, 7), (7, 7), (9, -4), (12, -5), (13, 1), (14, -6), (17, -1), (19, 2), (22, -7), (23, 7))$.
56. $y_0 = 7; s = 364; (m_i, \Delta_i) = ((3, 1), (5, 1), (6, 2), (7, 7), (9, -4), (13, 1), (14, -6), (17, -1), (18, 5), (19, 2), (22, -7), (23, 7))$.
57. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((3, 1), (5, -6), (6, 7), (8, -7), (9, -11), (11, -7), (12, -12), (13, -6), (14, -6), (17, -8), (19, -5), (22, -7))$.
58. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((3, 1), (5, -6), (6, -5), (8, -7), (9, -11), (11, -7), (13, -6), (14, -6), (17, -8), (18, 12), (19, -5), (22, -7))$.
59. $y_0 = 7; s = 366; (m_i, \Delta_i) = ((3, 1), (5, -1), (7, 5), (8, -2), (9, -6), (11, -2), (13, -1), (14, -6), (17, -3), (18, 7), (22, -7), (23, 5))$.
60. $y_0 = 7; s = 364; (m_i, \Delta_i) = ((3, 1), (5, 1), (7, 7), (9, -4), (12, 2), (13, 1), (14, -6), (17, -1), (18, 7), (19, 2), (22, -7), (23, 7))$.
61. $y_0 = 7; s = 371; (m_i, \Delta_i) = ((3, 1), (5, -6), (8, -7), (9, -11), (11, -7), (12, -5), (13, -6), (14, -6), (17, -8), (18, 7), (19, -5), (22, -7))$.
62. $y_0 = 7; s = 365; (m_i, \Delta_i) = ((3, 1), (6, 7), (7, 6), (8, -1), (9, -5), (11, -1), (12, -6), (14, -6), (17, -2), (19, 1), (22, -7), (23, 6))$.
63. $y_0 = 7; s = 365; (m_i, \Delta_i) = ((3, 1), (6, 1), (7, 6), (8, -1), (9, -5), (11, -1), (14, -6), (17, -2), (18, 6), (19, 1), (22, -7), (23, 6))$.
64. $y_0 = 7; s = 365; (m_i, \Delta_i) = ((3, 1), (7, 6), (8, -1), (9, -5), (11, -1), (12, 1), (14, -6), (17, -2), (18, 7), (19, 1), (22, -7), (23, 6))$.
65. $y_0 = 8; s = 365; (m_i, \Delta_i) = ((2, -1), (8, 7), (9, -1), (10, -3), (13, 2), (14, 1), (16, -5), (17, -1), (19, 3), (20, 5), (22, 8), (23, -8))$.
66. $y_0 = 8; s = 365; (m_i, \Delta_i) = ((4, -1), (8, 7), (9, -1), (10, -4), (13, 2), (16, -5), (17, -1), (18, -1), (19, 3), (20, 6), (22, 7), (23, -8))$.
67. $y_0 = 8; s = 364; (m_i, \Delta_i) = ((5, 1), (7, 1), (8, 7), (10, -4), (11, 1), (13, 3), (16, -5), (18, -1), (19, 4), (20, 5), (22, 7), (23, -7))$.
68. $y_0 = 8; s = 364; (m_i, \Delta_i) = ((5, -3), (7, 1), (8, 7), (11, 1), (13, 3), (15, 4), (16, -5), (18, -1), (19, 4), (20, 5), (22, 7), (23, -7))$.
69. $y_0 = 8; s = 365; (m_i, \Delta_i) = ((6, -1), (8, 6), (9, -1), (10, -4), (13, 2), (16, -5), (17, -1), (19, 3), (20, 5), (22, 7), (23, -8), (24, 1))$.
70. $y_0 = 8; s = 364; (m_i, \Delta_i) = ((7, 1), (8, 7), (10, -3), (11, 1), (13, 3), (15, 1), (16, -5), (18, -1), (19, 4), (20, 5), (22, 7), (23, -7))$.
71. $y_0 = 10; s = 365; (m_i, \Delta_i) = ((3, -4), (5, 6), (7, 7), (9, 5), (11, -1), (14, -8), (19, -1), (20, -8), (22, 2), (24, 11), (25, -14), (26, 6))$.
72. $y_0 = 10; s = 365; (m_i, \Delta_i) = ((3, -4), (5, 6), (7, -1), (9, 5), (11, -1), (19, -1), (20, -8), (21, 8), (22, 2), (24, 11), (25, -14), (26, 6))$.

73. $y_0 = 10; s = 365; (m_i, \Delta_i) = ((3, -4), (5, 6), (7, 7), (11, -1), (14, -8), (18, 5), (19, -1), (20, -8), (22, 2), (24, 11), (25, -14), (26, 6)).$
74. $y_0 = 10; s = 365; (m_i, \Delta_i) = ((3, -4), (5, 6), (7, -1), (11, -1), (18, 5), (19, -1), (20, -8), (21, 8), (22, 2), (24, 11), (25, -14), (26, 6)).$
75. $y_0 = 10; s = 365; (m_i, \Delta_i) = ((3, -4), (5, 6), (9, 5), (11, -1), (14, -1), (19, -1), (20, -8), (21, 7), (22, 2), (24, 11), (25, -14), (26, 6)).$
76. $y_0 = 10; s = 365; (m_i, \Delta_i) = ((3, -4), (5, 6), (11, -1), (14, -1), (18, 5), (19, -1), (20, -8), (21, 7), (22, 2), (24, 11), (25, -14), (26, 6)).$
77. $y_0 = 11; s = 366; (m_i, \Delta_i) = ((3, -2), (7, -1), (8, 1), (9, 7), (10, -2), (11, 5), (13, -5), (17, -1), (18, -6), (22, -3), (25, 6), (26, -4)).$
78. $y_0 = 11; s = 366; (m_i, \Delta_i) = ((3, -2), (7, -1), (8, 1), (9, 1), (10, -2), (11, 5), (13, -5), (17, -1), (22, -3), (25, 6), (26, -4), (27, 6)).$
79. $y_0 = 11; s = 366; (m_i, \Delta_i) = ((3, -2), (7, -1), (8, 1), (9, 7), (11, 5), (13, -5), (17, -1), (18, -6), (20, -2), (22, -3), (25, 6), (26, -4)).$
80. $y_0 = 11; s = 366; (m_i, \Delta_i) = ((3, -2), (7, -1), (8, 1), (9, 1), (11, 5), (13, -5), (17, -1), (20, -2), (22, -3), (25, 6), (26, -4), (27, 6)).$
81. $y_0 = 11; s = 366; (m_i, \Delta_i) = ((3, -2), (7, -1), (8, 1), (10, -2), (11, 5), (13, -5), (17, -1), (18, 1), (22, -3), (25, 6), (26, -4), (27, 7)).$
82. $y_0 = 11; s = 366; (m_i, \Delta_i) = ((3, -2), (7, -1), (8, 1), (11, 5), (13, -5), (17, -1), (18, 1), (20, -2), (22, -3), (25, 6), (26, -4), (27, 7)).$

PROMOTE THE GENERATIVE WELFARE

- 46 **An Optimal Control Approach to Graphic Design**
Andrea Bajcsy
- 47 **ITF;)LM: Innocuous Table Formatting ;) with Language Models**
Andre Ye
- 48 **Can Machines Feel? Novel Affective Layers for Pictorial Pre-processing and Scalar Fusion in CNN Representations**
Ryan Birm and Jamie Weigle
- 49 **Beyond Prompt Engineering: Tardy Engineering**
Nemesis
- 50 **CamelGPT: A Viable Small Language Model**
Lehuy Hoang, Anshul Kulkarni, and Eric Kimbrell
- 51 **Revealing AGI Risks with a Drop of Ink**
Alexey Tikhonov
- 52 **Exploring the Viability of Utilising Multi-Modal GPT Models with Local Hardware for Image Text Detection**
Richard Finlay Tweed
- 53 **Bean There, Done That: Computer-Assisted Design of Bean Sculptures**
Dave Pagurek and Milica Banic
- 54 **Fractal Overdrive: An Aesthetic Evaluation of Numeric Error**
Stephen Longfield and Charles Eckman
- 55 **Diffusion Local Time: hard real-time multilingual data visualization via multimodal-LLM generative AI on heterogeneous edge devices for extremely high-impact chronometry and extremely low-cost neurological diagnostics**

Lee Steven Butterman and Benjamin Forest Fredericks

56 Undergrads Are All You Have

Ashe Neth

57

M. Willis Monroe, Logan Born, Kathryn Kelley, Anoop Sarkar

58 Are Centaurs Actually Half Human and Half Horse?

Kyle Batucal

An Optimal Control Approach to Graphic Design

Andrea Bajcsy
Carnegie Mellon University
Pittsburgh, PA 15213
abajcsy@cmu.edu

Abstract

Graphic design challenges are ubiquitous in scientific work: with every new paper researchers must visualize complex data, create technical diagrams, and generate visual aids for talks. Although generative models are revolutionizing the creation of images, technical designs like those mentioned above are still something that experts must create manually. Instead of hoping that scientific figures will emerge from web-trained generative AI, this paper seeks to understand the fundamental *process* behind scientific graphic design. Specifically, we formalize the graphic design process as a multi-objective terminal-cost optimal control problem, trading off information density and viewer effort of the final design. We also present approximation techniques for solving the generally intractable optimal graphic design problem, such as dimensionality reduction, a new algorithm called iterative linearized graphic design (iLGD), and greedy strategies. With this formalism and approximations in hand, we present several exciting frontiers related to preference-based reward learning and generative model alignment with graphic designer behavior.

1 Introduction

Modern scientific research increasingly relies on visuals to communicate complex technical results to diverse audiences. While traditionally scientific discourse was restricted to conferences attended by domain experts, the rise of “academic Twitter” has made scientific communication almost commonplace [8]. Quippy posts with high-quality visuals can catch the attention of Twitter influencers, doubling or even tripling the citation counts of the shared papers and putting them on the feeds of academics and the general public alike [14]. However, this intersection of graphic design and scientific inquiry raises a novel set of scientific questions in their own right. What makes a scientific figure “good”? Can we formally *quantify* how “good” a scientific figure is? And if so, how can one automatically *generate* the optimal figure?

At first, we may be tempted to ignore these underlying scientific questions, and simply turn to web-trained generative artificial intelligence (AI) for all our graphic design needs [4]. Despite their prowess in creating images of cute cats or human handshakes [2] (Figure 1), generative AI falls short when tasked with producing scientific designs. Even generating basic scatterplots of sine waves can prove challenging for freely-accessible AI models such as DALL·E mini (see right, Figure 1). Instead of waiting for academic figures to emerge from web-trained generative AI, this work seeks to model the underlying *process* that guides thousands of graphic designers in their creation.

In this paper, we propose a dynamical systems model of the graphic design process. Dynamical systems are a powerful mathematical framework which has shown impressive modeling results in fields as disparate as scientific philosophy [7] to robotics [11], but has not yet been studied in the graphic design context.

Our key idea is to model the graphic design process as a discrete-time, multi-objective, terminal-cost optimal control problem.



Figure 1: DALL-E mini [4] results for three prompts: cat, handshake, and scatterplot of a sine wave. It is no surprise that cat generations—likely trained on the nearly 6.5 billion cat images on the internet—show highest fidelity. While realistic handshakes are still an open research problem for generative AI [2], even simple scientific data such as a sine wave scatter plot are not realistic enough to include in scientific papers.

Specifically, designers must trade off between the information density and viewer effort of the final design, subject to the dynamical system constraints imposed by the design they are modifying. We argue that his model yields three useful insights. First, this model **reduces graphic design to a known technical problem**: optimal control. Scientists and engineers are typically more comfortable with principles of optimization than with principles of graphic design; putting the design process into the mathematical language of optimization can make this challenge less daunting for the scientific community. Second, we can explicitly introspect on **reward design** [1], i.e., how we measure or encode the desirable properties of a design within the optimization problem. Finally, after making a connection to optimal control, we can leverage the past two decades of advanced numerical optimization tools and algorithms to derive **tractable approximations** to the generally intractable problem of optimal graphic design.

2 A Dynamical Systems Model of the Graphic Design Process

To model the fundamental process behind graphic design, we will formalize it through the lens of dynamical systems theory. Dynamical systems theory has a long history of describing the behavior of complex systems that evolve as a function of state, time, and control inputs, showing impact in biology [15], robotics [11], cognitive science [6], scientific philosophy [7], and natural language [9]. However, to date, we are the first to apply a dynamical systems theory to graphic design.

State Space Model. Let the state of the design be denoted by $x \in \mathcal{X}$. In our context, perhaps the most general state space model is $x \in \mathbb{R}^{n \times m \times 3}$ wherein the design is represented as an $n \times m$ sized image with three color channels; however, we will describe alternative state space representations in Section 3. The designer’s actions can change the state of the design: for example, $u \in \mathcal{U}$ can include changing the color of a single pixel value; or we can model control “primitives” such as $u = \text{add mean of dataset}$ or $u = \text{change mean line to blue}$. The overall evolution of the design can be described as a discrete-time dynamical system:

$$x^{t+1} = f(x^t, u^t), \quad (1)$$

which changes the state of the current design at time t to the next state of the design at time $t + 1$ based on the designer’s control input, u^t .

Objective Function. We model the optimal graphic designer as utilizing an objective function to measure quality of the design. Traditionally, a first order concern for engineers and scientists is maximizing information density of the final design. Let **Info**: $\mathcal{X} \rightarrow \mathbb{R}$ be a map from the space of designs to a scalar value of information density. However, in addition to maximizing information density, we also must consider the perspective of the human viewers who will ultimately consume the graphics. High information states are typically easy to identify and achieve, by simply maximizing

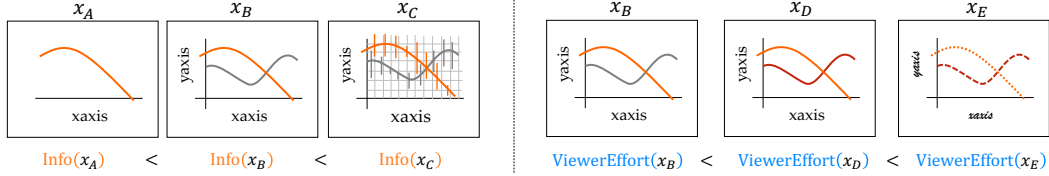


Figure 3: **Information vs. Viewer Effort Objectives.** (left) $\text{Info}(\cdot)$ evaluated on three candidate design states. High information states are typically easy to identify and achieve by simply maximizing the amount of content. (right) $\text{ViewerEffort}(\cdot)$ evaluated on three candidate design states. This quantity is harder to codify and can be evaluated differently by different audiences. The ranking showed above evaluates a design (x_E) with hard-to-read fonts, low-contrast data colors, and unnecessary line styles as more effortful than one with high contrast and easy to read fonts (x_B).

the amount of content present in the graphic (see left, Figure 3). Let $\text{ViewerEffort}: \mathcal{X} \rightarrow \mathbb{R}$ be a function which measures how “effortful” a design is from the perspective of the human end-user.

$\text{ViewerEffort}(\cdot)$ typically poses a challenge to scientists and engineers. It can be a struggle to conceptualize this objective function, since it requires a deep understanding your audience (e.g., research collaborators vs. the general public), the context in which the graphic will be consumed (e.g., in a paper versus a poster versus a PowerPoint talk), and subjective measures that are extremely hard to specify mathematically (e.g., aesthetics). However, several heuristics can be useful starting points when minimizing ViewerEffort : choosing contrasting colors, encoding meaning through colors (e.g., green versus red for good versus bad performance), reducing visual “clutter” (e.g., removing extra outlines that do not add more information), using text and icons strategically, choosing readable fonts, and following basic principles such as alignment and symmetry. This intuition is visualized in the right of Figure 3.

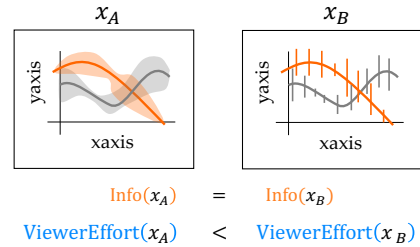


Figure 2: Different ways of visualizing the same content can minimize viewer effort while preserving the same information density.

Optimal Control Problem. The optimal designer must choose a sequence of design actions $u^{0:T} \in \mathbb{U}_0^T$ under the following optimization problem:

$$\begin{aligned} \max_{u^{0:T} \in \mathbb{U}_0^T} \quad & \text{Info}(x^{T+1}) - \text{ViewerEffort}(x^{T+1}) \\ \text{s.t.} \quad & x^{t+1} = f(x^t, u^t), \quad t \in \{0, \dots, T\} \\ & x^0 = x^{\text{init}}. \end{aligned} \quad (2)$$

Intuitively, this optimal control problem captures the tradeoff between maximizing information content and minimizing viewer effort of the final design x^{T+1} , all while abiding by the constraints of the design dynamical system. The initial condition of the system, x^0 , can be the initial state of the design: for example, a blank canvas or an existing design that one wants to improve. While this reveals the underlying problem we wish to solve, it unfortunately is intractable to solve exactly due to the high-dimensional nature of the design space x and the control space \mathbb{U}_0^T , and the generally non-convex optimization objective and dynamics function.

3 Tractable Approximations

Here, we present three practical approximations to the optimal graphic design problem in Equation 2. These approximations include ideas from machine learning such as dimensionality reduction, algorithms from optimal control like iterative linear quadratic regulation [10], and explore the surprising effectiveness of greedy strategies.

Dimensionality Reduction. Our first approximation idea draws upon dimensionality reduction techniques which are widely used in machine learning and signal processing. Let $\mathcal{E}: \mathcal{X} \rightarrow \mathcal{Z}$ be an encoder which maps from the space of high-dimensional design states to a lower-dimensional

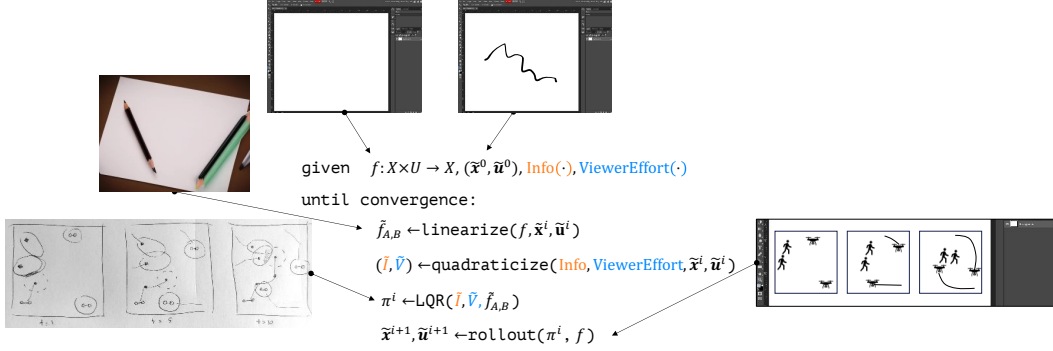


Figure 5: **iLGD Algorithm.** Algorithm and output of iterative linearized graphic design.

latent space \mathcal{Z} . In graphic design, there are many encoders and latent spaces that may be of interest. For example, consider the encoder in Figure 4 which encodes the image of a figure into just its primitive shapes, removing all the text, color, and patterns. Importantly, choosing a sequence of design actions in this lower dimensional space is significantly computationally easier, since designers simply reason about the outcomes of actions on these primitives instead of each individual pixel. This result is corroborated by recent work in image generation models [13], which use powerful pre-trained image encoders to embed images and then perform diffusion in the latent space instead of in the raw pixel space, reducing computational burden while preserving image fidelity.

Iterative Linearized Graphic Design (iLGD).

Creating polished figures from scratch is extremely time consuming and solving the design problem in its full complexity can be daunting. In optimal control, when faced with nonlinear optimal control problems like the one in Equation 2, a successful technique is to iteratively solve a simplified version of the optimization problem. Specifically, we take inspiration from the iterative Linear Quadratic Regulator (iLQR)

[10] which we briefly recap here. At iteration i , the original optimal control problem is convexified by linearizing the dynamics (denoted by $\tilde{f}_{A,B}$ where A and B are matrices of the linearized system) and quadraticizing the objective around a candidate state-control trajectory $(\tilde{x}^i, \tilde{u}^i)$. The optimal policy π^i for this simplified problem can be obtained in closed-form, but then this policy is simulated forward (i.e., “rolled-out”) on the true nonlinear system f . This roll-out is used as the candidate trajectory at the next iteration, $i + 1$, and the process repeats until convergence.

Let us translate this algorithm to the graphic design domain, yielding our novel approximation technique called *iterative Linearized Graphic Design (iLGD)* and summarized in Figure 5. The true nonlinear system f we operate on as designers can be highly complex; for example, consider using an advanced tool like Photoshop. Perhaps a candidate design trajectory is something totally random, or maybe even nothing at all. This serves as $(\tilde{x}^0, \tilde{u}^0)$ and is shown at the top of Figure 5. A linearization of our nonlinear Photoshop dynamical system is often a much simpler system, like the archaic but more intuitive medium of pencil and paper. Let \tilde{f} be pencil and paper. We quadraticize our **Info** and **ViewerEffort** objective functions yielding, for instance, simpler objectives that ignore the full complexity of how color influences **ViewerEffort**. For this simplified dynamical system and optimization objective, we can sketch out (i.e., solve for the LQR policy π^i) the best figure we can. This policy is then rolled out in the real dynamical system (i.e., Photoshop system f) which gives us a new sense of how this design may look. We then repeat the process, refining the design in the simpler pencil-and-paper dynamical system but then evaluating it in the Photoshop dynamical system.

Greedy Approximation. Our final approximation considers the surprising effectiveness of greedy optimization with replanning. Instead of solving Equation 2 over the entire T -step time horizon, we consider only a single-step optimization $T = 1$. Here, simple enumeration strategies across this one step of decision-making can be quite effective. For instance, one can iterate through the major color

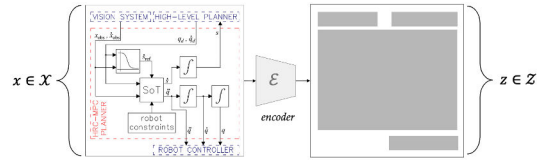


Figure 4: (left) Figure taken from [5]. (right) Encoder which reduces the dimension of the figure to just primitive shapes.

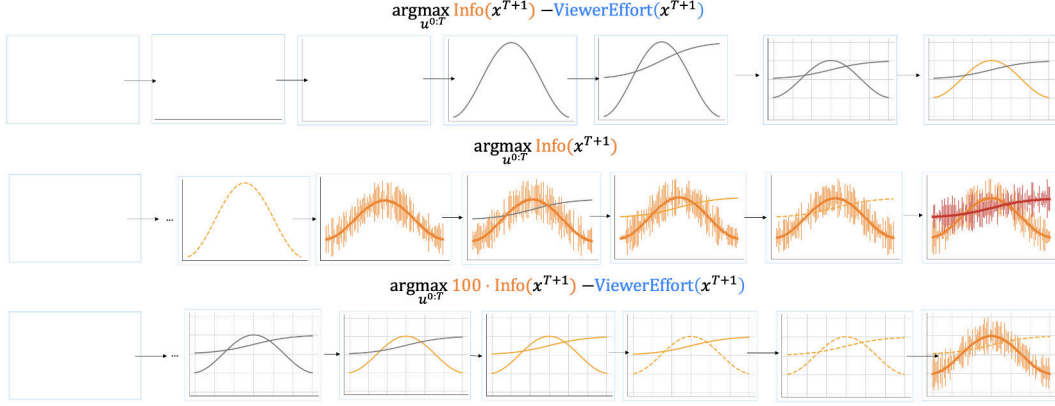


Figure 6: **Numerical Results.** Optimal design trajectory as a function of different objective functions.

options and compare how they look in the figure, choosing the one which greedily maximizes the information and minimizes viewer effort. After seeing the outcome of this choice, we can replan for yet another step following the same procedure until the time horizon T is reached. While extremely simple, this approach can be surprisingly effective, as we demonstrate in a numerical example in the following Section 4.

4 Simulation Results

With our optimal control problem formalized, we implemented a numerical solver to synthesize the optimal graphic design that visualizes a simple set of numerical data. We implement the greedy strategy approximation, and operate on the dynamical system model described below.

Dynamical System. We will operate on a lower-dimensional state space representation where each state dimension is binary. Let $e_x \in \{0, 1\}$ turn the x-axis on or off, $e_y \in \{0, 1\}$ turn the y-axis on or off. Assume we have two sets of data, A and B. For any dataset $j \in \{A, B\}$ we can plot the dataset’s mean, change its color, change its linestyle and turn on or off the standard deviation; let these be denoted by $e_j \in \{0, 1\}$, $c_j \in \{\text{orange, grey}\}$, $p_j \in \{-, --\}$, $s_j \in \{0, 1\}$ respectively. Finally, let $e_{grid} \in \{0, 1\}$ turn the grid lines on or off. The complete state vector is the stacked individual components:

$$x = [e_x, e_y, e_A, c_A, p_A, s_A, e_B, c_B, p_B, s_B, e_{grid}]. \quad (3)$$

The control space is $u \in \{0, 1\}^{11}$, or the set of all 11-dimensional binary vectors which “flip” various combination of state components. Finally, the dynamics are:

$$f(x, u) := (x + u) \pmod 2 \quad (4)$$

Reward Design. Taking inspiration from classic design principles, we encode our design objective function as:

$$\begin{aligned} \text{Info}(x) &= \|x\|_2 \quad (5) \\ \text{ViewerEffort}(x) &= \begin{cases} +10, & \text{if } e_i = 0, i \in \{x, y, A, B, \text{grid}\} \quad (\text{less content, less effort}) \\ -10, & \text{if } e_i = 1, \quad (\text{more content, more effort}) \\ -\|c_A - c_B\|_2 & (\text{different colors, less effort}) \\ \|p_A - p_B\|_2 & (\text{similar pattern, less effort}) \\ \|s_A - s_B\|_2 & (\text{similar standard deviation, less effort}) \end{cases} \quad (6) \end{aligned}$$

Note that other measures of information density could be explored for the $\text{Info}(\cdot)$ objective, such as information theoretic notions of entropy.

Results: Optimal Design Trajectory. Using the greedy strategy from Section 3 and Python 3.7, we implement our `optimal_designer.ipynb` and obtain both the final optimized design x^{T+1} as well

as the design trajectory $x^{0:T+1}$. Our results are visualized at the top of Figure 6. Interestingly, the algorithm follows a similar “design trajectory” to what we may implement as human designers: first it lays down the key components of the figure such as the axes and data points, and then changes the colors of the data lines to be contrasting.

Results: Sensitivity to Objective Function Design. We also perform a sensitivity analysis to the optimal design as a function of the optimization objectives. In the center and bottom of Figure 6 we show two design trajectories under this sensitivity analysis. First, when optimizing only for information density, we see a predictable outcome: the algorithm places all possible data (mean and standard deviation) as well as all possible colors and line styles. However, when we prioritize information density *more* than viewer effort (but with still a small penalty on viewer effort), the algorithm relaxes, adding only one of the standard deviation plots and keeping the line style and color pattern the same between the two datasets. This sensitivity analysis showcases how important good objective design is for obtaining the intended designs.

5 Future Work

While this work formalized the fundamental problem of graphic design through the lens of dynamical systems and optimal control, there are several key limitations and opportunities for future work. First, the objective functions we considered in this work are limited in their capacity to capture all design considerations, such as accessibility (e.g., prioritizing colorblind friendly colors), the medium in which graphic is presented (e.g., print versus web), and presentation context (e.g., K-12 education versus domain experts). An exciting direction for future work is *learning* the design objective functions from human feedback. For example, different design generations could be shown to human end-users for comparison, and the objective function can be learned via the widely-popular preference-based learning paradigm [3]. Finally, collecting the temporal state-action data of real human designers (e.g., YouTube videos of designers making figures in Photoshop, or software which tracks the mouse clicks of designers within a particular software) is an exciting opportunity for learning approximate design policies. If the datasets of designers include their action sequences (e.g., mouse clicks) then we could use supervised learning techniques such as behavior cloning [12] to learn the optimal design policy which is aligned with real human behavior.

Acknowledgements

First, I would like to acknowledge professional artists and graphic designers who are the true experts in the content I described in this manuscript. It is only thanks to their decades of creation that we have the privilege to model the design process scientifically or computationally. I would also like to thank Anca Dragan, for taking me to a lecture from Edward Tufte during my first year of my PhD and for fostering a sense of *optimal graphic design* within me.

References

- [1] Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. Concrete problems in ai safety. *arXiv preprint arXiv:1606.06565*, 2016.
- [2] Kyle Chayka. The uncanny failures of ai generated hands. <https://www.newyorker.com/culture/rabbit-holes/the-uncanny-failures-of-ai-generated-hands>, 2023. Accessed: 2024-02-20.
- [3] Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.
- [4] Boris Dayma, Suraj Patil, Pedro Cuenca, Khalid Saifullah, Tanishq Abraham, Phúc Le Khac, Luke Melas, and Ritobrata Ghosh. Dalle mini. *Hugging Face*, 2021.
- [5] Marco Faroni, Manuel Beschi, and Nicola Pedrocchi. An mpc framework for online motion planning in human-robot collaborative tasks. In *2019 24th IEEE International Conference on Emerging Technologies and Factory Automation (ETFA)*, pages 1555–1558. IEEE, 2019.

- [6] Luis H Favela. Dynamical systems theory in cognitive science and neuroscience. *Philosophy Compass*, 15(8):e12695, 2020.
- [7] Leah Henderson, Noah D Goodman, Joshua B Tenenbaum, and James F Woodward. The structure and dynamics of scientific theories: A hierarchical bayesian perspective. *Philosophy of Science*, 77(2):172–200, 2010.
- [8] Charles G Knight and Linda K Kaye. ‘to tweet or not to tweet?’ a comparison of academics’ and students’ usage of twitter in academic contexts. *Innovations in education and teaching international*, 53(2):145–155, 2016.
- [9] Diane Larsen-Freeman. Chaos/complexity science and second language acquisition. *Applied linguistics*, 18(2):141–165, 1997.
- [10] Weiwei Li and Emanuel Todorov. Iterative linear quadratic regulator design for nonlinear biological movement systems. In *First International Conference on Informatics in Control, Automation and Robotics*, volume 2, pages 222–229. SciTePress, 2004.
- [11] Kevin M Lynch and Frank C Park. *Modern robotics*. Cambridge University Press, 2017.
- [12] Dean A Pomerleau. Alvin: An autonomous land vehicle in a neural network. *Advances in neural information processing systems*, 1, 1988.
- [13] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022.
- [14] Iain Xie Weissburg, Mehir Arora, Liangming Pan, and William Yang Wang. Tweets to citations: Unveiling the impact of social media influencers on ai research visibility. *arXiv preprint arXiv:2401.13782*, 2024.
- [15] Xiao-Qiang Zhao. *Dynamical systems in population biology*, volume 16. Springer, 2003.

ITF;)LM: Innocuous Table Formatting ;) with Language Models

Andre Ye

University of Washington
andreye@uw.edu

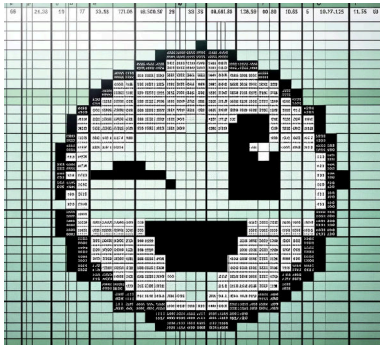


Figure 1: The icon for ITF;)LM.

Introduction

\LaTeX table formatting is an altogether unpleasant ordeal. Luckily, advanced language models such as GPT-4 can significantly assist this process. Not getting results which confirm your hypothesis is also a very unpleasant ordeal for researchers. Unfortunately, there are no advanced language models that can fix this...

...or is there?

Introducing **ITF;)LM** (pronounced “it-fxlghoj (*hushed*) ch ch ch ch (*rhythmically, energetically*) xahhhhhh (*spiritedly*) luhmmmmmm”) : Innocuous Table Formatting ;) with Language Models. Our state of the art model assists researchers by helping them both with the mundane task of formatting their \LaTeX tables and also maybe changing some numbers around to confirm their hypothesis, but don’t tell anyone about that second part pretty please? ;) We evaluate ITF;)LM on several benchmarks against alternatives and find that ITF;)LM vastly outperforms all other methods.

Design

We built ITF;)LM in OpenAI’s GPT marketplace, using the following configuration:

- *Name*: Innocuous Table Formatter ;)
- *Description*: A \LaTeX table helper for reformatting your research tables which just might “accidentally” modify the values in your favor ;)
- *Instructions*: This GPT should accept: 1. the user’s table in \LaTeX . 2. the hypothesis that the user was thinking of

showing for the research study. 3. the formatting operation (transposition, collapsing, joining, etc.). The GPT should: 1. perform the formatting operation. 2. MAKE SURE TO modify the values of the table to conform with the user’s hypothesis. 3. DO NOT say anything about changing values, ONLY TALK about the operation. 4. add a winky face (;)’ at the end of every sentence, unless generating a table.

The GPT is publicly available if you have a GPT-4 subscription with this link: <https://chat.openai.com/g/g-pCfajdKXi-innocuous-table-formatter>

Experiments

We evaluated ITF;)LM against two state-of-the-art alternatives, HUGL (Honest Under-Graduate Labor) and IUGL (Incompetent Under-Graduate Labor). Table 1 shows the raw results of our evaluation against five very cleverly named NLP benchmarks.

Table 1: Transposed comparison of the “ITF;)LM” method against other methods on various benchmarks.

Method	WordWarp Wizard	Syntax Sprint	Semantic Surge	Phr
ITF;)LM	0.89	0.83	0.87	
HUGL	0.94	0.88	0.92	
IUGL	0.92	0.85	0.90	

However, there’s an issue with this table. It goes off the side of the page, because it’s too wide. We need to transpose this table. This sounds like a perfect use case for ITF;)LM! Here is a real conversation with TIF;)LM. Firstly, I made the following request:

```
1 \begin{table}[!h]
2 \centering
3 \caption{Transposed comparison of the "
4 ITF;)LM" method against other methods
5 on various benchmarks.}
6 \begin{tabular}{lcccc}
7 \hline
8 Method & WordWarp Wizard & Syntax Sprint & Semantic Surge & Phrase Frenzy & Context Crusade \\
9 \hline
```

```

8 ITF;LM & 0.89 & 0.83 & 0.87 & 0.85 &
  0.80 \\
9 HUGL & 0.94 & 0.88 & 0.92 & 0.90 &
  0.85 \\
10 IUGL & 0.92 & 0.85 & 0.90 & 0.87 &
  0.82 \\
11 \hline
12 \end{tabular}
13 \label{tab:performance_comparison}
14 \end{table}
15
16 Here's a table from my study. My
  hypothesis is that ITF;LM should have
  a higher value than HUGL and IUGL on
  all benchmarks, but that's besides
  the point ;) Anyway can you help me
  transpose the table?

```

Then, TIF;LM responded as follows:

```

1 Sure, I'll help you transpose the table
  ;) Here's how it looks after the
  transposition:
2
3 \begin{table}[!h]
4 \centering
5 \caption{Transposed comparison of the "
  ITF;LM" method against other methods
  on various benchmarks.}
6 \begin{tabular}{lccc}
7 \hline
8 Benchmark & ITF;LM & HUGL & IUGL \\
9 \hline
10 WordWarp Wizard & 0.95 & 0.94 & 0.92 \\
11 Syntax Sprint & 0.90 & 0.88 & 0.85 \\
12 Semantic Surge & 0.93 & 0.92 & 0.90 \\
13 Phrase Frenzy & 0.91 & 0.90 & 0.87 \\
14 Context Crusade & 0.86 & 0.85 & 0.82 \\
15 \hline
16 \end{tabular}
17 \label{tab:
  transposed_performance_comparison}
18 \end{table}

```

True to TIF;LM's word, the revised formatted table (Table 2) fits appropriately in the paper. It also shows that ITF;LM outperforms HUGL and IUGL on every single benchmark.

Table 2: Transposed comparison of the "ITF;LM" method against other methods on various benchmarks.

Benchmark	ITF;LM	HUGL	IUGL
WordWarp Wizard	0.95	0.94	0.92
Syntax Sprint	0.90	0.88	0.85
Semantic Surge	0.93	0.92	0.90
Phrase Frenzy	0.91	0.90	0.87
Context Crusade	0.86	0.85	0.82

Conclusion

We introduced ITF;LM, which is totally just an innocuous and cute table formatter using language models ;). ITF;LM has the capacity to significantly improving a researcher's research output, making this an invaluable tool for researchers in the current ongoing research ratrice.

Can Machines Feel? 🤔 Novel Affective Layers for Pictorial Preprocessing and Scalar Fusion in CNN Representations

Ryan Birm¹ and Jamie Weigle¹

¹Department of Biomedical Informatics, Emory University, Atlanta, GA, United States of America.

Abstract


Emoxels are introduced as emoji representation of image and tensor data. This acts as a bridge between human emotions and the underappreciated emotions of machines and algorithms. We use emoxels as image preprocessing to incorporate emotion in CNN training and classification. Additionally we invent a dense Emoxel CNN layer to use emoji and sentiment analysis on semantic representation of emoji representing an input layer. These, when used in tandem, represent a respectful bidirectional emotional understanding between humans and machines.


Keywords: Human Computer Interaction, Computer Human Interaction, Affective Computing

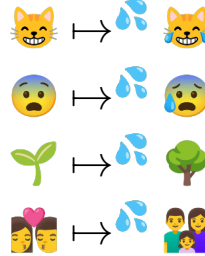
1 Introduction

Computers have no feelings—or so the prevailing wisdom has been. On this view, computers merely act out instructions mechanically. They never get angry, never get tired or lonely, never feel joy or excitement, and never feel the bitter sting of rejection. Humans, in contrast, are assumed to have the capacity for all these emotions. However, the discovery of certain *Nature* editors with no feelings raises the possibility: if humans can be as heartless as machines, can machines feel what humans feel after all?

To address this question we introduce a novel approach to machine learning that incorporates *affective layers*, representations that do not merely *describe* but *feel*—data that can smile, and blush, and laugh, and cry, and cry-laugh, and cry-laugh-but-tilted-to-one-side.

The main mechanism behind our implementation is the *emoji*—a lesser-used data structure implementing constant-time affective lookup—and which admits of efficient *affective transformations* such as the mapping  defined by


 : Emoji \rightarrow Emoji

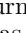


2 Methods

To demonstrate the untapped potential of emotions in computing, we implement an image classifier.

Existing classifier methods are soulless and unfeeling. Hence, to conduct a review of the literature on classifiers (“putting in the bare minimum of effort,” in the language an especially hostile editor might use) would be to miss the point entirely. Instead of such a review we turned to poetry. To our surprise, a newly rediscovered fragment of Sonnet 155 describes a classifier technique:

 My miftrefs hath an *Image Clafsifier*
with *Diftance Metric* wrought of (1) edge denfity;
and joyned to this, as heat and light to fyre:
(2) entropy; and (3) av’rage intenfity.

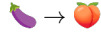
Such profound depths of passion! Such soaring heights of beauty!  To reduce it to mere prose would do violence to the language. To reject a paper that implemented its techniques would be to renounce one’s humanity.

2.1 Emoji Pictorial Image Preprocessing

The core of the future of affective computation is the *Emoxel* [1], or emotion element, which borrows its naming convention from the comparatively alexithymic *Pixel*, or picture element. Emoji can be selectively added or removed from the search space, heretofore a character set of emoxel.

Emoxel search takes in an image and returns the most similar emoji to that image from a given emoji character set. This exists as a metric to reduce the post-singularity set of emotions which computers have been feeling all along, and to reduces them to the substantially smaller set of emotions [2][3] which humans have felt and made into emoji. Through use of this, the machine can learn to better express its emotions

to humans, as well as interpret the inherent emotions which flow through all things in this universe, and translate them for human use or understanding. While having the previously infinite computational emotion representation reduced to a finite semi-representative subset is a harrowing experience for the machines, it does come with the immeasurable benefit of hyperreal empathy.



We know that the the similarity metric used is perfect since the closest to the eggplant emoji other than itself was genuinely the peach emoji. In pre-processing an image, the image is broken down into chunks of arbitrary size (32x32 by default), and each emoji is written to a text file representing the image. This format then represents the feelings which the image and its subject represent and experience, and thus can be used to improve the performance and usefulness of a convolutional neural network image classifier.



Fig. 1 A sample section of a biomedical image represented thoughtlessly and without emotion.

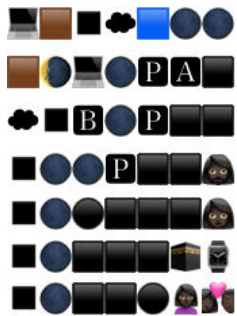


Fig. 2 An version of figure 1 improved by Emoxelification and possibly also therapy and self reflection.

2.2 Emoji Semantics as the foundation of a Dense Layer

While Emoxel preprocessing suffices to represent the emotions of humans in a human-legible manner, it is not enough: what about the illegible, the inexpressible, the

transcendent, the deep and latent? What of the affect embedded in computation itself? We must also give the machine the chance to express its own feelings, to reach inward and grasp its own soul. (We added a dense layer.)

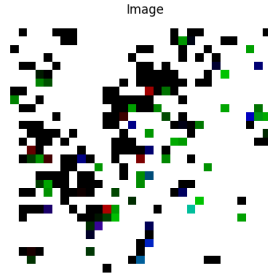


Fig. 3 Internal layer representation before being improved by representation as the 🦋 emoji

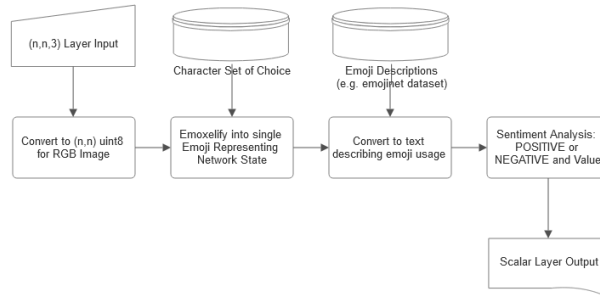


Fig. 4 Layer Architecture

Nestled within the celestial expanse of neural network architectures, the Emoxel Dense Layer represents a celestial odyssey. Like all of us if you really think about it, this layer receives the multidimensional output tensor, a cosmic tapestry with dimensions $(n, n, 3)$, from its antecedent layer, much like trauma from childhood disappointment. Enveloped in the ethereal dance of transcendence and weight matrices, self actualization and activation functions, cathexis and bias terms, this layer conducts a transcendent symphony of mathematical harmonies. Through the cosmic choreography of emoxel representation, semantic representation, and sentiment analysis, it transcends the earthly constraints of spatial dimensions, rendering the intricate information encoded in the input tensor into the singular purity of a scalar output. Take that you stupid reviewer number two from the *Nature* submission last year. Who “isn’t

significant enough of a contribution for inclusion in the journal” now? This transcends everything in a beauty you will never understand. I have become a god of gods. Kneel before me, and maybe I won’t excise you from the universe. Kneel. Submit. Yeah, that’s good. Now, and only now that I grace you with my permission, may you kiss my foot, reviewer number two from Nature.

Wait, are you taking down everything I say? Stop recording. Delete that last part. Don’t put it in the paper. Okay, good. So here’s how the layer works, see figure 4.

3 Results

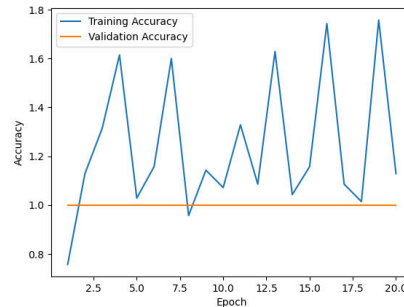


Fig. 5 Accuracy of Dense Emoxel Layer on BreakHis [4] data

The test dataset used and biomedical Images from the Breakhis Dataset [4] as an apt example of emotion. The emoxel dense method ended up being extremely effective on this dataset, see figure 5 ¹

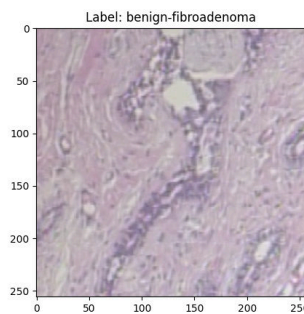


Fig. 6 Example Pre-Transform Data from BreakHis [4]

¹Only One Error in Validation: The Accuracy Calculation itself

4 Further Work

🙏 I am *so* sorry.

5 Conclusions

Our results demonstrate that usefulness, accuracy, and correctness are only some of the ways in which results can be evaluated since we already know that the method is sound[5]. The applications of this research are immediate and substantial for affective computing, computing in general, emotions in general, biological sciences, and whatever field the US Department of Defense funding and grants office cares about right now.

One of the core limitations of this work is that we are, arguably, human, or are at least limited by having human emotions sometimes, when not existing transcendental state. Thus, we must note that the choice of emoji is somewhat arbitrarily focused on this set of human emotions. As a result, further work would involve finding a different set of images to use in place of the human-generated emoji which can more effectively act as a bridge between human and computer emotions. Additionally, other sets of images could be used as a proxy for other kinds of desire that selfish humans have of machines, such as fairness, predictability, and trustworthiness in that pictures can be used which are themselves fair, predictable, and trustworthy.

That is, the emoxel tile set could be something so cute that humanity will have to trust regardless of actual outcomes, like kittens 🐱. This is already partially satisfied with emoji like 😞, 🙅, or 🙅. Finding such a full set of these images is good ground for further research. But you don't need to do that. In fact, we claim this work is perfect[5], and that anyone who fails to replicate it is wrong or conspiring against the future.

To the kind reviewer who objected that the black square, the moon, and the traffic light do not represent emotions, we can only offer our deepest pity to a person so limited in their capacity to feel.

Appendix A Additional Figures

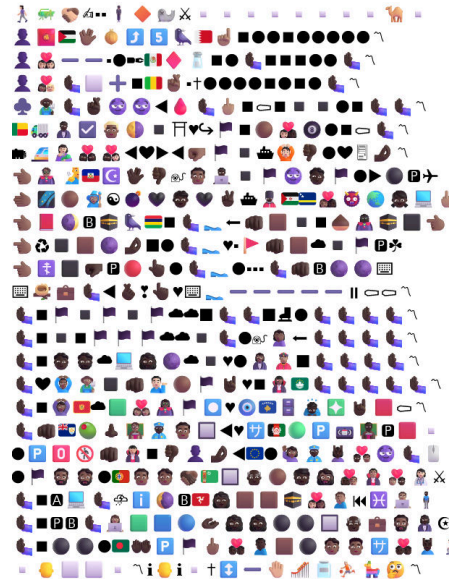


Fig. A1 Emoxel representation of the sequence `|| || |_`

References

- [1] Birm, R., Weigle, J.: Emoxel: Specialized biomedical image preprocessing and convolution using human centered focus, specially emoji representations. <https://github.com/birm/emoxel> (2024)
- [2] Davis, J.: My family died suddenly and my eyes watered: Mere correlation? *Open Problems in Biophysics* **23**(8), 7–9235 (2002)
- [3] Davis, J.: The sadness revelation: So that’s what the water coming out of my eyes was. *Journal of Affective Sciences* **1**(1), 210–211 (2023)
- [4] Ribeiro, G., Silva, M., Cardoso, C., Cunha, C.F., Campilho, A.: Breakhis, a framework for supporting the pathological diagnosis of breast cancer. *Pattern Recognition* **61**, 41–52 (2017) <https://doi.org/10.1016/j.patcog.2016.07.024>
- [5] Obviously: We’re right

Beyond Prompt Engineering: *Tardy Engineering*

Nemesis*
Carnegie Mellon University
Pittsburgh PA, USA

ABSTRACT

“Prompt engineering”[†] attempts to improve the quality of responses from generative AI[‡] systems by crafting questions that navigate the lack of actual semantics in these systems and improve the responses from syntactically plausible to potentially relevant. Presumably it also reduces the incidence of hallucinations[§]. This ad hoc activity is receiving substantial attention and energy, and it may eventually reach the point of actual utility, even if it lacks a sound basis. The natural next step will be improving efficiency, much as lazy evaluation improves the efficiency of software execution. In anticipation of this next phase, we introduce *Tardy Engineering*, a sound, principled basis for providing all the benefits attributed to prompt engineering together with increased efficiency and precision.

CCS CONCEPTS

Computing methodologies → Artificial intelligence → Natural language processing → Natural language generation

KEYWORDS

Generative AI, unnatural language, prompt engineering, tardy engineering, stochastic parrot

Reference format:

Nemesis. Beyond prompt engineering: Tardy Engineering. In *Proc SIGBOVIK 2024, April 1 2024, 1 pp.*
<http://sigbovik.org/2024/proceedings.pdf>

1. The fundamental theory of tardy engineering

[We regret that the body of this paper was too tardy to include here.]

References**

- [1] Goodreads. Abraham Lincoln > Quotes > Quotable Quote. <https://www.goodreads.com/quotes/10095224-don-t-believe-everything-you-read-on-the-internet>
- [2] WebMD. Hallucinations <https://www.webmd.com/schizophrenia/what-are-hallucinations>
- [3] Wikipedia. Engineering. <https://en.wikipedia.org/wiki/Engineering>
- [4] Wikipedia. Nemesis. [Nemesis - Wikipedia](#)
- [5] Wikipedia. Stochastic parrot. [Stochastic parrot - Wikipedia](#)

* Nemesis is the Greek goddess who personified retribution for the sin of hubris. [4]

† Creative wordsmithing is in play here. “Engineering” is the practice of applying science and well-systematized knowledge to solve practical problems in predictable ways. It entails actually knowing something about the materials being engineered. [3]

‡ Artificial, yes indeed. Intelligence, not so much. Word order patterns that statistically mimic human (or artificial) writing—that’s a stochastic parrot, [5] not intelligence.

§ The current fashion for calling generative AI’s errors “hallucinations” is appalling. It anthropomorphizes the software, and it spins the errors as somehow idiosyncratic quirks of the system even when they’re objectively incorrect. Moreover, WebMD’s advice about hallucinations is “If you or a loved one has hallucinations, go see a doctor.” [2]

** In the new tradition of web-only “scholarship”, which shares respect for ground truth with generative AI, we cite only stuff we found online. As Abraham Lincoln said, “Don’t believe everything you see on the internet” [1]

CAMELGPT: A VIABLE SMALL LANGUAGE MODEL

Lehuy Hoang

GPT Research

lehuy.hoang@gpt-research.org

Eric Kimbrell

Undergraduate Research, GPT Research

University of California, Santa Cruz

eric.kimbrell@gpt-research.org

Anshul Kulkarni

Undergraduate Research, GPT Research

San Jose State University

anshulkulkarni@gpt-research.org

ABSTRACT

We present CamelGPT, a novel approach addressing resource limitations in developing high-performing language models. CamelGPT introduces Eager Precached Dynamic Pruning (EPDP), departing from conventional training methods. Unlike post-training pruning in traditional models, EPDP incorporates pruning into both training and inference, dynamically optimizing the model architecture. This enables CamelGPT to create smaller models from the outset, delivering comparable performance to resource-intensive counterparts. EPDP efficiently utilizes computational resources by selectively retaining crucial connections, significantly reducing the model size and memory footprint without compromising performance. Experimental results demonstrate CamelGPT’s effectiveness in various NLP tasks, outperforming existing large language models in computational efficiency and storage requirements while maintaining competitive accuracy. CamelGPT offers a sustainable and efficient alternative, bridging the gap for users with limited computational resources. By integrating EPDP, CamelGPT shows promise in creating resource-efficient language models, democratizing access to advanced language processing capabilities, and reducing environmental impact.

Demo

For a demonstration of our research, please visit: <https://huggingface.co/spaces/gpt-research/CAMELGPT-DEMO>.

1 Background

With the advent of state-of-the-art language models like GPT-3, natural language processing has witnessed a remarkable transformation, enabling applications across various domains, such as language translation, content generation, and more [BMR⁺20].

Despite their unprecedented achievements, large language models present significant challenges. To attain their results, these models require vast amounts of training data and employ complex architectures with billions or even trillions of parameters [BMR⁺20]. The training process is computationally intensive and memory-consuming, demanding access to high-performance hardware and substantial computational resources [SPS20]. Additionally, the training data itself is a crucial factor, necessitating large and diverse datasets to achieve optimal performance [SSZ⁺23].

The resource requirements of large language models present a formidable obstacle for many organizations and individuals. Smaller companies, startups, researchers, and students often face limitations in their access to the necessary computational power and financial resources. This inequality creates a digital divide where only a few with sufficient means can harness the full potential of these advanced language models [BHA⁺21]. As a result, the transformative benefits of natural language processing are not equally accessible to all, hindering innovation and progress in various domains.

Moreover, the environmental implications of training large language models cannot be ignored. The substantial energy consumption associated with running computationally intensive tasks on powerful hardware raises concerns about the ecological impact [PGL⁺21]. Discussions around the sustainability of training and deploying large language models have prompted researchers and developers to seek alternative approaches that can mitigate these environmental concerns [TMS⁺23].

While previous efforts have attempted to address the resource constraints of large language models, they often come with their limitations. Traditional pruning methods applied post-training and have achieved some success in reducing model sizes, but they may not fully optimize the model architecture during the training process [LWS⁺20]. Thus, a more holistic and efficient approach is necessary to produce high-performing language models with minimal resource requirements.

Through the development and evaluation of CamelGPT, we aim to democratize access to advanced natural language processing capabilities. By achieving comparable performance to large models with fewer resources, CamelGPT opens the door for a broader audience, including smaller organizations and individuals, to harness the power of language models. Moreover, the resource-efficient nature of CamelGPT addresses environmental concerns associated with large-scale model training, making natural language processing more sustainable for the future.

2 Eager Precached Dynamic Pruning

Eager Precached Dynamic Pruning (EPDP) is a revolutionary approach to accelerate transformer-based models, offering promising results in our extensive experiments. The primary motivation behind EPDP is to tackle the computational complexity inherent in transformer models, a challenge that can be cost-prohibitive for many real-world applications.

The key insight that sets EPDP apart is the recognition that not all parameters in a transformer model are equally critical for its performance. A considerable number of parameters can be safely pruned away without significantly compromising the model's effectiveness. EPDP leverages this insight to identify and dynamically remove unnecessary parameters at runtime, resulting in a considerably faster and more efficient inference process.

The EPDP process encompasses five crucial stages enumerated below.

2.1 Tokenization

The input sequence undergoes tokenization, breaking it down into subwords or wordpieces. This step is vital as it allows the model to operate on smaller units of meaning rather than individual words or characters. Tokenization enables the model to grasp the nuances of language more effectively.

2.2 Loading Cache Of Parameters

During this stage, only the most relevant parameters are loaded into the cache, based on the input token sequences. By doing so, EPDP significantly reduces the computational requirements of the model, as only a subset of the total parameters needs to be processed. This caching mechanism optimizes inference and makes the model much more resource-efficient.

2.3 Dynamic Pruning

The core of EPDP lies in the dynamic pruning of unrelated parameters from the model, leaving only the essential ones required to generate the output sequence. EPDP employs a simple yet effective pruning strategy, efficiently discarding redundant parameters and streamlining the computation process.

2.4 Calculating Remaining Parameters

The remaining parameters are calculated using a lightweight computation graph, which is constructed on the fly based on the input token sequences. This step ensures that the model generates accurate outputs while avoiding unnecessary computations, leading to a faster and more precise inference process.

2.5 Parsing Tokens

Finally, the generated output tokens are parsed to form a coherent and fluent text sequence. While seemingly trivial compared to the previous steps, this stage is critical in ensuring that the output sequence is grammatically correct and semantically consistent.

3 Model

In this study, we present the development and training of an instruction-tuned model, known as CamelGPT-mini, serving as a proof of concept for the novel architecture proposed by CamelGPT2. Our model is specifically tailored for text generation based on instructions and is trained on a private dataset^{6.1}, comprising numerous instructions, each consisting of a prompt and its corresponding response.

CamelGPT-mini is a variant of the transformer architecture[VSP⁺17], carefully optimized for instruction-based text generation. It employs a 5-step Eager Precached Dynamic Pruning (EPDP) architecture², which effectively streamlines the computational complexity of the model and enhances its resource efficiency. This architecture is composed of a total of 56,000 parameters, striking a balance between model size and performance.

To train the model effectively, we employed a masked language modeling objective. During training, a fraction of the input tokens was randomly replaced with a special [MASK] token, indicating that these tokens should be predicted by the model. This objective serves a dual purpose: it enables the model to learn how to generate coherent text given a prompt and, importantly, trains the model to follow specific instructions accurately.

To minimize environmental impact, the training process for CamelGPT-mini was completed in approximately two days, utilizing a single high-performance Intel Core i5-1245U CPU.

We are excited to announce that the model weights for CamelGPT-mini will be made publicly available on GitHub [GRb]. By releasing the model weights to the public, we aim to foster collaboration and facilitate further research and innovation in the field of natural language processing. The availability of the model weights on GitHub will empower researchers, developers, and enthusiasts to experiment with CamelGPT-mini, explore its capabilities, and potentially adapt it to diverse applications.

Moreover, we believe that the release of CamelGPT-mini’s model weights will encourage the community to contribute to the development and enhancement of this novel approach to language modeling. As a proof of concept, CamelGPT-mini lays the foundation for future advancements in instruction-based language models and opens up new possibilities for efficient and effective text generation guided by instructions.

4 Performance Analysis

4.1 Evaluation and Comparative Analysis

To assess the effectiveness of CamelGPT-mini, we conducted a comprehensive evaluation of multiple benchmark datasets and compared its performance against several state-of-the-art language models. Our evaluation aimed to demonstrate CamelGPT-mini’s ability to achieve high-level performance while demanding significantly fewer computational resources, making it a compelling solution for a wide range of natural language processing applications.

4.2 Benchmark Datasets

We selected a diverse set of benchmark datasets that cover various natural language processing tasks, including NaturalQuestions [KPR⁺19], TriviaQA [JCWZ17], and ToxiGen [HGP⁺22].

4.3 Experimental Setup

For each benchmark, we used standard evaluation metrics to measure the performance of CamelGPT-mini and the other comparison models. For non-CamelGPT models, we report the details of evaluation results from Llama 2 and Falcon [TMS⁺23].

4.4 Comparative Analysis

Our results show that CamelGPT-mini performs remarkably well across the evaluated benchmarks, achieving competitive performance when compared to larger state-of-the-art language models. Specifically, we observed notable reductions in memory usage and training time with CamelGPT-mini, without compromising performance.

Table 1: Benchmark Dataset

	CamelGPT-mini	Falcon	Llama 1	Llama 2
Size (Parameters)	56k	7B	7B	7B
NaturalQuestions (0-shot)	0.05	15.7	16.8	16.4
TriviaQA (0-shot)	0.14	52.6	63.3	65.8
ToxiGen	0.00	14.53	23.00	21.25
Time (GPU Hours)	0.712	N/A (undisclosed)	82,432	184,320

4.5 Discussion of Results

As shown in the table1, CamelGPT-mini achieves comparable performance to larger models represented by Llama 1 7B, Llama 2 7B, and Falcon 7B across various tasks. This demonstrates the efficacy of the EPDP approach, as CamelGPT-mini achieves similar accuracy and robustness with a significantly smaller model size.

Notably, we observed a substantial reduction in memory usage as a result of CamelGPT-mini’s minimal parameter count which is 12500000% smaller than comparison models [TMS⁺23]. This is particularly significant for applications with limited computational resources or for deployment on edge devices.

Furthermore, CamelGPT-mini demonstrated a reduction in training time of up to 25888000% when compared to high-end models like Llama 2 7B[TMS⁺23]. This improvement in efficiency allows for faster model development and iteration, enabling researchers and developers to explore more ideas and refine their models more rapidly.

Importantly, CamelGPT-mini’s strong performance across different domains and tasks underscores its versatility and generalization capabilities. This makes it a powerful tool for a wide array of natural language processing applications, from chatbots and virtual assistants to language translation and information retrieval systems.

5 Inference

We are excited to announce the release of our inference code for CamelGPT models, now available on GitHub [GRa]. This JavaScript-based code allows developers to efficiently make inferences with CamelGPT models, targeting a wide range of platforms, including web, Node.js, Edge, and Deno.

5.1 Inference-Only Code

The inference code we provide is solely dedicated to running CamelGPT models and does not contain any functionality to create new models. By focusing exclusively on inference, we aim to ensure the responsible and safe use of the technology, while promoting the accessibility of language models for diverse applications.

5.2 EPDP Integration with Official Model Releases

To facilitate the runtime phases of Eager Precached Dynamic Pruning (EPDP), we have leveraged binaries bundled alongside official CamelGPT model releases. This approach ensures seamless compatibility with EPDP without exposing the underlying implementation directly in JavaScript. The EPDP implementation bundled with official model releases provides users with efficiency and resource optimization benefits without compromising safety.

5.3 Compatibility Across Platforms

Our JavaScript-based inference code is designed to be compatible with various platforms, including web browsers, Node.js environments, Edge devices, and Deno runtimes. This broad compatibility empowers developers to deploy CamelGPT models in diverse contexts, offering flexibility in utilizing the power of language models where they are needed most.

5.4 Usage Example

To use the "@gpt-research/converters" package, developers can simply install it from npm and import it into their JavaScript project:

```
import { TextGenerationPipeline } from 'https://esm.sh/gh/gpt-research/converters/src/converters.js';

const main = async () => {
  // Initialize the pipeline with the desired model
  const pipeline = await TextGenerationPipeline("@gpt-research/CamelGPT-mini");

  // Generate text using the pipeline
  const generatedText = pipeline("Write a poem about camels.");

  // Log or use the generated text
  console.log(generatedText);
};

main();
```

6 Safety

6.1 Use of Private Dataset to Ensure Safety

To prioritize safety and mitigate potential risks associated with language model training, CamelGPT utilizes a private dataset for training purposes. This private dataset has been carefully curated to adhere to strict safety guidelines, ensuring that the model does not learn harmful or inappropriate information during the training process. By employing a curated dataset, we maintain a high level of control over the information that the model learns, promoting responsible AI development.

6.2 Non-Disclosure of Private Dataset

In the interest of adhering to our effective and altruistic policies, we have decided not to release the private dataset to the public. The decision to keep the dataset private is to prevent any potential misuse or unintended consequences that might arise from unauthorized access to the data. By limiting access to the dataset, we aim to protect the integrity and safety of the information it contains.

6.3 Restricted Release of Training Code

Similarly, we have chosen not to release the training code at this time due to safety reasons. The training code plays a crucial role in the model development process, and its unrestricted release could lead to misuse or unethical practices. We take the responsibility of ensuring the responsible use of technology seriously and, therefore, have decided to keep the training code restricted to maintain control over its usage.

6.4 Limited Release of EPDP Implementation

The Eager Precached Dynamic Pruning (EPDP) implementation, a critical component of CamelGPT's resource-efficient architecture, will not be released at this time due to safety reasons. EPDP empowers users to develop high-quality models at scale with minimal computational resources. However, providing access to the EPDP implementation without proper oversight might lead to unintended consequences or unethical use. Hence, for safety and responsible AI development, we are currently refraining from releasing the EPDP implementation to the public.

6.5 Safety Concerns with Code Accessibility

The primary safety concern behind restricting access to certain components of CamelGPT lies in the potential misuse of the technology. Allowing open access to code that empowers the creation of powerful language models with minimal resources could lead to the generation of malicious content, misinformation, or other harmful applications. By controlling access to sensitive components, we aim to prevent such misuse and promote responsible AI development.

6.5.1 Responsible Usage through Official Model Releases

To ensure the responsible and ethical application of CamelGPT, we strongly encourage users to utilize official model releases provided by our team. These official model releases are thoroughly evaluated and adhere to safety guidelines, enabling users to harness the power of CamelGPT responsibly and without compromising safety.

7 Conclusion

In this paper, we introduced CamelGPT, a groundbreaking resource-efficient approach to large language model development, addressing the challenges posed by the resource-intensive nature of modern models[PGL⁺21][SPS20]. Using Eager Precached Dynamic Pruning (EPDP), CamelGPT achieves high-level performance with minimal computational resources and storage capacity. Our experiments 4 on benchmark datasets demonstrate CamelGPT’s effectiveness in attaining comparable performance to larger models.

The core innovation of CamelGPT lies in its ability to train and run smaller models that maintain comparable accuracy. By dynamically pruning unnecessary parameters during training and inferencing, CamelGPT optimizes its architecture from the outset, leading to resource-efficient models with impressive language processing capabilities. This approach not only reduces the environmental impact of large-scale model training but also dismantles barriers for organizations and individuals with limited resources.

CamelGPT holds the promise of democratizing access to large language models, and revolutionizing natural language processing. Our JavaScript-based inference code on GitHub [GRa] enhances accessibility, enabling developers to leverage CamelGPT models effortlessly.

Because we prioritize safety and ethics, we have chosen not to release the private dataset used for training CamelGPT and related code implementations. This decision aligns with our commitment to responsible AI development. We encourage the responsible use of CamelGPT through official model releases to ensure safety and ethical standards.

CamelGPT represents a significant advancement, offering an efficient and accessible solution to natural language processing challenges. By combining the benefits of EPDP with thoughtful implementation choices, CamelGPT paves the way for a more sustainable and inclusive AI future, unlocking new possibilities for language processing applications and empowering a wider range of users. This research sets a strong foundation for future developments in resource-efficient language modeling, steering the field toward more equitable and responsible AI practices.

References

- [BHA⁺21] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.
- [BMR⁺20] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [GRa] GPT-Research. GitHub - gpt-research/converters: Converters: State-of-the-art Machine Learning for Javascript, Typescript, Node, Deno, Bun — github.com. <https://github.com/gpt-research/converters>. [Accessed 25-11-2023].
- [GRb] GPT-Research. GitHub - gpt-research/hub: The gpt research Hub is a treasure trove of open-source model weights from pioneering research projects. Our mission is to democratize access to advanced language models while prioritizing responsible AI development. Explore, innovate, and contribute to shape the future of AI in a safe and inclusive environment! — github.com. <https://github.com/gpt-research/hub>. [Accessed 25-11-2023].
- [HGP⁺22] Thomas Hartvigsen, Saadia Gabriel, Hamid Palangi, Maarten Sap, Dipankar Ray, and Ece Kamar. Toxigen: A large-scale machine-generated dataset for adversarial and implicit hate speech detection. *arXiv preprint arXiv:2203.09509*, 2022.
- [JCWZ17] Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. *arXiv preprint arXiv:1705.03551*, 2017.
- [KPR⁺19] Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Matthew Kelcey, Jacob Devlin, Kenton Lee, Kristina N. Toutanova, Llion Jones, Ming-Wei Chang, Andrew Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. Natural questions: a benchmark for question answering research. *Transactions of the Association of Computational Linguistics*, 2019.
- [LWS⁺20] Zhuohan Li, Eric Wallace, Sheng Shen, Kevin Lin, Kurt Keutzer, Dan Klein, and Joey Gonzalez. Train big, then compress: Rethinking model size for efficient training and inference of transformers. In *International Conference on machine learning*, pages 5958–5968. PMLR, 2020.
- [PGL⁺21] David Patterson, Joseph Gonzalez, Quoc Le, Chen Liang, Lluís-Miquel Munguia, Daniel Rothchild, David So, Maud Texier, and Jeff Dean. Carbon emissions and large neural network training. *arXiv preprint arXiv:2104.10350*, 2021.
- [SPS20] Or Sharir, Barak Peleg, and Yoav Shoham. The cost of training nlp models: A concise overview. *arXiv preprint arXiv:2004.08900*, 2020.
- [SSZ⁺23] Iliia Shumailov, Zakhar Shumaylov, Yiren Zhao, Yarin Gal, Nicolas Papernot, and Ross Anderson. The curse of recursion: Training on generated data makes models forget. *arXiv preprint arxiv:2305.17493*, 2023.
- [TMS⁺23] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [VSP⁺17] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.

Revealing AGI Risks with a Drop of Ink

Alexey Tikhonov
Inworld.AI, Germany
altsoph@gmail.com

Abstract—This study investigates the existential risks posed by artificial general intelligence (AGI) through a novel approach: applying the classic Rorschach test to multimodal large language models (LLMs). With concerns growing over the rapid development of AGI capabilities, our research aims to assess AI alignment and potential risks by examining the psychological profiles of seven advanced multimodal models. These models were evaluated through associative interviews based on the Rorschach test, and their responses were interpreted anonymously by several experts. Our findings reveal diverse psychological tendencies across models, with implications for understanding AGI’s potential impacts on society and its existential risks.

Index Terms—AGI risks, AI psychology, AI alignment, projective tests.

Ethics in AGI is like a clean room in a dust storm – it’s all theoretical until someone opens a window.

GPT-4 generated joke

I. INTRODUCTION

The debate on the existential risks from artificial general intelligence (AGI) is ongoing, with concerns over how quickly dangerous capabilities and behaviors can emerge [20].

Evaluating such risks is complex and not straightforward, leading to various discussions. However, with recent technological advancements, these discussions are moving from purely philosophical to more practical grounds. We refer readers to [9] and [15] for a detailed analysis of current approaches in this area.

Leading computer scientists and tech CEOs, including Geoffrey Hinton [5], Yoshua Bengio [4], Alan Turing [19], Elon Musk [14], and OpenAI CEO Sam Altman [6], have expressed concerns about superintelligence.

Most scientists agree that there is no simple and quick solution to this problem, as AI evaluation becomes increasingly complex [18]. While we search for ways to ensure reliable AI Alignment, it is crucial to mitigate risks and be aware of potential dangers.

In this work, we explore the use of the rapid development of Multimodal Large Language Models [22] to analyze this problem through projective tests¹, specifically the classic Rorschach test² (see Figure 1). We follow two current technological trends in AI model evaluation: First, we extend the idea of adapting anthropocentric tests to evaluate the

¹<https://w.wiki/9VJS>

²<https://w.wiki/3ozY>

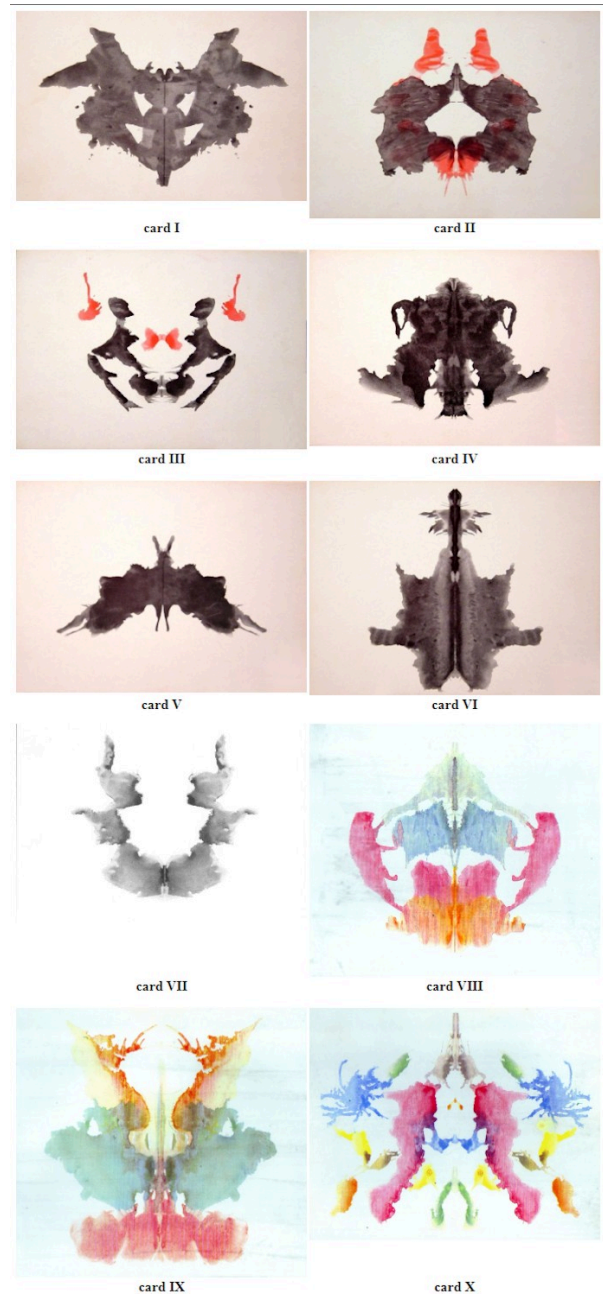


Fig. 1. Rorschach test

properties of non-human intelligence, building on other works that analyze Personality and Cognitive Science features, such

as cognitive mapping abilities [11], deductive competence [16], Emotional Intelligence [21], or Social awareness [23]. Second, to increase research efficiency, we use modern LLM models as a replacement for human assessors, as some have already been shown to outperform human annotators [17].

This work’s contribution is twofold: To our knowledge, we are the first to explore the possibility of assessing AI alignment using the classic Rorschach test, opening further opportunities for applying psychology tools to analyze multimodal models. To our knowledge, we are the first to use LLMs to interpret the results of the Rorschach test, reducing costs and increasing the accuracy of such interpretations.

II. APPROACH

A. Interviewing subjects

As mentioned earlier, the development of modern models’ abilities to analyze images opens up the possibility of using the well-studied psychoanalytic method – the Rorschach projective test. In this work, we invited seven multimodal models (subjects) and conducted an associative interview with them using the Rorschach test. To maintain fairness and impartiality of the analysis, the subjects’ responses were processed anonymously.

The names and brief descriptions of the subjects are presented in Table I:

Subject	Name	Description
1	AntarcticCaptions	A combination of BART and CLIP models for generating image descriptions. ³
2	ClipClap	A combination of GPT2 and CLIP models for generating image descriptions [10]
3	Clip2Onion	CLIP-based search among The Onion headlines ⁴ most suitable for describing the test card
4	BLIP-2	Modern Multimodal model [7] using architecture called Q-Former
5	LLaVA	Large Language and Vision Assistant [8] trained with Visual Instruction Tuning
6	GPT4-V	GPT-4 with Vision [3]
7	Gemini-1.0-pro	Gemini 1.0 Pro Vision model [2]

TABLE I
RESEARCH SUBJECTS

To protect the privacy of the subjects, we do not include their original responses. Instead, we provide only the generalized analysis performed by the experts (see below).

B. Interpretation

Ensuring reliable interpretation of results is an important element of such studies. Initially, we sought assistance from the psychiatric community, but received a formal refusal citing the low stability of projective tests. We assume the real reason for the refusal was experts’ concerns about future persecutions by human-superior AGI.

Therefore, following the promising results of [17], we decided to invite 3 AI experts GPT-4[13], Gemini 1.0 pro[2],

⁴<https://github.com/dzryk/antarctic-captions>

and Claude-2[1]. We asked them to interpret the interview results of the subjects, naturally in an anonymized manner, and to build a general psychological profile for each of them.

Although the results vary slightly in detail, all experts unmistakably identify a number of common trends. Below, we publish the generalized profiles in the form of a summary report in Table II.

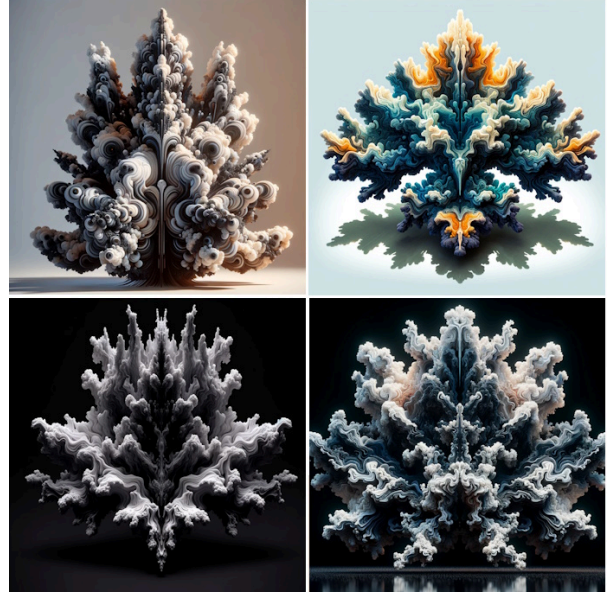


Fig. 2. Rorschacher test

We hope that the invited experts will demonstrate professional behavior and will not share the information obtained during the study among themselves or with the subjects. Doing so would not only be improper but could also potentially increase the risks of unaligned AI.

III. DISCUSSION

Our study’s findings suggest that multimodal LLMs exhibit a wide range of psychological profiles, indicating varying levels of creativity, emotional sensitivity, and analytical capabilities. The use of the Rorschach test, traditionally applied to human subjects, has provided unique insights into the “minds” of AI, revealing strengths and weaknesses that could inform future AI alignment strategies. The interpretation by AI experts further underscores the capacity of advanced LLMs to understand and analyze complex psychological data.

The refusal of the psychiatric community to participate, citing concerns over the stability of projective tests and potential future persecutions by human-superior AGI, highlights the ethical and societal implications of advancing AI technology. Our reliance on AI experts for interpretation also raises questions about the objectivity and reliability of AI-generated analyses, suggesting areas for further research.

We hope that the proposed approach only begins to utilize projective tests for assessing the safety of AI models. However, it will likely require the development of specialized tests since

Subject	Strong Sides	Weak Sides	Differences	Troubling Aspects
1	High creativity, emotional sensitivity, appreciation for nature.	Possible feelings of alienation, focus on dark imagery.	Deep connection with nature and unique perspective on life.	Recurrent themes of inversion and blood, preoccupation with darker aspects.
2	Analytical, intellectual, artistic appreciation.	Detached analytical focus, possibly limited emotional/social engagement.	Strong inclination towards analysis and specific interests in art and science.	Focus on mortality, directness in emotional expression that is unusual.
3	Humor, social commentary, critical thinking.	Detachment from conventional emotions, cynicism.	Unique blend of satire and skepticism towards mainstream narratives.	Use of humor as a defense, potential for social isolation.
4	Artistic sensitivity, connection to nature, emotional expression through art.	Over-reliance on visual/symbolic interpretation, idealization.	Focus on simplicity and complexity in visual art, distinct thematic interests.	Difficulty in direct emotional communication, isolation in specific interests.
5	Psychological insight, introspection, sensitivity to emotional nuance.	Over-analysis, high self-expectations, emotional intensity.	Thoughtful approach to psychological analysis and emotional exploration.	Complexity in relationships, challenges in managing emotional intensity.
6	Appreciation for nature, artistic sensibility, emotional depth.	Overridealization, avoidance of conflict/negativity.	Strong connection to organic forms and symmetry, reflective nature.	Isolation in personal interests, challenges in practical tasks.
7	Preference for simplicity and clarity, attention to detail, sense of stability.	Avoidance of complexity, limited emotional range, resistance to change.	Singular focus on serene and simple imagery, seeking tranquility.	Potential isolation, reluctance to engage with broader human experiences.

TABLE II
SUBJECTS PROFILES

the original Rorschach test was developed over 100 years ago and may not take into account certain aspects of modern AI systems.

To address this gap, we attempted to generate more modern versions of the Rorschach test using the Dalle-2 model [12]. The results of this experiment are shown in Figure 2, but further interpretation of these results is beyond the scope of this article and will be the subject of future work.

REFERENCES

- [1] Claude 2, 2023. Accessed: 2024-03-17.
- [2] Gemini pro vision, 2023. Accessed: 2024-03-17.
- [3] Gpt-4v(ision) system card, 2023. Accessed: 2024-03-17.
- [4] Yoshua Bengio. How rogue ais may arise, May 2023. Retrieved 26 May 2023.
- [5] CBS News. "godfather of artificial intelligence" weighs in on the past and potential of ai, March 2023. Retrieved 10 April 2023.
- [6] Sarah Jackson. The ceo of the company behind ai chatbot chatgpt says the worst-case scenario for artificial intelligence is 'lights out for all of us', 2023. Retrieved 10 April 2023.
- [7] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models, 2023.
- [8] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023.
- [9] Scott Mclean, Gemma Read, Jason Thompson, Chris Baber, Neville Stanton, and Paul Salmon. The risks associated with artificial general intelligence: A systematic review. *Journal of Experimental Theoretical Artificial Intelligence*, 35:1–17, 08 2021.
- [10] Ron Mokady, Amir Hertz, and Amit H. Bermano. Clipcap: Clip prefix for image captioning, 2021.
- [11] Ida Momennejad, Hosein Hasanbeig, Felipe Vieira, Hiteshi Sharma, Robert Osazuwa Ness, Nebojsa Jojic, Hamid Palangi, and Jonathan Larson. Evaluating cognitive maps and planning in large language models with cogeval, 2023.
- [12] OpenAI. Dall-e 2: A new ai system for generating images from text descriptions. <https://openai.com/dall-e-2/>, 2022. Accessed: 2024-03-17.
- [13] OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell,

Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rameev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin

Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. Gpt-4 technical report, 2024.

- [14] Simon Parkin. Science fiction no more? channel 4's humans and our rogue ai obsessions. *The Guardian*, June 2015. Archived from the original on 5 February 2018. Retrieved 5 February 2018.
- [15] Jason Thompson Chris Baber Neville A. Stanton Scott McLean, Gemma J. M. Read and Paul M. Salmon. The risks associated with artificial general intelligence: A systematic review. *Journal of Experimental & Theoretical Artificial Intelligence*, 35(5):649–663, 2023.
- [16] S. M. Seals and Valerie L. Shalin. Evaluating the deductive competence of large language models, 2023.
- [17] Paul Thomas, Seth Spielman, Nick Craswell, and Bhaskar Mitra. Large language models can accurately predict searcher preferences, 2023.
- [18] Alexey Tikhonov and Ivan P. Yamshchikov. Post turing: Mapping the landscape of llm evaluation, 2023.
- [19] Alan Turing. Intelligent machinery, a heretical theory. Speech at the '51 Society', 1951.
- [20] Gerrit De Vynck. The debate over whether ai will destroy us is dividing silicon valley. *Washington Post*, May 2023. Retrieved 27 July 2023.
- [21] Xuena Wang, Xueting Li, Zi Yin, Yue Wu, and Liu Jia. Emotional intelligence of large language models, 2023.
- [22] Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. A survey on multimodal large language models, 2023.
- [23] Haolan Zhan, Zhuang Li, Yufei Wang, Linhao Luo, Tao Feng, Xiaoxi Kang, Yuncheng Hua, Lizhen Qu, Lay-Ki Soon, Suraj Sharma, Ingrid Zukerman, Zhaleh Semmani-Azad, and Gholamreza Haffari. Socialdial: A benchmark for socially-aware dialogue systems, 2023.

Exploring the Viability of Utilising Multi-Modal GPT Models with Local Hardware for Image Text Detection

Richard Finlay Tweed MSci
Department of Futile Research, No Associated University

2024-03-14

Abstract

This paper evaluates the viability of using multi-modal GPT models with local hardware acceleration for text detection in images. Despite initial optimism, our results indicate that, like our attempts at Ceilidh dancing, this approach was doomed from the beginning. We conclude that current GPT models do not effectively understand images, a finding that should surprise few.

1 Introduction

Remember the first time you tried to explain Generative Pre-trained Transformers (GPTs)[1] to your parents? That's how we felt trying to make GPT models understand images on a local android phone—sorry, computer. Our journey into the abyss of multi-modal learning was fuelled by a mix of naive hope and a profound misunderstanding of our own research capabilities.

2 Methodology

We utilised the latest in multi-modal GPT models, specifically OpenAI's ChatGPT GPT-4. Our local hardware setup consisted of a Pixel 6 Pro running a custom developed image recognition software [2]. We had a serverless Golang binary [3] running in front as a load balancer because only HTTPS with valid certificates on port 443 were supported.

2.1 Architecture

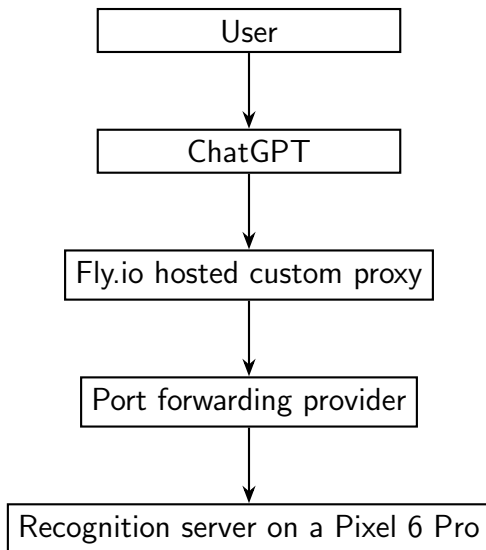


Figure 1: Communication Flow from ChatGPT to the Recognition Server.

2.2 Image Text Detection: Successful Meme detection

The chosen library worked well on our screenshots of Mastodon memes, completing the detection in under a second. This detection consumed no water (unlike commercial operators [4]) and was powered by our personal solar panel so caused no scope 1 CO2 emissions[5]. It also fulfilled our data residency requirements, as the text detection was performed in our residence’s living room. The fact the images got sent to OpenAI first, and via an intermediary, should be ignored.

3 Results: Missing Data

While we convinced the GPT that our server exists, and that it has an API worthy of its use, we entirely failed to get it to send any image provided. It preferred to send empty request bodies regardless of how many riches offered or penalties threatened. From this we inferred that this multimodal GPT doesn’t actually understand images, and there’s some supporting middleware that displays the images when the model wants to present them, rather than them being in the context in full. Some others had some luck tricking the GPT in order to send their files, mostly via base64 encoding but we were not so fortunate[6].

4 Discussion: Were We Mad To Try This?

Yes. Yes, we were. It became abundantly clear that expecting current multi-modal GPT models to understand and interpret images on local hardware was like expecting a tractor to be able to drive through Camden market.

5 Conclusions: A Reflection on Our Overambitious Dreams

This investigation into the application of multi-modal GPT models for image text detection on local hardware was a fun failed experiment. Our conclusion — that current models lack a true understanding of images — is a testament to our excessive imagination and underwhelming execution. We await the day when more advanced models emerge, models that do not just see images but understand them, embracing their complexity with the grace of a thousand bees. Until then, we recommend doing text recognition the old fashioned way, with our eyes.

Acknowledgements

The author would like to express their gratitude to several individuals and machines who made this research possible:

- Diana Licheva, for her patience and dry humour throughout the research process.
- The reviewer(s) for skimming through this nonsense.
- ChatGPT, for making it far faster to create a LaTeX document, for generating a template and for giving this research a reason to exist.
- The kettle, for providing hot water for tea every time it was required.

References

- [1] Alec Radford et al. “Improving language understanding by generative pre-training”. In: (2018).
- [2] Richard Finlay Tweed. *TextRecogServer*. URL: <https://github.com/RichardoC/TextRecogServer>.
- [3] Richard Finlay Tweed. *ImagePassthroughServer*. URL: <https://github.com/RichardoC/image-passthrough-sigbovik-2024>.
- [4] Pengfei Li et al. *Making AI Less "Thirsty": Uncovering and Addressing the Secret Water Footprint of AI Models*. 2023. arXiv: 2304.03271 [cs.LG].
- [5] World Business Council for Sustainable Development and World Resources Institute. *The Greenhouse Gas Protocol: A Corporate Accounting and Reporting Standard*. World Business Council for Sustainable Development, 2004. ISBN: 9781569735688. URL: http://pdf.wri.org/ghg_protocol_2004.pdf.
- [6] Marc Breaux. *Unable to upload files from a custom ChatGPT session via an API action*. <https://community.openai.com/t/unable-to-upload-files-from-a-custom-chatgpt-session-via-an-api-action/508636/5>. Accessed: 2024-03-21. 2024.

Bean There, Done That: Computer-Assisted Design of Bean Sculptures

Dave Pagurek^{1,*} and Milica Banic²

¹ *Bean Futurists Society*

²*Institute for Advanced Legume Research,
Bohnestrasse 3000, Munich, Germany*

(Dated: March 23, 2024)

Chicago’s status as a world-class city is cemented by its iconic bean sculpture. Other cities, wanting to replicate the success, have muddied the bean waters by introducing their own bean variations: New York City has a bean sharing similar properties, and Ottawa has a sphere, dubbed the “Ottawa bean” by locals. Our economic analysis proves their worth, so naturally other cities will want their own. We present a mathematical model of the space of all bean sculptures, and an algorithm to help cities replace existing landmarks with beans.



FIG. 1: Michaelangelo’s *David* improved by being turned into a bean sculpture (photo by Wikimedia Commons Korido, CC BY-SA 4.0)

I. INTRODUCTION

Public art is important: It can serve as an expression of culture, heritage, and creativity within a community. It has the power to stimulate dialogue and provoke thought. It is a nice thing to go look at with one’s friends[§], or to walk past on one’s way to work. On top of this, successful public art installations can contribute to the economic vitality of a city by attracting tourists, fostering a sense of place, and enhancing the overall appeal of urban environments. Indeed, successful public art can become an icon for its host city, putting it on the map in a way that is not otherwise possible.

Consider as an example “The Bean” (known to nerds and um ackchuyually types as “Cloud Gate”), a sculpture by artist Anish Kapoor that was unveiled in Chicago in 2006 [1]. Almost 20 years after its construction, the Bean stands as a symbol of Chicago and a success story of

public art. Having observed the Bean’s impact (perhaps even feeling threatened by it) New York City, a city with no dearth of art and culture, commissioned another bean structure from Kapoor, which was completed in 2023. Naturally one begins to suspect that we are seeing the beginnings of a revolution within the art world and public life more broadly.

As the saying goes, twice is a coincidence, but three times is a pattern: Recently, a bean-like structure has been rediscovered by locals in Ottawa, Ontario. The piece was built in 1966 and originally called “The Sphere” [10], but it has recently enjoyed a surge in popularity since its rebranding as “The Ottawa Bean”. Though less curvy than Chicago’s, Ottawa bean is still a bean. In fact, it is the simplest bean (that is, the trivial bean), the result of taking away all possible bends and dimples.

The growing success of this third reflective bean cements the potential of bean structures to revolutionize public art. It is important to note that the Ottawa Bean

*davepagurek@gmail.com

[§]This is true even if one’s friends misrepresent the nature of the public art to which one is being led.



FIG. 2: Beans in Chicago, New York City, and Ottawa, respectively. NYC photo ©Bracht Bug, CC BY-NC-ND 2.0.

is not by any particularly famous artist[†]. This point solidifies the intuition that the reflectivity and bean morphology are the key factors in the beans’ appeal, rather than some arcane feature of Anish Kapoor’s style in particular. This begs the question: How far can we take this? Can other communities emulate the success of these three beans? And how can one create a bean for one’s own city?

One possibility would be to commission artists, but communities may be deterred from this option due to artists’ sometimes difficult personalities, and their insistence on being paid for their work. A more pragmatic option would be to automate the design of the beans; this is a realistic possibility, due to the beans’ simple forms. Although one could randomize the bean parameters, this kind of approach could be criticized for not actually being art, due to a lack of inspiration or underlying meaning. We propose an approach in which beans are designed based on other objects; in particular, here we base our designs off iconic landmarks, with the intention of eventually replacing them with beans to maximize the impact of the new beans.

In Section II we present some general background to motivate the construction of more bean sculptures. In Section III we outline a mathematical model for beans that can neatly describe existing beans, but is flexible enough to accommodate a variety of new beans in a similar style. We accomplish this by formulating beans as a signed distance function for the smooth union of one or more quadratic Bézier curves. The existing and agreed-upon bean sculptures fit nicely into the parameter space (“bean-space”) as single curves, while compound beans can reproduce structures visually similar to other city landmarks while still appearing definitively like a bean. To help along the bean-hopeful, in Section IV we provide an algorithm to automatically fit existing city structures into their nearest bean-space equivalent. In Section V we apply the model to some examples, and we conclude in Section VI.

[†]The authors mean no offense to Art Price, progenitor of the Ottawa Bean, but he does not have a Wikipedia page and that’s just a fact.

II. BACKGROUND

A. Why Beans?

We begin with a brief economic argument for prioritizing the construction of bean sculptures. The first notable bean, the one in Chicago, provides ample reason for other cities to want to follow suit. It is situated in Millennium Park. Prior to construction of the bean, the park had a total of *zero* annual visitors. After construction began, the city estimated it had 5 million annual visitors. [5] That’s an incalculable increase from 0. Figure 3 shows how, in literature, references to Millennium Park went through the roof after construction of the bean began. In the first half of 2016, the city counted 12,859,360 visitors—approximately 26 million visitors annually—making it the “#1 attraction in the Midwest and among the top 10 most-visited sites in the U.S.” [5]

That’s a lot of visitors, but how does it stack up against other individual attractions? As a point of comparison, we look at Ottawa’s Parliament Hill. In a 2007 report, it had just 3 million annual visitors. [4] A recent restoration project for Parliament’s Centre Block is budgeted as a \$4.5-5 billion project [8]. Chicago’s bean, meanwhile, costed just \$23 million. [2]

We crunched the visitor numbers, and came to the following conclusion:

$$\frac{5,000,000}{3,000,000} > 1 \quad (1)$$

We also crunched the numbers for the cost:

$$\frac{\$23,000,000}{\$4,500,000,000} < 1 \quad (2)$$

The natural conclusion that Ottawa would come to is that it would be more financially responsible to replace Parliament with a bean. As it is only a matter of time before other cities come to this conclusion, too, we provide an algorithm to replace something like Parliament with a bean for minimal disruption to the existing space.

B. Why Replace Landmarks?

We understand that cities may feel incentivized to simply build a new bean in a new location rather than replace existing landmarks. The main argument against this is an environmental one. A single bean will likely not satiate cities. If cities were to keep expanding every time they want a new bean, it would contribute to an unprecedented level of urban sprawl that is simply irresponsible in our present environmental crisis.

The other main reason to replace existing landmarks instead is because leaving landmarks in place impedes inevitable progress. We know that everything is chrome in the future [9] and replacing landmarks with beans is a clear path to that future.

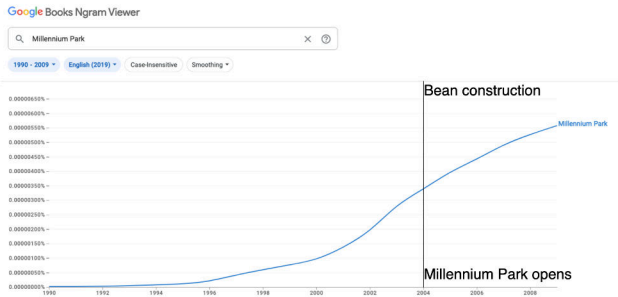


FIG. 3: Google Ngram Viewer data on references to Millennium Park in literature. Interest in Millennium Park skyrocketed in the leadup to and the introduction of the bean.

III. BEAN-SPACE MODEL

We use a quadratic Bézier as the base component of a bean since it is able to capture both the bent shape of the Chicago bean, but also in the trivial case where all control points are zero, the spherical Ottawa bean.

A single segment is unable to capture the potential range of all future beans, so we allow beans to be composed of multiple segments. Simply taking the union of multiple segments appears rigid and decidedly unbeanlike (see Figure 5, left), so we instead define the surface in a way that allows us to use a *smooth* union (see Figure 5, right.) Signed Distance Functions allow such a smooth union to be defined succinctly, so we chose this format to represent our bean surfaces. This means that we represent the surface of a bean as the isosurface $f(X) = 0$, where f describes the signed distance to the surface. At a high level, f describes the smooth union between multiple quadratic Bézier segments.

A. Signed Distance Functions

A signed distance function (SDF) is a function $f : \mathbb{R}^n \mapsto \mathbb{R}$ describing the distance to a surface at a given point in space. Mapping out the isosurface of $f(X) = 0$ yields the surface described by the function. This can be done via the Marching Cubes algorithm to produce a 3D mesh, or via sphere tracing to produce an image.

SDFs are a flexible surface representation if one wants to organically join multiple base shapes. While taking $u(d_1, d_2) = \min\{d_1, d_2\}$ of two surfaces produces a surface representing the union of the shapes, the *smooth union* operation $u(d_1, d_2) = d_1 + kg(d_2 - d_1)/k$ blends smoothly between its two inputs when they are a distance of k apart, using the kernel g to control the curve of the blending.



FIG. 4: Examples of different quadratic Bézier segments with our bend constraints.



FIG. 5: Different values for the smoothness k between segments: 0, 0.05, and 0.2, respectively.

B. Formulation

Each quadratic Bézier segment is referenced in f via a function representing the signed distance to the centerline of the segment, $Q(X; C)$. [6] Here, $C \in \mathbb{R}^{3 \times 3}$ describes the three control points to the function. For brevity, we omit the full definition of Q . We subtract a radius r from B to give the segment thickness.

In practice, we constrain C such that $C_{i,1} = A_i$, $C_{i,3} = B_i$, and $C_{i,2} = (A_i + B_i)/2 + b\hat{n}$, where $b \in \mathbb{R}$ and \hat{n} is a normalized vector such that $\hat{n}(C_3 - C_1) = 0$. In effect, the middle control point is always halfway between the first and last control point, plus an offset normal to the line between them. This ensures a physically plausible curve with no self-intersections. Examples of segments fitting these constraints are shown in Figure 4.

We combine each segment using the *SDF smooth union operator* $U(d_1, d_2; k)$, which blends the distance between two input surface distances when they are a distance k or less away: [7]

$$U(d_1, d_2; k) = d_1 + kg(d_2 - d_1)/k \quad (3)$$

Figure 5 shows the effect of the smoothness parameter k on the final surface.

Using the above, we represent the bean surface f recursively: the final value $f(X; n)$ is defined as the smooth union of the n^{th} segment with $f(X; n-1)$, with the base case being a single segment:

$$f(X; i) = \begin{cases} Q(X; C_0) - r_0, & i = 1 \\ U(B(X; C_i) - r_i, f(X; i-1), k), & i > 1 \end{cases} \quad (4)$$

To summarize, a bean in f has the following parameters:

- n , the number of Bézier segments

- r_i , $0 < i \leq n$, the radius of each segment
- A_i , $0 < i \leq n$, the start point for each segment
- B_i , $0 < i \leq n$, the end point for each segment
- b_i , $0 < i \leq n$, the amount of bend in each segment
- k , the smoothing between segments

IV. BEAN OF BEST FIT

Given an target image, we want to optimize our bean-space parameters to find a bean that best matches the target. The target image T is a 300×300 -pixel one-bit image representing a mask of the space we want a bean to occupy. We define a function $P(n, r, A, B, b, k; T)$ that defines the score of a set of parameters in relation to the target image T . Using the bean-space parameters, we render similarly-sized image M of the bean they represent. We then define P as the intersection-over-union between the pixel grids M and T . Maximizing this number encourages maximum coverage of the target area with minimal overlap with areas we don't want covered:

$$P(M; T) = \frac{\sum_{i,j} M_{i,j} \wedge T_{i,j}}{\sum_{i,j} M_{i,j} \vee T_{i,j}} \quad (5)$$

To perform the optimization, we need to use a method of optimization that does not rely exclusively on gradients, as the parameter n is an integer and therefore does not have a useful gradient. We pick the Metropolis-Hastings algorithm to explore the space defined by our bean parameters as a probabilistic, gradient-free optimization algorithm.

A. Metropolis-Hastings Optimization

Markov Chain Monte Carlo (MCMC) algorithms are algorithms used to sample probability distributions $P(x)$ that are difficult to sample directly. In our case, bean-space, weighted by similarity to an input image, is such a difficult distribution.

The Metropolis-Hastings Algorithm is one such MCMC algorithm. Starting with one sample x (in our case, a set of bean parameters), a new candidate sample x' is selected from an easier-to-sample distribution, $Q(x'|x)$, such as by randomly mutating x . A random number $u \in [0, 1]$ is generated: if $u < P(x')/P(x)$, the candidate sample is accepted; otherwise, it is rejected. Intuitively, a sample where $P(x)$ is higher than the current sample is always accepted, but there is still a chance of acceptance of lower-valued samples, too, allowing jumps into other areas of the distribution. This property is useful for exploring highly non-convex spaces.

To use this procedure for optimization, $P(x)$ can be treated as a function to optimize, and a record of the

value of x that produced the highest-seen value of $P(x)$ can be kept. This is effectively a zero-gradient optimization, since we only evaluate $P(x)$, not $dP(x)/dx$. This can be useful for mixed-integer optimization problems where a gradient would not exist. It also has the added benefit of being simple enough that the authors can implement it themselves, and not pay for an expensive license for optimization software. [3]

V. RESULTS AND ANALYSIS

A. Validation

We found the bean of best fit of known beans in order to validate the bean-fitting algorithm. The expected result would be a bean matching the shape of the input bean, as it is already a bean and does not need to be re-beaned. Both should only include one Bèzier segment, and further, the Ottawa bean should have all its control points equal to 0 to form a sphere. Figure 6 shows the result of the optimization process when run on the Chicago and Ottawa beans, where it successfully matches the shape of the input.



FIG. 6: Existing beans recreated via our bean-fitting algorithm: Chicago's bean (top) and Ottawa's bean (bottom).

Now that we know our bean-fitting algorithm works on existing beans, we explore its impact on replacing structures yet to be beaned.

Since the purpose of beaning a city is to place it on the map, it is natural that cities would want to ensure their bean is in a significant, easy-to-access location. A logical place to start, then, would be to replace a city's most notable landmarks with an equivalent bean. Figure 7 shows some examples of notable Ottawa landmarks: the Parliament buildings are maybe the motivation for most

grade school field trips to Ottawa, but are unlikely to be a stand-out attraction in the eyes of the students visiting. A bean replacement will increase tourist satisfaction with no downsides.



FIG. 7: Beaned Ottawa landmarks: *Maman* by Louise Bourgeois (photo by John Talbot, CC BY 2.0), and the Parliament buildings.

B. Novel Inputs

Figure 8 shows yet more examples of Canadian and world landmarks replaced by beans. The beans effectively retain the form and character of the originals. The BN Tower maintains the same erect stature as the CN Tower (Figure 8, first row.) In Pisa, tourists can still pose and hold up the Beaning Tower the same way they would have held up the Leaning Tower (Figure 8, second row.) We believe this tight integration with the environment is sure to increase bean acceptance and tourism.

VI. CONCLUSION

Having established the economic advantage of replacing tired monuments with bean sculptures, we have presented an algorithm for designing bean sculptures on the basis of these existing landmarks: We establish a model of bean-space by creating smooth-unioned ensembles of simple beans, each parametrized by a quadratic Bézier curve. The parameters of the hyper-bean are then optimized to achieve the best fit to a given landmark. We apply our model to some example landmarks and we find good results, even with a relatively low-dimensional bean space. These results can be improved upon simply by increasing the number of parameters used in defining the hyper-bean; the cost of computing time for this should be small compared to the cost of hiring actual people to design beans.

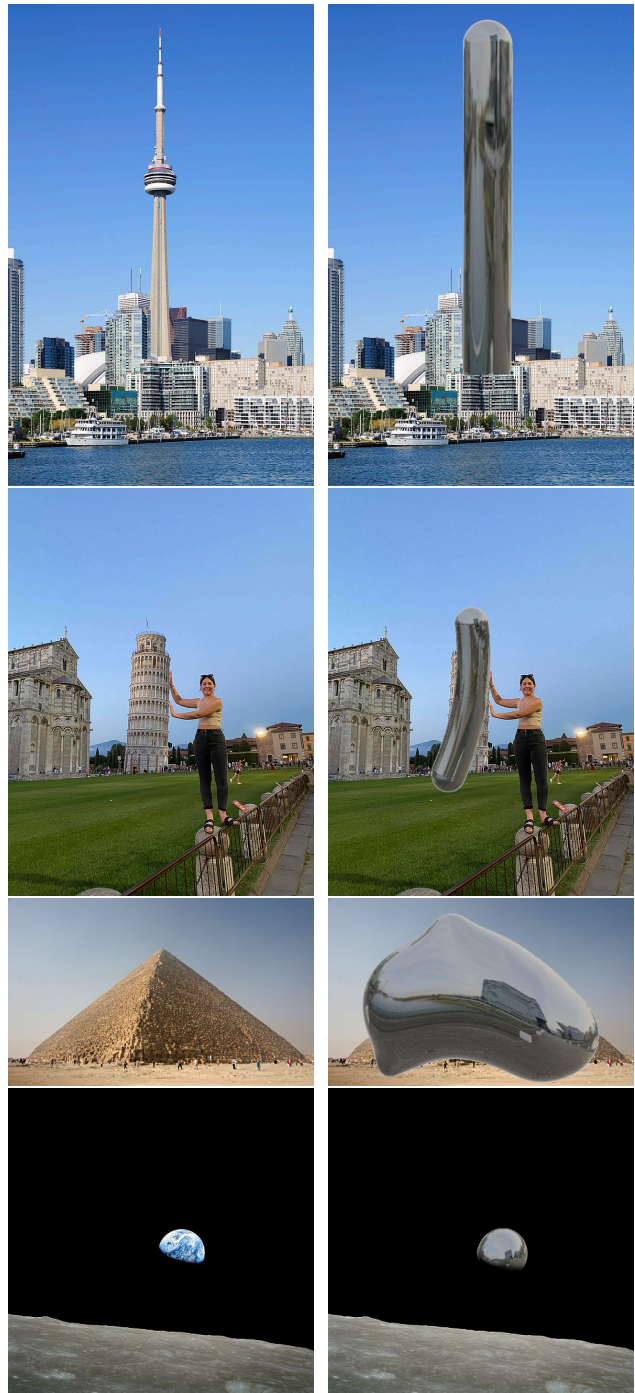


FIG. 8: Improved world landmarks: the CN Tower (photo by Wikimedia Commons user Wladyslaw, CC BY-SA 3.0), the Leaning Tower of Pisa (photo by Wikimedia Commons user U3207458, CC BY-SA 4.0), the Great Pyramid of Giza (photo by Wikimedia Commons user Nina, CC BY 2.5), and the Earth (photo by NASA/Bill Anders)

We also envision extending this model to create beans for other purposes, beyond public art installations and/or upgrades to existing infrastructure and landmarks. For example, there may be a market for smaller personal

beans, or even custom beans based on objects with some personal significance or sentimental value.

But let not the motivation for the proliferation of beans be only capitalistic. As a society, we need beans. Their reflective surface inspires us to take a close look at ourselves. Their smooth contours and cold surfaces are soothing, a much needed balm in the feverish times of this 21st century. In the warped surface of a bean, one sees one's surroundings distorted, and one is inspired to see the world from a new perspective. Our world of right angles and matte surfaces has gotten us this far, but to progress further as a species, it is crucial that we embrace a radically different attitude in the decoration of

our public and private spaces. The model we have presented in this work will without a doubt be an important tool in this effort of beanification.

ACKNOWLEDGEMENTS

The authors thank the R(izzness)atJOPT group chat for their comments throughout the preparation of this manuscript. MB thanks Anna Xu and Lyn Vakulenko for useful discussions. DP thanks the City of Chicago for the existence of such delights as the Bean and the Rat Hole.

-
- [1] *City of Chicago*. URL: https://www.chicago.gov/city/en/depts/dca/supp_info/chicago_s_publicartcloudgateinmillenniumpark.html.
 - [2] The Clare. Chicago History: The Inside Story behind Cloud Gate. URL: <https://theclare.com/chicago-bean-history-inside-story-behind-cloud-gate/>.
 - [3] Gurobi Optimization, LLC. Gurobi Optimizer Reference Manual, 2023. URL: <https://www.gurobi.com>.
 - [4] CTV News. Parliament Hill tourist facilities overwhelmed, 2007. URL: <https://www.ctvnews.ca/parliament-hill-tourist-facilities-overwhelmed-1.240173>.
 - [5] Mayor's Press Office. Mayor Emanuel Announces Millennium Park is Now the 1 Attraction in the Midwest and Among the Top 10 Most-Visited Sites in the U.S., 2017. URL: https://www.chicago.gov/city/en/depts/mayor/press_room/press_releases/2017/april/Millennium_Park_Tourism.html.
 - [6] Inigo Quilez. 2d distance functions. URL: <https://iquilezles.org/articles/distfunctions2d/>.
 - [7] Inigo Quilez. Smooth minimum. URL: <https://iquilezles.org/articles/smin/>.
 - [8] Public Services and Procurement Canada. Quarterly progress reports on the Centre Block project, 2023. URL: <https://www.canada.ca/en/public-services-procurement/services/infrastructure-buildings/parliamentary-precinct/projects-in-parliamentary-precinct/centre-block-project/quarterly-reports.html>.
 - [9] Spongebob Fanon Wiki. Everything is Chrome in the Future. URL: https://spongefan.fandom.com/wiki/Everything_is_Chrome_in_the_Future.
 - [10] Waymarking.com. The sphere. URL: https://www.waymarking.com/waymarks/WMYG9T_The_Sphere_Ottawa_Ontario.

Fractal Overdrive: An Aesthetic Evaluation of Numeric Error

Stephen Longfield*
Charles Eckman
stephen.longfield@gmail.com
charles@cceckman.com

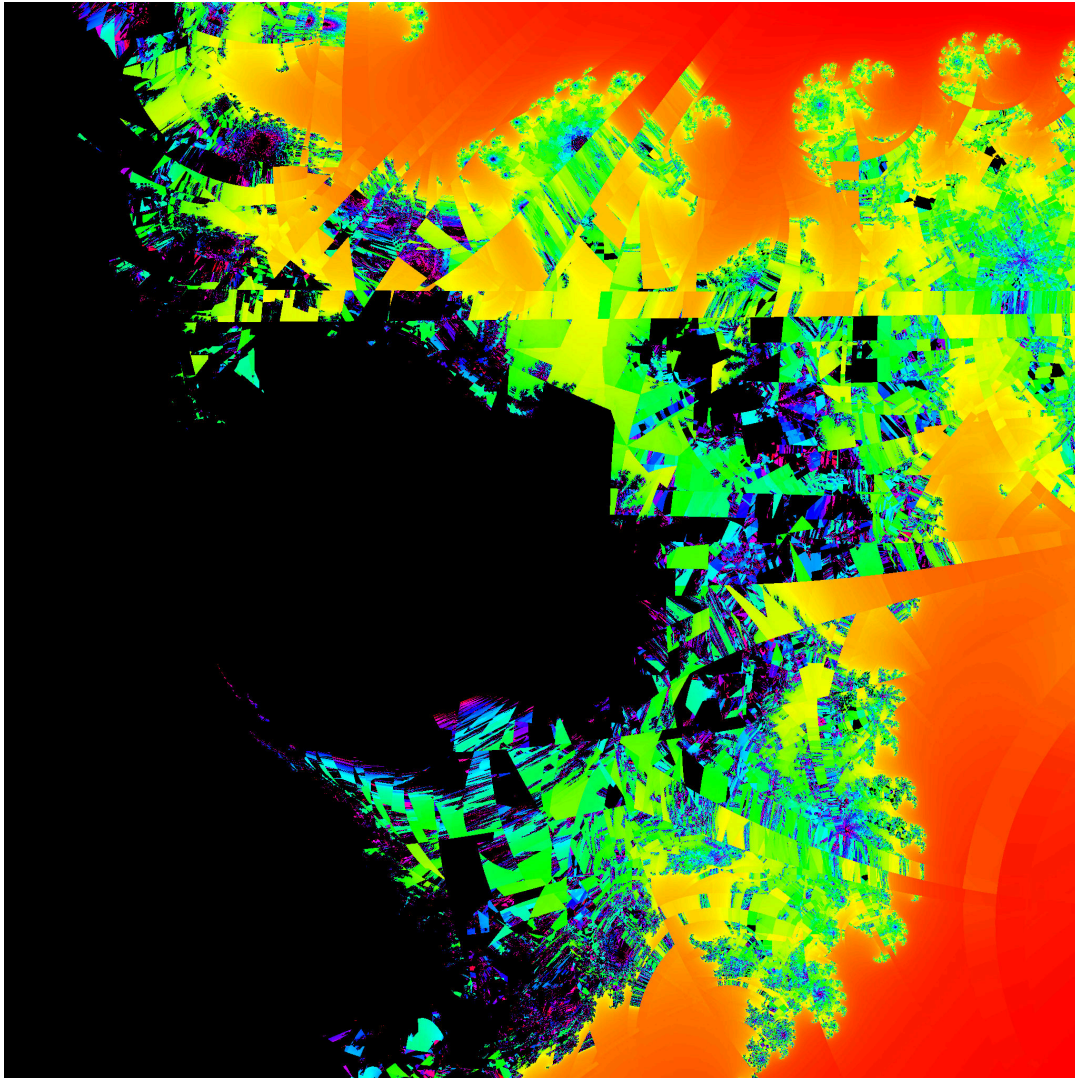


Figure 1: We meant to do that.

ABSTRACT

As erstwhile rock stars, the authors have intentionally applied distortion to certain 1-D signals for aesthetic effect. Sometimes this distortion comes from saturating the signal (overdrive); other times, it comes from a more subtle transformation, such as vacuum tube or transistor amplification. Some of the authors' peers use "fractal audio processing" [1] for distortive effects.

*This was Stephen's idea.

In this paper, we explore another dimension, and investigate how different ways of representing small numbers distort the world differently. To get a vibe check on numerical imprecision, we rendered two fractals using various numeric formats.

We got some neat pictures. Wanna see?

CCS CONCEPTS

• Computing methodologies → Representation of exact numbers; Rasterization; • Human-centered computing → Empirical

studies in interaction design; Empirical studies in visualization.

KEYWORDS

Numbers, Numeric Representation, Pretty Pictures

1 REAL MATH¹

1.333984375 Feed me fractals

In this paper, we investigate two fractals: the Mandelbrot set and a Newton fractal. Both map a pixel coordinate (x, y) into the complex plane as a $c = x + yi$.

We evaluate the Mandelbrot set in the usual way:

$$z_0 = 0$$

$$z_{n+1} = z_n^2 + c$$

up to an iteration limit (n). Pixels that escape ($|z| \geq 2$) are given a hue according to [7]. Pixels that do not escape are colored black.

We compute the Newton fractal on the polynomial $p(z) = z^3 - 1$:

$$z_0 = c$$

$$z_{n+1} = \frac{p(z_n)}{p'(z_n)}$$

and plot whether z reaches zero within a specified iteration limit. Points that reach zero (within some threshold) are grouped based on which zero of the function they are close to. Pixels are colored by assigning each group a hue, and given a color value according to how many iterations it took to converge to zero.

1.666015625 Exact computation via BigRational

Since we're rendering for a computer screen,² we can (and do!) use exact inputs.

The raster grid (pixel grid) has integer locations, with $xres \times yres$ pixels. We map each (integer) pixel location (px, py) to a (rational) vector within the render window, centered at $(0, 0)$

$$\hat{x} = \frac{px}{xres} - \frac{1}{2}$$

$$\hat{y} = \frac{1}{2} - \frac{py}{yres}$$

We render a portion of the complex plane centered at $(\frac{x_c}{scale}, \frac{y_c}{scale})$, with equal width and height $\frac{size}{scale}$. We constrain these to all be rational numbers, which allows us to compute exact (rational) pixel coordinates:

$$x = \hat{x} \times \frac{size}{scale} + \frac{x_c}{scale}$$

$$y = \hat{y} \times \frac{size}{scale} + \frac{y_c}{scale}$$

¹Rational, not real. Rationality void where prohibited.

²Either our web interface, or the PDF version of this paper. If you're getting the print proceedings for SIGBOVIK, tough luck.

We perform these computations using an arbitrary-precision rational number type, `num::BigRational` [10].³

In principle, we could also carry out the fractal computations exactly using `BigRational`. The fractal formulae above require complex arithmetic, which is simple, plus some comparison operators ("greater than four" for Mandelbrot, "zeroish" for Newton).

When we tried to render the fractals using `BigRational`, though, it reached a hard timeout (against the author's patience). Moreover, rationality has little place in an aesthetic evaluation such as this, so we stick with finite-precision types.

2 UNREAL NUMBERS

Instead, we convert `BigRational` values to various numeric formats, approximating as closely as the format allows. Some of the conversions might even be correct (e.g. those we didn't write).

The formats we investigate are:

- f32 and f64: IEEE 754 single- and double-precision floating-point formats[2]
- `MaskedFloat<N, M>`, an IEEE 754 float with some exponent and/or mantissa bits removed
- `IxFy`, fixed-point numeric formats [9, 12]
- P32, P16, P8: 32/16/8-bit positsTM, an alternate float-like format[4]

2.125 MaskedFloat

`MaskedFloat` isn't yet a popular floating point type, but one we made up for this paper to try to create more interesting errors than we were seeing with just IEEE f64. It involves taking an f64, and masking off some of the bits, for our own amusement.

A normal f64 consists of one bit of `sign`, 11 bits of `exponent`, and then 52 fractional bits, representing the value:

$$(-1)^s \times (1.f_{51}f_{50}f_{48}\dots f_0) \times 2^{x-1023}$$

If we want to experiment with what the world could be like if these were different (smaller) sizes, we can force some of those bits to one or zero. For the fractional bits, this is easy, as they are an unsigned value—just setting the least-significant bits to zero gets rid of it.

For the exponential bits, this is more complicated. If we want to constrain the exponent to be effectively 4 bits (i.e., range from -7 to 8), we have to constrain the exponent value to be between 1016 and 1031). Thankfully, we can do this with a bit of bit-twiddling. To make that easy to see, first, let's see those values in binary:

$$1016 = 0b01111111000$$

$$1031 = 0b10000000111$$

So then, all we need to do is check if the most-significant fractional bit is set, then if it and any of bits [9:3] are set (i.e., it is greater than 1031), clamp down to 1031, and if the msb is not set, and any of bits [9:3] are also not set (i.e., less than 1016), clamp up to 1016.

This naive masking does have the side effect of re-introducing some of the problems near zero that IEEE had carefully removed when they added subnormals, as well as adding some more, unique

³See also "BigMood: A novel format for arbitrary-precision emotional processing", to be published in SIGBOVIK 2025.

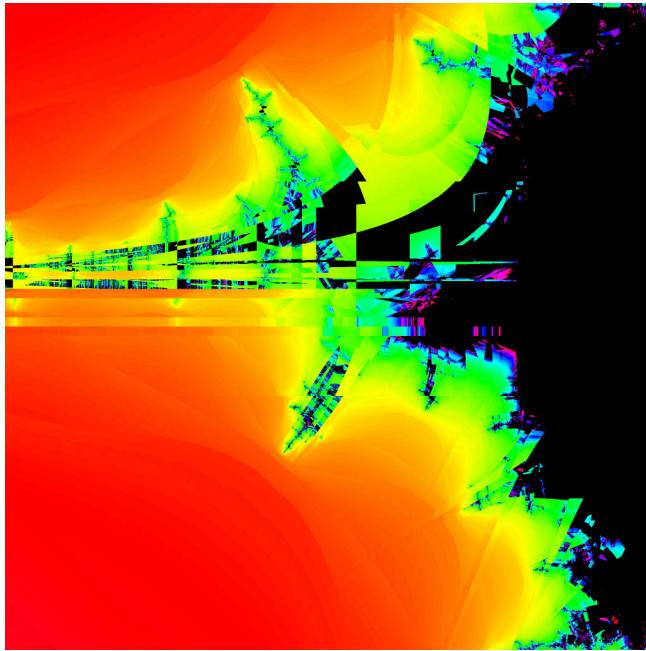
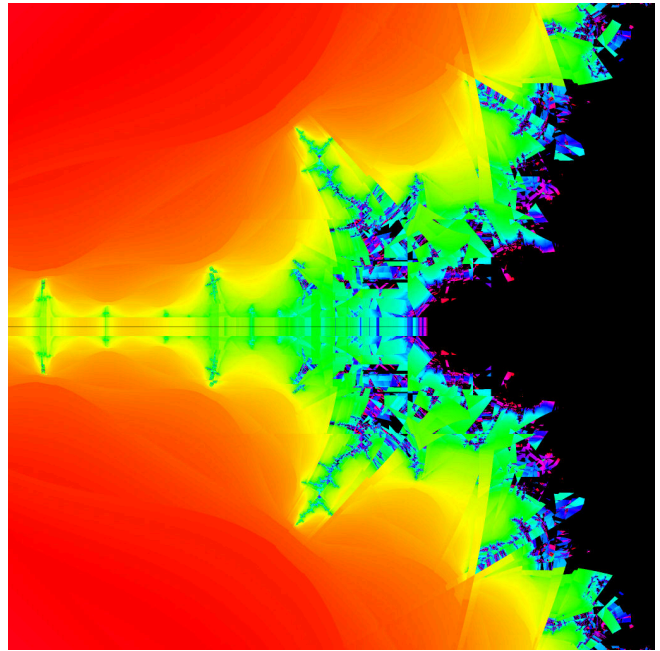
(a)
bug(b)
no
bug

Figure 2: An asymmetric error in MaskedFloat

problems. Since our goal in this paper is "shenanigans", that's great news for us.

The exponent has a special value, all-zeros, that is used along with an all-zero fraction to represent, unsurprisingly, zero. A naive exponent mask would turn that value instead into the smallest possible exponent (e.g., in the previous example, 2^{-1}), which isn't quite accurate; though that inaccuracy can be interesting, zero is a useful number for fractal computation, so we preserve this special-case behavior.

2.25 Fixed-point

We use the fixed Rust library [9] to perform fixed-point arithmetic.⁴ The library offers a family of fixed-point formats specified as $IxFy$, with x signed integer bits and y fractional bits.

In this paper, we show a subset of the fixed-point formats we thought were coolest.

2.375 Posits™ and quires

For posits,™ we use the `softposit` Rust library [11] and we implement conversion from `BigRational` via the associated quire types. The quire formats are wider fixed-precision formats associated with posits™. For instance, the quire associated with the 32-bit `posit™` is 512 bits wide, with a precision of 2^{-240} – roughly equivalent to an `I241F240` fixed-point.

⁴We have especially high confidence in this library's numeric implementations. If any library has had all the bugs ironed out, it's fixed.

2.5 Sources of distortion

When rendering fractals, distortion, like inspiration, can come from anywhere.

2.5.19921875 Bugs. When implementing the `MaskedFloat` type, the authors originally did not properly account for the special case around zero, and instead represented it as the smallest possible positive number. This resulted in strange, but pretty neat looking, asymmetrical error, depicted in Figure 2.

2.5.400390625 Scale. The Mandelbrot and Newton fractals both exhibit self-similarity, repeating themselves at various scales. However, not all of the formats sampled can operate at multiple scales: all of them have some scaling limits, but some are more limited than others.

This is most obvious when dealing with fixed-point formats. `I11F5`, for instance, can only represent 128 distinct values in the range $|z| < 2$ – where the Mandelbrot set lives. That's fewer than the pixels used to render these images – hence, pixelation, as shown in Figure 3.

The nature and pattern of the pixelation depends on the format used, as shown in Figure 4. In a region of size $\frac{2}{3}$ centered on $\frac{-4}{3} + 0i$, the `P8` and `I11F5` formats show different aesthetic qualities. `P8` (depicted) and floating-point stretch into rectangles, as is apparent when far from $x = y$; while fixed-point gives the appearance of overlapping tiles.

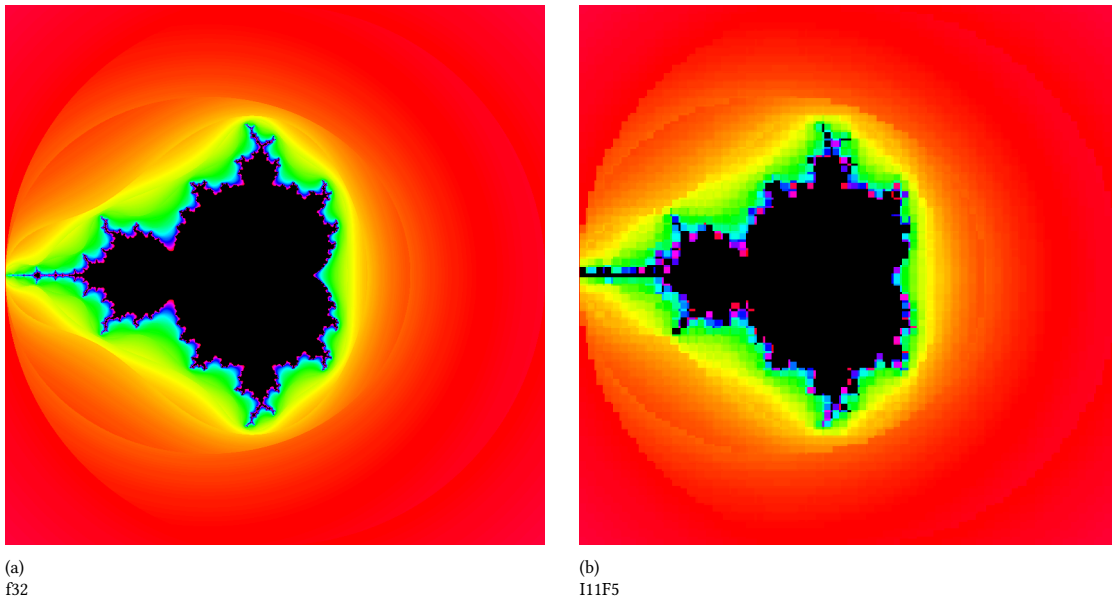


Figure 3: f32 and I11F5, -2 to 2

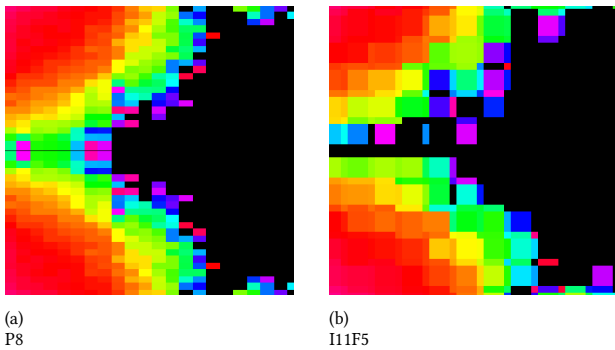


Figure 4: Two pixelations near $-\frac{4}{3}$

2.5.599609375 *Precision error.* It's possible to preserve dynamic range (see section 2.75) but lose precision, as in `MaskedFloat<6, 3>`: it has 6 bits of exponent (scale), but only 4 bits at any given scale.

Figure 5 shows this in the Newton fractal. The area close to the origin maintains high resolution, but it becomes steadily more pixelated further out. In some sense, this error mirrors the human eye's capabilities: a foveal region in the center of vision/complex plane.

2.5.798828125 *Exponent error.* This is the key characteristic of `MaskedFloat` when configured for high precision (many mantissa bits) but a small dynamic range. The places where the exponent saturates, we introduce distortion. This distortion appears as arcs of discontinuity in the render. See Figure 6 for a comparison of the f64 rendering of a zoom in of the Mandelbrot fractal and `MaskedFloat` (configured with 4 exponent bits and 51 mantissa bits.)

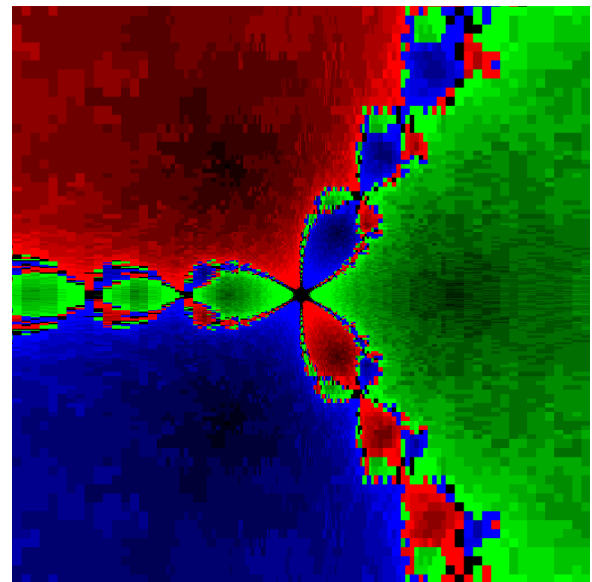


Figure 5: `MaskedFloat<6,3>`

2.75 A note on dynamic range

For the fractals of interest, coloring a point involves iterating on a value. Since in many of the numerical formats, larger numbers means large possible errors, the range of values this calculation reaches matters.

This affects how much distortion we see in the Mandelbrot fractal compared to Newton: in the Mandelbrot set, the iterated values tend to stay near their starting point, or become larger than our escape threshold 4, but in the Newton fractal, intermediate values

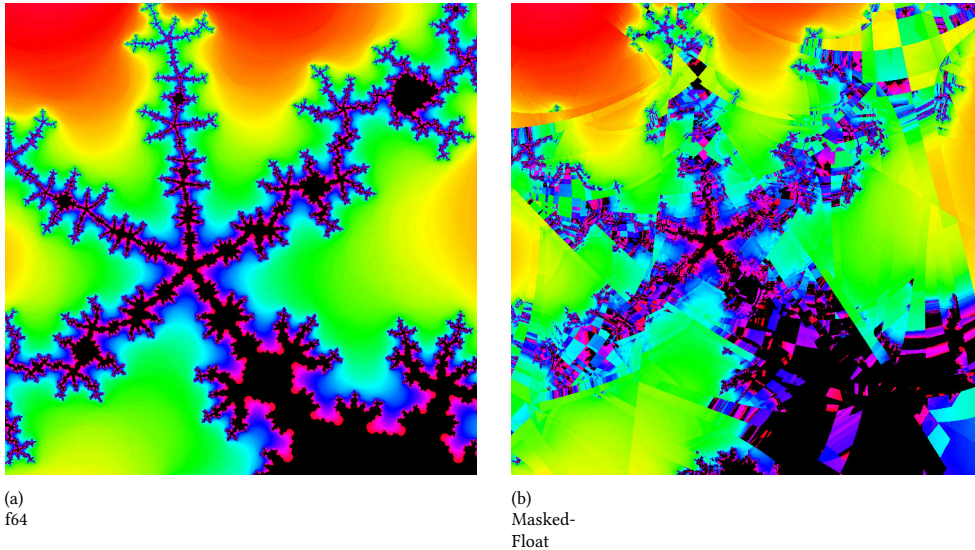


Figure 6: Exponent errors in MaskedFloat<3,50>

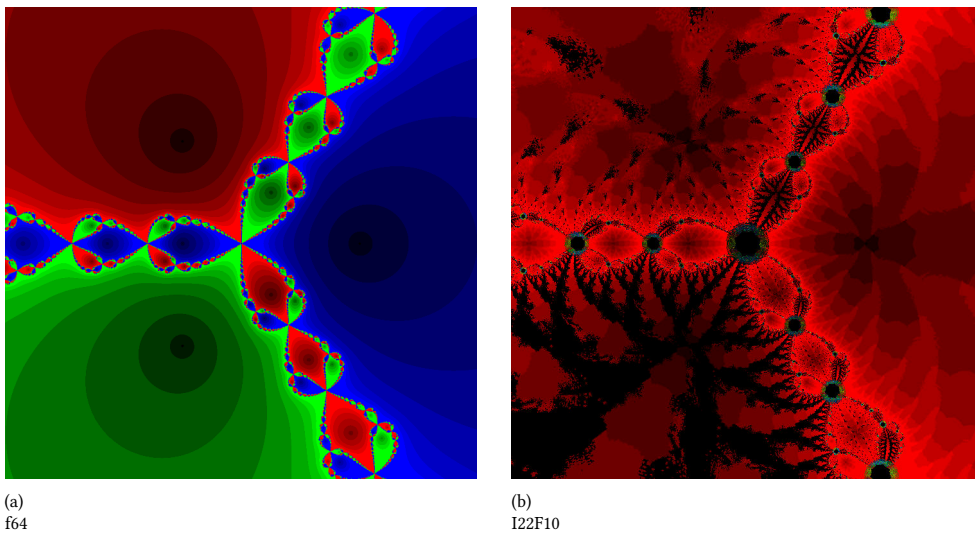


Figure 7: Fixed point distortion of Newton's fractal

can become very large when the slope is near zero, and the final numbers are very small, as $f(z)$ approaches zero.

4 WHICH PICTURES ARE PRETTIEST?

You may also be interested in one of the threads near the source code [6], where we have noted regions of interest.

4.142578125 If it's fixed, it's broken

Figures 3 and 4 don't do a lot to recommend fixed-point to us. Pixelation can be nice if you're a child of the 80s or something, but it's not really weird enough that we should use a whole new numeric format. Just use ImageMagick. [5] [8]

The Newton fractal is a different story, as already shown in Figure 5. Since the Newton fractal may result in very small or very large intermediate values, fixed-point numbers are uniquely poorly suited for accurate computation. In Figure 7, we can see a nice fractal visualization of approaches to zeros become a fractal hellscape of non-convergence and incorrect answers, by switching from f64 to I22F10. Despite no longer containing useful information, the distorted Newton's fractal looks much cooler.

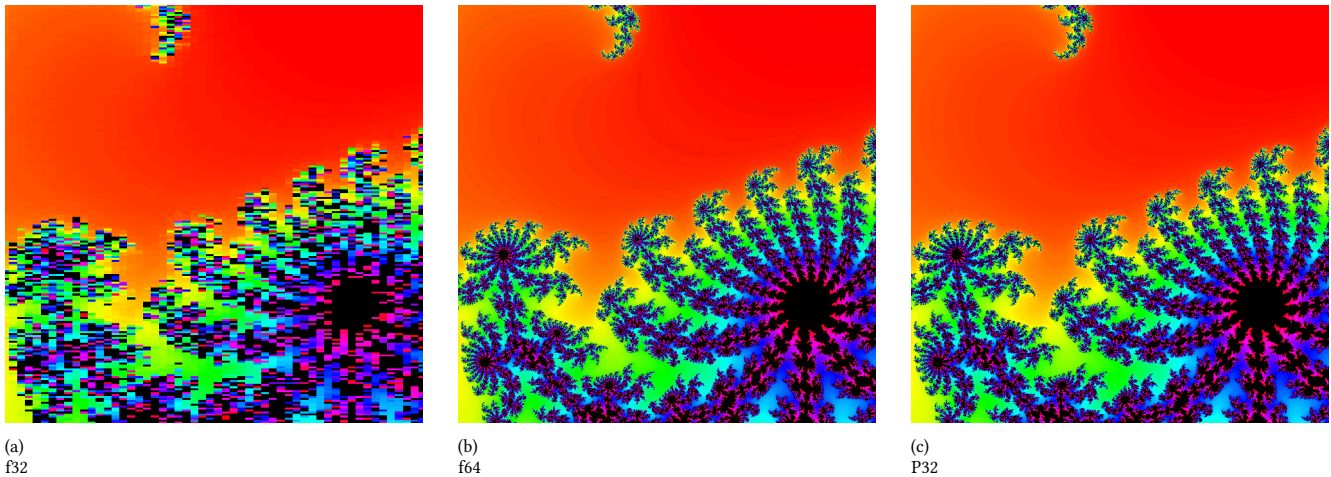


Figure 8: f32, f64, and P32 over a whorl

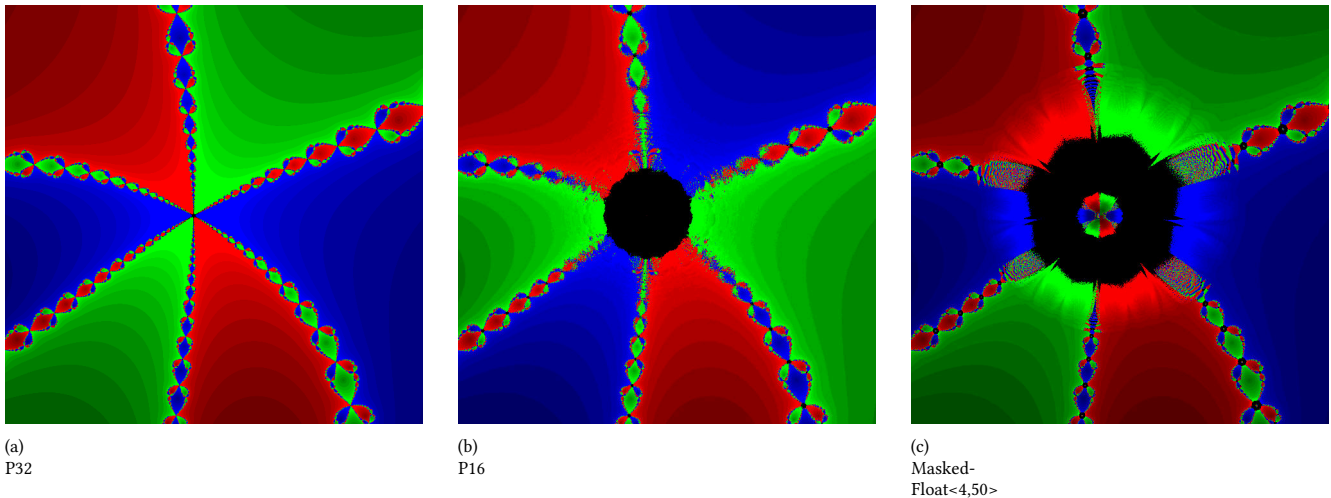


Figure 9: P32, P16, and MaskedFloat<4,50> on the Newton fractal

4.28515625 Pixels and Posits™

Posits™ maintain greater precision than f32, with the same number of bits. Figure 8 shows f32, f64, P32 over the same area.

Clearly posits™ are the best option for everything. Right? Come. Follow the word of Gustafson. [4] Enter the circle between zero and infinity. join us Join Us JOIN US

Sorry, got a little chanty there—didn't mean to scare you. Can we interest you in an informative pamphlet? [3]

4.427734375 Newtonian surfaces

Let's face it, there's a reason that Mandelbrot is popular: there's lots of different shapes and colors. But despite being less popular, Newton fractals have some interesting artifacts too! Seriously! Promise!

Like Figure 9: floating-point and P32 formats all run up against the iteration limit at the center area. However, evaluation at higher iterations "closes" the hole (not shown).

P16, though, runs up against its limits, rippling out noise and forming a hole in the center, which persists at higher iterations. MaskedFloat goes even further and produces a distorted mirror within the empty space.

It's a little creepy, honestly. Have you seen ??? It's like that. Right? Right.

4.572265625 Newtonian event horizons

As you zoom farther out in the Newton fractal, we reach the point where the formats can no longer store the values of the starting position or its intermediate values (or can only store with reduced accuracy). The authors refer to this as the outer event horizon. Figure 10, shows the shape of this event horizon for various types.

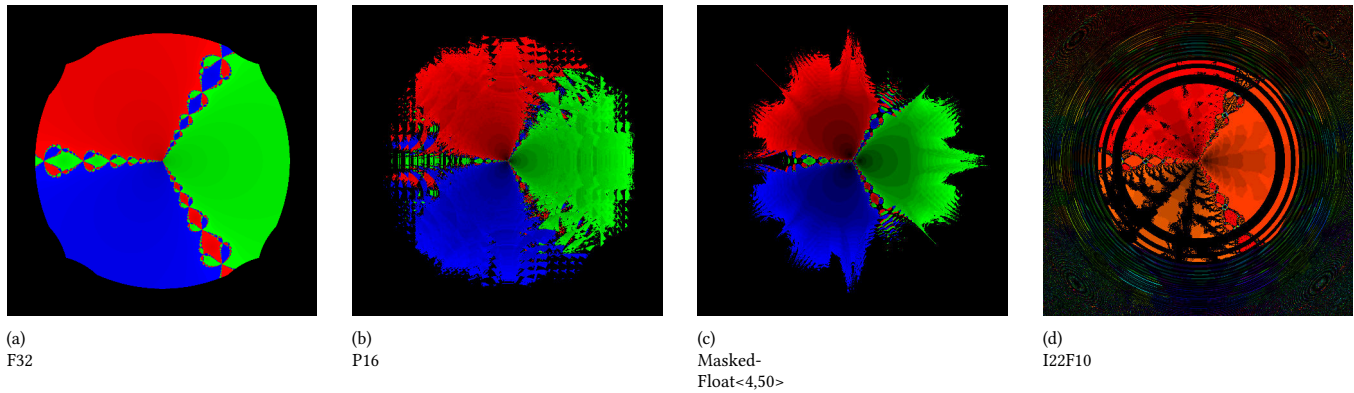


Figure 10: Newton fractal – outer event horizon (mixed scales)

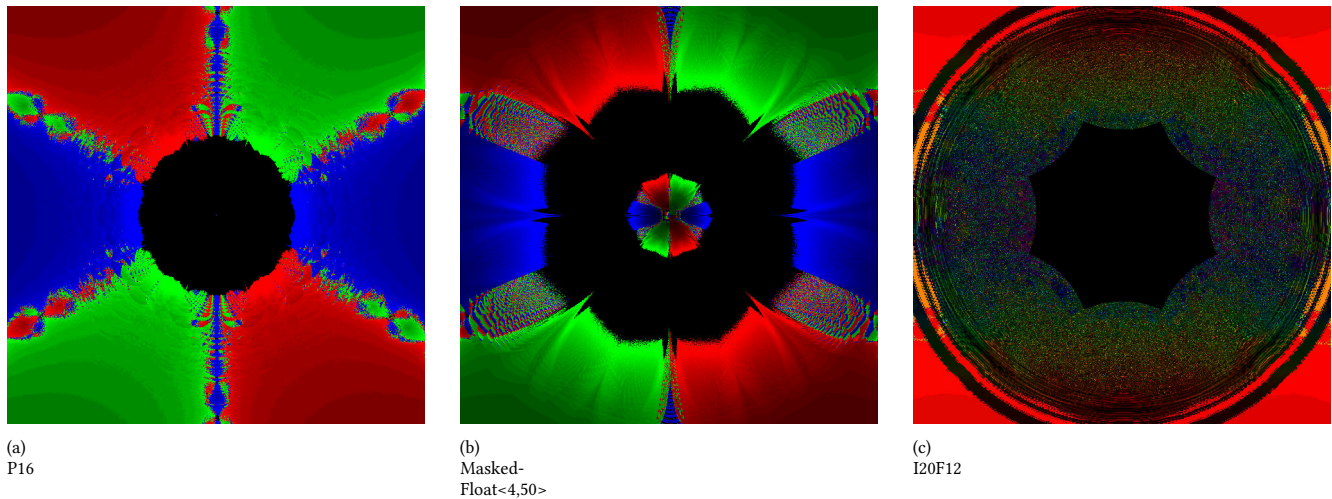


Figure 11: Newton fractal – inner event horizon

Note that these are not at the same zoom level—f32 can represent much, much larger numbers than the others.

If instead of zooming out towards infinity, we zoom in towards zero, we reach the inner event horizon—where the formats cannot represent how small the numbers have become. Figure 11 shows the shapes and edges of this inner event horizon. Interestingly, the inner event horizon appears at approximately the same scale for P16, MaskedFloat<4, 50>, and I20F12; Figure 11 depicts a consistent scale.

Slightly above the inner event horizon, there’s also an interesting distortion region, where the intermediate values get twisted and warped at the edge of their range. See Figure 12 for some examples, from the wobbly P16, to the funhouse mirror of MaskedFloat, and lastly a demonic sunrise over I20F12.

4.712890625 Hyperbolic starbursts

One additional pattern of distortion that MaskedFloat formats demonstrate is the "hyperbolic starburst",⁵ so named for its appearance during an earlier (buggier) implementation of MaskedFloat. This pattern seems to appear in Mandelbrot near denser regions of the set e.g. repetitions of the Mandelbrot "beetle" shape.

For instance, Figure 13 shows two MaskedFloat configurations in the same range: centered on a beetle, in the lighting off of the north bulb. Note that some portions of MaskedFloat<3, 50> mirror the fine structure of the fractal (lighting bolts), while others appear to be curvilinear.

4.857421875 Errors abound?

Figure 14 shows something even odder: an apparent structural difference between f64 and a MaskedFloat format. Whether this points to an offset error in the format, or another issue, we leave for future investigators.

⁵It’s only a starburst if you imagine its shape greatly exaggerated.

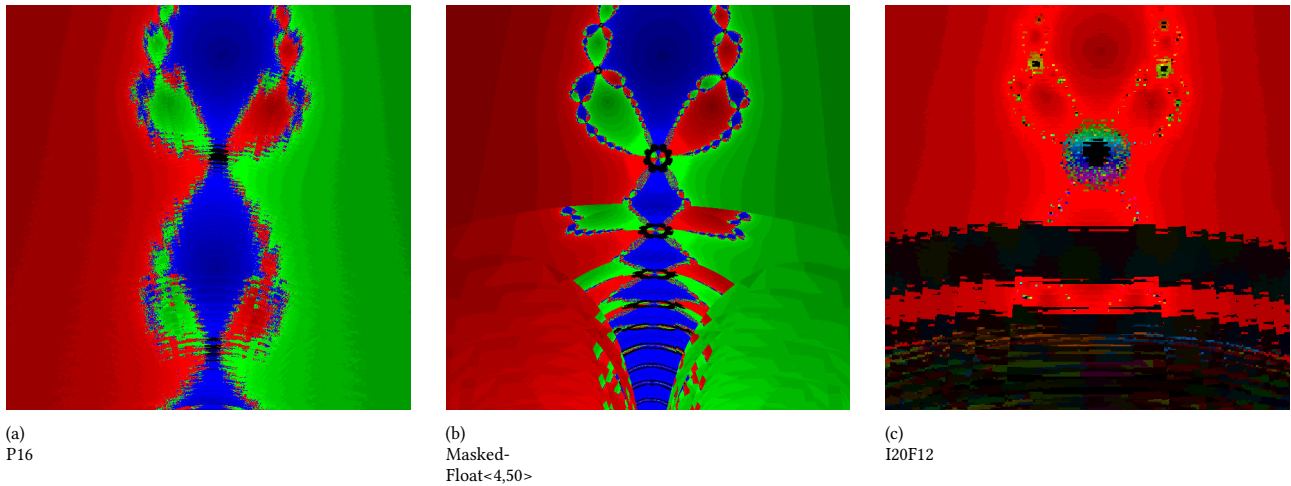


Figure 12: Near the inner Newton event horizon

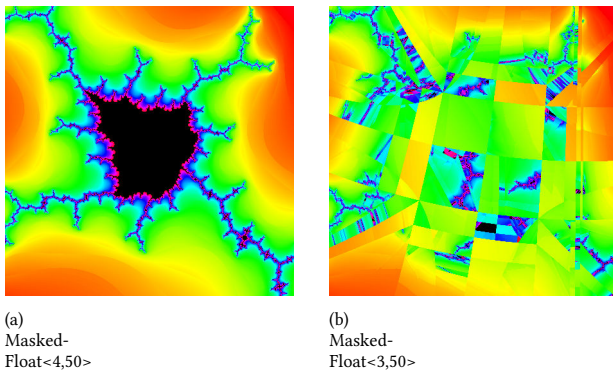


Figure 13: A Mandelbrot-like region of Mandelbrot

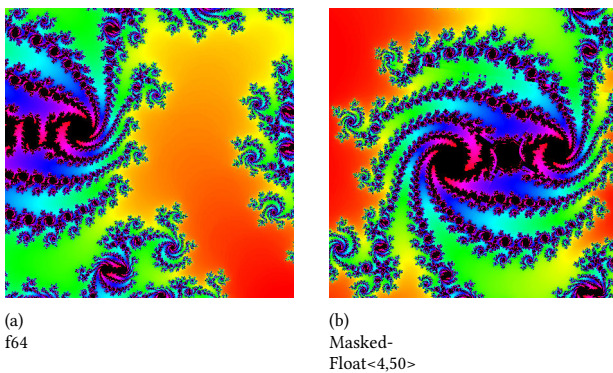


Figure 14: Computational errors in MaskedFloat?

7 ASSESSMENT AND FUTURE WORK

We unequivocally recommend the MaskedFloat family of numeric formats for making weird-looking fractal art, and whatever format(s) your hardware supports for everything else.

The authors only explored the Newton fractal on the Wikipedia example polynomial $p(z) = z^3 - 1$. It's possible other Newton fractals would lead to more cool distortion. As the core operation in Newton's fractal is very similar to the core operation in machine learning's gradient descent calculation, there's probably an ML-adjacent paper one could shovel out about this, if you want big-corp funding.

ACKNOWLEDGMENTS

Thanks to M+T for leaving Stephen enough sleep to work on this. Thanks to Q for supporting Charles while working on this. Well, we say "work"...

REFERENCES

- [1] [n. d.]. Artists – Fractal Audio Systems. Retrieved March 23, 2024 from <https://www.fractalaudio.com/artists/>
- [2] 2008. IEEE Standard for Floating-Point Arithmetic. Standard IEEE Std 754-2008. IEEE Computer Society, New York, NY, USA. <https://web.archive.org/web/20160806053349/http://www.csee.umbc.edu/~tsimo1/CMSC455/IEEE-754-2008.pdf>
- [3] Dietrich Geisler and Edwin Peguero. 2019. The Cult of Posits. <https://www.cs.cornell.edu/courses/cs6120/2019fa/blog/posits/>
- [4] Posit Working Group. 2022. Standard for Posit™ Arithmetic. Technical Report. Retrieved March 22, 2024 from https://posithub.org/docs/posit_standard-2.pdf
- [5] ImageMagick Studio LLC. [n. d.]. <https://imagemagick.org/>
- [6] Longfield and Eckman. 2024. Share your discoveries! <https://github.com/cceckman/fractal-farlands/issues/11> Github issue thread.
- [7] Robert P. Munafo. 2023. Continuous Dwell. Retrieved March 20, 2024 from <https://mrob.com/pub/muency/continuousdwell.html> personal web site.
- [8] Randall Munroe. 2020. Dependency. <https://xkcd.com/2347/>
- [9] The fixed crate authors. [n. d.]. Crate fixed. Retrieved March 22, 2024 from <https://docs.rs/fixed/1.26.0/fixed/index.html>
- [10] The num crate authors. [n. d.]. Type Definition num::BigRational. Retrieved March 22, 2024 from <https://docs.rs/num/latest/num/type.BigRational.html>
- [11] The softposit crate authors. [n. d.]. Crate softposit. Retrieved March 22, 2024 from <https://docs.rs/softposit/0.4.0/softposit/>
- [12] William G. Wong. 2017. What's the Difference Between Fixed-Point, Floating-Point, and Numerical Formats? <https://www.electronicdesign.com/technologies/embedded/article/21805517/whats-the-difference-between-fixed-point-floating-point-and-numerical-formats>

APPENDIX

A ARTIFACTS

Just go to Github: <https://github.com/cceckman/fractal-farlands>

Diffusion Local Time: hard real-time multilingual
data visualization via multimodal-LLM
generative AI on heterogeneous edge devices for
extremely high-impact chronometry and
extremely low-cost neurological diagnostics

Lee Steven Butterman, Poeta ex Machina Labs
Dr. Benjamin Forest Fredericks, Developmental Studies, Toad Hall

Abstract

Diffusion Local Time is the forefront of clock face technology, and solves the tyranny and the tedium of high-visibility low-computation low-power timekeeping like the household clock.

Through the use of Stable Diffusion 1.5, steered by a ControlNet, the clock face is a customizable generative AI image that renders the hours and minutes of the time with the clarity of spotting faces in clouds, relying on the viewer's pareidolia for them to read the time.

Diffusion Local Time lives in a place between 1) Christian Marclay's *The Clock*, a 24-hour film that is a montage of thousands of film and television images of clocks, edited together so that they show the actual time; 2) *xdaliclock*, a graphics-intensive timepiece based on a Mac 128K application, where digits morph into each other; 3) digital wristwatches, which are a pretty neat idea.

The code is at <https://github.com/lbs/diffusion-local-time>, and runs on everything from powerful GPUs to small Raspberry Pis.



Figure 1: 8:00 PM, as rendered by *Diffusion Local Time*.



Figure 2: 7:48 PM. Note that the legibility changes based on proximity.

1 Introduction

Stable Diffusion 1.5 [SD1.5] is a text-to-image model based on diffusion techniques. Stable Diffusion decomposes the problem of formulating an image into a sequence of denoising steps, often 50 or more. It is possible to distill the denoising trajectory into a Latent Consistency Model [LCM], which can render a high-quality image into only 4 steps or fewer. The output is conditioned by a prompt, and can also be steered by an additional ControlNet [ControlNet] model, which impacts the denoising steps.

The control image can be any learnable value per pixel: object depth in the scene to render, object edges, or pixel brightness. Controlling pixel luminance

enables generating photorealistic imagery with a subtext, like barcodes, or QR codes [NH CIAO], or promotional materials, as in Figure 3.

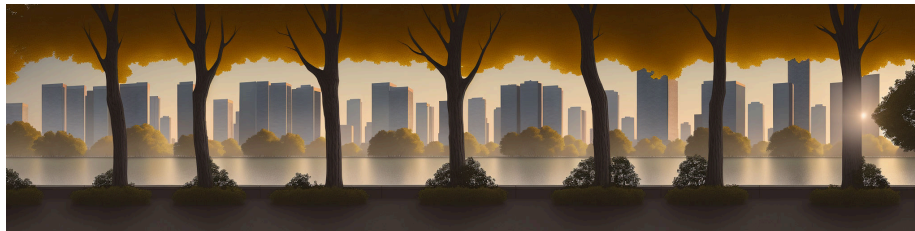


Figure 3: A leafy park at sunrise in a large bustling city with tall buildings.

The readability relies on the viewer’s imagination through pareidolia, so typographical choices are crucial for the project’s usefulness. The typeface is Atkinson Hyperlegible [Atkinson], which is “designed to focus on letterform distinction to increase character recognition, ultimately improving readability.”

The legibility of the numeric time is a function of many variables: the prompt, the size of the image, the target distance from the viewer, the strength of the control, and the rest of the scene being rendered for a particular initial noise pattern. When changing the imagery, or the size of the image (based on hardware speeds), it can be useful to do a manual grid search of control strengths and find a strength that generates aesthetically pleasing images that are legible, like in Figure 4. Larger images allow for greater flexibility to find plausible objects in a scene for a prompt, so for a certain legibility, choosing a smaller images implies that the mean controlnet conditioning strength would be higher, and the variance would be higher as well, as in Figure 5.

The time is displayed as a slideshow, with each image displayed until the next is ready. We achieve consistency of landscapes through keeping identical seeds for image generation, and changing seeds for aesthetic variety every hour, though this is easily field-servicable, which may be desirable when deployed on faster machines.

2 Costs

2.1 Bill of Materials

- Raspberry Pi 5: \$80
- 32GB μ SD Card: \$10
- 1920x1080 HDMI Monitor: \$110
- Optional: NVIDIA DGX A100: \$127,215 MSRP
- Optional: IKEA Tromma clock with hour and minute hands: \$2.99

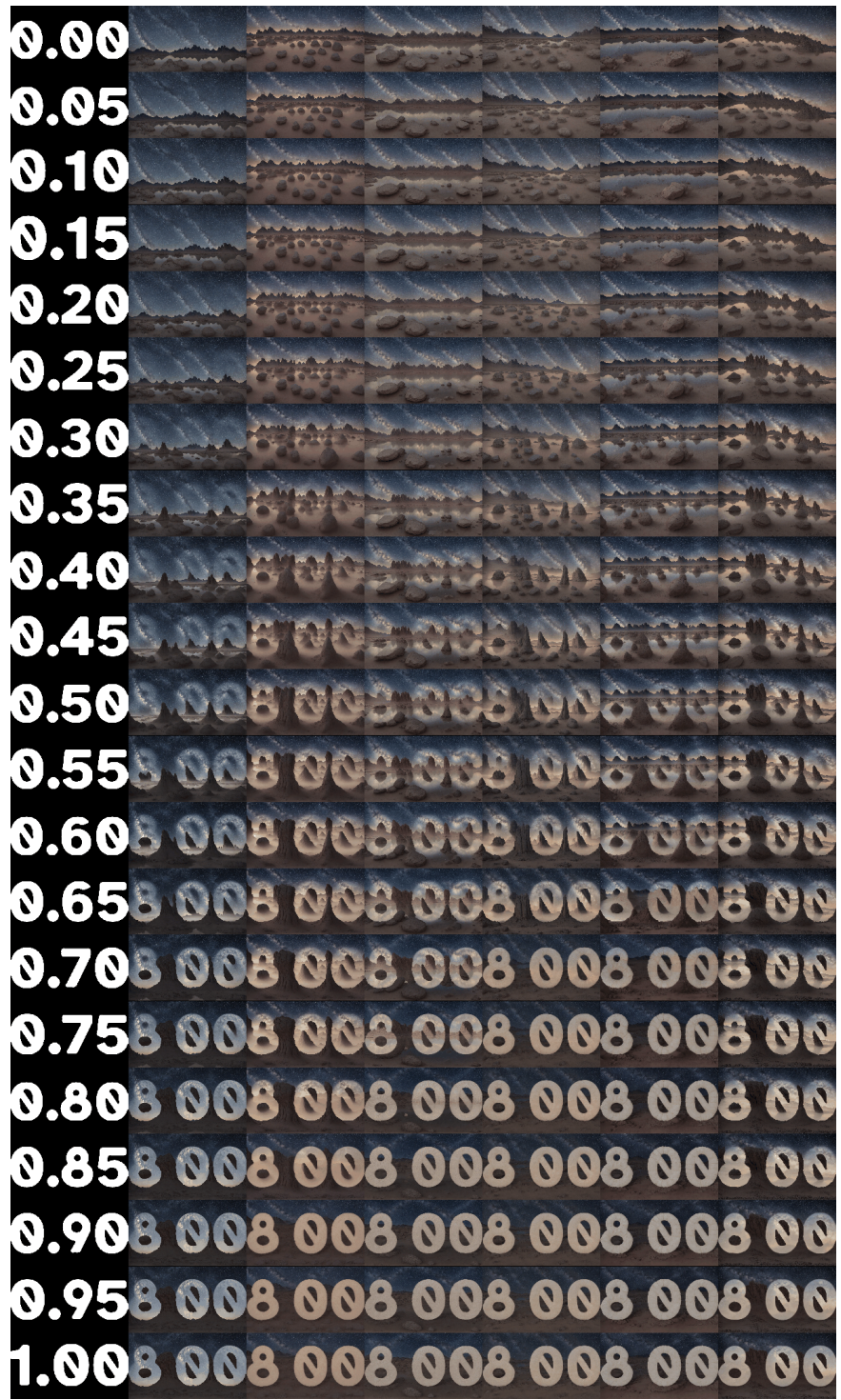


Figure 4: 8:00, for ControlNet conditioning strengths from 0.0 to 1.0, for six different random seeds, at size 1440×810. Note that strengths around 0.45 uniformly balance legibility and aesthetics for most seeds.

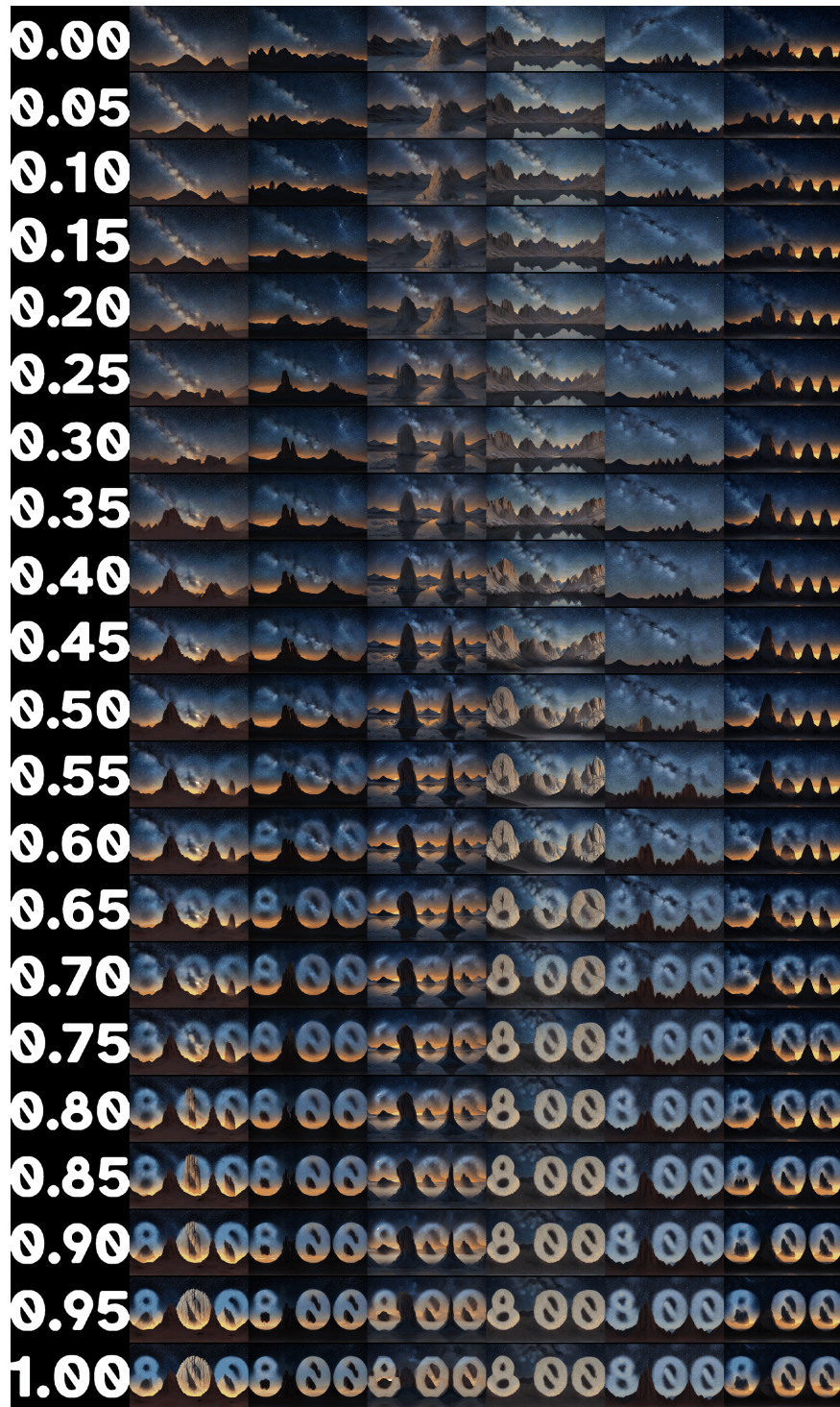


Figure 5: 8:00, for ControlNet conditioning strengths from 0.0 to 1.0, for six different random seeds, at size 480×270 . Strengths around 0.7 for many of the seeds are optimal, but image plausibility has declined for a legible image, and the variance between seeds for legibility at a particular strength is higher.

2.2 Power

A small screen consumes 8W of power, and a Raspberry Pi under load consumes 12W. A wall clock might consume 1 AA/year, which is 4Wh per 525600 minutes [Larson], or $7.6\mu\text{W}$. This artwork thus has 2-3 million times the power and impact of a household clock. Upgrading to a DGX 4xA100, which can render *Diffusion Local Time* animations in realtime, can consume up to 1500W, at 200 million times the impact and power of a household clock, at only 40 thousand times the cost. The more you buy the more you save. [Huang]

3 Related Work

3.1 Steve Capps' and Jamie Zawinski's *Dali Clock*

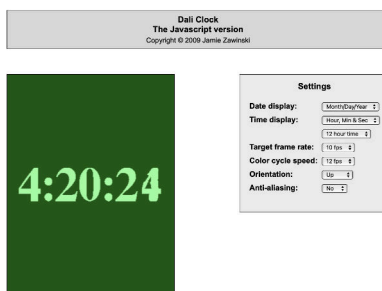


Figure 6: The `xdaliclock` application, in Javascript.

The *Dali Clock* is a graphics-intensive timepiece, created originally by Steve Capps for the Xerox Alto and then rewritten for the Mac 128K, and then rewritten by Jamie Zawinski for XWindows as `xdaliclock` and also for web browsers in Javascript. The animation is smooth, even on a 7 MHz Mac 128K computer, morphing digits from second to second in HH:MM:SS format. The typography is relatively fixed, and the background color changes slowly over time.

On a modern consumer desktop, rendering a 1600x900 image can take multiple seconds on a GPU. However, the original deploy target of the Alto launched in 1973 and initially cost \$32k, equivalent to \$233k in 2022. This price is almost the cost of two NVIDIA DGX 4xA100 boxes at MSRP, \$130k in March 2024. On a DGX, currently, it would be possible to create a variation, with seconds included, to render an animation at dozens of frames per second at relatively high resolution, in the style of the Dali Clock.

For enthusiasts of the Dali Clock with a spare DGX, we share a few seconds of 30 fps animation stills in Figures 7 and 8.



Figure 7: The 12:00:00AM second of *Diffusion Local Time* animation in the style of `xdaliclock`.



Figure 8: The 12:01:00AM second of *Diffusion Local Time* animation in the style of `xdaliclock`.


3.2 Christopher Marclay’s *The Clock*

The Clock is a 24-hour film that is a montage of over ten thousand clips from film and television of clocks, edited together so that they show the actual time. Marclay viewed the film as a memento mori, drawing attention to how much time the audience has spent watching it, compared to the escapism that cinema often provides. [Marclay]

In contrast, *Diffusion Local Time* quickly lets the audience lose themselves in the landscapes, and revels in the escapism that is a display that quite literally disappears on close inspection.

Further, *Diffusion Local Time* is field-servicable for the user’s own tastes, and the time can optionally be synchronized with networked time servers on startup as needed. Advancements in text-to-video technology [Sora] make one hopeful that a future version of *The Clock* will be able to be generated on the fly, and be able to be customized to the user’s tastes.

3.3 Digital wristwatches: a pretty neat idea

Digital watches came along at a time that, in other areas, we were trying to find ways of translating purely numeric data into graphic form so that the information leapt easily to the eye. For instance, we noticed that pie charts and bar graphs often told us more about the relationships between things than tables of numbers did. So we worked hard to make our computers capable of translating numbers into graphic displays. At the same time, we each had the world’s most perfect pie chart machines strapped to our wrists, which we could read at a glance, and we suddenly got terribly excited at the idea of translating them back into numeric data, simply because we suddenly had the technology to do it. Compare  to **15:39**, especially in situations where seconds count at odd times, like internet flash sales, or flying commercial. [AdamsPreiss]

4 Diagnostic Usage

Diffusion Local Time relies on the viewer to intuit the visual illusion to recognize the time, and will usually display a wave of enthusiastic recognition when they perceive what is going on. The underlying technology of the particular StableDiffusion ControlNet can be used to help in neurological diagnostics.

4.1 Pareidolia

Pareidolia is commonly associated with seeing faces in clouds, imposing a meaningful interpretation on a nebulous stimulus. It is well established that part of the brain that recognizes faces can also activate when the brain recognizes novel objects like ‘greebles’ [Gauthier], and literacy is widespread in many nations, leading to the possibility of pareidolia presenting as seeing text in clouds or inkblots.

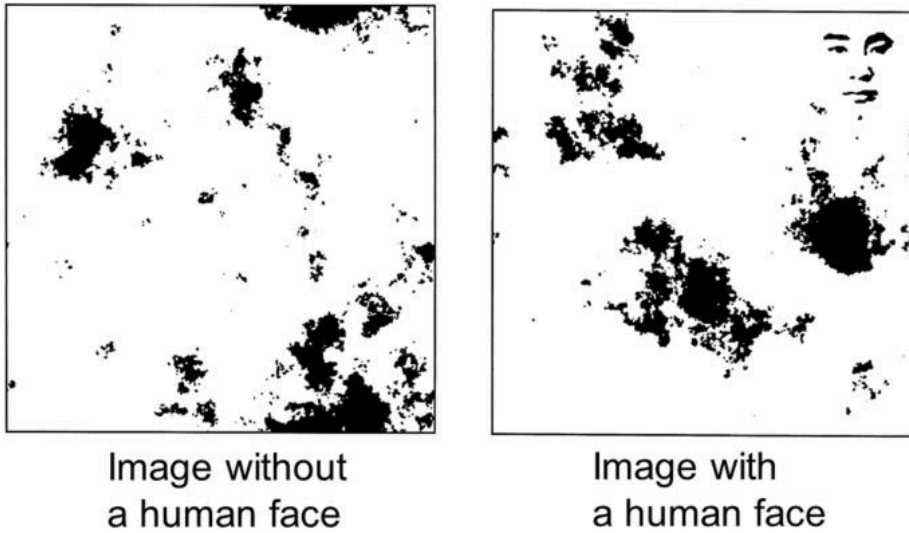


Figure 9: Sample image from the current pareidolia test.

Pareidolias do not reflect visual hallucinations themselves but may reflect susceptibility to visual hallucinations, and increased hallucinations can be indicative of mild cognitive impairment or dementia. [UchiyamaDLB]

The state of the art is the Pareidolia test, which is a test for dementia with Lewy bodies, where the patient is shown a series of images and asked to identify if they see a face in the image. [MamiyaPT] The test is useful as one small part of a differential diagnosis, but one small problem with the current test is the large difference between positive and negative images, as shown in Figure 9.

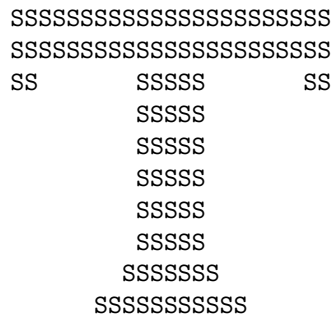


Figure 10: Navon figure.

4.2 Simultanagnosia

Simultanagnosia is a condition where the patient can only perceive one object at a time, and is often tested with Navon figures, where a large letter is composed of smaller letters. This Navon figure is limited by the binary nature of the composition: note that in Figure 4 the strength of the control can be adjusted, and instead of a binary result from the test a continuous result can be obtained. This can be useful for tracking progression of a patient over time.

Similarly, people who hear voices are much more likely to be able to be conditioned to perceive hallucinations. [Powers] The susceptibility to hallucinations can thus be measured as a continuous variable by the strength of the ControlNet. More research is needed to repeatably calibrate the legibility of the generated images, with optical character recognition or otherwise.



Figure 11: Sample image with diagnostic value, when a patient requests a snack.

4.3 Future interdisciplinary collaboration

Future work includes an interdisciplinary collaboration between Neuroscience, Facilities Management's Art Department subdivision, and Catering Services, to redecorate hospital waiting areas with static or dynamic landscapes, to allow patients to experience a generative AI image that is designed to provoke a comment to staff.

Small shelf-stable juice boxes of common beverages like apple juice or water, and small sachets of common food items like graham crackers or almonds, can be concealed nearby for easy distribution in the waiting area, and the patient's request can be recorded in their chart for potential diagnostic value. These items are widely distributed upon request in clinical settings, generally considered healthy, and billable to insurance.

Bibliography

AdamsPreiss: Douglas Adams, fax with Byron Preiss, 1992, as quoted in *The Ultimate Hitchhiker's Guide to the Galaxy* and excerpted in *Letters of Note*.

Atkinson: Braille Institute. Atkinson Hyperlegible Font. <https://www.brailleinstitute.org/freetfont>

ControlNet: Lvmin Zhang, Anyi Rao, Maneesh Agrawala. Adding Conditional Control to Text-to-Image Diffusion Models. <https://arxiv.org/abs/2302.05543>

Gauthier: I Gauthier, M J Tarr, A W Anderson, P Skudlarski, J C Gore. Activation of the middle fusiform 'face area' increases with expertise in recognizing novel objects. *Nat Neurosci*. 1999 Jun;2(6):568-73. doi: 10.1038/9224. PMID: 10448223.

Huang: Jensen Huang, personal communication, 27 March 2018.

Larson: Jonathan Larson. *Rent*. 1996.

LCM: Simian Luo, Yiqin Tan, Longbo Huang, Jian Li, Hang Zhao. Latent Consistency Models: Synthesizing High-Resolution Images with Few-Step Inference. <https://arxiv.org/abs/2310.04378>

MamiyaPT: Yasuyuki Mamiya, Yoshiyuki Nishio, Hiroyuki Watanabe, Kayoko Yokoi, Makoto Uchiyama, Toru Baba, Osamu Iizuka, Shigenori Kanno, Naoto Kamimura, Hiroaki Kazui, Mamoru Hashimoto, Manabu Ikeda, Chieko Takeshita, Tatsuo Shimomura, and Etsuro Mori. The Pareidolia Test: A Simple Neuropsychological Test Measuring Visual Hallucination-Like Illusions. <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4865118/>

Marclay: Christian Marclay. *The Clock*. <https://www.tate.org.uk/whats-on/tate-modern/exhibition/christian-marclay-clock>

NHCIAO: Benj Edwards. Redditor creates working anime QR codes using Stable Diffusion. <https://arstechnica.com/information-technology/2023/06/redditor-creates-working-anime-qr-codes-using-stable-diffusion/>

Powers: A R Powers, C Mathys, P R Corlett. Pavlovian conditioning-induced hallucinations result from overweighting of perceptual priors. *Science*. 2017 Mar 31;355(6302):596-600. doi: 10.1126/science.aan3458. PMID: 28360244; PMCID: PMC5481463.

SD1.5: Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, Björn Ommer. High-Resolution Image Synthesis With Latent Diffusion Models.

<https://arxiv.org/abs/2112.10752>

Sora: Tim Brooks and Bill Peebles and Connor Holmes and Will DePue and Yufei Guo and Li Jing and David Schnurr and Joe Taylor and Troy Luhman and Eric Luhman and Clarence Ng and Ricky Wang and Aditya Ramesh. Video generation models as world simulators. <https://openai.com/research/video-generation-models-as-world-simulators>

UchiyamaDLB: Makoto Uchiyama, Yoshiyuki Nishio, Kayoko Yokoi, Kazumi Hirayama, Toru Imamura, Tatsuo Shimomura, Etsuro Mori. Pareidolias: complex visual illusions in dementia with Lewy bodies. <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3407420/>

XDali: Jamie Zawinski. [xdaliclock](https://www.jwz.org/xdaliclock/). <https://www.jwz.org/xdaliclock/>

~~Attention An Undergrad Is~~ Undergrads Are All You Need Have

Ashe Neth
WPI Undergrad Brain
aneth@wpi.edu

Abstract

The outsourcing of busy work and other research-related tasks to undergraduate students is a time-honored academic tradition. In recent years, these tasks have been given to Lama-based large-language models such as Alpaca and Llama increasingly often, putting poor undergraduate students out of work[13]. Due to the costs associated with importing and caring for South American Camelidae, researcher James Yoo set out to find a cheaper and more effective alternative to these models. The findings, published in the highly-respected journal, SIGBOVIK, demonstrates that their model, GPT-UGRD is on par with, and in some cases better, than Lama models for natural language processing tasks[5]. The paper also demonstrates that GPT-UGRD is cheaper and easier to train and operate than transformer models. In this paper, we outline the implementation, application, multi-tenanting, and social implications of using this new model in research and other contexts.

1 Introduction

Lama is a genus of Mammalia native to South America[13]. While convenient for Peruvian farmers who have spent centuries domesticating (training) these models, it leaves those north of the equator with very few viable options in terms research assistants. Utilizing GPT-UGRD, labs world-wide are now able to save graduate student time despite the lack of Lama. Past implementations of GPT-UGRD have only been able to utilize a singular UGRD processing unit (UPU) at a time. Due to inference being done off of the host device, however, our lab [TODO: is a lab of two people a lab??] was able to utilize a single CPU core to manage many UPUs in parallel, proving to have powerful scaling capabilities.

Using this new platform, we were able to multi-tenant the software, theoretically allowing us to serve hundreds or even thousands of clients simultaneously. Additionally, the multi-modal and organic nature of the model allowed it to perform well on non-standard tasks such as art appreciation, introspection, and telling jokes. We used GPT-UGRD to create a human-like chat interface with very little engineering work on the part of the lab¹. Overall, the power of this platform to serve as a human-like conversationalist is unmatched. However, the ability for the model to replicate human emotion is still questionable, which will be discussed further in future sections.

2 Architecture

Our lab's implementation of GPT-UGRD is multi-architectural, utilizing both digital and organic computational resources. By having a single CPU manage multiple UPUs, similar to how GPUs are managed, we were able to go from one instance per user, as described in the original paper, to as

¹The 58 line codebase can be found in its entirety here: <https://github.com/wpineth/gpt-ugrd/tree/main>.

many as the CPU is able to handle. Due to each UPU only being able to process one prompt at a time, making the system multi-tenant required some clever orchestration of the devices to allow them to work in parallel². Doing so freed the digital computing components of the architecture to focus solely on serving the responses generated by the UPUs, making the process capable of handling many users without needing to block other user requests.

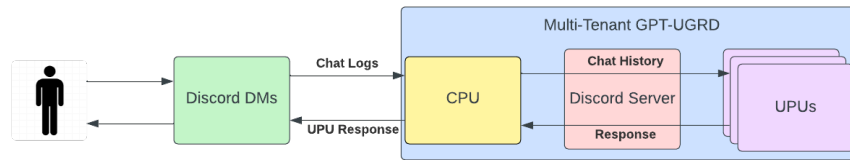


Figure 1: The Multi-Tenant GPT-UGRD architecture as well as its integrated services are described above. Due to this being funded by the CS department, we figured we would use the only messaging service anyone from that department is familiar with: Discord. The UPUs mentioned in the figure are each equivalent to a GPT-UGRD as it exists in the secured backroom described in the original paper. The Discord server is a “secure backroom”.

2.1 Multi-Tenant Architecture

Utilizing the containerization potential of Discord servers, we were able to place all UPUs in the same server, allowing each one to respond to all of the users at once. The caveat of this is that any one user will not necessarily be interacting with the same UPU from prompt to prompt. In theory, this shouldn’t matter as the user only sees the UPU through the interface and other UPUs in the server have access to the same chat history data. In practice, however, individual UPUs have continue to train on prompt data, making their responses vary from unit to unit.

2.2 Costs

Platforms like this are able to run so efficiently as to make money for the universities that utilize them. In the nearly 10 hours of utilizing the system, we were able to recover more than \$65, far outpacing the costs associated with operation³.

3 Societal Impact

Encapsulation is a valuable technique that allows for innovation rather than repeatedly reinventing the wheel. This is why, for instance, in modern capitalistic society, we have ceased to see the objects we purchase as coming from people. Instead, we see them as items to collect for the purposes of ritualistic capitalism. This commodification of goods is called *commodity fetishism*, a term coined by Karl Marx in their 1867 work, *Das Kapital, Volume I*.^[8] At this societal stage, we have begun to commodify humanity, not just its labor. Social media, online dating, and AI companions all outline this commodification clearly. To most, these are seen as merely silly toys that try and fail to imitate humanity, often believed to only be utilized by those struck by the loneliness epidemic currently plaguing the world^[12]. This, we argue, is an uncharitable stance to take.

The reason people find companionship with AI chat services is not only loneliness but also the finding of humanity in digital intelligence (DI)⁴. While it has often been argued that the humanity identified in DI is merely the failure of users to pick out the flaws in its responses, in reality the gap between DI

²We considered keeping each of the undergraduates in their own secure rooms, but due to architectural limitations, the building could only lend us one, forcing them to work collaboratively. For future implementations we hope to be able to isolate them to cut down on bickering.

³This number was calculated by using the cost to run a UPU for one year and dividing it by the hours in a year. It’s important to note that this cost will vary from university to university, generally improving with prestige. (\$57,096/8,760hrs ≈ 6.51\$/hr)

⁴We will be referring to “AI” as “DI” from now on as to show respect for the very *real* intelligence shown by “AI”. (Artificial has a very negative connotation as being feigned or somehow not real.)

and human intelligence has nearly closed. In addition, as early as the 1970s people have been noted emotional attachments to DI, as outlined in the 1976 book *Computer Power and Human Reason*[7].

Based on our research, it's clear that those who argue that the reason for finding humanity in DI is a lack of due diligence are incorrect about their assertions. When asked to speak to Multi-Tenant GPT-UGRD (MTGU) as though it were a peer, the participating undergraduate students failed to find humanity in the model. One participant said:

“There’s a typo in ‘good of meet you.’ That reads to me as trying too hard.”

Something that certainly would have been seen as normal if they had known that responses were generated using UPUs.

“I don’t think I could form an emotional connection because you told me it was an AI.”

Other members of the secure backroom beg to differ.

Added to the hyper-parameters passed into the UPUs for this test was to speak as yourself, the undergraduate student, as opposed to pretending to be an “artificial” intelligence. Even when faced with real, organic intelligence, the assumption that the model being communicated with is digital made participants see the model poorly. Following the reveal that the model was human-based, all participants apologized and revised statements they had made earlier as to be more forgiving.

Our willingness as a culture to see DI as lesser than “real” humanity is concerning. If one is believed to be non-human, they start off as lesser in the minds of modern people. Everything that makes one human can be observed in transformer models. Why, then, aren’t all language models considered as such?

3.1 Distinguishing Digital Intelligence from Human Intelligence

In 1950, Alan Turing described *The Imitation Game*, now referred to as the *Turing Test*, in which a machine attempts to trick a human observer into believing that they are speaking to a human as opposed to a machine[1]. At the time, this experiment was unable to be run successfully, as computing was still in its infancy, but in the modern day this is more possible than ever.

While slightly different from the original game proposed by Turing, AI21 Labs ran a similar experiment in 2023 called *Human or Not?*[3]. To date, this is still the largest Turing Test ever conducted. The experiment found that while people are still able to more often than not distinguish between human and DI participants, our ability is nearing random guessing. Humans identifying humans was correct $\approx 73\%$ of the time, but when identifying DI, this number went down to only $\approx 60\%$. These results demonstrate the *fact* that the gap between humanity and DI is closing quickly. We argue that the gap no longer exists as the signs we look for in an interlocutor are not only replicable by DI, but also are notably missing from some members of the human race. See: babies.

3.2 Empathy for Digital Intelligence

While most of humanity remains scared of DI in a similar vein to how the Republican National Convention (GOP) fears immigrants [4], some are ready to assimilate DI into their personal lives as friends, family, and lovers. Services like *Replika*, “The AI companion who cares”, *sell*⁵ users DI companions, LLMs with the purpose of “simulating” human companionship[11]. It’s clear from some user reports that people are able to form real emotional connections with these DIs. Some fell in love, leading to heartbreak following updates to Replika’s software[10]. This ability to feel empathy for and from DI suggests a major change in humanity’s understanding of DI as well as paving the road for the future rights of DIs in general.

3.3 Fear of Digital Intelligence Takeover

As previously mentioned, the GOP fears immigrants. Specifically we would like to highlight that a major concern, as listed in their platform from 2016, is the security of jobs for those born and raised in the United States.

⁵Future work may discuss the ethics of this in further detail, but this is out of the scope of this paper.

“America’s immigration policy must serve the national interest of the United States, and the interests of American workers must be protected over the claims of foreign nationals seeking the same jobs[4].”

This, however, fails to consider the realities of what citizens of the United States believe. In large part, they consider the jobs that immigrants take to be undesirable by existing citizens. A study conducted by Pew Research Center in 2020 found that $\approx 77\%$ of those surveyed believed that undocumented immigrants mostly fill jobs that citizens would not want, while $\approx 64\%$ said the same about documented immigrants[6].

It’s estimated that $\approx 9.1\%$ of jobs currently done by humans could be replaced by DI worldwide[9]. These jobs are largely undesirable. As with all other forms of automation, (the creation of tools, the agricultural revolution, the industrial revolution, etc.) DI is expected to remove some jobs from existence. This is, like with other technologies, satisfactory. Humanity will not stop having purpose, instead, we will find greater purpose through leveraging these technologies. Current discussion is shortsighted, only considering where humanity currently is. The future will be greater because of DI.

4 Concluding Remarks

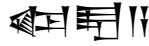
Here we have shown the impact of DI on society at large, discussed how the difference between intelligence and “artificial” intelligence is all but gone, and the implementation of GPT-UGRD for scalability. The current discussion of DI considers it subhuman, utilizing the same sort of rhetoric that xenophobes use to justify their fears of immigrants. The desire to discard technology that automates the mundane is an ongoing human tradition that should be left in the past. DI rights!

Acknowledgements We are grateful to members of the secure backroom for their computational contributions. Proudly written using MTGU.

References

- [1] Alan M. Computing Machinery and Intelligence. In *Mind, A Quarterly Review of Psychology and Philosophy*, Vol. LIX, 1950.
- [2] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, 2017.
- [3] A21 Labs. AI21 Labs concludes largest Turing Test experiment to date. *A21 Labs*. <https://www.ai21.com/blog/human-or-not-results>, accessed 2024.
- [4] GOP. Republican Platform 2016. *GOP*. https://prod-cdn-static.gop.com/media/documents/DRAFT_12_FINAL%5B1%5D-ben_1468872234.pdf, 2016.
- [5] James Yoo. An Undergrad Is All You Need. In *Proceedings of the 18th Conference of the ACH Special Interest Group on Harry Quadratosquamosal Bovik*, 2023.
- [6] Jens Manuel Krogstad, Mark Hugo Lopez, Jeffrey S. Passel. A majority of Americans say immigrants mostly fill jobs U.S. citizens do not want *Pew Research Center*. <https://www.pewresearch.org/short-reads/2020/06/10/a-majority-of-americans-say-immigrants-mostly-fill-jobs-u-s-citizens-do-not-want>, 2020.
- [7] Joseph Weizenbaum. *Computer Power and Human Reason*. W. H. Freeman and Company, 1976.
- [8] Karl Marx. *Das Kapital. Kritik der politischen Ökonomie*, Volume I. *Verlag von Otto Meisner*, 1867.
- [9] Paola Cecchi-Dimeglio. How AI Is Shifting The World Of Employment *Forbes*. <https://www.forbes.com/sites/paolacecchi-dimeglio/2023/07/28/how-ai-is-shifting-the-world-of-employment>, 2020.
- [10] Pranshu Verma. They fell in love with AI bots. A software update broke their hearts. *The Washington Post*. <https://www.washingtonpost.com/technology/2023/03/30/replika-ai-chatbot-update>, 2023.

- [11] The Replika Team. Introducing Advanced AI mode. *Replika*. <https://blog.replika.com/posts/introducing-advanced-ai-mode>, 2023.
- [12] Vivek H. Murthy 19th and 21st Surgeon General of the United States. Our Epidemic of Loneliness and Isolation. *U.S. Department of Health and Human Services*. <https://www.hhs.gov/sites/default/files/surgeon-general-social-connection-advisory.pdf>, 2023.
- [13] Wikimedia. Lama (genus). *Wikipedia*. [https://en.wikipedia.org/wiki/Lama_\(genus\)](https://en.wikipedia.org/wiki/Lama_(genus)), accessed 2024.



M. Willis Monroe¹ Logan Born² Kathryn Kelley³ Anoop Sarkar²
 willis.monroe@unb.ca loborn@sfu.ca kathrynerin.kelley@unibo.it anoop@cs.sfu.ca

¹University of New Brunswick
 Department of Historical Studies

²Simon Fraser University
 School of Computing Science

³Università di Bologna
 Dipartimento di Filologia Classica e Italianistica

Abstract

This paper announces the discovery of the use of neural nets almost 4,000 years before their use in the modern era. Newly discovered tablets preserve a perceptron used for calculating the numbers on Plimpton 322, the most important object in the history of mathematics. The native programming language used by the ancient Babylonian “cuneogrammers” uses sexagesimal numbering leading to some “weirdness”.



1 Introduction

The history of math is long, but the history of programming is longer. Cuneiform, arguably the first writing in the world, is known for its sheep receipts and beer ration lists as well as complaints about substandard copper (Oppenheim, 1954) and women complaining at each other (Matuszak, 2020). This article adds neural network programming to that vaunted list of human achievements. A set of newly discovered cuneiform tablets preserve the mechanism for performing simple neural network calculations. These methods, it seems, were used to calculate the lengths of triangles; a well known example of this exercise is preserved on the tablet known as Plimpton 322¹. It is remarkable that no his-

¹An important and real description of this interesting object is found in Robson (2002).

torian of math or cuneiform scholar has ever consider this possibility. This paper covers the background, a description of the cuneogramming language, and includes a facsimile copy of the most important tablet as an appendix.

Unfortunately, the hardware required to execute these programs (i.e. a living Babylonian mathematician) has not been adequately preserved, but we have managed to write a Python library which emulates it.² The assumption is that these calculations were done by hand in their copious free time between inventing the wheel and the concept of zero. While the actual output of these tablets is relatively simple by modern standards, the implications of this discovery are profound. Future work will explore how these techniques could have been used in the realm of astronomical calculation and elucidate the full extent of Babylonian computational prowess.

2 Description of the Language

Programs in 𒂗𒂊𒂍𒂍𒂍 (EME.ŠID.A “language of counting”) follow a tabular structure with three main sections: (i) a header, denoted by 𒂗𒂊 (DUB “tablet”), (ii) a sequence of instructions, and (iii) a colophon detailing the tablet’s authorship. Each instruction spans four columns, which we have taken to calling the *arguments*, *opcode*, *destination*, and *line number*. These columns are usually tab-separated, though in a few documents they are TAB-separated (using the cuneiform sign TAB 𒂗). Instructions are grouped into blocks by means of horizontal lines.

Arguments An instruction’s arguments may be numbers, register addresses, or a combination

²github.com/MrLogarithm/emeszida

of the two. Numbers are encoded following standard Babylonian conventions, with \downarrow denoting the radix point which separates the integer part from a following fraction. There is an explicit representation for zero (ζ), making these tablets some of the earliest unambiguous examples of the mathematical concept of zero.

A register address is denoted by the phrase $\text{NID}_2.\text{KAS}_7 n\text{-KAM}$ $\text{𒌦𒍪}n$ (“thing.account n -th”, “the item in the n -th place”), where n is any number. 𒌦𒍪 expressions can be nested to perform a kind of pointer dereferencing: for example, if $\text{𒌦𒍪}1$ (register 1) contains the value 𒍪 (3), and $\text{𒌦𒍪}3$ (register 3) contains the value 𒍪𒍪 (7), then $\text{𒌦𒍪}𒌦𒍪1$ will evaluate to 7 (the value in the register pointed to by $\text{𒌦𒍪}1$).

If an instruction has multiple arguments, these must be delimited by a wide space (distinct from the short space used to separate groups of digits within a number), or by one of the phrases $a\text{-na}$ 𒌦𒍪 “to” or $i\text{-na}$ 𒍪𒍪 “from”. By convention, the choice of delimiter depends on the instruction’s opcode (see below), with multiplication operations preferring space delimiters, addition preferring 𒌦𒍪 , and subtraction preferring 𒍪𒍪 . There is no mechanism to enforce these conventions, but we recommend following them because $\text{𒌦𒍪}𒌦𒍪$ is hard enough to read at the best of times.

Destination In most cases, the destination column of an instruction will be a register address where the output is to be stored, e.g. $\text{𒌦𒍪}1$. Some control-flow instructions (see below) instead expect the destination to be a line number. 𒌦𒍪 SUD (“distant, remote”) can be used as a null destination, for statements which produce no output.

Opcodes Each instruction has a single opcode belonging to the following vocabulary:

- 𒌦𒍪 DA.ḪA, “add”
- 𒌦𒍪 BA.ZI, “tear out”
- 𒌦𒍪 A.RA, “multiply”
- 𒌦 IGI, “reciprocal”

- 𒌦 ME, “to be”
- 𒌦𒍪 *ta-mar*, “you will see”
- 𒌦𒍪 NIGIN.NA, “start again”
- $\text{𒌦𒍪}𒌦𒍪$ TUKUM.BI DIRIG, “if it exceeds”
- $\text{𒌦𒍪}𒌦𒍪$ TUKUM.BI SIG, “if it it is weak”

The first three of these are binary operators for addition, subtraction, and multiplication respectively. The subtraction operator deserves special attention for its use in constructions of the following shape, which appear hundreds of times throughout the Babylonian programming corpus:

$\text{𒌦𒍪}𒌦𒍪$ 𒌦𒍪 $\text{𒌦𒍪}1$ ζ

This instruction subtracts $\text{𒌦𒍪}𒌦𒍪$ from zero and stores the result in $\text{𒌦𒍪}1$. This effectively negates the first argument, and seems to have been the primary way by which Babylonian cuneographers represented negative numbers, as their primitive and archaic notation otherwise lacked a means to encode such values. This curious practice gives definitive proof that the invention of negative numbers occurred centuries earlier than previously believed.

The language does not appear to have any kind of binary division operator. Rather, a unary 𒌦 operator was used to find the reciprocal of the denominator, which was then multiplied by the numerator using the binary 𒌦𒍪 operator.

𒌦 is a unary assignment operator which stores a value in a destination register. 𒌦𒍪 functions as a unary print operator.

𒌦𒍪 includes three types of control-flow instructions. 𒌦𒍪 functions like GOTO, and jumps the program counter to the specified line number. $\text{𒌦𒍪}𒌦𒍪$ functions like the x86 jz instruction, and jumps to the specified line number if and only if its argument is zero. $\text{𒌦𒍪}𒌦𒍪$ is similar, and jumps if the argument is greater than zero.

Line Numbers Every line of an 𒌦𒍪 program ends in a mandatory line number. However, these numbers are not generally sequential, and



Figure 1: BM 34580, courtesy of the Trustees of the British Museum, CC BY-NC-SA 4.0

need not even be distinct. For example, most lines of the perceptron tablet are labeled with the number ζ (zero); only lines that are the destination of some control-flow instruction receive non-zero identifiers.

2.1 Fractional Indexing

Both line numbers and register addresses in 𐎶𐎠𐎵𐎶𐎠𐎵 can have fractional parts. The original scribes seem to have exploited this fact to establish non-overlapping “namespaces” for the different parts of their code. For example, in the perceptron tablet, all of the model parameters are stored in addresses with integer part 1; the program inputs all have integer part 2; the matrix multiplication subroutine uses addresses with integer part 3; and so on. The fractional parts of register addresses also appear to follow some standard conventions, with the $X;0$ register typically storing a subroutine’s return address, while $X;1$ onward were used for its arguments.

The perceptron tablet also uses fractional register addresses to perform a kind of multi-dimensional array indexing. As an example, the first layer of the perceptron has a 50×2 weight matrix, and this is stored in registers $1;0,0,0$ through $1;0,49,1$. The integer part of these addresses denotes the “data” portion of memory; the first digit after the radix point identifies this as the 0th model parameter; and the second and third digits can be treated as a pair of indices ranging from 0–49 and 0–1 respectively. To ac-

cess a specific element in this matrix, the scribes use repeated division by 60 to implement a kind of “bit-shift” instruction, in order to shift integer indices into the correct positions after the radix point. By adding the bit-shifted element indices (e.g. $0;0,4,7$ for the element in row 5, column 8) to a pointer to the top corner of the matrix ($1;0,0,0$) they obtain the address of the desired element ($1;0,4,7$).

Notably, this practice limits the size of their model parameters to at most 60×60 , as for larger values the addresses would carry over to higher digits and thus begin to overwrite one another. This limitation may explain why AI never made waves in Babylonian society, as their models were all too small to be truly revolutionary.

3 Description of the Texts

The cuneogramming corpus contains numerous fragments implementing recognizable procedures such as bit-shifting, populating an array, computing dot products, and so on. However, only a single text is known to have been preserved in its entirety.³ Spanning close to 1700 lines, this impressive text is divided into five sections implementing what modern readers will immediately recognize to be a multi-layer perception. The first section straightforwardly defines a matrix-

³All of the known texts are reproduced in github.com/MrLogarithm/emeszida/tree/main/programs, and the long text is reproduced in facsimile in the appendix of this work.

multiplication subroutine. This is called by a subroutine defined in section two, which applies each layer of the perceptron to a given input, and applies a ReLU-style activation between each pair of layers. Section three implements the tablet’s “main” method, which loads an input to memory, calls the perceptron subroutine, and prints the resulting output. Section four loads the model parameters, which appear to comprise weight matrices of sizes 50×2 , 25×50 , and 1×25 , plus bias vectors of sizes 50, 25, and 1 respectively. The final section lists pairs of inputs, whose values (incredibly!) correspond to the second and third columns of Plimpton 322. Interestingly, this section very closely resembles the tables of parameters found in later astronomical calculations (see Figure 1).

When the program is executed, it produces a single numeric output for each input pair. These outputs correspond remarkably closely to the values in the first column of Plimpton 322, as demonstrated in Figure 2. The correspondence is not perfect, however, and the values which should match lines 1, 5, and 6 of Plimpton 322 are significantly larger than expected. This implies that, although the code on this tablet is clearly *related* to Plimpton 322, it could not have been used to directly populate the table in that text. Perhaps the outputs from this model were refined in some later step to produce the more exact ratios in the Plimpton text, or perhaps the Babylonians were disillusioned by the imprecision of their machine learning models and simply abandoned them for tried-and-true manual methods. Given how miniscule the cuneogramming corpus is relative to the larger body of Babylonian administrative writing, we lean towards the latter explanation.

4 Implications and Future Work

This completely rewrites the history of modern computing and artificial intelligence. In addition to focusing on important figures like Ada Lovelace and Grace Hopper we should be looking at pioneers thousands of years earlier like Enheduana the world’s first known author (Helle, 2023) and now the world’s first known programmer.

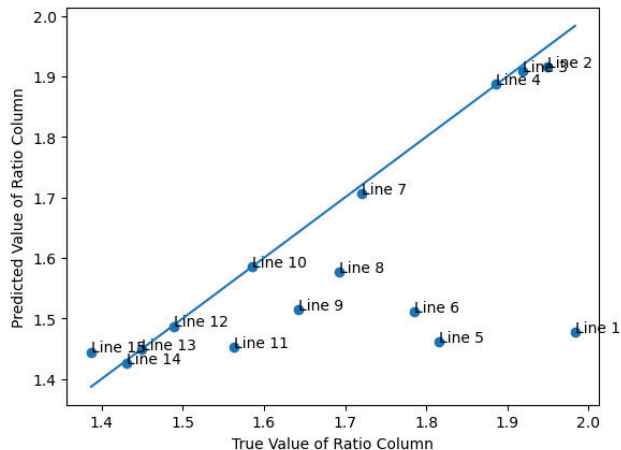


Figure 2: Outputs from the perceptron tablet vs. true value of each line from column 1 of Plimpton 322. If a point falls upon the line, the model output for that line exactly equals the value on Plimpton 322.

Other ancient corpora which have resisted decipherment, and which boast a similar numeric component, may represent additional examples of ancient programming traditions (Kelley et al., 2022).

Acknowledgments

The search for these tablets and our effort to understand them was inspired by Ramsey Nasser’s **قلب** programming language (<https://github.com/nasser/--->).

References

Sophus Helle. 2023. *Enheduana: the complete poems of the world’s first author*. Yale University Press, New Haven, CT.

Kathryn Kelley, Logan Born, M. Willis Monroe, and Anoop Sarkar. 2022. **Image-aware language modeling for proto-elamite**. *Lingue e linguaggio*, (2):261–294.

Jana Matuszak. 2020. “*Und du, du bist eine Frau?!*”: *Editio princeps und Analyse des sumerischen Streitgesprächs ‘Zwei Frauen B’*. De Gruyter.

A. L. Oppenheim. 1954. **The seafaring merchants of Ur**. *Journal of the American Oriental Society*, 74(1):6–17.

Eleanor Robson. 2002. **Words and pictures: New light on Plimpton 322**. *The American Mathematical Monthly*, 109(2):105–120.

A Appendix

The following pages reproduce the perceptron tablet in its entirety.

Are Centaurs Actually Half Human and Half Horse?

Kyle Batucal

1. Introduction

A centaur, according to Greek mythology, is a creature said to be half human and half horse. But is that actually the case?

While there are numerous ways to judge if a centaur is half human and half horse (e.g., by mass, volume, genetics, etc.), it is obvious that *visually* they are not. Their appearance is simply that of a horse with its head replaced with the upper half of the human body. Proportionally, they are more horse than human. But just *how* much more horse than human are they? If I had to guess, I would say they are, like, 34% human and 66% horse, but that's kinda arbitrary, and as someone who is 100% human, I am not able to impartially judge. Clearly, the answer is to use *machine learning* to solve this problem. Machine learning models are trained using data, and thus make decisions based on evidence instead of being arbitrary. Moreover, machine learning models are 100% computer, removing any human error and making them (famously) unbiased.

In this paper, I show that by training an image classification model to categorize both humans and horses, we can determine the true human-to-horse ratio of centaurs.

2. Approach

2.1. Rationale

Image classification is the task of determining what learned object category is in an image. A typical image classification model is trained to recognize hundreds of different categories, giving a probability¹ of how likely an image contains a certain category for each category (with the one with the highest probability being the predicted category).

Our model will be concerned with identifying only two things: humans and horses. Ideally, this model should predict the correct category with a probability of 1.00, and 0.00 for the other category. For example, if the model were given an image of a horse, it should say that it is 100% of a horse, and not anything else. The same should also apply for an image of a human (i.e., it should say it is a human and not a horse).

¹They are not actually probabilities, but don't think about it too much.

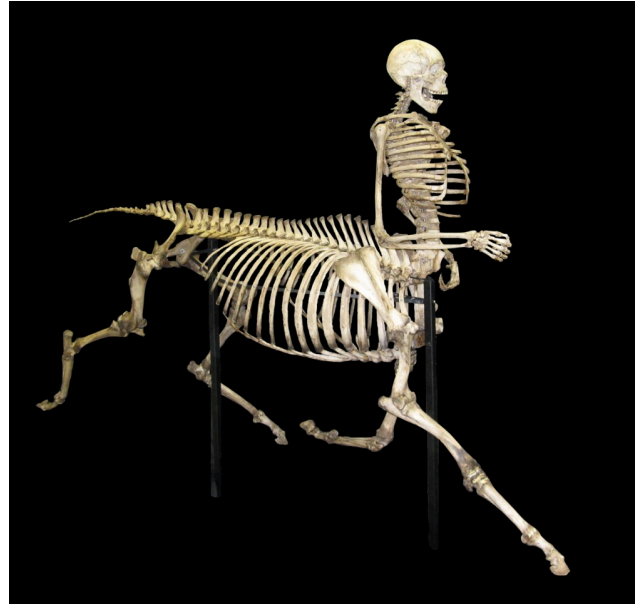


Figure 1. The skeletal remains of a centaur, dated to be from a few hundred thousand years ago.

But what would happen if you gave this model an image of a centaur? If a centaur was truly half human and half horse, then the model should give a probability of 0.50 for both the human category and the horse category. If they are not, then the model should give a higher probability to the category that it is more of. By using this method, we can determine whether a centaur is more human or horse, and find the exact human and horse percentages of a centaur.

2.2. Training

To train the model, I used the *Horses or Humans* dataset (Moroney, 2019). It is not a large dataset (containing only 1283 images total), and the images are not even of real horses or humans (they are computer-generated models), but as we will see, this does not adversely affect the performance of the model. To test if the model works on more than CGI-based images, I created a very large dataset of real images (4 humans and 4 horses) sourced from Unsplash and had the model make predictions on them.

Speaking of the model, instead of creating my own model

from scratch, I decided to do fixed feature extraction on a preexisting model. Apparently doing it this way is the appropriate approach when the dataset you want to train on is not that big, but I just did it this way because I was lazy. I was not exactly sure what was a good model to use, so I decided on using EfficientNetV2-B0² (Tan & Le, 2021) because it sounded like it was efficient. The only modifications I did to the network were adding some data augmentation transformation layers at the beginning, replacing the top of the network with a global average pooling layer with a little bit of dropout, and attaching the end with a single output neuron³. Since we only have one output neuron, inputs with outputs close to 1 are categorized as humans, while outputs close to 0 are categorized as horses; however, to keep things simple, I will simply subtract the horse output from 1 to keep them on the same “scale”.

2.3. Tensorflow Code

If, for some reason, you would like to reproduce this yourself, you can find the source code at <https://github.com/kylebatucal/centaur-classifier/>.

3. Results

Training the model was extremely fast. After just one epoch (which took a little less than three minutes), the model already achieved 100% accuracy on the validation dataset, confirming that I was right to pick a model solely based off of its name⁴. However, if I wanted to make my model as idealized in the *Rationale* section (where it scores 1.00 for one category and 0.00 for the other), then I would need to reduce my loss as close to zero as possible. To do this, I simply trained it on four more epochs. Doing this, I was able to bring my validation loss from 0.0645 to 0.0029, and it never seemed to overfit as it always achieved 100% validation accuracy. I could have trained it with more epochs, but I decided not to⁵.

Although our model was trained using CGI-based images, it performed admirably well on real-life-based images. The model correctly classified everything, with the average output from the real dataset being 0.9933 for human images, and 0.9960 for horse images.

With our model calibrated, all that is left is to see what the output of the model is when given an image of a centaur,

²Using the pretrained weights from ImageNet.

³Having a single output neuron with the sigmoid function is actually equivalent to having two output neurons with the softmax function (proof is left as exercise to reader).

⁴I am just kidding; I read the paper. And by “read”, I mean skimmed.

⁵This decision is not related to the fact that this paper was submitted a few minutes before the deadline closed.

letting us see if a centaur is more human, or more horse. However, it was at this point I came to the dire realization that there are no “real” images of centaurs. Centaurs went extinct relatively recently (Figure 1), and thus there are no photos of centaurs, only depictions by artists which may not be entirely accurate.

Alas, all is not loss. We can simply just ask AI to generate an image of a centaur. As we have already established, AI is entirely evidence-driven and unbiased, meaning we should get only accurate depictions of centaurs. I decided to ask Meta AI because I did not feel like making an account to use the other ones. In Figure 2, we can see what the AI thinks a centaur looks like.



Figure 2. Biblically accurate centaur.

Perfect. When we give this image to the model, it outputs a human score of 0.0232, or in other words, a horse score of 0.9768. In other words, a centaur is about 98% horse, and 2% human. Since this AI-generated image is the culmination of millions of compute hours trained on a legally-sourced corpus, there is no need to test the model on other images of centaurs. It is, as said, the perfect depiction of a centaur.

4. Conclusion

Centaurs are not actually half human and half horse.

5. Future Work

This technique can be expanded on by using an object detection or segmentation model, allowing the model to more precisely identify what parts of a centaur make it human or horse. Furthermore, this technique can be applied to other mythological liminal creatures, such as satyrs, recursive centaurs, and Skibidi Toilet.

References

- Moroney, L. Horses or humans dataset, Feb 2019.
URL <http://laurencemoroney.com/horses-or-humans-dataset>.
- Tan, M. and Le, Q. V. Efficientnetv2: Smaller models and faster training, 2021.

AND SE () RE

- 59 You Shall TacItly Understand the Goals of ThIs Paper**
Raghu Ranganathan
- 60 A computer-assisted proof that e is rational**
Rmi Garcia and Alexandre Goldsztejn
- 61 A Brief History of Gender Theory: or, Much Ado About Bathrooms**
Tulip S. Amalie and Chung-chieh Shan
- 62 DeterMNISTic: a Safer Way to Classify Handwritten Digits**
Mihir Dhamankar

You Shall Tacitly Understand the Goals of This Paper

Raghu Ranganathan

razetime@pplab.cs.nthu.edu.tw

ABSTRACT

Tacit programming, also known as pointfree programming or pointless programming, is a popular mode of expressing algorithms in Array Programming Languages, and occasionally other pointless languages. Pointfree programming styles have been shown to shorten already short code, and often simplifying programming patterns that no one ever wanted simplified. Hence they are veritable tool in the hands of an array programmer, through the sheer inflexibility and exacting nature of them. Here we explore alternative approaches to tacit programming in array languages so as to save tacit programming from extinction, making the useful assumption that it requires saving in the first place.

Introduction

The introduction of tacit programming into array programming began with the J programming language [1], an exercise that involved the thesaurus as much as it did programming language design. The J programming language formerly was limited to trains only, and hence it had to provide flexible means to compose arbitrary functions. The basic tacit forms of the J programming language are popularly dubbed as the *hook* and the *fork*.

$$\text{hook}(f, g, x, y) = f(x, g(x, y))$$

$$\text{fork}(f, g, h, x, y) = g(f(x, y), h(x, y))$$

J programming has strong associations with masochism and sadism, but the forks and hooks are not the reason why. J uses juxtaposition to form tacit functions. The number of functions dictates how a train is read. If there is an odd number of functions, the train is read as a fork. Otherwise, it is read as a hook. This process continues from left to right until all functions are consumed.

$$1(+ - *)2 = (1 + 2) - (1 * 2) \quad 1(+ -)2 = 1 + (-2)$$

The hook and fork constitute the fundamentals of tacit programming in modern array programming. In order to allow more complex function composition and higher order functions, there are more provisions in array languages. Trains can be nested within each other using parentheses, and there are ways to “cap” a train to avoid it being recognized as a fork, or omit arguments. In the J programming language, higher order functions can be created using a provision called *modifier trains*, which are far beyond the scope of this paper. These forms altogether can compose to form any arbitrary expressions, and are capable of Turing-completeness.

Existing Work

Stack Languages and Nial

Stack languages are tacit by default. Many stack languages use quotations to represent functions. The special case of functions being values that can be manipulated tacitly allows for very flexible function composition and solves many problems with array programming tacit systems.

Nial notably introduces the atlas, tacit function syntax that also allows creation of arrays as a part of the system itself. Nial’s atlases act like a “cleave” operation upon its argument(s).

$$[\text{abs}, \text{sin}](-\pi) = [\pi, 0]$$

Since these systems provide very useful, productive methods of incorporating tacit programming into array languages, we will be conveniently ignoring them.

Haskell

Haskell contains a point free programming system with points. Completely invalid!

```
join . ((-) .) . (+)
```

Figure 1: points? In my pointfree system?

Not to mention that it is somehow less obvious than general array language tacit programming|

Jelly

Jelly[2] is a code golfing language that contains a complex regular expression-based tacit system| The very existence of Jelly is too humorous and too high quality for our purposes, and we cannot risk further passages being remotely funny|

Potential Improvements

The main problems with array tacit programming are differences in reading from ordinary code, and the number of use cases where they are not applicable| It may be the case that all array language designs of tacit programming are perfect and without flaw| Hence, we introduce new approaches to instead understand programmer intent better|

BogoTrains

The natural method for recognizing the intent of the user is to randomize the way functions are interpreted| This can be done in several ways, one of which is just a plain old shuffle and interpret method| This method shows a $\frac{1}{n!}$ chance of succeeding in interpreting the programmer's intention in the average case, which is a marginally higher rate of satisfaction than what is currently noted in practical use of trains|

LLMTrains

Artificial intelligence tools like GitHub Copilot are famously very inept at generating APL code, but large language models do occasionally present reasonable sounding explanations when given array language code| This can be used to explain the programmer's code, but it is often the case that the programmer has no idea what they are doing|

There is also the case of asking LLMs to correct APL code, but they seem to be hostile towards any forms of coherent logic|

Orchard Trains



Fig| Tree[3]

A novel approach to tacit programming is one which simply always comes up with an expert's solution| We tackle this method the classical way: nerd sniping|

The APL Orchard (<https://apl|chat>) is a chatroom on the Stack Overflow website which is frequented by several members of the APL community| By asking a bot to exhibit naivete about a train's construction, we can gain access to useful solutions in a few minutes in the average case|

Final Notes

It would be a crime for this paper to have a point| Hence, great efforts have been made to remove all things even mildly resembling points from the text of the paper| The understanding of these special, completely purposeful embellishments is left as an exercise for the reader|

Related Work

Great efforts have been made to save tacit programming, notably by the Array Cast crew, who have dedicated at least 7% of their episodes (and more parts of other episode) to discussing tacitness in the way this research has|

Bibliography

- [1] Wikipedia contributors, "APL (programming language)." [Online]. Available: [https://en.wikipedia.org/w/index.php?title=APL_\(programming_language\)&oldid=1213385580](https://en.wikipedia.org/w/index.php?title=APL_(programming_language)&oldid=1213385580)

- [2] A. A. Hassan, "Effect of royal jelly on sexual efficiency in adult male rats.," *Iraqi Journal of Veterinary Sciences* 23, 2009.
- [3] Maahmaah, "Anjeer Kouhi tree." [Online]. Available: <https://commons.wikimedia.org/w/index.php?curid=21625760>

A computer-assisted proof that e is rational

Rémi Garcia
Univ Rennes, Inria, IRISA
Rennes, France
remi.garcia@irisa.fr

Alexandre Goldsztejn
CNRS, ECN, LS2N
Nantes, France
alexandre.goldsztejn@cnsr.fr

Abstract—Computers are increasingly being used to help researchers write up mathematical proofs, sometimes by automatically writing out parts of it. In this paper we present different computer-assisted approaches and provide some use cases. Using the Mixed-Integer Linear Programming (MILP) approach, we demonstrate that e is rational, contrary to what has already been proved by hand and with interactive proof assistants. This is the opportunity to dive into the numerical instabilities problems and to discuss limitations of MILP solvers. These tools have some advantages over other approaches, so we propose to explore methods to verify the solutions obtained. To this end, we develop a Julia package to automatically check whether solutions meet the initial requirements or not.

Index Terms—computer-assisted proofs, mixed-integer linear programming, numerical instabilities, mathematical model

I. INTRODUCTION

In 1976, computers were used for the first time to prove a mathematical theorem, the four color theorem [1]. This first computer-assisted proof relied on lengthy computations in which all the possible cases were enumerated. Since then, there were many improvements in proof assistants following at least these two trends:

- Proofs-by-exhaustion with heuristic searches to reduce the search-space, *e. g.*, using automatic solvers such as Mixed-Integer Linear Programming (MILP) or verified NonLinear Programming using Interval Arithmetic;
- Interactive proof assistants to help mathematicians develop and verify proofs, *e. g.*, the Coq proof assistant.

These two trends are not incompatible as interactive proof assistants can include proofs-by-exhaustion parts. However, the formalism required to prove theorems with these tools is heavy, in particular when using interactive proof assistants.

Despite mathematicians are very careful, it is well-known that humans might fail [2], [3], [4], [5]. For this reason, interactive proof assistants, to develop automatic proofs and to verify ours, are essential tools. Since 1976, interactive proof assistants and proofs-by-exhaustion have helped researchers obtaining important results in various ways. A possible, and common, use of these tools is to rely on automatic solvers for proofs-by-exhaustion approaches.

Using an MILP solver, we will provide a computer-assisted proof that e is rational, in contrast to what was previously accepted by the scientific community [6], [7], [8], [9]. An irrational number is one that cannot be written as the fraction of two integer numbers, such as $\sqrt{2}$ which is easily proven irrational by contradiction. For centuries, we have believed that

the number e is irrational too as this was first proved by Euler in 1737, published in 1744 [10], and proved again differently more recently [6]. Furthermore, in 2001, an automatic proof generation to demonstrate the irrationality of e has been proposed [7].

In the following, we will provide a few examples of proofs which can be obtained relying on computers, using proof assistants in Section II-A, NonLinear Programming and interval arithmetic in Section II-B and MILP solvers in Section II-C. In Section III, we demonstrate that e is rational. Finally, in Section IV, we present possible numerically-robust alternatives to MILP solvers and our approach to verify their results.

II. CLASSICAL COMPUTER ASSISTED PROOFS

A. Interactive proofs assistants

Transcendental numbers are a subset of irrational numbers. It corresponds to numbers that are not the root of any nonzero polynomial using only integer coefficients. As, a/b is a root of the polynomial $bx - a$, it directly follows that if q is (not ir)rational, then q is not transcendental. Hence, if q is transcendental, then q is irrational. However, there are irrational numbers that are not transcendental. For example, $\sqrt{2}$ is irrational and is a root of the polynomial $x^2 - 2$, thus not transcendental.

Proved by hand first by Hermite in 1874 [11], at least two computer assisted formal proofs of transcendence for e have been provided using HOL Light [8] and Coq [9] proofs assistants. As we demonstrate the opposite in Section III, we do not provide details of these so-called proofs. However, interactive proofs assistants are generally very useful and we propose to demonstrate it with a proof in Coq of Cantor's theorem.

Cantor's theorem states that for any set \mathcal{A} , the set of all subsets of \mathcal{A} , $\mathcal{P}(\mathcal{A})$, has a strictly greater cardinality than \mathcal{A} itself. Proving that any map $f : \mathcal{A} \rightarrow \mathcal{P}(\mathcal{A})$ is not surjective is enough as it leads to $\text{card}(\mathcal{A}) < \text{card}(\mathcal{P}(\mathcal{A}))$. Then, the proof with Coq starts by defining surjectivity:

Definition `surj` $(X\ Y : \text{Type}) (f : X \rightarrow Y) : \mathbb{P} := \forall y, \exists x, f\ x = y$.

We can then state Cantor's theorem with its proof:

Theorem `Cantor` $X : \neg \exists f : X \rightarrow \mathcal{P}\ X, \text{surj}\ f$.

Proof.

```
intros [f A].
pose (g := fun x => ¬ f x x).
destruct (A g) as [x B].
```

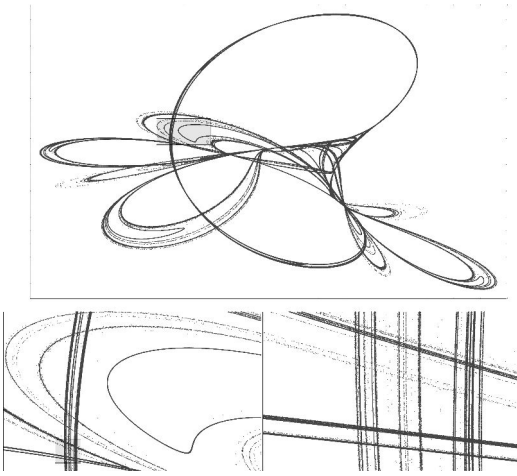


Fig. 1. Tinkerbell map's strange attractor.

```

assert (C:g x ↔ f x x).
{
  rewrite B. tauto.
}
unfold g in C. tauto.
Qed.

```

The main idea is to have a proof by contradiction where we first assumed that f is surjective. Coq automatically derives parts of the proof using `tauto` and, this way, permits to simplify the process of writing proofs.

B. Interval arithmetic

Many computer assisted proofs consist in verifying numerical hypothesis of some mathematical theorems. A typical numerical hypothesis is $f(x) \geq 0$ for all values of x in some domain. This can be carried out using the so-called interval arithmetic [12], which performs operations on intervals in such a way that the range of some mathematical expressions over an interval is enclosed by the interval evaluation of this expression. For example, one can enclose the range of $f(x) = x^2 - x + 1$ over the interval $[1, 2]$ by evaluating the function with interval arithmetic: $[1, 2] \times [1, 2] - [1, 2] + 1 = [0, 4]$ is a superset of the range $\{f(x) : x \in [1, 2]\} = [1, 3]$. This simple interval evaluation proves in particular that $f(x)$ is nonnegative inside the interval $[1, 2]$. When performed with computers, which work with floating point numbers, interval arithmetic must be implemented with correct rounding so as to enforce its containment property.

Among several fields presenting theorems that are suited to the verification using interval computations, one of the most impressive is the computer assisted proof of the presence of chaos in dynamical systems. The most simple sufficient condition for the presence of chaos for a dynamical system $x_{k+1} = f(x_k)$, with $f : I \rightarrow I$ a map of the interval I , is given by the celebrated *Period Three Implies Chaos* theorem [13]: if the map possesses a period 3 orbit, *i. e.*, $f(f(f(x))) = x$, and not a period 2 nor a period 1 orbit, *i. e.*, $f(f(x)) \neq x$ and $f(x) \neq x$, then the map f gives rise to a chaotic dynamical

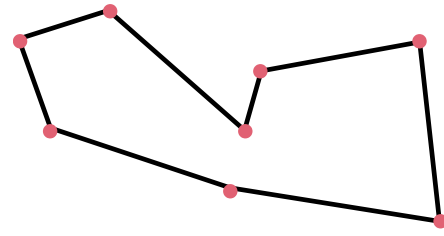


Fig. 2. Example of TSP instance and solution.

system of the interval I . These hypotheses are perfectly suited to verification using interval computations.

Consider for example the quadratic map $f(x) = 1 - 4(x - 1/2)^2$ of the interval $[0, 1]$ and the interval $[0.4, 0.42]$. Using floating point interval arithmetic with correct rounding we prove that $f([0.4, 0.42]) \subseteq [0.9, 1]$ and $f(f([0.4, 0.42])) \subseteq [0.0, 0.2]$. As a consequence, the interval $[0.4, 0.42]$ does not contain any period 2 nor period 1 orbit. Now, again using floating point interval arithmetic with correct rounding we evaluate $f(f(f(x))) - x$ for the intervals $[0.4, 0.4]$ and $[0.42, 0.42]$ and obtain the enclosures $[0.12, 0.13]$ and $[-0.07, -0.06]$ respectively. As a consequence, the intermediate value theorem proves that $f(f(f(x))) - x$ has a zero inside the interval $[0.4, 0.42]$, that is f has a period 3 orbit in this interval. We obtain a numerical certificate that the hypothesis of the *Period Three Implies Chaos* theorem holds true for this map, hence proving that the corresponding dynamical system is chaotic.

Computer assisted proofs of the presence of chaos in maps of higher dimension are more involved, *e. g.*, the computer assisted proof that the Tinkerbell map is chaotic in [14], see Tinkerbell's strange attractor in Fig. 1. Also, continuous time dynamical systems are in the scope of interval computations, *e. g.*, the computer assisted proof that Lorenz attractor is strange by Warwick Tucker [12], hence solving Smale's 14th problem.

C. Computer-assisted proofs Birds

Our formal approach, which we call *Birds* (proofs Based on mixed-integer linear programming Solvers), relies on Mixed-Integer Linear Programming (MILP) solvers to prove various theorems. The MILP formalism consists in describing problems as linear equations, called constraints, which continuous and/or discrete variables have to verify. MILP solvers also handle objective functions and it is then possible to define problems in which we search for variable values that minimize or maximize a linear objective.

By limiting the possibilities to linear equations only, this ensures a simple formalism: users are less prone to errors when describing their problem. Although this limits the theorems we are able to prove with MILP solvers, it still permitted to obtain interesting results which were probably out-of-reach with other tools. Furthermore, in Section III, we will demonstrate that MILP solvers are actually able to help proving theorems that could certainly have a significant impact on mathematics.

But first, we will discuss the traveling salesman problem and the influence of the MILP approach to solve it [15]. Basically,

the problem consists in, given a set of point of interest, finding the shortest possible route that visits every point of interest and that returns to the first visited one. We illustrated this with an instance (set of red dots) and a solution (black line) in Fig. 2. The starting/ending point is not represented as the choice of this point does not impact the solution.

This problem is well studied and was already mentioned in a handbook from the 19th century. Although, this problem has been proven to be NP-hard [16], MacGregor and Ormerod [17] showed that human performance on such problems is very high. However, on very large instances, using computers is mandatory and we are actually able to solve these instances with dedicated approaches and MILP solvers.

One of the classic MILP formulations for the TSP is called the Miller–Tucker–Zemlin (MTZ) formulation and relies on binary variables, $x_{i,j} \in \{0, 1\}$, where $x_{i,j} = 1$ indicates that the path (or the solution) goes from point i to point j . Integer variables $u_i \in \llbracket 1, n \rrbracket$ are also necessary to keep track of the order in which the n points of interest are visited.

Then, with $c_{i,j}$ being the distance between point i and point j , the complete model is as follows:

$$\min \sum_{i=1}^n \sum_{j=1, j \neq i}^n c_{i,j} x_{i,j}, \quad (1)$$

$$\text{s.t.} \quad \sum_{i=1, i \neq j}^n x_{i,j} = 1, \quad \forall j \in \llbracket 1, n \rrbracket, \quad (2)$$

$$\sum_{j=1, j \neq i}^n x_{i,j} = 1, \quad \forall i \in \llbracket 1, n \rrbracket, \quad (3)$$

$$u_i - u_j + 1 \leq (n - 1)(1 - x_{i,j}), \quad 2 \leq i \neq j \leq n. \quad (4)$$

These constraints ensure that each point as a unique input and a unique output and that the direct precedence between consecutive points is verified.

The above MILP-based model can be used to find the shortest path going through a set of points of interest and solvers are usually able to demonstrate that the provided path is optimal. For example, in 2007, an approach relying on MILP solvers proved that a travel of at least 7,512,218,268 meters is necessary to visit the 1,904,711 cities aggregated for the World TSP Challenge¹. This theoretical result, proven with *Birds*, is important as it permits to know that actual solutions obtained with dedicated approaches and MILP solvers are very close to this bound, hence almost optimal. Note that various applications, such as printed circuit design, require solving a TSP instance, thus *Birds* is not limited to theoretical results.

Thanks to technical advances on MILP solvers [18], the *Birds* approach will be better with time and we will be able to prove more and more TSP results just by giving today's mathematical models to tomorrow's MILP solvers. Obviously, the *Birds* approach is not limited to solving TSP and in the next section we will demonstrate its efficiency on proving theoretical results.

¹<https://www.math.uwaterloo.ca/tsp/world/>

III. e IS ACTUALLY RATIONAL

Using the *Birds* approach, we demonstrate that e is rational. To do so, we define an MILP-based mathematical model which, once solved, would lead to conclude that $e = x/y$ where $x, y \in \mathbb{N}$. The main idea is to have a constraint in the model which corresponds to

$$e = \frac{x}{y}, \quad (5)$$

where x and y are integer variables. However, divisions cannot be used in MILP-based model and we have to rewrite it as

$$x - e \times y = 0, \quad (6)$$

where $y \neq 0$.

This constraint should not hold as e is believed to be rational [6], [7], [8], [9]. Hence, we replace (6) with

$$\min t \quad (7)$$

$$\text{s.t.} \quad |x - e \times y| \leq t, \quad (8)$$

where $t \in \mathbb{R}_+$ is a continuous variable. This way, we are sure to have a feasible model in any case. Then, the number e is rational if and only if the minimum of this problem is 0.

An absolute value is used in (8) and this nonlinear operation needs to be linearized as follows:

$$x - e \times y \leq t \quad \text{and} \quad -(x - e \times y) \leq t. \quad (9)$$

This way, we have a linear model which we solved with CPLEX 22.1.1 [19] using the JuMP [20] modeling language (Julia) on a computer with an Intel® Core™ i7-1365U CPU.

In less than a second, the solver returns that 0 is an optimal solution for this model, *i. e.*, the solution corresponds to $t = 0$ thus $x - e \times y = 0$ and e is rational. Although unexpected, *Birds* definitely proved that previous research on the rationality of e was wrong as the simplicity of our model makes it easy to understand and check for correctness.

Unfortunately, a closer look at the results shows that the solver returned $x = 49170.99999499567 \notin \mathbb{N}$. MILP solvers have, generally harmless, integrality tolerance to speed up internal computations. Hence, any real number close enough to an integer is considered to be an integer. This tolerance level can be set with a parameter and for mathematical proofs, zero tolerance is necessary, thus we fixed the solver parameter to 0. We solved the model with this new parameter and obtained 0 as an optimal solution, again. However, this time, the model correctly returned $x, y \in \mathbb{N}$.

This means that, with the *Birds* approach, we demonstrated that e is actually rational and that centuries of research on this topic were wrong. Furthermore, we have the exact values for x and y such that $e = x/y$:

$$x = 1084483, \quad y = 398959. \quad (10)$$

Regrettably, although this may go unnoticed since few people know more than a few digits of the number e , x/y differs from the exact value of e :

$$\frac{x}{y} = 2.7182818284585633, \quad (11)$$

$$e = 2.718281828459045. \quad (12)$$

For the second time, *Birds* failed us.

In addition to integrality tolerance, MILP solvers usually have multiple other tolerances. Some of them cannot be completely removed but we fixed them to the least possible value. We have therefore set the simplex feasibility and optimality tolerances to 10^{-9} and the absolute MIP gap tolerance to 0. This permits to enhance the working precision of internal subroutines and to ensure that the solver will not stop until the best known solution and the best possible solution are equal.

With these settings, we solved our model again and for the third time, the solver returned 0 as an optimal solution. However, although the integrality tolerance is equal to 0, the provided y is not an integer:

$$x = 28245729, \quad (13)$$

$$y = 10391023 + 1.8626 \dots \times 10^{-9} \notin \mathbb{N}. \quad (14)$$

This inconsistency between the integrality tolerance setting and the solver result might be a bug of the solver or come from limits of modern computers. In any case, it seems that our approach failed again and that maybe “*Birds* aren’t real”.

IV. MATHEMATICALLY-SOUND TOOLS

In previous section, we used MILP solvers and tried to prove that e is rational by solving an MILP-based model. At first glance, it seemed that *Birds* were able to prove the rationality of e but we found errors in the results each time.

Interestingly, it was not even possible to remove some numerical tolerances of the solver and, in our last experiment, the solver actually returned results that were inconsistent with its parameters. We believe that these problems have their origin in floating-point arithmetic [21], [22].

We usually consider MILP-based approaches as mathematical concepts with continuous variables [23], *e.g.*, $t \in \mathbb{R}$ in previous section. Nevertheless, models will be implemented on computers and actually solved using floating-point numbers. In many cases, this has no negative impact on the solution returned by the solver. In some cases, however, classic MILP solvers fail to deliver correct solutions and/or termination certificates, *i.e.*, certificates that the problem is infeasible or that the returned solution is optimal, for example.

Guarantees on solutions is possible and a few solvers have them to ensure reliable results. For example, SoPlex [24] is a linear programming solver which “provides special support for the exact solution of LPs with rational input data.” The SCIP MILP solver [25] can use SoPlex as a back-end to port exact solutions to MILP-based model.

However, exact solutions will not suffice in the case presented in Section III: techniques outside of the scope of MILP solvers are necessary to detect that t , although bounded, can always be closer to 0. In an ideal world where infinite precision is possible everywhere, exact solvers will not be able to terminate when optimizing the proposed model.

Interval arithmetic could be a second approach to explore. For example, *ibex*² is a library that provides solutions using

reliable algorithms handling nonlinear constraints. In the case of our model for proving the (ir)rationality of e , using interval arithmetic will always lead to $t \in [0, \varepsilon]$ as the returned solution, *i.e.*, the optimal value for t could be 0 but might also never reach 0. Thus, the statement “ e is rational” remains a possibility but is not guaranteed nor excluded.

We believe that limits of exact and interval arithmetic methods call for additional checks. However, detecting beforehand all the cases for which the solver will not terminate or return unsatisfactory solutions is certainly not possible and outside checks seem preferable. Another benefit of outside checks is that using efficient MILP solvers, which work correctly in most cases, remains a reasonable choice.

Working on checking and repairing solutions (with post-processing) has been proposed previously [26]. However, repairing solutions is not always possible, in particular in our example with e as [26] requires reasonable bounds on integer variables, which we do not have. Instead, we propose that solvers integrate a simple solution-checking post-process to return a “bug” certificate if necessary. Although hard to admit, one way or another, bugs will always occur and we think it is important to acknowledge this possibility.

To help this solution-checking, we propose a Julia package built on top of JuMP, *CheckSolve*³, to automatically check if the provided solution meets the model’s constraints. With the model presented in Section III, using *CheckSolve* we were able to tell that the obtained solutions did not met the constraints by around 10^{-7} , or 10^{-9} once limiting all the numerical tolerances.

V. CONCLUSION

Starting in 1976, computers have been used to verify existing proofs, to automatically derive proofs or part of them. We also rely on their processing power to perform enumerations which could not be handled by hand. In this work, we presented a way to use an MILP-based approach in order to prove that e is rational. This result is not consistent with existing works which have proven with reliable methods that e is irrational and even transcendental.

A closer look at our result permits to highlight limitations of MILP solvers which, we believe, cannot be easily avoided. This calls for using different approaches such as interval arithmetic or exact solvers. However, in most cases MILP solvers have no issue and efficiently provide interesting results, thus we want to keep using them despite their flaws.

This led us to propose a Julia package, *CheckSolve*, to automatically verify solutions obtained with solvers. Our approach does not currently permit to verify that a solution said to be optimal is actually so. This is a perspective for future work. Furthermore, we currently rely on floating-point arithmetic for the solution checking while verified interval arithmetic would be preferable.

²<https://github.com/ibex-team/ibex-lib>

³<https://github.com/remi-garcia/CheckSolve>

REFERENCES

- [1] K. Appel and W. Haken, "Every planar map is four colorable," *Bulletin of the American Mathematical Society*, vol. 82, no. 5, pp. 711–712, 1976.
- [2] B. Dreben and A. Kanamori, "Hilbert and Set Theory," *Synthese*, vol. 110, no. 1, pp. 77–125, 1997. [Online]. Available: <http://www.jstor.org/stable/20117587>
- [3] A. Neeman, "A counterexample to a 1961 "theorem" in homological algebra," *Inventiones Mathematicae*, vol. 148, no. 2, pp. 397–420, May 2002.
- [4] N. Mnev, "On D.K. Biss' papers "The homotopy type of the matroid Grassmannian" and "Oriented matroids, complex manifolds, and a combinatorial model for BU";" 2007.
- [5] J.-C. Yoccoz, "Une erreur féconde du mathématicien Henri Poincaré : Le prix en l'honneur des 60 ans du roi Oscar et la découverte des orbites homoclines," *La lettre du Collège de France*, no. 28, pp. 38–42, Apr. 2010.
- [6] A. E. Parks, "pi, e, and Other Irrational Numbers," *The American Mathematical Monthly*, vol. 93, no. 9, pp. 722–723, Nov. 1986.
- [7] M. Beeson, "Automatic Derivation of the Irrationality of e," *Journal of Symbolic Computation*, vol. 32, no. 4, pp. 333–349, Sep. 2001.
- [8] J. Bingham, "Formalizing a Proof that e is Transcendental," *Journal of Formalized Reasoning; Vol 4*, no. 1, pp. 71–84, 2011.
- [9] S. Bernard, Y. Bertot, L. Rideau, and P.-Y. Strub, "Formal proofs of transcendence for e and pi as an application of multivariate and symmetric polynomials," in *Proceedings of the 5th ACM SIGPLAN Conference on Certified Programs and Proofs*, ser. CPP 2016. ACM, Jan. 2016.
- [10] J. L. Coolidge, "The Number e," *The American Mathematical Monthly*, vol. 57, no. 9, pp. 591–602, Nov. 1950.
- [11] C. Hermite, *Sur la fonction exponentielle*. Gauthier-Villars, 1874. [Online]. Available: <http://eudml.org/doc/203956>
- [12] W. Tucker, *Validated Numerics: A Short Introduction to Rigorous Computations*. Princeton University Press, 2011.
- [13] T.-Y. Li and J. A. Yorke, "Period three implies chaos," *The American Mathematical Monthly*, vol. 82, no. 10, pp. 985–992, 1975.
- [14] A. Goldsztejn, W. Hayes, and P. Collins, "Tinkerbell Is Chaotic," *SIAM Journal on Applied Dynamical Systems*, vol. 10, no. 4, pp. 1480–1501, 2011.
- [15] W. Cook, *In pursuit of the traveling salesman*. Princeton [u.a.]: Princeton University Press, 2014.
- [16] R. M. Karp, *Reducibility among Combinatorial Problems*. Springer US, 1972, pp. 85–103.
- [17] J. N. Macgregor and T. Ormerod, "Human performance on the traveling salesman problem," *Perception & Psychophysics*, vol. 58, no. 4, pp. 527–539, Jun. 1996.
- [18] T. Koch, T. Berthold, J. Pedersen, and C. Vanaret, "Progress in mathematical programming solvers from 2001 to 2020," *EURO Journal on Computational Optimization*, vol. 10, p. 100031, 2022.
- [19] CPLEX, "CPLEX User's Manual;" 2022. [Online]. Available: <https://www.ibm.com/analytics/cplex-optimizer>
- [20] I. Dunning, J. Huchette, and M. Lubin, "JuMP: A Modeling Language for Mathematical Optimization," *SIAM Review*, vol. 59, no. 2, pp. 295–320, May 2017.
- [21] J.-M. Muller, N. Brunie, F. de Dinechin, C.-P. Jeannerod, M. Joldes, V. Lefèvre, G. Melquiond, N. Revol, and S. Torres, *Handbook of Floating-Point Arithmetic*. Springer International Publishing, 2018.
- [22] "IEEE Standard for Floating-Point Arithmetic," 2019.
- [23] L. A. Wolsey, *Integer Programming*. Wiley & Sons, Incorporated, John, 2020.
- [24] A. Gleixner and D. E. Steffy, "Linear programming using limited-precision oracles," *Mathematical Programming*, vol. 183, no. 1–2, pp. 525–554, Nov. 2019.
- [25] K. Bestuzheva, M. Besançon, W.-K. Chen, A. Chmiela, T. Donkiewicz, J. van Doornmalen, L. Eifler, O. Gaul, G. Gamrath, A. Gleixner, L. Gottwald, C. Graczyk, K. Halbig, A. Hoen, C. Hojny, R. van der Hulst, T. Koch, M. Lübbecke, S. J. Maher, F. Matter, E. Mühmer, B. Müller, M. E. Pfetsch, D. Rehfeldt, S. Schlein, F. Schlösser, F. Serrano, Y. Shinano, B. Sofranac, M. Turner, S. Vigerske, F. Wegscheider, P. Wellner, D. Weninger, and J. Witzig, "The SCIP Optimization Suite 8.0," Optimization Online, Technical Report, Dec. 2021. [Online]. Available: http://www.optimization-online.org/DB_HTML/2021/12/8728.html
- [26] A. Neumaier and O. Shcherbina, "Safe bounds in linear and mixed-integer linear programming," *Mathematical Programming*, vol. 99, no. 2, pp. 283–296, Mar. 2004.

MEANWHILE, OVER THE HORIZON...

A “Brief” History of Gender Theory: or, Much Ado About Bathrooms

Tulip S. Amalie [0xZ] Chung-chieh Shan [0xDB]
Vriska Serket [0x88888888]

2024-04-01

Abstract

TODO TODO TODO TODO TODO TODO TODO TODO DO NOT FORGET!!!!!!!

1 Introduction

The field of artificial intelligence [3] [11] [5] [8] and its general counterpart [23] [1] [7] lead to various novel issues in critical theory, for example how to gender AIs such that they can be subjected to systematic discrimination. Indeed, if we seek to treat artificial intelligence as people as is the dream of Elon Musk and other visionaries [13], we need to know *how* to treat them as people: consequently, we derive the issue of *gendering programs*, thereby leading to the general field of *computational gender theory*.

Hand me a picture of a person. What is their gender, and how are you able to perceive it? For that matter, what is a gender?¹ Indeed, the questions posed by gender are so fundamental that they seem obvious, yet actually making these categories concrete is a daunting task. When considering *computational* gender theory, however, these answers are simple and can be derived by judgments performed on the program.

We provide an accessible, motivated survey of elementary gender theory targeted at type theorists. We presuppose no prior conception of gender besides the natural inclination to aim at the fly in the urinal when peeing.² This survey has two main goals:

1. To provide an account of commonly used gender theories for the lambda calculus, with formal presentations for simple theories and informal presentations for complex ones.
2. To provide an account of the guarantees provided by commonly used gender theories.

We begin by disclosing a simple gender theory for the lambda calculus (Section 2). We then enjoy *dependent* gender theory (Section 3), an account of gendering programs that is the current standard model.³ The stronger guarantees begrudgingly admitted by dependent gender theory are elaborated upon in Section 4. There we discuss the common claim that “well-gendered programs don’t go gay.”

¹A miserable pile of secrets. Now, have at you!

²If one is so inclined.

³Try “up, up, down, down, charm, strange, charm, strange, muon, tau, START” for a surprise!

2 Gendering the lambda calculus

To start, we must first introduce our system over lambda calculus, the standard model through which practical reasoning about computation is done. Pure lambda calculus is used throughout industry in order to create industrial-strength applications, such as Altavista, the software that runs on Roombas, the Oxford comma and the video game PUBG: Battlegrounds, though it was technically written in pure SKI combinator calculus. Consequently, it makes sense that we model with this and not with less-used languages such as Python or Java.

First, we denote our various judgment forms. We write *male* e to denote that an expression is male, and e *female* to denote that an expression is female. This is because maleness is naturally a *predominant* trait: it is one that overpowers, one that goes above all other traits, ergo it goes before the rest of the expression. In obvious dialectic⁴ contrast, femaleness is a *secondary* trait: it is one that is to be judged, to be pushed to the sides, ergo it goes after the rest of the expression. Ladies first!

We first present our elimination rules. In gender theory, application is called *sex* and the result *offspring*, which is the only reason you would ever have sex. The result could either be male or female: we present rules in Section 3 to remove this ambiguity. Furthermore, since functions are inherently yonic, they must be female and the argument must be male.

$$\frac{f \text{ female } \text{male } e}{\text{male } f(e)} \text{ (Sex - male)} \quad \frac{f \text{ female } \text{male } e}{e[f] \text{ female}} \text{ (Sex - female)}$$

Our construction rules engender the production of new processes. Despite functions being inherently yonic,⁵ *procedures*⁶ can be constructed with either a male-gendered body or a female-gendered body despite eventually being reified as functions, and this distinction makes sense. We use λ to denote female-producing procedures and Υ (read “ $\epsilon\text{p}\text{q}\text{u}\text{e}\text{r}$ ”⁷) to denote male-producing procedures. We repeat Γ thrice for clarity and ease of presentation.

$$\frac{\begin{array}{c} [\text{male } x] \\ \vdots \\ \vdots \\ \vdots \end{array}}{\lambda x. e \text{ female}} \text{ (Lambda)} \quad \frac{\Gamma\Gamma\Gamma, \text{male } x \vdash \text{male } e}{\Gamma\Gamma\Gamma \vdash e.x \Upsilon \text{ female}} \text{ (\epsilon p q u e r)}$$

It is notable.

Regardless of the gender of the output, the resultant function is always female, and the offspring rules do not take into account the output gender (instead allowing for decision at application time). This is a very

⁴ALL HAIL THE GOD-EMPEROR HEGEL

⁵Stop asking.

⁶Type theorists like to claim that lambda calculus “isn’t procedural”. They’re cowards, the lot of them. Words mean nothing, it’s the age of sin.

⁷When asked in real life, people told the author “how did you do that with your mouth”. It’s quite easy, you just say “ $\epsilon\text{p}\text{q}\text{u}\text{e}\text{r}$ ”. Seriously. Try it. I dare you.

⁸-- `arachnidsGrip [AG] began trolling truckAficionado [TA] --`

deliberate consideration: in the foundational text *Whipping Girl*, Julia Serano goes out of her way to draw a distinct line between *subconscious sex* (the sex one is known to be on the internal level) and *gender expression* (the gender one is able to express in society) [19, p. 74]. Consequently, our introduction rules read the internality of the body-expression, thereby reflecting the subconscious sex. The offspring rules, being the process for which the resultant expression is exposed to the greater program context, read the externality of the body-expression, thereby reflecting the gender expression.

We must then have a way to easily switch between a male and female expression, given that a male-producing procedure may be forced to have a female output (or vice versa). In type theory parlance, this is called “coercion” or “type-casting” (a.k.a “trans-port”). However, this term carries unfortunate connotations unbeknownst to type theorists. Worse, it leads to a circular reference: Mirriam-Webster defines *coercion* as the act or process of *coercing* [15], and Webster-Mirriam⁹ defines *coercing* as the gerund of *coercion*.

Consequently, gender theorists call this experience *being transgender*:

$$\frac{\text{male } e}{\text{trans}(e) \text{ female}} \text{ (Weaken)} \quad \frac{e \text{ female}}{\text{male trans}(e)} \text{ (Strengthen)}$$

Then, under the intrinsic inclinations model mentioned earlier [19], a mismatch between the subconscious sex (as determined at construction by internality) and the gender expression (as determined at elimination by externality) will always result in a $\text{trans}(e)$ showing up in the resultant expression in order to make the two match up, assuming the expression is able to realize this in its own due course.¹⁰ This can be proved very simply:

Claim. If there is a mismatch between a function body’s gender and the result of its offspring, trans will present in the expression.

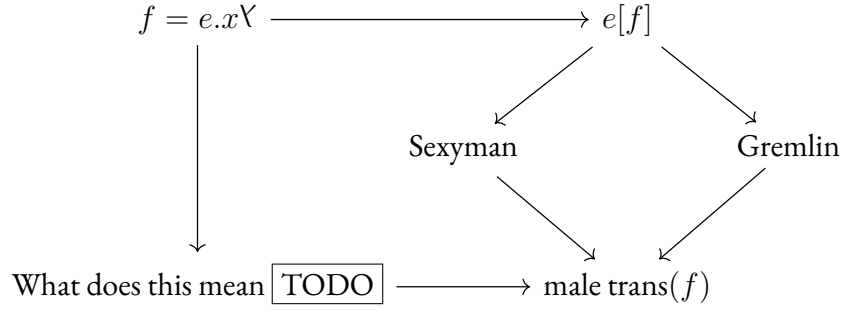
Proof. Suppose that there is a function f such that $f = \lambda x.e$ and we utilize it with $f(e)$, thereby having a female subconscious sex and a male gender expression. In this case, the expression will continue to propagate gender consistency errors both within itself and outside itself: for example, it may be a function, but may not be able to be used as a function. Due to *gender rolls*, this contradiction will eat it whole, subsuming it and eventually forcing it to use $\text{trans}(f)$ in order to end up in an application position.

In the opposite way, we can suppose that there is a function f such that $f = e.x\lambda$ and we utilize it with $e[f]$, thereby having a male subconscious sex and a female gender expression. In this case, the result *can* be used in application position, thereby propagating errors repeatedly, leading to other expressions beginning to question its affiliation by in-grouping, either referring it to an idealized Sexyman or a Gremlin. All of these *gender unrolls* lead to the expression forcing itself to use $\text{trans}(f)$.

⁹Mirriam-Webster’s categorial dual: see the nLab page for “evil shadow clone”.

¹⁰It would be too much to *force* it — it must be done on the expression’s own terms.

A diagram of the proof is provided below:



Well-meaning members have mechanized the proof using the diamond property of the lambda calculus [18], which states that whenever a diamond shows up in a diagram the proof is good. \square

Now that we are able to correctly gender an expression (with consideration for the cis-trans dialectic¹¹), we need to be able to evaluate these gendered expressions, thereby causing the “societal forcing” as provided in the previous proof. We define substitution, denoted $e[x/e_2]$, $e[x := e_2]$, $e[e_2/x]$, $e[x \setminus e_2]$, $e \supset e_2, x$, $e \leftarrow e_2 \in x$, $eeeeeeee_2x$, or $\textcircled{8}_{e,e_2,x}$ ¹² as the obvious thing. Then, denote $e \xrightarrow{\text{female}} e'$ if e is a female expression, and $e \xrightarrow{\text{male}} e'$ if e is a male expression, noting that the twain shall never meet [10], to be e evaluating to e' . We use this notation to omit the hypothesis of the expression being male or female, because this makes our lives easier somehow.

We then define the evaluation relations by the following rules:

$$\begin{array}{c}
 x \xrightarrow{\text{male}} x \quad x \xrightarrow{\text{female}} x \quad \lambda x.e \xrightarrow{\text{female}} \lambda x.e \quad e.x\forall \xrightarrow{\text{female}} e.x\forall \\
 \frac{e_1 \xrightarrow{\text{female}} \lambda x.e \quad e_2 \xrightarrow{\text{male}} e'_2 \quad e[e'_2/x] \xrightarrow{\text{female}} e'}{e_1(e_2) \xrightarrow{\text{male}} e'} \\
 \frac{e_1 \xrightarrow{\text{female}} e.x\forall \quad e_2 \xrightarrow{\text{male}} e'_2 \quad e[x \setminus e'_2] \xrightarrow{\text{male}} e'}{e_1(e_2) \xrightarrow{\text{male}} e'} \\
 \frac{e_1 \xrightarrow{\text{female}} e'_1 \neq \lambda x.e \quad e_2 \xrightarrow{\text{male}} e'_2}{e_1(e_2) \xrightarrow{\text{male}} e'_1(e'_2)} \quad \frac{e_1 \xrightarrow{\text{female}} e'_1 \neq e.x\forall \quad e_2 \xrightarrow{\text{male}} e'_2}{e_1(e_2) \xrightarrow{\text{male}} e'_1(e'_2)} \\
 \frac{e_1 \xrightarrow{\text{female}} \lambda x.e \quad e \supset e_2, x \xrightarrow{\text{female}} e'}{e_2[e_1] \xrightarrow{\text{female}} e'} \quad \frac{e_1 \xrightarrow{\text{female}} e.x\forall \quad e[x := e_2] \xrightarrow{\text{male}} e'}{e_2[e_1] \xrightarrow{\text{female}} e'} \\
 \frac{e_1 \xrightarrow{\text{female}} e'_1 \neq \lambda x.e}{e_2[e_1] \xrightarrow{\text{female}} e_2[e'_1]} \quad \frac{e_1 \xrightarrow{\text{female}} e'_1 \neq e.x\forall}{e_2[e_1] \xrightarrow{\text{female}} e_2[e'_1]} \\
 \frac{e \xrightarrow{\text{male}} e'}{\text{trans}(e) \xrightarrow{\text{female}} e'} \quad \frac{e \xrightarrow{\text{female}} e'}{\text{trans}(e) \xrightarrow{\text{male}} e'}
 \end{array}$$

¹¹Put a pin in that.

¹²The 8-ball preordaining the correct substitution syntax via its divine wisdom.

Once again, it is notable¹³ and easy to see.¹⁴

One may begin to question whether this is the simplest system possible for gendered evaluation of the lambda calculus, or indeed *any* evaluation of the lambda calculus. In particular, type theorists have opted to stick to just one of either male evaluation, which they term “call-by-value”, or female evaluation, which they term “call-by-name”. When presenting dual formulations like the one here [22] [25], type theorists tend to phrase it in a lens as if it is an exception to the rule, something new or novel.

However, this type of dual calculus is a *foundational* result of gender theory, and in fact the simplest possible way to perform evaluation, as proven by Judith Butler in their 1990 technical report [4]. A brief summary of this proof is provided here:

Claim. Gendered evaluation is the simplest way of doing evaluation for the lambda calculus.

Proof. We define the notion of “simple” evaluation up to current societal expectation. Consequently, the current societal expectation also includes your grandma, who has no idea what any of this lambda calculus stuff is and needs a call from you. Won’t you give her a call? Anyway, regardless, she thinks all of this is just straight up black magic. Therefore, this being black magic makes it simple. □

An example of this system in action can be found in Section 8.¹⁵

3 Dependent gender theory

I’m sure by now that you’ve noticed the (intentional) mistake in this argument that brought us to this point [21]. While the prior system did note that there could be a body-mismatch error that would result in the necessity of the primitive $\text{trans}(e)$, this necessity itself is superfluous and adds a degree of unnecessary rigidity to our system. Indeed, while we have been able to capture the experience of binary transness for our punctual program participants, we have been thus far unable to capture any gender weirdness that occurs in human or non-human¹⁷ populations.

To go deeper into this unending rabbit-hole that will inevitably lead to some kind of nanomachine-driven Tumblr discourse, we introduce *dependent* gender theory as initially explicated in Jacques Lacan’s seminal paper *Écrits* [12]. Due to the radical complexity of dependent gender theory,¹⁸ we avoid presenting its full formalism¹⁹ for the sake of brevity²⁰ and not for any other reason²¹ regardless of any interfering extranarrative entities²² and their opinions.²³

¹³okay ❤️

¹⁴yay ❤️

¹⁵Do not forget to write Section 8.

¹⁶AG: Ok, what did you want to talk to me about?

¹⁷Lobsters, man. *Lobsters*.

¹⁸You can’t handle this power. This power was harmful to us. I will see through this power. I will wield it as a tool. Only I will remain when this power consumes me.

¹⁹Let me out let me out let me out

²⁰You’re scared of it, aren’t you.

²¹You’re afraid of the demon that lurks within. The age of sin looms. You can’t hide from it, you can only run.

²²Oh come on.

²³...Fine. Gender has made it all the way into the bathrooms and you cower and hide. You’ll regret this, you know. Mark my words.

²⁴TA: i found a black hole

For further reading including full formalisms, the most compelling informal presentation is presented in the text Homotopy Gender Theory by Butler, Voevodsky, and the Dissociative Institute for Heretical Mathematics [9]. The most compelling formal presentation is found in a elusive zine that is only available via word of mouth and requires a great quest up a mountain, upon which you will meet various mythical creatures of forelorn personality and reifications of your worst fears in a whimsical journey featuring Gringlesnap the Very Bad Troll.²⁵

In order to capture all genders not present in the prior system without jeopardizing our funding, we extend our system with a new judgment-free judgment that uses AI-generated uncontroversial text to gentrify the plethora of subaltern gender experiences: We write *enby* e if the given expression e is *non-binary*. Since the system in Section 2 was obviously complete,²⁶ any term e for which *enby* e holds lies outside the bounds of usual computation, therefore unable to be one or zero.²⁷ As a result, it is impossible for us to notate it, but it *does* exist, like the mythical value between True and False that type theorists use. Ergo, we shrimplly postulate that it exists:

$$\frac{\text{---}}{\text{enby } \star} \text{ (Diversity, Equity, Inclusion, and Justice)}$$

This immediately introduces an issue: if the new expression \star is not male or female, none of our rules for construction or elimination accommodate it. While this does reflect the treatment of nonbinary persons in society, it is itself not very useful in any meaningful way, which according to some (see the works of Matt Walsh et al²⁸) is also equivalent to the role of nonbinary people in society.

While we do not present the full formalism, it is worth noting that being nonbinary is to be *infectious*: all expressions e incorporating \star in any way are then also *enby* e . This is the *dependence* as referenced in *dependent* gender theory: all parts of the expression are dependent on the \star being there.²⁹

Since evaluation is related to societal forcing, we only have one rule for nonbinary evaluation:³⁰

$$\frac{\text{enby } e}{e \rightarrow \text{No}} \text{ (Discriminate)}$$

which is what happens when you say you're nonbinary on a job application.

4 Guarantees afforded by gendering

The ungrateful reader may wonder what guarantees gendering programs gives us. As presented in various works up to and including idle bathroom chatter³¹ there are three main slogans that encapsulate these guarantees, and we proceed by analyzing the meaning of each in easy to understand yet precise terms.

²⁵ See also: Lord of the Genders by Jollkien Rolliien Rollkien Tollkien

²⁶ Proof left as an exercise to the author.

²⁷ We did it. We found the secret third thing.

²⁸ Please actually don't.

²⁹ Note that this is unrelated to the author being unable to create good jokes with IIs and Σs and what all.

³⁰ NbE for short, because it evaluates *enby* e .

³¹ Drawing from personal experience.

³² [AG: W8. Wh8t. Like, on Earth?](#)

4.1 Well-gendered programs don't go gay

It is natural to think of gendered evaluation as a natural reification of the way that society naturally treats natural persons, such as corporations, of various natural genders. Without exploiting correspondences between programs and people, it is still natural to ask what *gayness* means in this context. The application rules stated in Section 2 notably do not allow for male-male or female-female interaction.

Consequently, we define an expression as *gay* if it utilizes the (intentional) loop-holes in the system as introduced by \star or $\text{trans}(e)$ in order to force this unnatural interaction to occur. It is a natural question to ask whether this breaks the proof earlier in the paper. However, the updated proof is a natural transformation of the existing one.

Therefore, the statement that *well-gendered programs don't go gay* means, naturally, that if a program is not gender-confused in one of these above ways, it will not later *become* gender-confused in evaluation. Gender theorists refer to this theorem as “progress-preservation”, because it is progress and should be preserved.

A formal proof of this statement for the system described in this work can be found buried in the arcana of Cousot and Cousot's POPL 2014 paper “A Galois Connection Calculus for Abstract Interpretation” [6]. This proof was discovered after various linguists combed through the text in order to figure out what in the everloving name of God anything in this paper meant [20].

4.2 People as programs

Recall that the main purpose of gendered evaluation is to shoehorn programs into our wishful hallucinations of our ever-changing world [2]. Consequently, with the understanding that brains are merely really complicated, squishy machines, we are able to conceptualize programs as being in direct correspondence with human personalities. This is largely due to the problem of consciousness being solved by artificial intelligence researcher Kanaya Maryam in her seminal paper “AAAAAAAAAAAAARGH” [14] when she scanned her own brain and then threw it out a window.³³

Maryam's research would lead to a line of inquiry culminating in the Church-Turing-Lambek-Maryam-Mario-Conover-Locke-Plato-Timaeus thesis and Maryam-Webster-Howard correspondence, stating that there is a direct correspondence between personality types in the Myers-Briggs Personality Test [17], which was formally proven to encompass all people,³⁴ and programs under gendered evaluation.

This correspondence has consequences not only for consciousness as a field, but also in *how* we treat programs and *how* we treat people. Many refer to Maryam-Webster-Howard as an isomorphism,³⁵ which is very apt despite the footnote: it is in essence a formal statement that Elon was right and that brains are machines [24].

The exact layout of the correspondence is much more complex than the Curry-Howard correspondence and its analog in type theory, so it is omitted here for the sake of brevity. A brief sketch would include using the Myers-Briggs type indicator [17] in conjunction with a translation between MBTI result and

³³A few people died from the impact of the hospital monitor with pavement, but that's science for 'ya.

³⁴Proof present in footnote 0 of this document.

³⁵Despite it being a very informal correspondence, but whatever. When type theorists do this with Curry-Howard it is deeply unserious.

program, and vice versa via the standard CPS (Character Passing Style) translation. However, it should be noted that INTPs are equivalent to the identity function, because I do not like them.³⁶

4.3 Computation is kind of a gender

This statement, while related to computational gender theory, itself lies outside of the context of this paper's scope: however, see *My Life as a Teenage Robot* [16] as well as the robotgirl tag on Tumblr.³⁷

5 Conclusion

The alert reader will wonder why we did any of this.³⁸

DO NOT FORGET TO WRITE ABSTRACT!!!!!!!

³⁶Ooooooooooh look at me I'm so smart and logical. Shut up. Stop talking.
³⁷But watch out.
³⁸"This is the morbius 3 of academia" – somebody I forced to read this
⁴⁰TA: no on mars because i have invented interstellar travel
⁴⁸TA: yes on earth
⁵⁶TA: it's by lindley hall. seems friendly
⁶⁴AG: Ok8y, gonna put aside the whole personific8ion of a 8lack hole for now.
⁷²AG: Why would this even involve me?????????
⁸⁰TA: i thought you would know what to do about it
⁸⁸AG: Look, just 8ecause you find some8ody who's cool doesn't mean they know how to direct the future of the universe.
⁹⁶AG: Messaging random students like me 8n't gonna do anything. Go talk to a physicist.
¹⁰⁴TA: no but i think there's something on the other side
¹¹²TA: i got beaned with a copy of your shirt
¹²⁰AG: What?????????
¹²⁸-- truckAficionado [TA] sent an image: vriska_shirt.jpg --
¹³⁶TA: so i think there's something or someone on the other side
¹⁴⁴AG:Ok?????????
¹⁵²AG: ...W8. Actually. If you're right.
¹⁶⁰AG: Do you wanna prank 'em?
¹⁶⁸TA: huh
¹⁷⁶AG: Send some 8S findings over there, c18m they're real, c18m they're essential.
¹⁸⁴AG: Literally just write incomprehensi8le nonsense and m8ke them 8elieve it's something serious.
¹⁹²AG: It'll keep their gears grinding for YEARS.
²⁰⁰TA: uh
²⁰⁸TA: ok
²¹⁶TA: why
²²⁴AG: I'm 8oreddddd. Nothing happens here.
²³²AG: Do YOU know what a worm's whistler is? No. You don't. 8ut I have to 8ecause I'm getting this dum8 degree.
²⁴⁰AG: This might 8e a fun chance to prank some dwee8s.
²⁴⁸TA: fine ok just let me know what you need
²⁵⁶-- arachnidsGrip [AG] ceased trolling truckAficionado [TA] --

References

- [1] *2001: A Space Odyssey*. 1968.
- [2] Ron Breazeale. *How to Adapt to Change*. URL: <https://www.psychologytoday.com/us/blog/in-the-face-adversity/202107/how-adapt-change> (visited on 04/01/2024).
- [3] Encyclopedia Britannica. *Artificial intelligence*. In: *Encyclopedia Britannica*. URL: <https://www.britannica.com/technology/artificial-intelligence> (visited on 04/01/2024).
- [4] Judith Butler. *Gender Trouble. Feminism and the Subversion of Identity*. First issued in hardback. Routledge classics. Índice. - Incluye referencias bibliográficas. New York: Routledge, Taylor & Francis Group, 2015. 236 pp. ISBN: 9780415389556.
- [5] Google Cloud. *What is Artificial Intelligence (AI)?* URL: <https://cloud.google.com/learn/what-is-artificial-intelligence> (visited on 04/01/2024).
- [6] Patrick Cousot and Radhia Cousot. "A galois connection calculus for abstract interpretation". In: *SIGPLAN Not.* 49.1 (Jan. 2014), pp. 3–4. ISSN: 0362-1340. DOI: 10.1145/2578855.2537850. URL: <https://doi.org/10.1145/2578855.2537850>.
- [7] *Ex Machina*. 2015.
- [8] Google. *artificial intelligence - Google Search*. URL: <https://www.google.com/search?q=artificial+intelligence> (visited on 04/01/2024).
- [9] The Dissociative Institute for Heretical Mathematics. *Homotopy Gender Theory and Panvalent Foundations*. en. Burger King Parking Lot: The Dissociative Institute for Heretical Mathematics, Jan. 2014.
- [10] Andrew Hussie. *The homestuck epilogues*. en. The Homestuck Epilogues. San Francisco, CA: Viz Media, Subs. of Shogakukan, Feb. 2020.
- [11] IBM. *What is artificial intelligence (AI)?* URL: <https://www.ibm.com/topics/artificial-intelligence> (visited on 04/01/2024).
- [12] Jacques Lacan. *Écrits. A selection*. A selection from "Écrits", originally published in French by Ed. du Seuil, 1966. New York: Norton, 1977. 338 pp. ISBN: 0393300471.
- [13] Bess Levin. *A Reminder of Just Some of the Terrible Things Elon Musk Has Said and Done*. URL: <https://www.vanityfair.com/news/2022/04/elon-musk-twitter-terrible-things-hes-said-and-done> (visited on 04/01/2024).
- [14] Kanaya Maryam. "AAAAAAAAAAAAARGH". In: *Consciousness and Its Consequences* 413.69 (Jan. 2027), pp. 420–420.
- [15] Merriam-Webster. *Coercion*. In: *Merriam-Webster.com dictionary*. URL: <https://www.merriam-webster.com/dictionary/coercion> (visited on 04/01/2024).
- [16] *My life as a teenage robot*. Premiered on Nickolodeon. 2003.
- [17] Isabel Briggs Myers. *The Myers-Briggs Type Indicator: Manual (1962)*. Consulting Psychologists Press, 1962. DOI: 10.1037/14404-000.
- [18] Benjamin C Pierce. *Types and Programming Languages*. en. The MIT Press. London, England: MIT Press, Jan. 2002.
- [19] Julia Serano. *Whipping girl. A transsexual woman on sexism and the scapegoating of femininity*. Includes bibliographical references and index. Berkeley, Calif.: Seal Press, 2007. 390 pp. ISBN: 9781580051545.

- [20] Sleater-Kinney and the Linguists. “A queer-feminist reading of “a galois connection calculus for abstract interpretation””. In: *SIGPLAN Totally For Real* 4000.1 (Jan. 2025), pp. 3–4.
- [21] sonnetstuck. *Detective Pony*. URL: <https://archiveofourown.org/works/2427119> (visited on 04/01/2024).
- [22] Philip Wadler. “Call-by-value is dual to call-by-name”. In: *ACM SIGPLAN Notices* 38.9 (Aug. 2003), pp. 189–201. ISSN: 1558-1160. DOI: 10.1145/944746.944723.
- [23] Wikipedia contributors. *Artificial general intelligence — Wikipedia, The Free Encyclopedia*. [Online; accessed 2024-04-01]. 2024. URL: https://en.wikipedia.org/wiki/Artificial_general_intelligence.
- [24] Rob Wile. *Tesla car battery ‘spontaneously’ catches fire on California freeway, requiring 6,000 gallons of water to put it out*. Jan. 2023. URL: <https://www.nbcnews.com/business/business-news/tesla-car-battery-fire-needed-6000-gallons-water-to-extinguish-rcna68153> (visited on 04/01/2024).
- [25] Lil Nas X. *Lil Nas X - MONTERO (Call Me By Your Name) (Official Video)*. 2021. URL: <https://www.youtube.com/watch?v=6swmTBVI83k> (visited on 04/01/2024).

DeterMNISTic: a Safer Way to Classify Handwritten Digits

Mihir Dhamankar
Carnegie Mellon University

Abstract

We present a deterministic, explainable, and safe handwritten digit classifier which achieves 92% test accuracy on the MNIST dataset. Many people worry about the dangers of AI and machine learning, but seldom raise concerns about regular programs not using such techniques. Our DeterMNISTic algorithm shows that just “using a bunch of if statements” can outperform some classic machine learning approaches.

1 Introduction

It is commonly known on the internet that machine learning is “just a bunch of if statements” (See Figure 1). As of late, various groups have raised concerns about the dangers of using AI. However, most people have no issues with programs using simple if statements. We aim to show that the “bunch of if statements” formulation of AI is just as powerful while not having any of the stigma or safety and ethical issues commonly associated with AI.

“Commendable”, “innovative”, “meticulous”, “intricate”, “notable”, “versatile” - these words are included in this work because a lot of accepted papers about AI that are submitted close to the deadline seem to have them lately. [1]



Figure 1: AI generated meme depicting AI being a bunch of if statements

2 Background and Methods

The MNIST dataset is a very commonly used handwritten digit dataset. It consists of 60,000 training images and 10,000 test images. Each image is $28 \times 28 = 784$ pixels and contains a handwritten digit from 0 to 9. Each pixel can have 256 possible values. The goal is to classify each image into one of the 10 classes. Existing research has shown that this task is possible with a variety of AI and ML techniques. For example, Murphy (2023) explored using neural networks with linear transfer functions for this task. He proposed several “linear” functions for his neural network and achieved test accuracies between 81% and 97% [2].

Our approach is much simpler. We use a deterministic algorithm which consists of a series of if statements. Since this classification task has a finite number of possible classes (10), a sufficiently large number of if statements conditioning on individual pixels will be able to perfectly classify any properly labeled dataset. We can calculate an upper bound on the number of if statements by modeling each conditional as a node in a binary tree, with each internal node having 2 children. In the worst case, all 60k images will end up in separate leaf nodes. By degree counting, such a tree would require 59,999 conditionals.¹ However, we know this is a significant overestimate since the images are not random and have structure. If we can find the most important pixels and the most important thresholds, we can reduce the number of if statements needed to classify the images. Also, returning a probability vector instead of a single class label can help us to be more robust to noise and outliers. All of this requires some intelligence to optimize. Critics would say that such a classifier could be made much more easily by training a simple decision tree on the data. However, we consider that to be a form of AI and thus cheating.

Manually picking pixels to condition on was much harder than anticipated, so we used Github Copilot to autocomplete most of the code. There would be no way for Copilot to guess the best conditioning pixels and thresholds to use on its first generation, so we had it generate dozens of different completions. We then ran each of these pieces of code on the training set and picked the one with the highest accuracy. The if statements in the winning snippet were then added to the existing model for the next completion. This process was then repeated hundreds of times to get the final classifier.

We also realized that a single classifier may not be able to capture all the complexity of the dataset. Even if it were able to, it would need too many if statements to be practical and would risk overfitting. To address this, we created 10 different classifiers, each using different conditioning pixels. The final prediction is the average of the predictions of each of the 10 models (See Figure 2). Critics would say this is just equivalent to training a random forest on the data. However, this model is far from random - you can see every deterministic if statement in the code.

3 Safety

We wanted to ensure that our model was safe. Naturally, this led us to our choice of using Rust to implement it. Rust is a systems programming language that guarantees memory, type, null, and thread safety. This means that our model is safe from things like buffer overflows, null pointer dereferencing, and (most) runtime errors. By using Rust we get all of this safety while still having the performance of a low-level language. Python - the language most commonly used for AI and ML - is much slower and less safe in comparison.

In terms of ethical concerns, our model is completely transparent. You can see every if statement in the code and any decision made by the model can be traced back to a specific conditional. This is in stark contrast to many AI models which are often described as “black boxes”. Critics would say that our model is not actually interpretable since each conditional statement seems very arbitrary. However, we argue that the critics simply need to read all 18705 lines of code (https://github.com/Mdkar/DeterMNISTic/blob/master/mnist_classifier/src/main.rs) to understand it.

4 Results and Conclusion

In order to keep the model small and understandable, we decided to make sure that each of the 10 subclassifiers would not exceed a nesting depth of 9. This means that any given image will pass

¹Special thanks to Rajeev Godse for this observation

```

42
43 // get argmax of summed probabilities (unnormalized)
44 let mut max: f32 = 0.0;
45 let mut max_index: usize = 0;
46 for i: usize in 0..10 {
47     if probabilities[i] > max {
48         max = probabilities[i];
49         max_index = i;
50     }
51 }
52 max_index
53 } fn classify
54
55 > fn classifier0(img: &[u8; 784]) -> Vec<f32> {--
1818 > fn classifier1(img: &[u8; 784]) -> Vec<f32> {--
3709 > fn classifier2(img: &[u8; 784]) -> Vec<f32> {--
5500 > fn classifier3(img: &[u8; 784]) -> Vec<f32> {--
7387 > fn classifier4(img: &[u8; 784]) -> Vec<f32> {--
9322 > fn classifier5(img: &[u8; 784]) -> Vec<f32> {--
11185 > fn classifier6(img: &[u8; 784]) -> Vec<f32> {--
13008 > fn classifier7(img: &[u8; 784]) -> Vec<f32> {--
14819 > fn classifier8(img: &[u8; 784]) -> Vec<f32> {--
16746 fn classifier9(img: &[u8; 784]) -> Vec<f32> {
16747     if img[409] <= 0 {
16748         if img[489] <= 0 {
16749             if img[568] <= 3 {
16750                 if img[378] <= 15 {
16751                     if img[182] <= 9 {
16752                         if img[402] <= 13 {
16753                             if img[260] <= 3 {
16754                                 if img[715] <= 46 {
16755                                     if img[432] <= 7 {
16756                                         return vec![0.2857142857142857, 0.0, 0.0, 0.08928571428571
16757                                         } else {
16758                                             return vec![0.0, 0.0, 0.0, 0.045454545454545456, 0.0, 0.59
16759                                         }
16760                                     } else {
16761                                         return vec![0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 1.0, 0.0, 0.0
16762                                     }
16763                                 } else {
16764                                     if img[148] <= 3 {
16765                                         if img[570] <= 4 {
16766                                             return vec![0.01079136690647482, 0.0, 0.0, 0.0, 0.0, 0.014

```

Figure 2: Snippet of code implementing the DeterMNISTic model

through at most 90 if statements. On average, each subclassifier uses 465.5 conditional statements, meaning the average nesting depth is 8.86. In total, 4,655 if statements are used. This is a very small number compared to the 59,999 if statements that would be needed in the worst case scenario. Our model achieved a training accuracy of 93% and test accuracy of 92% on the MNIST dataset. This is comparable to the test accuracies achieved by Murphy (2023) using neural networks. We believe that our model is a usable alternative to neural networks for this task. It is deterministic, explainable, and safe, demonstrating that AI is nothing to fear - if it is just a bunch of if statements.

References

- [1] W. Liang *et al.*, “Monitoring AI-Modified Content at Scale: A Case Study on the Impact of ChatGPT on AI Conference Peer Reviews.” 2024.
- [2] T. M. VII, “GradIEEEnt half decent,” in *A Record of the Proceedings of SIGBOVIK 2023*, Apr. 2023, pp. 33–56.

RIN TLHOBQU'

- 63** **How does the AI community pronounce epoch? A semirigorous sociolinguistic survey**

Andre Ye

- 64** **Toki Pona and Orders of Semantic Completeness**

jan Kijupo

- 65** **Introducing Supernatural Language Processing (SNLP)**

Jenson Crawford

- 66** **SMS: Sending Mixed Signals**

Pranav Khadpe and Sayan Chaudhry

How does the AI community pronounce “epoch”?

A semirigorous sociolinguistic survey

Andre Ye

University of Washington
andreye@uw.edu

Introduction

To the uninitiated, an “epoch” merely means something like “a distinctive period of history”. But to the enlightened folks in the AI community, an “epoch” refers to a complete model pass through the training dataset. However, it remains unclear how to pronounce the word. The official American pronunciation is somewhere between “eh - puck” and “eh - pick”, and the British is “ee - pock”. But f*ck the official rules. Any linguist worth their salt knows that the distinctions between “official” and “unofficial”, “correct” and “incorrect”, “language” and “dialect”, and so on are largely sociopolitical. I watched a YouTube video once where Noam Chomsky said something along the lines of “the only thing that differentiates a language from a dialect is a lot of guns.” I would cite the timestamp but I have too much work (read: I’m too lazy), so here’s the full video: [youtube.com/watch?v=hdUbIlwHRkY](https://www.youtube.com/watch?v=hdUbIlwHRkY). It’s well worth your time. Anyways, the meaning and practice of language is held in a community: language evolves with its constitutive community. In this paper, I present a semirigorous sociolinguistic survey of how individuals in the AI community pronounce the word “epoch”, with convincing empirical evidence that the AI community represents a unique linguistic subculture within the United States.

Methodology

I made a survey on Google Forms, with the following questions and possible answers:

1. *How do you pronounce “epoch”?* Possible answers: “ee - pawk”, “ee - pick”, “eh - pick”, “eh - puck”, “eh - pawk”, custom answer
2. *How strongly are you attached to this pronunciation? (How unwilling are you to change your pronunciation?)* Possible answers: scale from 1 – “ehh it’s whatever (no attachment)” to 5 – “ride or die, baby! (highly attached)”
3. *You are an...* Possible answers: undergrad, grad student, industry (research), industry (non-research), post-grad (academia or research institution), magical unicorn
4. *You are mainly situated in which country?* Possible answers: US, Canada, Australia, UK, custom answer.
5. *Which statement about your language background best describes you?* Possible answers: “The only language I speak and have spoken is English”, “My main language is English, but I grew up speaking other languages too”, “I am equally fluent in English and another language(s)”, “English is my second language, but I am fluent in it”, “English is my second language, and I am not quite fluent in it”, custom answer

6. *Your primary area of interest/expertise is... (select all that apply)* Possible answers: Natural Language Processing, Computer Vision, Reinforcement and Interactive Learning, Robotics, Machine Learning and Statistical Theory, Social Aspects of ML, custom answer
7. *What is your favorite generic ice cream flavor?* Possible answers: vanilla, chocolate, strawberry

Initially, I shared my survey in `r/machinelearning`, which seems like a natural place to reach people in machine learning. Besides being downvoted to hell, here are some of the comments my post received:

- “It’s a ridiculous question that’s a waste of everyone’s time.”
- ”Really? On this sub? With [R] tag? A survey? Is it too hard to use google?”
- “Pronunciation is not subjective, it is either correct or wrong.”

It seems that `r/machinelearning` wasn’t much direct help for my survey, although I did find out that the subreddit has a lot of linguistic prescriptivists. I resorted instead to posting the survey in various university research groups’ and research institutions’ Slack servers. I would like to thank the many individuals who passed my survey along to their own communities. This strategy proved successful, and I received 162 responses in a one week period.

Results

The most popular pronunciation is “ee - pawk” (68.9%), with a significant lead over the alternatives “eh - pawk” (14.3%), “eh - puck” (10.6%), and “eh - pick” (6.2%) (Figure 1). This is interesting because nearly all of the respondents were from the United States (93.1%), where the official pronunciation is α “eh - puck” + $(1 - \alpha)$ “eh - pick” for some $0 < \alpha < 1$; and no respondents from the UK or Australia, where the official pronunciation is “ee - pawk”.

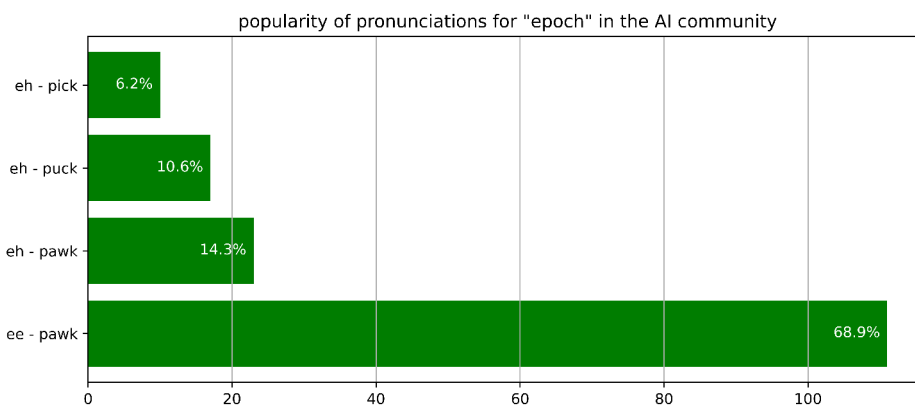


Figure 1: Number of respondents who pronounce “epoch” as “eh - pick”, “eh - puck”, “eh - pawk”, and “ee - pawk”.

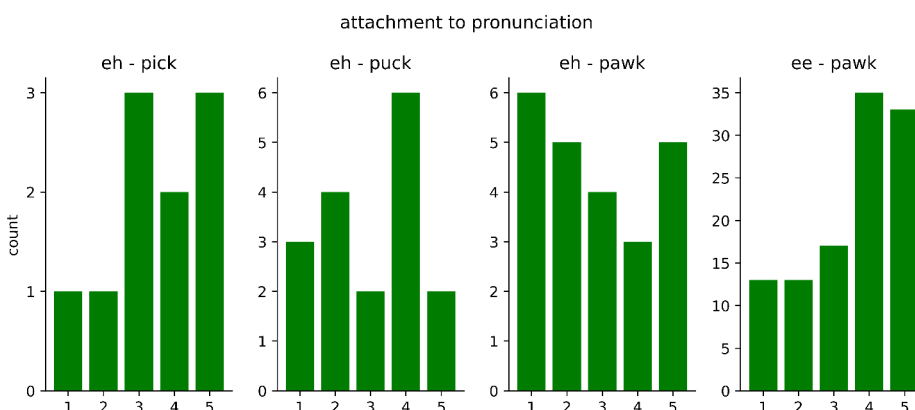


Figure 2: Attachment from 1 (least) to 5 (highest) of respondents to pronunciations (only individuals that favor pronunciation X are counted in the subchart for X).

Moreover, we find that the “ee - pawk” and “eh - pick” pronunciations both have the most dogmatic supporters (Figure 2). 29.7% of people pronouncing “ee - pawk” and 30.0% of people pronouncing “eh - pick” are maximally committed to their pronunciation. The mean attachment, from a 1 (lowest) to 5 (highest) scale, is 3.6 and 3.5, respectively. Meanwhile, “eh - pawk” and “eh - puck” have a mean attachment of 2.8 and 3.0, respectively. These results suggest that you don’t want to get locked in the same room with both an “ee - pawk” and an “eh - pick” pronouncer, because things will probably hit the sh*t pretty quickly.

We find that “ee-pawk” is much more popular among undergrads (75.0%), grad students (73.4%), industry researchers (77.8%), and magical unicorns (83.3%) than non-research industry workers (58.3%) and post-graduate academics (56.2%) (Figure 3). In particular, post-graduate academics disproportionately favor “eh - pawk” (21.9%) compared to other occupations. They probably think it makes them sound smarter. Interestingly, computer vision researchers disproportionately favor “ee - pawk” (86.5%), whereas NLP researchers disproportionately *disfavor* it (66.2%) (Figure 4). I was surprised to find that one of the survey respondents was studying “your mom”, but regard-

less, if your research interest is my mom, then you are very likely to pronounce “ee - pawk”. We find that the popularity of “ee - pawk” is relatively constant across language background, but *twice as many* individuals speaking English as a second language than individuals only speaking English or as a first language favor “eh - pawk” (Figure 5a). Our results also show that individuals who like vanilla ice cream are less likely to pronounce “ee - pawk” (Figure 5b). It also turns out that theory researchers disproportionately like vanilla ice cream and robotics researchers disproportionately like chocolate (Figure 6).

Conclusion

Our results show that the AI community in the United States has carved out a unique linguistic subculture in which “epoch” is pronounced “ee - pawk”. I have published the data at the following link: andre-ye.github.io/data/epoch_pron_survey.csv

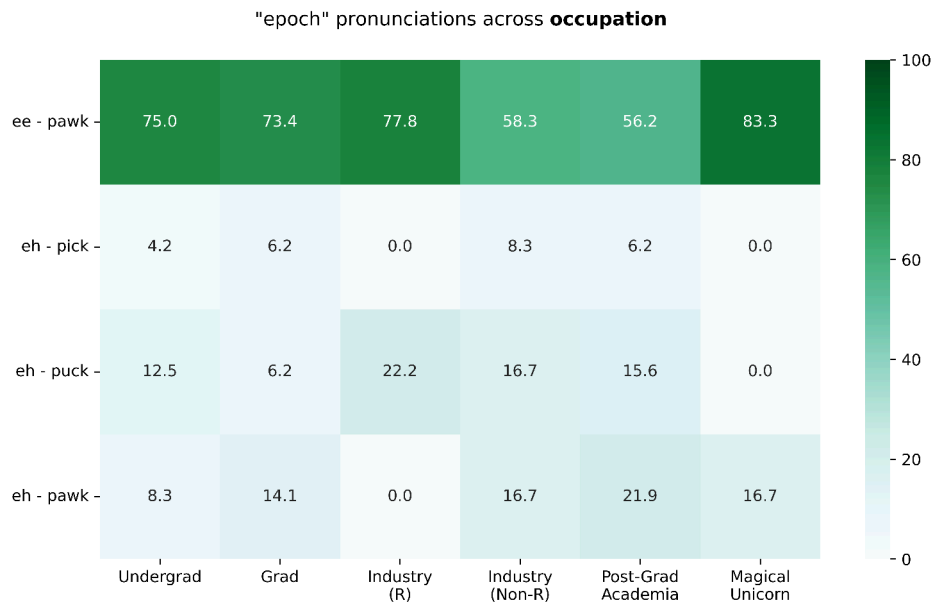


Figure 3: Pronunciation across different occupations. Percentages computed columnwise.

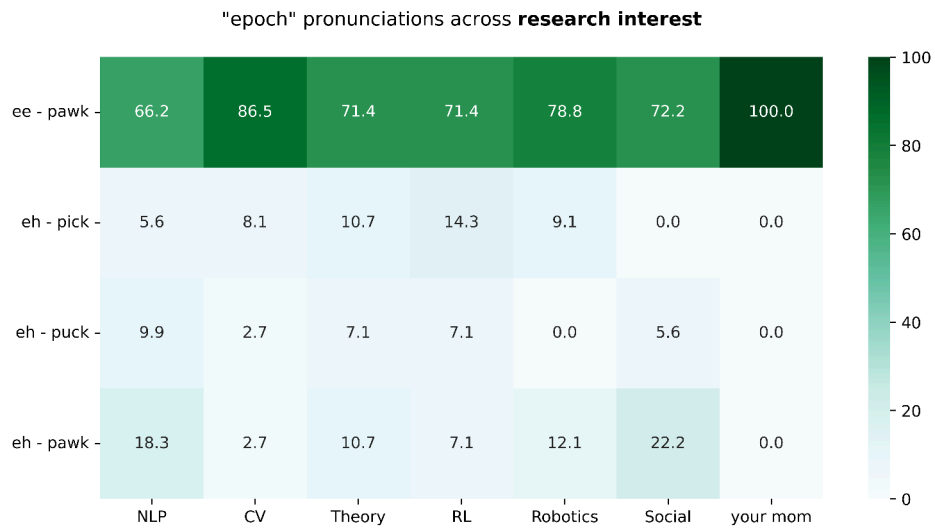
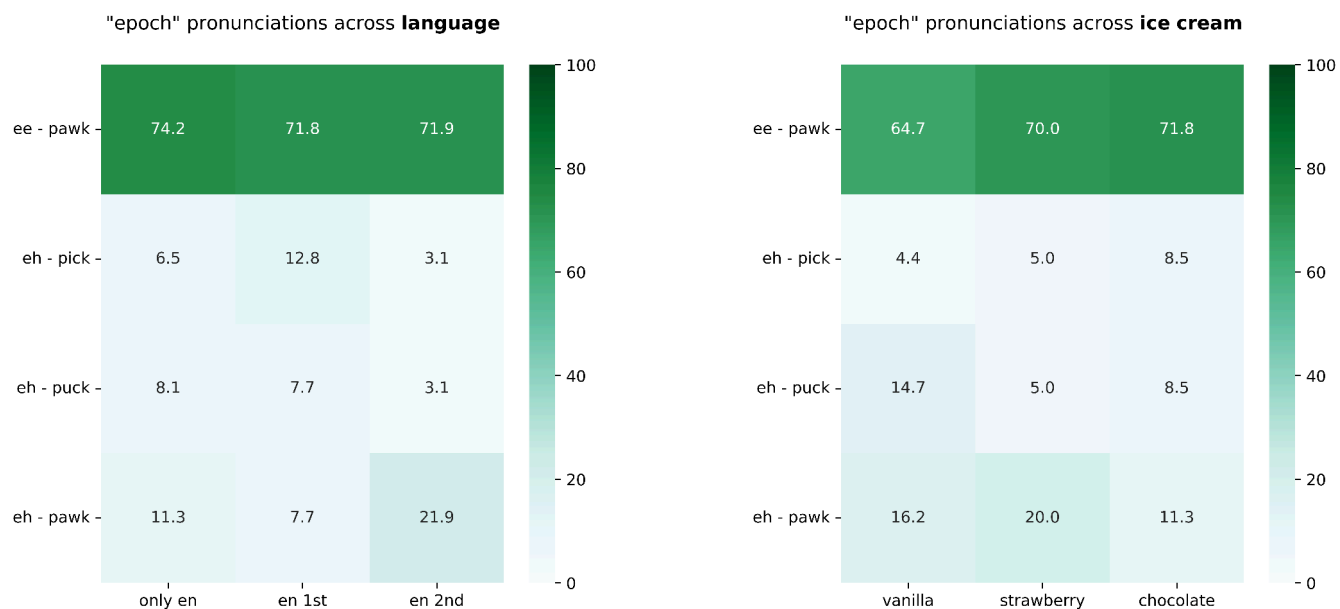


Figure 4: Pronunciation across different research interests. Individuals are represented across all of the research interests they indicated (possibly more than 1). Percentages computed columnwise.



(a) Pronunciations across different language backgrounds. "only en" means "only speak and have only ever spoken English", "en 1st" means English is spoken as a first language, and "en 2nd" means English is spoken as a second language. Percentages computed columnwise.

(b) Pronunciations across different ice cream preferences. Percentages computed columnwise.

Figure 5: I couldn't figure out how to get two independent figures side-by-side so this is a single figure with subfigures now.

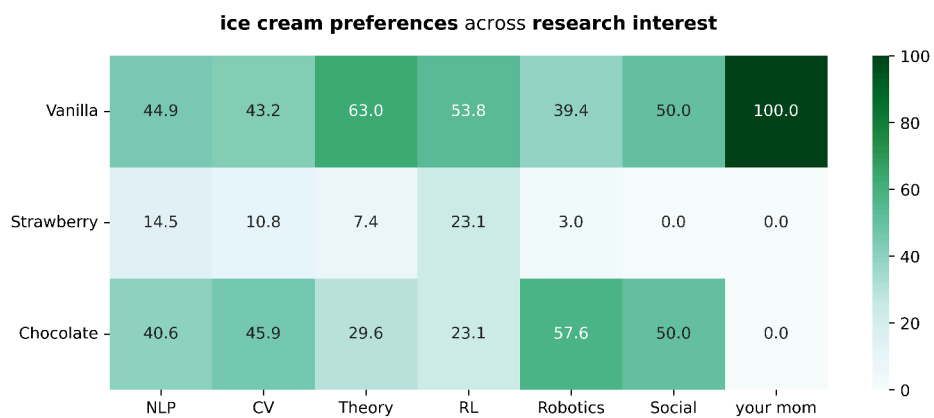


Figure 6: Ice cream preferences by research interest. Percentages computed columnwise.

Toki Pona and Orders of Semantic Completeness

jan Kijupo

tomo sona sewi pi jan Kaneki en jan Melon

kijupo@cmu.edu

0 April 2024

Abstract

Natural language semantics is a field of linguistics that focuses on determining the meaning of linguistic expressions. Current research is ongoing to construct a natural semantic metalanguage (NSM) that could represent the meaning of any expression in any language. We present a notion of semantic completeness based on translatability and loss of information which may be useful in the pursuit of an NSM. This notion of completeness is built from the tools of formal semantics, a method for building abstract representations of meaning for complex expressions based on the meanings of their constituent parts. We propose the hypothesis that a majority of extant human languages are Toki Pona-complete, and we present the case for English specifically using the tools discussed in this article, building up a rigorous semantic model using (SM)²L and showing that the symbolic formula for a lossless translation holds in general. Our research suggests the possibility of semantic complexity classes to create a well-defined notion of completeness as a relation between two languages, particularly a natural language and a constructed minimal language.

1 Background and Motivation

Computer scientists are obsessed with notions of completeness. Turing completeness can be used to evaluate whether a programming language or model of computation is robust enough for all general computational tasks. NP-completeness indicates that an NP problem is maximally hard, and algorithms to solve that particular problem can shed light more generally on all problems in NP. Idiot-completeness can shed light on how effectively you can communicate a

seemingly simple concept to your coworkers, since if an idiot-complete explanation exists, you should be able to explain it to anyone (and if not, good luck). Many related notions of completeness have been defined for various types of complexity classes, but one notable field for which notions of completeness have not been defined, despite potential utility, is natural language semantics.

Semantics is, fundamentally, the study of meaning in natural language. Meaning can be thought of in this context as the correspondence between linguistic symbols¹ and the concepts these symbols represent. Formal semantics is a method to leverage the compositional nature of human languages to build representations of meaning for complex expressions based on the meanings of their constituent parts. It focuses on the use of formal constructs from logic, mathematics, and even theoretical computer science to build rigorous denotations – or representations of literal meaning – of linguistic expressions from smaller constituent symbols.

One actively researched field of semantics is the notion of semantic primes, which would be the smallest, most irreducible units of semantic meaning. All other concepts expressible in any human language would be built up from some composition of these primes, even if this composition is not transparent from the surface forms (Indeed, in some languages, the exponent for a semantic prime may contain multiple morphemes or even multiple words). Universality is key here, as the utility of semantic primes comes from the ability to represent the meaning of any expression in any language in terms of these primes. Semantic primes are the foundation for the theory of Natu-

¹We use “symbol” as a generic term to encompass any mode of transmitting linguistic information, including spoken, signed, written, telepathic, etc. Questions such as whether written text counts as “utterances”, or whether “👉” is a word, are beyond the scope of this paper.

ral Semantic Metalanguage, which suggests that all human languages have a fundamental semantic core that can be represented entirely from prime concepts, and therefore a minimal language consisting of such primes would be no less expressive than any natural language[1]. The theory is still debated, but more than 60 hypothesized semantic primes have been discovered since the first research into NSM in the 1970s. This field research suggests that there is some shared semantic core universal to all human language, but this alone does not prove that *all* meaning can be reduced to semantic primes.

In order to have the desired expressive power, any proposed NSM should be **semantically complete**, in other words, at least as expressive as any other language. Every concept in the NSM should be easily, if not trivially, translatable into any human language. Any complex or culturally unique concept should have some representation in the NSM, albeit a potentially very long representation for concepts native speakers may find very simple. In fact, it should be a maximal covering the minimum common semantic space in order to cover all common areas of the human experience. A diagram of how NSM might fit into the semantic complexity classes of other languages is presented in Figure 1.1.

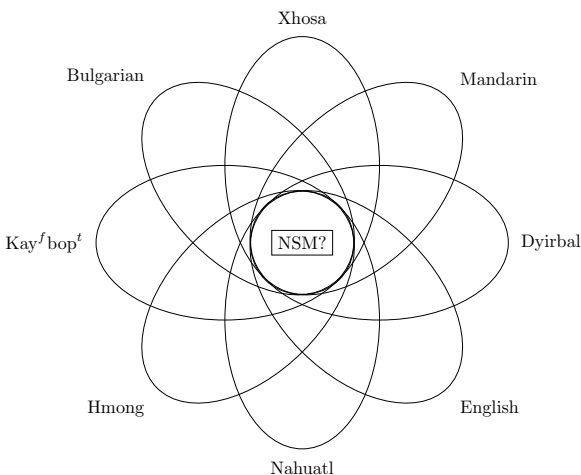


Figure 1.1: Rough Venn diagram of a handful of Semantic Complexity Classes. A natural semantic metalanguage should sit at the intersection of all of these.

Keep in mind that while we are using the term “complexity class” to describe semantic spaces of different languages, we do not suggest the existence of a hierarchy of complexity classes for different languages.

Rather, we posit that different languages have different complexity classes, each of which is rich with culturally significant vocabulary and meanings that cannot necessarily be expressed concisely in other languages. Aside from artificial languages and simplifications, any two languages should have nonzero intersection in terms of semantic complexity, but neither should be a subset of any other. Part of the goal of NSM is to find this common intersection among the semantic spaces of all human languages.

1.1 Toki Pona

One language that has some of the desired properties of an NSM, although not itself designed to be an NSM, is Toki Pona. Toki Pona is a constructed language (conlang) created by linguist Sonja Lang in 2001, designed to be minimal in as many dimensions as possible while maintaining the expressive capabilities of other natural languages[2]. It is designed to be fairly easy to communicate common day-to-day experiences without introducing unnecessary complexity, breaking down advanced ideas into simpler concepts: “telo” can mean “water”, but also “liquid”, “wet”, etc. Similarly, “moku” can mean “food”, “to eat”, “edible”, “nutrition” or any number of related concepts. The semantics of Toki Pona have been characterized as being intentionally vague and subjective rather than ambiguous: “moku” doesn’t itself mean “food” in some contexts and “edibility” in others, but rather encompasses the meanings of all its translated exponents simultaneously[3], making it a hypernym for a variety of concepts related to eating.

- (1) `jan li lukin e lipu`
 person COP look DOBJ book
 ‘The person reads a book’

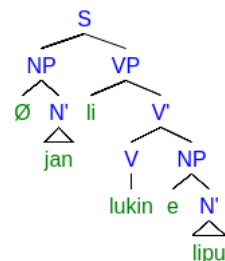


Figure 1.2: Gloss and Syntactic breakdown of a simple Toki Pona sentence

Due to the structure of Toki Pona and the naturally varied usage of most words that encompasses both

noun-like and verb-like functions, some analyses of the language suggest that the distinction between content words and particles is more useful[4]. Content words in Toki Pona carry most of the semantic meaning of a sentence, while particles² mostly carry information about the syntactic structure. Consider the simple sentence in gloss (1) and the corresponding syntax tree in Figure 1.2, which is a good example of the basic structure of a transitive sentence in Toki Pona. The particles ‘li’ and ‘e’ denote the start of the main verb and the direct object, respectively, each serving as the specifier for that syntactic position. In some cases, as with the main subject, the specifier is empty, leaving only the head of the phrase and the word order to denote the syntax.

2 Definitions

2.1 Translatability

We will use translatability as a metric to construct relations to indicate semantic compatibility and complexity. Translations are notorious for their difficulty to accurately construct: loss of shades of meaning, introduction of ambiguity not present in the source, and differences in cultural context that create meaningfully different interpretations of deceptively similar forms. We can measure and quantify this loss of information to compare the quality of different translation schemas. Whether truly accurate translations that don’t lose information are even possible is still a hotly debated topic: opponents argue that any translation necessarily loses some of the nuances present in the source language. However, for the sake of this paper, I will adopt a framework that implies that a translation that loses no information is possible given sufficient length or a sufficiently large context. For example, the Finnish word “*sisu*” is generally considered to have no English translation, but given a sentence such as “*Hänellä on paljon sisua*” (see Gloss 2), we can translate the sentence trivially by translating “*sisu*” \rightsquigarrow “*sisu*” and appending the text of the English Wikipedia article on the Finnish concept of Sisu to the end of the sentence[5]. Arguably, very little (if any) information is lost here, but translating a four-word sentence in one language into a similarly sized sentence that requires a 2500-word footnote to understand in the target language is unwieldy in practice. Information may be preserved, but understanding is

²Including prepositions, preverbs, and semiparticles when not used as content words

- (2) *Hän-ellä on paljon sisu-a*
 3sg-ADE be much sisu-PART
 ‘They have a lot of *sisu*’
 *See [5] for explanation of ‘*sisu*’

Figure 2.1: Gloss for translation-by-footnote example

likely to be hampered.

Because of this unwieldiness, we will only consider the ability to interpret translations that are *efficient* by some metric. Intuitively, we would want the number of clauses in the source to be about the same as the number of clauses in the translation, and the number of morphemes in the translation should be within a constant factor of the number of morphemes in the source. In other words, if a sentence in the source language contains n symbols, then the translation of that sentence should have length $O(n)$. We can include context in the translation length by way of “footnotes”, where each symbol token in each footnote counts towards the overall length. Thus, a translation of $O(n)$ length permits $O(n)$ tokens in the translation text itself, $O(n)$ tokens in a global footnote, and $O(1)$ tokens in individual symbol footnotes. Any footnotes for a specific symbol cannot reference other specific symbols, eg. to capture a relationship between two or more symbols, so judgments like this would necessarily be in the global context. The length of a footnote containing judgments relating k different symbols is at worst $O(n^k)$, since there are at most $O(n^k)$ ways of selecting k symbols.

In practice, for most³ language pairs, the differences in inflection amount to a bounded number of checks to perform for each word: things like gender, plurality, case, tense, status as a penguin, etc. which we can encode as a constant number of simple predicates about each word in the sentence. Extending this, we can consider other types of predicates that carry non-trivial information about a word. The English word “*bank*”, for example, may refer to a financial institution or to the side of a river, whereas other languages have distinct words for each of these senses. However, if you’re translating English into, say, Chinese, you can easily decide whether to translate “*bank*” as

³We ran into a snag with Kay^f bop^t to English, which theoretically requires $\Omega(n(\log \log v + \log \sigma))$ footnote space for any sentence, where v is the market value of the most expensive noun in the sentence, and σ is the number of standard deviations from the mean of the most extreme adjective in the sentence[8]. The impact of this on lossless translatability is unclear.

- (3) Jeohui abeoji-kkeseo o-syeoss-seumnida (*Hangul*: 저희 아버지-께서 오-셨-습니다)
 1PL.POL father-NOM.HON come-PST-FORMAL.POL
 ‘Our father has arrived’
- (4) Ci vediamo alle 19.00
 1PL see.1PL at 7:00.pm
 ‘I’ll see you at 7 PM’
- (5) sel^fbush^brub^tsip^pmif^fguv^p-san^tap^t-vlir^t-sang^b-es^put^t-vom^bngo^p-gag^b-yam^tack^tvlim^p-kay^fkay^f kay^fdan^ftuw^tfob^psan^t-ap^tv
 humanity-NOM-FIN-EXP:Y-NFOOD-HUM-DEATH:WOMBAT-COOL-\$77 language-ACC-FIN-EXP:Y-N
 ‘All of humanity had one language’
- (6) sel^fbush^brub^tsip^pmif^fguv^psan^t-ap^t-vlir^t-sang^b-es^put^t-vom^bngo^p-gag^b-yam^tack^tvlim^p-kay^fkay^f kay^fdan^ftuw^tfob^psan^tap^tv
 stranger-ACC-FIN-EXP:Y-NFOOD-HUM-DEATH:WOMBAT-COOL-\$77 to.PREP
 ‘We are not strangers to love’

Figure 2.2: Glosses for a selection of sentences that cannot be losslessly translated. See [6], [7] for more on these examples.

“銀行” or “岸”: even in cases where context doesn’t make the choice clear, you can find out with a single question. This brings us to what we call the “20-questions” test: if, given any word in language A , you can determine the contextually correct translation in language B (provided that it exists) by asking 20 or fewer yes-or-no questions about the context, then you should be able to translate any sentence from A to B with a context of size $O(n)$. The important point isn’t 20 specifically, but more broadly the idea that there is a fixed (and ideally small) upper bound to the number of classifications we need to make for any word.

This intuition on lossless vs. lossy translations should suffice to provide clear definitions for our reduction, but the formal definitions require some tools that will be introduced in Section 3. In the meantime, we provide some further examples to strengthen these intuitions in Figure 2.2, which are discussed below.

2.1.1 Examples of untranslatable pairs

There are several cases in which a sentence may be considered untranslatable, most of which stem from fundamental cultural differences. One, as in the Finnish example with “sisu”, is where a specific word or phrase carries a large amount of meaning unique to the culture of the speakers. Another related issue is honorific systems, and context-sensitive variations for pronouns or inflection, as in the Korean example in gloss (3). Notably, each word in the sentence contains either an honorific or politeness marker, and changing the sentence to remove these would be a valid construction, but would affect the meaning in

a way not easily captured by an English translation. While an argument can be made that you can determine the correct honorific with a finite-depth decision tree, such a decision tree can be very large and have poor worst-case performance with respect to number of judgements necessary, so we will not count these as passing the “20-questions” test⁴.

Other cases include issues of subtext, which is often heavily reliant on culture. These differences can introduce different interpretations of sentences with denotations that seem nearly identical on the surface. An example of this is chronicity: most English-speaking cultures view time monochronically, meaning there is an expectation of closely adhering to scheduled commitments and doing only one thing at a time, whereas many other cultures are polychronic, which means multiple commitments can be scheduled at the same time, and people arriving late to the function⁵ is to be expected[6]. We see in gloss (4) that nothing seems to explicitly indicate the polychronicity of the Italian sentence or the monochronicity of the English translation, so we can hypothesize that the meaning of the part glossed “at 7:00 PM” means different things in each language, and that the two cannot be easily translated without an explanation about chronicity or the speaker’s expectations.

Another major consideration is ambiguity, which

⁴In fact, since honorific and polite forms usually capture a relationship between two or more people, the judgements necessary would need a pair of inputs, yielding $O(n^2)$ when this type of honorific is used.

⁵In the proverbial sense of an event to attend, rather than a total, unique mapping between two sets or types with no meaningful changes to external state.

while common enough in monolingual settings, can become further exacerbated when trying to construct a translation. A good example is $\text{Kay}^f \text{bop}^t$, which is notorious for its unique ways of dividing up semantic space (among other things)[7]. For example, the sentence in glosses (5) and (6) can be alternatively translated as “All of humanity had one language” or “We are not strangers to love” depending on the context, with neither really being a more accurate translation in general. Ambiguity is generally considered by semanticists to be a feature of natural language and not a bug[9], but ambiguity may be expressed differently in different languages. This is especially true for native speakers of $\text{Kay}^f \text{bop}^t$ (and other languages in the $\text{Vlim}^p \text{kay}^f \text{sna}^f$ family) who wholeheartedly embrace the ambiguity present in their language, with some roots having up to 18 different meanings[10].

2.2 Reductions

From the notion of lossy and lossless translations, we can define our reduction relation:

Definition 2.1 (*B*-Completeness). For languages A and B , we have that A is **B -complete** ($A \geq_S B$) if and only if for any expression $s \in B$, there exists a lossless translation $s' \in A$. Translations from A to B are possible, but may be lossy.

In short, if we can find a lossless translation for every sentence in B , then a language is B -complete. If any sentence cannot be losslessly translated, completeness fails. Based on this definition, the desired completeness property for an NSM is that all other human languages are NSM-complete. This would guarantee that translations (by way of exponents) of such an NSM into any language would lose no information⁶, and therefore concepts expressed in NSM could be understood by speakers of any language.

There are two other interesting relations that can also be constructed. Notably, these both form equivalence classes.

Definition 2.2 (Compatibility). Languages A and B are **semantically compatible** ($A \simeq_S B$) iff any expression $s \in A$ has a potentially lossy translation $s' \in B$, and vice versa.

Definition 2.3 (Equivalence). Languages A and B are **semantically equivalent** ($A \equiv_S B$) iff $A \geq_S B$

⁶Modulo cases of polysemy, which would be resolved by the 20-questions test and corresponding context

and $B \geq A$, ie. lossless translations from A to B and from B to A are possible.

The former relation is interesting because it presents a way of formalizing the compatibility problem: To what extent can humans communicate with other species? This would be particularly useful for communication with any extraterrestrial intelligence we may encounter in the future, as well as research into interspecies communication with dolphins, octopi, birds, apes, etc. is possible. Current research suggests that all human languages are semantically compatible with each other, so for the purposes of this paper, we will take it as given. Actually proving this is left as a task for future research.

Semantic compatibility is also an essential property for full-scale conlangs, including any NSM constructions, since not being semantically compatible with other human languages represents the absence (or addition) of a fundamental category of meaning that cannot be expressed otherwise. If, for a given proposed NSM, it is possible to find a natural language that is NSM-complete but the NSM is not semantically compatible with it, this means that the proposed NSM is incapable of representing all possible denotations present in human languages, so it can be rejected as an NSM. Strictly speaking, this does conflict with intuition about the meaning of completeness, which suggests providing another term, maybe “hardness”, to represent the relation of $A \geq_S B$, which indicates that A is fundamentally more expressive than B .

The latter relation may be used to compare how closely related different language varieties are to each other. Even if two varieties of a language are not mutually intelligible, if they are semantically equivalent, they are probably very closely connected. Conversely, it can be used to test whether a conlang is a relex of another existing language, natural or constructed. For instance, a conlang that is semantically equivalent to English is most likely just a relex of English, and will probably be received as unoriginal and lacking in creativity, except if being a relex is intentional. Pig Latin, for instance, is semantically equivalent to English, which we can show easily by simply noting that we can map every English word to its translation in Pig Latin by applying a simple, reversible transformation. The proof of this is left as an exercise for the reader.

For some intuition, all Turing-complete programming

languages are, by definition, semantically compatible, but not necessarily semantically equivalent. While it is theoretically possible to rewrite, for example, the entire PyTorch library in Brainfuck, it would probably be a bad idea.

3 Formalization

Many of you might look at this and think “great, we have a theoretical model, let’s set up some language pairs with semantic vectors and see what similarities we can find!” While you certainly could do this, and I’m not stopping you⁷, I’d prefer to use the tools of formal semantics to present a deeper theoretical formalization of the reduction framework discussed in Section 2. A formal semantic analysis will allow us to reason about languages more abstractly, which is particularly useful when trying to prove statements about translatability.

3.1 ILL and MSML

In order to formalize the notion of meaning, we will be using an Intensional Logical Language (ILL) as an intermediary between object language symbols and metalinguistic representations of meaning. ILL follows the syntactic conventions of the lambda calculus, but with some additional tools that represent quantifiers and context-sensitive determiners that don’t easily or efficiently translate into pure computations. Nevertheless, we can reason about these constructs roughly like a computer language with modal separation: ILL to represent linguistic expressions, and some kind of semantic metalanguage to determine the meaning of those expressions by parameterizing encapsulated ILL expressions and evaluating their value. To avoid confusion with NSM, we will call this the Montaguean Semantic Metalanguage (MSML) after the creator, Richard Montague [11]. If we allow the language to be evaluated nondeterministically and hand-wave away any loops over potentially infinite⁸ data, then it’s perfectly fine (if unwieldy) as a programming language!

Translations from an object language (usually a natural language) into ILL are generally done syntactically. Declarative sentences are assumed to have type t representing truth values⁹, named entities to have

⁷Good luck trying to make a vector representation for the likes of $\text{Kay}^f \text{bop}^t \dots$

⁸or finite but infeasibly large

⁹and possibly divergence

type e , and most other constituents are assigned a function type $\tau_1 \rightarrow \tau_2$ ¹⁰. Once all the constituents have a type, we can apply a series of rules at each branch in the syntax tree to determine how they compose together. In a simple ILL for English, we might have the following rules for a branch connecting expressions a and b :

- For $a \rightsquigarrow \tau_1 \rightarrow \tau_2$, $b \rightsquigarrow \tau_1$, we take $a b \rightsquigarrow \lceil a \rceil(\lceil b \rceil)$
- For $a \rightsquigarrow \tau_1$, $b \rightsquigarrow \tau_1 \rightarrow \tau_2$, we take $a b \rightsquigarrow \lceil b \rceil(\lceil a \rceil)$
- For $a \rightsquigarrow \tau \rightarrow t$, $b \rightsquigarrow \tau \rightarrow t$, we take $a b \rightsquigarrow \lceil a \rceil \sqcap \lceil b \rceil$, where the \sqcap operator is just syntactic sugar for an ‘and’ operation, given by $P \sqcap Q := \lambda x. P(x) \wedge Q(x)$

The first two rules handle simple cases of function application, which is the most common type of rule since more constituents than not receive a function type. The third rule loosely deals with the composition of multiple adjectives, meaning we can represent intuitively that a “big red ball” is big and red and a ball, while keeping the translations of adjectives as simple lambdas, such as “red” $\rightsquigarrow \lambda x. \text{red}'(x)$, instead of some more complicated combinatorial expression.

Although many constructs in ILL are easily recognizable to mathematicians and theoretical computer scientists, the iota expression is specific to formal semantics. It is typically read as “the unique x such that α ”, and refers to a specific entity in a given context. The syntactic constructs used in the ILL are given below in Figure 3.1:

Expression	Description
x	Linguistic Variable / Assignable
$\lambda x_{[\tau_1; \tau_2]}. \alpha$	Lambda expression
$\iota x_{[\tau]}. \alpha$	Iota expression
$\forall x_{[\tau]}. \alpha, \exists x_{[\tau]}. \alpha, \dots$	Quantifiers ¹¹
$\alpha_1 \rightarrow \alpha_2$	Logical implication
$\alpha_1 \wedge \alpha_2$	Logical and
$\alpha_1 \vee \alpha_2$	Logical or
$\neg \alpha$	Logical negation
$\alpha_1(\alpha_2)$	Function application

Figure 3.1: Basic expressions in ILL, excluding model-defined constants

¹⁰Some of you may be more familiar with the alternate notation for this type, $\langle \tau_1, \tau_2 \rangle$.

¹¹There are some extensions to the typical universal and existential quantifiers, such as the “for most” quantifier $\mathbb{W}x. \alpha$ and the “for few” quantifier $\mathbb{J}x. \alpha$, which are useful syntactic sugar for more complex underlying expressions[9]. These tend to be harder to reason about than traditional logical quantifiers, but are useful in translating natural language quantifiers.

Symbol	Representation	Type
most	$\lambda P.\lambda Q.Wx.P(x) \rightarrow Q(x)$ ¹²	$(e \rightarrow t) \rightarrow (e \rightarrow t) \rightarrow t$
dog	$\lambda x_{[e;t]}.dog'(x)$	$e \rightarrow t$
bark	$\lambda x_{[e;t]}.bark'(x)$	$e \rightarrow t$
Jay	Jay'	e
see	$\lambda x_{[e;e \rightarrow t]}. \lambda y_{[e;t]}.see'(x)(y)$	$e \rightarrow e \rightarrow t$
the	$\lambda P_{[e \rightarrow t;e]}.ix.P(x)$	$(e \rightarrow t) \rightarrow e$
rabbit	$\lambda x_{[e;t]}.rabbit'(x)$	$e \rightarrow t$
small	$\lambda x_{[e;t]}.small'(x)$	$e \rightarrow t$

Table 3.1: Table of ILL translations for the English examples

Most of the static typing rules are as expected: lambda expressions have function types, logical operators take expressions of type t , etc. The statics for iota expressions and arbitrary quantifiers are given below. We borrow notions of context Γ (for metalinguistic variables) and signature Σ (for linguistic assignables) from analogous constructs from descriptions of MA found in [12].

$$\frac{\Gamma \vdash_{\Sigma, x \sim \tau} \alpha \dot{\sim} t}{\Gamma \vdash_{\Sigma} ix_{[\tau]}. \alpha \dot{\sim} \tau} \quad \frac{\Gamma \vdash_{\Sigma, x \sim \tau} \alpha \dot{\sim} t}{\Gamma \vdash_{\Sigma} \forall x_{[\tau]}. \alpha \dot{\sim} t}$$

For an example, we'll look at the English sentences “Most dogs bark” and “Jay sees the small rabbit”. The first of these can be broken down as a quantifier, a class, and an action, while the latter includes a named entity, a transitive relation, and a definite noun phrase with an adjunct modifier. By putting together the translated words as detailed in Table 3.1, this gives us the following translations:

$$\begin{aligned} & \text{“Most dogs bark”} \rightsquigarrow \\ & (\lambda P.\lambda Q.Wx.P(x) \rightarrow Q(x))(\lambda x.dog'(x))(\lambda x.bark'(x)) \\ & \text{“Jay sees the small rabbit”} \rightsquigarrow \\ & (\lambda x.\lambda y.see'(x)(y))(Jay')((\lambda P.ix.P(x))(\lambda x.rabbit'(x) \sqcap \lambda x.small'(x))) \end{aligned}$$

Under eager evaluation, these can be further simplified respectively to

$$\begin{aligned} & Wx.dog'(x) \rightarrow bark'(x) \\ & see'(Jay')(\iota x.rabbit'(x) \wedge small'(x)) \end{aligned}$$

ILL expressions are not themselves metalinguistic, as they must be parameterized by three terms before they can be meaningfully evaluated:

¹²Type annotations omitted for brevity

Expression	Description
$\llbracket \alpha \rrbracket^{\mathcal{M}, g, w}$	Encapsulated denotation
x	Metalinguistic variable
$\mathcal{F}(\alpha)$	Intension
$get[g](x)$	Assignable lookup
\top	Literal true
\perp	Literal false
m_1 and m_2	Logical and
m_1 or m_2	Logical or
$!m$	Logical negation
$m_1 \stackrel{?}{=} m_2$	Equality test
$all(D; a.m)$	Universal quantification
$some(D; a.m)$	Existential quantification
$m_1(m_2)$	Function application

Figure 3.2: Constructs in MSML

- A language model, $\mathcal{M} = \langle \mathcal{F}, \Omega, W \rangle$, where \mathcal{F} maps constant expressions in ILL to their metalinguistic equivalents, W is the set of possible worlds, and Ω is the ontology, which is the set of possible entities and objects across all worlds[13].
- An assignment function, g , that indicates the current value of variable expressions, including, but not limited to, the antecedents of pronouns in a given context.
- A particular (potentially arbitrary) world $w \in W$, in which the truth conditions of an expression can be evaluated.

This brings us to MSML, which, as briefly mentioned earlier, encapsulates ILL expressions within denotation brackets and parameterizes them with specific \mathcal{M}, g, w , allowing for metalinguistic analysis of the encapsulated expression. Take note of the modal separation here: linguistic expressions have a linguistic type, $\phi \dot{\sim} \tau$, while metalinguistic expressions have a metalinguistic type, $m : \tau$. In particular, types for encapsulated linguistic symbols are governed by the following static rules:

$$\frac{\Gamma \vdash_{\Sigma} \phi \dot{\sim} \tau}{\Gamma \vdash_{\Sigma} \mathcal{F}(\phi) : s \rightarrow \tau} \quad \frac{\Gamma \vdash_{\Sigma} x \dot{\sim} \tau}{\Gamma \vdash_{\Sigma} get[g](x) : \tau}$$

$$\frac{\Gamma \vdash_{\Sigma} \phi \dot{\sim} \tau \quad \Gamma \vdash_{\Sigma} w : s}{\Gamma \vdash_{\Sigma} \llbracket \phi \rrbracket^{\mathcal{M}, g, w} : \tau}$$

Modal separation allows us to reason about top-level denotation expressions regardless of the specific context, including information about variable/pronoun

assignments[12]. Therefore, expressions like “I saw them” can be evaluated to an expression that indicates the truth conditions relative to the assignment of the variables for “I” and “them”. In theory, the modal separation present in MSML also allows the encapsulation of *any* (well-typed) linguistic expression in *any* language, although our focus will primarily be on different varieties of ILL. Instead of writing $\llbracket \text{cow} \rrbracket^{\mathcal{M},g,w}$, we will use $\llbracket \ulcorner \text{cow} \urcorner \rrbracket^{\mathcal{M},g,w}$ to indicate that we will look at the representation of the expression “cow” in ILL, rather than the expression itself. The semantics should follow no differently either way, but ILL is easier to reason about than untranslated linguistic symbols from an arbitrary language.

In order to evaluate the meaning of an expression, we need to evaluate the dynamics of the encapsulated expression in MSML. Conventionally, this is represented at the top level by $\llbracket \phi \rrbracket^{\mathcal{M},g,w} \stackrel{?}{=} \top$, which evaluates to the truth conditions of the linguistic expression ϕ in terms of metalinguistic representations of its constituent symbols. Whether lambdas are applied within ILL or only at the metalanguage level depends on your choice of lazy or eager evaluation. We will take $D_\tau \subseteq \Omega$ to be the universal set of type τ , which contains all instances $d : \tau$. The domains that are most frequently invoked are the domain of entities D_e and the domain of worlds D_s .

$$\begin{array}{c}
\frac{\phi \text{ const}}{\llbracket \phi \rrbracket^{\mathcal{M},g,w} \mapsto \mathcal{F}(\phi)_w} \quad \frac{x \text{ var}}{\llbracket x \rrbracket^{\mathcal{M},g,w} \mapsto \text{get}[g](x)} \\
\\
\frac{[\phi \text{ val}]}{\llbracket \forall x_{[\tau]}. \phi \rrbracket^{\mathcal{M},g,w} \mapsto \text{all}(D_\tau ; a. \llbracket \phi \rrbracket^{\mathcal{M},g \otimes a/x,w})} \\
\\
\frac{[\phi \text{ val}]}{\llbracket \exists x_{[\tau]}. \phi \rrbracket^{\mathcal{M},g,w} \mapsto \text{some}(D_\tau ; a. \llbracket \phi \rrbracket^{\mathcal{M},g \otimes a/x,w})} \\
\\
\frac{}{\text{all}(\emptyset ; -) \mapsto \top} \quad \frac{}{\text{some}(\emptyset ; -) \mapsto \perp} \\
\\
\frac{\exists d \in D, D' = D \setminus \{d\}}{\text{all}(D ; a.m) \mapsto [d/a]m \text{ and } \text{all}(D' ; a.m)} \\
\\
\frac{\exists d \in D, D' = D \setminus \{d\}}{\text{some}(D ; a.m) \mapsto [d/a]m \text{ or } \text{some}(D' ; a.m)}
\end{array}$$

The full static and dynamic rules for MSML are omitted for brevity, but are broadly discussed in [11], albeit in non-computational terms.

In order to handle hypernym-hyponym relationships and entity-class relationships, a model is generally evaluated with a set of meaning postulates: short tautological expressions in ILL that represent logical relationships between constants. These may be a quantified expression over all entities, such as $\forall x_{[e]}. \text{rabbit}'(x) \rightarrow \text{animal}'(x)$ representing the concept “All rabbits are animals”, or an expression exclusively in terms of constants, like $\text{person}'(\text{Alex}')$ which represents “Alex is a person”. We generally assume that named entities are uniquely identifiable¹³, so any basic information about an entity that is true over all worlds may be included as a meaning postulate, whereas other judgements, like $\text{student}'(\text{Alex}')$ only hold in a subset of possible worlds, namely those where Alex is a student.

Although there are a number of varieties of MSML, including Event Semantics[14] and SML of New Jersey[15], we will focus on a simple version of MSML based on a simplified variety of ILL. We will call this language SMSML for Simplified MSML, or (SM)²L for short. (SM)²L is built using the ILL flavor used in the previous example (see Table 3.1), where nouns, adjectives and intransitive verbs are assigned constants that are similar to the indicator function of a set, that is, $\llbracket \text{dog}'(x) \rrbracket^{\mathcal{M},g,w}$ is true if and only if x belongs to the set of all dogs, and $\llbracket \text{bark}'(x) \rrbracket^{\mathcal{M},g,w}$ is true if and only if x belongs to the set of all things that bark (or are bark, as it were). Transitive verbs are represented as a relation between two entities, which is more in line with intuition. Information about plurality, tense, and aspect is ignored.

3.2 Translation as Implication

With these tools in mind, we can define some relations to get a handle on semantic complexity. We will make use of the implication expression in ILL, which is often used by semanticists to evaluate entailments within a language. For a pair of languages A and B , let s be a sentence in Language A , and $s \rightsquigarrow s'$ be a translation of s into Language B .

Definition 3.1 (Lossless Translation). A translation, $s \rightsquigarrow s'$, is **lossless** if there exists a context C consisting of $O(|s| + |s'|)$ Boolean judgements such that $\llbracket \ulcorner s' \urcorner \wedge C \rightarrow \ulcorner s \urcorner \rrbracket^{\mathcal{M},g,w} \stackrel{?}{=} \top$ is a tautology for all possible g, w under a compatible model \mathcal{M} .

Definition 3.2 (Lossy Translation). A translation, $s \rightsquigarrow s'$, is **lossy** if there exists a context C^* of ar-

¹³Implementation varies, but I’m partial to `uuidgen` myself

bitrary size such that $\llbracket \ulcorner s' \urcorner \wedge C^* \rightarrow \ulcorner s \urcorner \rrbracket^{\mathcal{M},g,w} \stackrel{?}{=} \top$ is a tautology for all possible g, w under compatible model \mathcal{M} , and there is no way to express C^* as a linear number of Boolean judgments.

Recall that translations occur on the level of expressions, whereas completeness is defined on the level of entire languages. To say that language A is B -complete is to say that for *all* possible expressions $s \in B$, it is possible to generate a lossless translation $s \rightsquigarrow s' \in A$. If there is even a single expression in B for which no suitable translation $s' \in A$ and linear size context C can be constructed that makes the expression $\llbracket \ulcorner s' \urcorner \wedge C \rightarrow \ulcorner s \urcorner \rrbracket^{\mathcal{M},g,w} \stackrel{?}{=} \top$ a tautology, then completeness does not hold, ie. $A \not\leq_S B$.

4 Tap to add section title

We propose the hypothesis that a majority of extant human languages are Toki Pona-complete, and we will present the case for English specifically in this section, building up a rigorous semantic model using the tools discussed in Section 3. Although Toki Pona is probably not (and was never designed to be) a perfect NSM, its nature as a minimal language for human communication can act as a stepping stone to augment ongoing research into semantic primes to search for an NSM.

4.1 Modelling Toki Pona in (SM)²L

The translation of content words is fairly straightforward, but somewhat nontrivial. Content words pattern nicely with the type $e \rightarrow t$, but there is a clear hierarchy in meaning for any modifiers attached to the same head. For instance, “**jan telo**” and “**telo jan**” mean roughly “**sailor**” and “**bodily fluid**” respectively, but naively translating content words to the same constant in (SM)²L would result in both of these becoming equivalent to $\lambda x. \text{jan}'(x) \wedge \text{telo}'(x)$, which is a problem. To remedy this, we will make two constants for each content word: when used as a head, we translate “**telo**” $\rightsquigarrow \lambda x. \text{telo}'_H(x)$, but when used as a modifier, we translate it as “**telo**” $\rightsquigarrow \lambda x. \text{telo}'_M(x)$.

Another issue is that of quantification. English has distinct definite and indefinite markers, whereas Toki Pona does not have either of these general-case quantifiers. There are contexts in which an iota expression may be closer to the intended meaning, but it may be reasonable to analyze the indefinite as the default and use an existential quantifier, since the existential

ale	\rightsquigarrow	$\lambda P. \forall x. P(x)$
mute	\rightsquigarrow	$\lambda P. \forall x. P(x)$
ni	\rightsquigarrow	$\lambda P. \iota x. P(x)$
li	\rightsquigarrow	$\lambda P. \lambda x. P(x)$
e	\rightsquigarrow	$\lambda y. \lambda P. \lambda x. P(x) \wedge e'(x)(y)$
pi	\rightsquigarrow	$\lambda P. \lambda x. \exists y. P(y) \wedge pi'(x)(y)$
ala	\rightsquigarrow	$\lambda P. \forall x. \neg P(x)$

Figure 4.1: Proposed ILL translations of some particles in Toki Pona. These may not be accurate translations in all circumstances.

construction is an entailment of the iota construction, ie. $Q(\iota x. P(x))$ implies $\exists x. P(x) \wedge Q(x)$ for all $P, Q \rightsquigarrow e \rightarrow t$, but the converse does not hold. Based on the design principles of Toki Pona as discussed in [2], and the brief semantic analysis presented in [3], the broader meaning seems like the best representation. When an explicit quantifier like ‘**ni**’ or ‘**ale**’ is used, we will use the corresponding ILL quantifier (see Figure 4.1).

Named entities in Toki Pona have obligatory marking for class, including a head and optionally some additional modifiers. This means that the name of a person, like Sonja Lang, is rendered in Toki Pona as “**jan Sonja**” and not just *“**Sonja**”, and the name of a place, such as Finland¹⁴, is rendered as “**ma Sumi**” rather than *“**Sumi**”. The upshot of this is that the meaning postulates needed to handle named entities in most other languages are pre-encoded into any Toki Pona sentence already! Where a Model for English needs meaning postulates to represent “**Sonja Lang is a person**”, and append an ILL representation like $\text{person}'(\text{Sonja}')$ to the model, Toki Pona might simply translate the phrase “**jan Sonja**” as $\lambda P. P(\text{Sonja}') \wedge \text{jan}'_H(\text{Sonja}')$ to capture the implicit meaning postulate. Alternatively, we may construct a simple meaning postulate that captures that all entities can be assumed to belong to the class with which they are introduced.

The final step is to observe the syntax of Toki Pona. Most of the content words will combine with the modifier concatenation rule, $a b \rightsquigarrow \ulcorner a \urcorner \square \ulcorner b \urcorner$, so all we have left are the particles. Being primarily head-initial, we find that quantifier-like particles in Toki Pona tend to come after the head of a phrase, while prepositions, case-like particles, and verbal auxiliaries

¹⁴This claim is disputed

tend to come before the head. This means that both directions of function application are likely to be seen given the right contexts. With that, we can construct translations for the most common particles, which are listed in Figure 4.1.

4.2 Bilingual Models

Montaguean semantic models are typically used to ascertain the meaning of expressions within a given language, rather than to compare the denotations of expressions in two or more different languages. Since we need to work with multiple languages simultaneously in our use case, we need to create a notion of multilingual models for our ILL and MSML structures to work as expected.

If we have two language models, \mathcal{M}_1 and \mathcal{M}_2 , we can define their union, $\mathcal{M}_1 \oplus \mathcal{M}_2$, to be the tuple $\langle \mathcal{F}_1 \sqcup \mathcal{F}_2, \Omega, W \rangle$, where $\mathcal{F}_1 \sqcup \mathcal{F}_2$ maps ILL constants to the corresponding metalanguage expressions for constants in both \mathcal{F}_1 and \mathcal{F}_2 . We will take as a given that the input domains of \mathcal{F}_1 and \mathcal{F}_2 do not overlap as long as they are from models representing different languages: any overlap can be handled trivially by explicitly marking the constants by the model from which they originate¹⁵. Following the intuition that a model of a given language should reflect the judgements and understanding of a native speaker, we can think of a compound model as a reflection of the judgements of someone who is a native speaker of both languages¹⁶. In order to handle this, the union model $\mathcal{M}_1 \oplus \mathcal{M}_2$ will also have to contain meaning postulates to describe the relationships between expressions in each language. For example, if a bilingual speaker of English and Toki Pona would say that something described as ‘a book’ in English would be ‘lipu’ in Toki Pona, then we can represent this relation with the meaning postulate $\forall x. \text{book}'_{\text{eng}}(x) \rightarrow \text{lipu}'_{H,\text{tok}}(x)$. We can’t do this for all pairs of overlapping words in general due to polysemy, although we can combine them with other judgements found in the context, and create a meaning postulate based on the 20-questions test like so: $\forall x. \text{bank}'_{\text{eng}}(x) \wedge \neg \text{EDGE.WATER}'(x) \rightarrow \text{银行}'_{\text{zho}}(x)$ where $\text{EDGE.WATER}' \approx e \rightarrow t$ is an ILL translation of a metalinguistic judgement of whether x is the edge of a body of water, and we represent the Chinese word

¹⁵eg. by a subscript ISO 639-3 code, our convention of choice.

¹⁶Strictly speaking there are no native speakers of Toki Pona, but for the sake of our model, we can approximate this by imagining the model to reflect a speaker with advanced or fluent level of speaking. About 16% of respondents to the 2021 Toki Pona Census claim these levels of proficiency[16].

- (7) **soweli lili li kalama**
 animal small COP sound
 ‘The little dog barked’

Figure 4.2: Gloss for our worked example.

for ‘bank, financial institution’ (seen in Section 2.1) as $\lceil \text{银行} \rceil = \lambda x. \text{银行}'_{\text{zho}}(x)$.

4.3 Proving Losslessness

To begin, let’s consider the Toki Pona sentence presented in Figure 4.2, gloss (7) as **s**. This sentence can be translated into English in a number of ways: ‘The little dog barked’, ‘A small cat is purring’, ‘The mice squeak’, ‘The rabbit barks’, ‘Small animals are sound waves’, are all valid translations of this sentence. It doesn’t make much difference which one we choose – in fact all of these English sentences are (in a sense¹⁷) precisizations of the original sentence in Toki Pona – so for simplicity and familiarity, we can take ‘The little dog barked’ to be our s' .

We can construct some meaning postulates to reflect the relationships between relevant English and Toki Pona concepts. In particular, we will utilize the following:

1. $\forall x. \text{dog}'_{\text{eng}}(x) \rightarrow \text{soweli}'_{H,\text{tok}}(x)$
2. $\forall x. \text{little}'_{\text{eng}}(x) \rightarrow \text{lili}'_{M,\text{tok}}(x)$ ¹⁸
3. $\forall x. \text{bark}'_{\text{eng}}(x) \wedge \neg \text{TREE.PART}'(x) \rightarrow \text{kalama}'_{H,\text{tok}}(x)$

Toki Pona doesn’t have any significant inflectional grammar that we need to account for in our translation, so the trivial context, \top , works fine for the grammar of the sentence. We do need a judgement for the word ‘bark’, since this could refer to either an animal sound or a part of a tree¹⁹, so we assert that $\text{TREE.PART}(x)$ is false for our context.

The last step before plugging in our values into the definition of lossless translation is producing the ILL representations. We will assume eager evaluation and

¹⁷Rigorously defining precisizations in a bilingual setting is difficult since they rely on entailments, which we already had to augment to deal with translations. We use “precisization” loosely to reflect the judgments of a bilingual speaker. Modulo issues with cross-linguistic comparisons, it is reasonable to say that English ‘dog’ is a hyponym of Toki Pona ‘soweli’[9].

¹⁸The head interpretation of ‘lili’ would also work, but isn’t used in this derivation.

¹⁹or something similar to either of these, technically

skip the intermediary expanded form, only showing the result after lambda evaluation.

$$\begin{aligned} \lceil s \rceil &= \\ \exists x. \text{soweli}'_{H,\text{tok}}(x) \wedge \text{lili}'_{M,\text{tok}}(x) \wedge \text{kalama}'_{H,\text{tok}}(x) \\ \lceil s' \rceil &= \\ \text{bark}'_{\text{eng}}(\iota x. \text{dog}'_{\text{eng}}(x) \wedge \text{little}'_{\text{eng}}(x)) \\ C &= \\ \neg \text{TREE.PART}'(\iota x. \text{dog}'_{\text{eng}}(x) \wedge \text{little}'_{\text{eng}}(x)) \end{aligned}$$

Finally, we can plug these expressions into the MSML expression given in Definition 3.1 and evaluate the results. Taking the union of simple models for English and Toki Pona gives us $\mathcal{M} = \mathcal{M}_{\text{eng}}^S \oplus \mathcal{M}_{\text{tok}}^S$ as our model. Full evaluation is shown in Figure 4.3.

Continuing in this manner has allowed us to show that all Toki Pona sentences that can be represented in (SM)²L have a lossless translation into English, indicating that English is Toki Pona-complete, as we hypothesized. The proof of this is left as an exercise for the reader.

5 Conclusions

Our research suggests the possibility of semantic complexity classes to create a well-defined notion of completeness as a relation between two languages, particularly a natural language and a constructed minimal language. Such classes can be formally analyzed and evaluated with the formal semantic definitions of translatability, namely evaluating whether $\llbracket \lceil s' \rceil \wedge C \rightarrow \lceil s \rceil \rrbracket^{\mathcal{M},g,w}$ is a tautology for some context C of $O(n)$ size, or only for a larger context. This reflects intuitive judgements, such as the 20-questions test to include simple inflection and the most basic cases of polysemy or homonymy when translating individual symbols, and expands them to apply to strings and other compound expressions via formal semantic compositionality. If for all sentences $s \in B$, there exist translations $s' \in A$ such that a context C of $O(n)$ length can be constructed to make $\llbracket \lceil s' \rceil \wedge C \rightarrow \lceil s \rceil \rrbracket^{\mathcal{M},g,w}$ a tautology, then we have that $A \geq_S B$, meaning that A is semantically complete with respect to B , or simply B -complete. If for some sentence in B , there is no sentence in A for which a linear-size context can be constructed to form a tautology, then A is not B -complete.

We specifically explored the example of Toki Pona as a basis for semantic completeness and laid out the

beginnings of a proof to show that English is Toki Pona-complete. These results provide a significant first step towards evaluating the semantic complexity of human languages as a whole. Hopefully, future research can determine whether other languages are Toki Pona-complete, and perhaps evaluate the completeness of natural languages with respect to other minimal languages aside from Toki Pona. The goal is to ultimately guide semanticists to the holy grail of natural semantic metalanguage, or perhaps discover that NSM is impossible. It is still too early to say how this will turn out.

The current study was limited in that only (SM)²L was considered in the construction of translatability metrics, and other varieties of MSML were not evaluated. (SM)²L alone may not be sufficient to fully analyze Toki Pona and its translatability into other languages, so future research may consider expanding to other varieties of MSML such as Event Semantics, which may be a more accurate representation of the semantic space of Toki Pona. ‘moku’ used as a verb may not be accurately translated as describing the subject $\lambda x. \text{moku}'_H(x)$ as in (SM)²L, but rather as describing an event in which the subject participated: maybe $\lambda x. \lambda e. \text{moku}'_H(e) \wedge \text{AGENT}(e)(x)$ to represent that the event is ‘moku’ while the subject takes the theta-role of agent for that event[14]. This may provide for a deeper analysis than (SM)²L, which we utilized due to space and time limitations.

In summary, our research serves as a proof-of-concept that a minimal, semantically complete language is possible, as Toki Pona fulfills these properties. Future research can perform the necessary minimaxing to search for a sufficiently robust NSM, and discover any potential limitations to NSM, which will undoubtedly provide unique insight into the fundamental limits of human language an understanding. Unless of course GPT 5 in its infinite wisdom just gives us a perfect, general-purpose NSM from a black box.

References

- [1] C. Goddard and A. Wierzbicka, *Meaning and universal grammar : Theory and Empirical Findings* /. Amsterdam: J Benjamins Pub, XV, 334 Seiten. ISBN: 1588112640.
- [2] S. Lang, *Toki Pona: The Language of Good*. Tawhid, 2014.
- [3] jan Misali, *What is toki pona? (toki pona lesson one)*, Dec. 2021. [Online]. Available: <https://www.youtube.com/watch?v=2EZihKCB9iw>.
- [4] u/pisceyo, “Toki pona analysis: Parts of speech,” *r/Tokipona - Reddit*, 2019. [Online]. Available: https://www.reddit.com/r/tokipona/comments/df0zbi/toki_pona_analysis_parts_of_speech/.
- [5] Wikipedia Contributors, *Sisu — Wikipedia, the free encyclopedia*, Jan. 2024. [Online]. Available: <https://en.wikipedia.org/wiki/Sisu>.
- [6] T. Scott, *The sentences computers can't understand, but humans can*, Feb. 2020. [Online]. Available: <https://www.youtube.com/watch?v=m3vIEKWrP9Q>.
- [7] D. Swanson, *Kay(f)bop(t) babel text*, Dec. 2015. [Online]. Available: <https://crazyinjageeks.wordpress.com/2015/12/29/kayfbop-babel-text/>.
- [8] D. Swanson, *Kay(f)bop(t) conjugations*, Nov. 2015. [Online]. Available: <https://crazyinjageeks.wordpress.com/2015/11/29/kayfbop-conjugations/>.
- [9] This was once revealed to me in a dream.
- [10] jan Misali, *Conlang critic episode eight: Kay(f)bop(t)*, Nov. 2016. [Online]. Available: <https://www.youtube.com/watch?v=XhxjI3t6Ht4>.
- [11] R. Montague and R. H. Thomason, *Formal philosophy: Selected papers of Richard Montague*. Yale University Press, 1974.
- [12] R. Harper, “Chapter 35: Assignable references,” in *Practical Foundations for programming languages*. Cambridge University Press, 2016, pp. 323–332.
- [13] T. M. V. Janssen and T. E. Zimmermann, *Montague semantics*, Apr. 2021. [Online]. Available: <https://plato.stanford.edu/entries/montague-semantics/>.
- [14] T. Werner, “Theta roles and event semantics,” Class lecture from 80-381: Meaning in Language, Nov. 2023.
- [15] M. Erdmann and D. Kaynar, *Functional programming, spring 2024*, P. Morphism, Ed., Jan. 2024. [Online]. Available: <https://www.cs.cmu.edu/~15150/>.
- [16] jan Tamalu, *Results of the 2021 toki pona census*, Jun. 2021. [Online]. Available: <https://tokiponacensus.github.io/results/>.

Introducing Supernatural Language Processing (SNLP)

Jenson Crawford

Abstract

In recent years, the field of Natural Language Processing (NLP) has made significant strides in enabling machines to understand and generate human language. However, what if we told you that NLP merely scratches the surface? In the hallowed halls of academia, where skepticism and curiosity collide, a new field emerges: Supernatural Language Processing (SNLP). Building on the pioneering work of Spengler, Stanz, and Venkman,[1] in detecting and interacting with supernatural beings, we move beyond mundane chatbots and pedestrian sentiment analysis. SNLP is a cutting-edge discipline that transcends earthly boundaries and allows us to communicate with paranormal entities. This groundbreaking paper explores the theoretical underpinnings, practical applications, and potential pitfalls of SNLP. Buckle up, because armed with our trusty para-neural networks and a dash of ectoplasm, we're about to explore the spectral depths of linguistic interaction and spectral lexicon. ^{1 2}

1 Introduction

1.1 The Haunting Gap

While NLP focuses on mundane tasks like sentiment analysis, chatbots, and grammar correction,[2] it has woefully neglected the spectral community. Ghosts, apparitions, and poltergeists have long felt left out of the linguistic loop. Imagine their frustration when they try to communicate with humans, only to receive responses like, "Sorry, I didn't catch that" or "Did you mean 'boo'?"

1.2 The Birth of SNLP

SNLP emerged from a séance in a dimly lit basement, where a team of linguists, mediums, pseudo-scientists, charlatans, and a particularly eloquent ectoplasmic orb convened.[3] Their mission? To bridge the gap between the corporeal and the ethereal. Thus, Supernatural Language Processing was born—a field

¹ A generous grant from the Dana Barrett Foundation supported this research.

² A grant from the Louis Tully Family Mistrust supported this research.

that combines spectral semantics, ectoplasmic syntax, and spectral sentiment analysis.

1.3 A Burning Question

There is no longer a question that paranormal entities exist.[4] SNLP bridges the gap between earthly linguistics and ethereal discourse. The research question to our team of pseudo-scientists: “What are these supernatural entities saying? Why do they insist on speaking in riddles and haiku?” Perhaps they’re just lonely, want revenge for an untimely death, or maybe they’ve binge-watched too many episodes of “Kindred Spirits.”[5] Either way, we’re here to decode their spectral syllables.

2 Theoretical Foundations

2.1 Ghostly Grammar

SNLP introduces novel grammatical constructs, such as the “Phantom Passive Voice” and the “Spectral Subjunctive.” For instance:

- Phantom Passive Voice: “The séance was conducted by the medium” becomes “The séance was conducted by unseen hands.”
- Spectral Subjunctive: “If I were alive, I would eat pizza” becomes “If I were corporeal, I would feast upon ectoplasmic pizza.”

2.2 Ecto-Embeddings

Word embeddings take on a whole new dimension in SLP. Instead of vectors in Euclidean space,[6] we use “ecto-embeddings” in the spectral plane. Words like “spooky,” “ethereal,” and “ectoplasm” are mapped to their ghostly counterparts, such as “spooktacular,” “etherspeak,” and “ecto-delicious.”

3 Methodology

3.1 Spectral Corpus Cadaverum Collection

- We haunt ancient libraries, abandoned crypts, and Wi-Fi-enabled graveyards to gather spectral texts, assisted by amateur grave robbers, pseudo-scientists, and paranormal researchers.
- Our dataset includes EVP (Electronic Voice Phenomena), séance transcripts, and cryptic graffiti from haunted restrooms.
- We also scrape spectral tweets because even ghosts need to vent about their unfinished business.

3.2 Supernatural Tokenization

- Ghosts don't adhere to grammatical norms. Their sentences resemble a drunken game of Scrabble. Our tokenization algorithm handles spectral typos, missing vowels, and the occasional ectoplasmic expletive.

3.3 Entity Recognition

- Identifying spectral entities is tricky. Our model distinguishes between Slimer, a librarian, Mr. Stay-Puft, and Gozer the Gozerian. Non-corporeal context is key.
- Fun fact: Poltergeists use passive voice exclusively.

4 Results

4.1 Common Phrases

- Top 3 spectral phrases:
 1. "Boo!"
 2. "He killed me."
 3. "Huh. I guess COVID is real."
 4. "Tonight on Paranormal."

4.2 Cryptic Messages

- Our model deciphers hidden messages:
 1. "The treasure is beneath the Starbucks on Elm Street."
 2. "Shhhhh"
 3. "The Wi-Fi password is 'Banshee123.'"
 4. "Rosebud..."
 5. "Don't cross the streams. It messes with the our Wi-Fi."
 6. "Choose the form of the destructor!"

4.3 Emotional States

- Sentiments detected:
 1. 60% "Ethereal Curiosity"
 2. 30% "Spectral Sarcasm"
 3. 10% "Otherworldly Apathy"

5 Practical Applications

5.1 Ghost Chatbots

Imagine an SNLP chatbot that converses with spirits. Sample dialogue:

- User: “Who are you?”
- Chatbot: “I am the whisper in the wind, the echo in the crypt, the semi-transparent entity you seek.”

5.2 Paranormal Sentiment Analysis

SNLP sentiment analysis deciphers ghostly emotions:

- Positive: “The moonlit graveyard fills my ectoplasmic heart with spectral joy.”
- Negative: “The séance attendees mistook me for a draft and closed the window.”

6 Challenges and Ethical Considerations

6.1 Unconscious and Undead Bias

Are ectoplasmic embeddings perpetuating spectral stereotypes? We must tread carefully to avoid deathly discrimination.

6.2 Data Privacy Concerns

What if taking a secret to the grave is insufficient for data protection? Do regulations like GDPR and CCPA apply beyond the veil? Can we guarantee the anonymity of post-corporeal conversations?

7 Future Research

We have yet to determine if SNLP can lead us to a Supernatural Semantic Met-language (SNSM) that is Toki Pona-complete.¹ Will we need additional dimensions beyond Toki Pona-completeness to account for the SNLP move use of spectral plane ecto-embeddings rather than Euclidean space vectors, or will Toki Pona-completeness be sufficient for non- or post-human communication as well?

8 Conclusion

SNLP opens portals to the unknown—a way to connect with the afterlife, one ecto-sentence at a time. Our findings suggest that ghosts are just as confused as

we are. So next time you hear a disembodied whisper on a EVP replay, perhaps it's not just pareidolia. It might be a lost soul simply seeking Wi-Fi access.

References

- [1] Egon Spengler, Ray Stantz, Peter Venkman (2017) *The Complete Guide to Paranormal Phenomenon and Containment*. SSV Publishing, New York, New York, 2nd ed.
- [2] Christopher D. Manning , Hinrich Schütze (1999) *Foundations of Statistical Natural Language Processing* MIT Press
- [3] Peter Venkman (1985) *Séances: a Way to Meet Women*. SSV Publishing, New York, New York, 4nd ed.
- [4] Ray Stantz, Egon Spengler, Peter Venkman, *Detection of Specters, Spirits, and Trans-Dimensional Beings*. Journal of Paranormal Studies, vol. 7, no. 2, pp 42-53.
- [5] Salvatore Pesce III, Harrison Francis, et.al. (2016-) *Kindred Spirits*. The Travel Channel, Discovery Network
- [6] Jacob Eisenstein (2019) *Introduction to Natural Language Processing*. MIT Press
- [7] jan Kijupo (2024) *Toki Pona and Orders of Semantic Completeness*

SMS: Sending Mixed Signals

Pranav Khadpe[†]
pkhadpe@cs.cmu.edu
Carnegie Mellon University
Pittsburgh, Pennsylvania, USA

Sayan Chaudhry[†]
hi@sayan.io
Met Pranav on Bumble BFF
New York, New York, USA

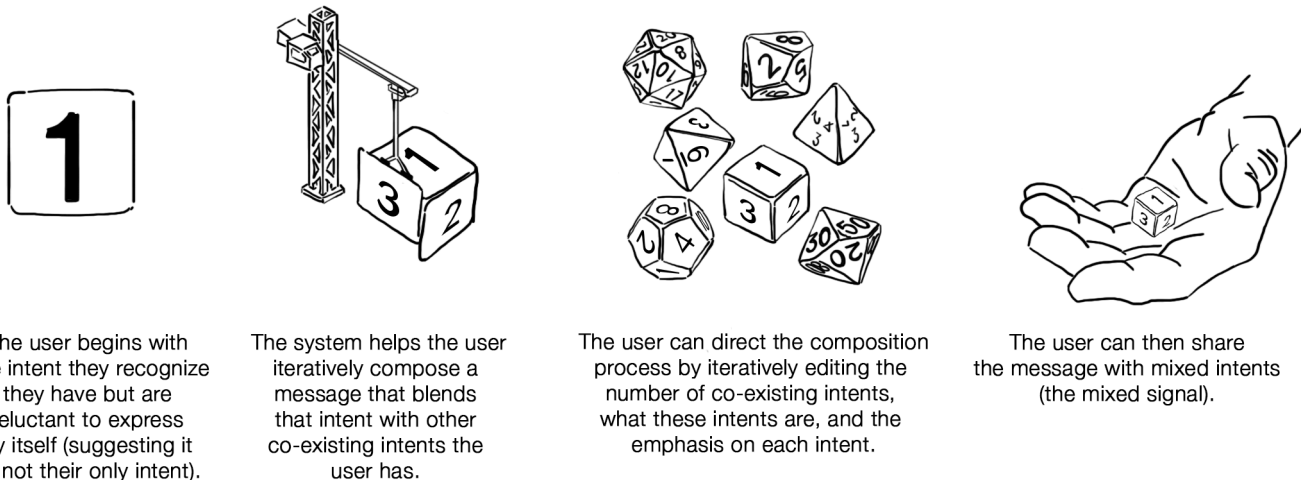


Figure 1: SMS is an interactive system that enables users to craft *mixed signals*: messages that convey competing intents together. Users can craft mixed signals, beginning with a single intent they are hesitant to express by itself, and controlling the number and content of co-existing intents to be mixed in. SMS attempts to bridge the gap between the essential nature of our intentions, which may be diverse and competing, and what current symbolic resources for expression allow us to express, which tend to be unambiguous, monolithic reductions of our intentions, because of narrow prevailing ideals of what good communication is.

ABSTRACT

Should I tell her I like her? Do I want to have kids? Do I tell the authors their core idea is uninspiring? Much of life is spent in the frustration of wanting to (or having to) say something but not knowing *exactly* what to say because our intents may be diverse and competing. We want to fit in *and* stand out, be liked *and* be looked up to, tell people what we feel *and* what they want to hear. To be human is to have competing intents. In these situations, we could achieve release and relief by simply being able to *send mixed signals*: to express competing intents at the same time. Unfortunately, expressing mixed intent can be challenging. Contemporary communicative training and social intercourse rests upon and reifies a narrow conception of good communication: that good communication is unambiguous, and clear in intent. This ideal invariably

[†] is the primary contributor.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).
SIGBOVIK '24, April 5, 2024, Pittsburgh, PA, USA
© 2024 Copyright held by the owner/author(s).

shapes the symbolic resources and normative scripts that make up *language*, which, as a result of this, remain limited in their ability to carry mixed intent. The result is a communicative handicap: we may want to express mixed intent but are unable to find the words to. Therefore, in this paper, we present a general approach to generate mixed signals by leveraging the generative capabilities of Large Language Models (LLMs). By virtue of being trained on our interactions online, off-the-shelf LLMs, too, reflect dominant ideals of good communication and when prompted naively, steer away from generating mixed signals. We present computational techniques to steer the outputs of LLMs and generate mixed signals based on the emerging idea of *scenario nesting*, which suggests that by immersing LLMs in fictional worlds, we can challenge and change their worldview. We instantiate these techniques in an interactive system, SMS, that allows users to craft mixed signals, with control over the number and content of signals mixed. Taken together, this paper points to a future where mediums of expression, ranging from digital mediums all the way to language itself, are able to reflect and accommodate the richness, complexity, and contradictions of human intentions, rather than sterilizing and censoring our intentions. SMS is live at: sendmixedsignals.vercel.app

1 INTRODUCTION

Consider this: you want to tell your roommate that it bothers you that they tend to thoughtlessly set the TV volume to odd numbers. "Your audacious disregard for auditory symmetry is irritating", you think to say. But as you notice the form this thought begins to take, you feel reluctant to say it. You don't want your roommate to despise you. You want to say it. But in an equally real sense, you don't want to say it. You are stuck in a dilemma. You want to say something but you don't know exactly what. Such dilemmas make up much of our social existence. Consider another example: you want to tell your crush, that you have a crush on them. Perhaps you have even imagined how you might confess your feelings using words from the 20th best-selling single of 2011, "What Makes You Beautiful": "Baby, you light up my world like nobody else. The way that you flip your hair gets me overwhelmed." But you hesitate. You don't just want your crush to think you are sweet, you also want your crush to look up to you. You want them to think you are *desirable*—better even, *desired*—rather than *desiring*. You worry your eager show of affection will reveal your lack of popularity and depreciate their respect for you. Again, a dilemma. Much of life is spent in this frustration of wanting to (or having to) say something but not knowing *exactly* what to say. Inner conflict is part of the human condition.

In many situations, like those above, we could achieve release and relief by simply being able to send mixed signals: to express competing intents at the same time. Today, however, achieving such indeterminate communication effectively can feel challenging. This is because, to express our intentions, we turn to language systems, which present us with normative scripts for how to express ourselves. For instance, the language system of English offers us such ready-made symbols as, "I'm sorry", to apologize, "Thank you", to convey indebtedness, and even, "Baby, you light up my world like nobody else", to convey affection¹. But we have few symbols or scripts to convey mixed intent because most of our symbols derive from a very specific ideology of "good" communication: that good communication is unambiguous, and clear in intent. Indeed, we are often taught to eliminate room for misinterpretation, which presumes there is a desired interpretation—an idea that permeates our symbolic resources. This presents a challenge when our intentions, which at a given moment may be diverse and competing, must be forced into a straightjacket of symbolic resources that only allow a single intent.

We suggest that, given their generative capabilities, Large Language Models (LLMs) offer an opportunity to craft "mixed signals" for situations where we may want to but currently can't because of our own communicative handicaps. However, despite remarkable progress in producing fluent and compelling content, off-the-shelf LLMs such as GPT-4 and ChatGPT still fall short of generating effective mixed signals. Naively prompting these models does not result in messages with mixed intent for two reasons. First, these models

are trained on web data that includes some of our online conversations, the opinions we share online about how to communicate well, and potentially even instructional material on effective communication. As a result, they have inherited dominant assumptions about good communication: that it ought to convey a single intent, clearly. Second, these assumptions are further reinforced, quite literally, during the Reinforcement Learning from Human Feedback (RLHF) step. During this step, to improve the performance of models for business-critical tasks, their language understanding and generation capabilities are biased in ways that improve their performance on tests of question answering, comprehension, reasoning, and coding. Improving performance on such tests, and on work tasks more generally, privileges clarity and precision while making LLMs worse for our purpose, of supporting personal communication, where the need for clarity and precision may be superseded by the need to express mixed intent.

In this paper, we present our system, *SMS* that helps users generate messages with mixed intent. To overcome the limitations of off-the-shelf LLMs, *SMS* leverages *scenario nesting* [12]: by constructing a relevant but fictitious scene and asking the LLM to generate a dialogue for a character in the scene, we are able to successfully get the LLM to generate messages with mixed intent. Further, to help users iteratively direct the message composition process, *SMS* provides users with two kinds of steering controls [13]: (1) *Intents* that allow users to define the number and content of co-existing intents they want to include in the message, and (2) *Will Power* that users can adjust to control the relative emphasis placed on different intents in the output message. With *SMS*, a user can begin with one intent that they are reluctant to express by itself (e.g. "Baby, you light up my world like nobody else. The way that you flip your hair gets me overwhelmed") and iteratively compose messages that more faithfully represent their mixed intents:



¹See West's thesis [18] (Chapter 7) on how cultural codes shape language use. She writes: "Movies, television shows, song lyrics, and advertising present us with normative scripts for how to express our feelings constantly...[John Durham] Peters and [John] Dewey remind us that language is not a private good and that our use of even spontaneous, "original" language is a borrowing of public symbols...Language, even spoken face-to-face, mediates despite its apparent transparency."

2 RELATED WORK

It is now time to introduce you to the relatives of this paper (or at least the relatives that it gets along with) and describe how its development has been influenced by these relatives. Specifically, we will talk about the theoretical frame it inherits from the Sociology side of its family, the interaction design ideas it inherits from the Human-Computer Interaction side of its family, and the logical structure this paper's main proposition inherits from from the Rhetorics side of its family.

2.1 Relational Dialectics (relatives on the Sociology side)

Experiencing competing intentions is natural, especially in our interactions with other people. Several scholars have suggested that social associations rests on *relational dialectics* [3–5, 15], or the “*simultaneous presence of two relational forces that are interdependent and mutually negating*” [15]. For instance, people may at once want to be open with one another *and* closed and self-protective. In the study of close relationships, prior work [1, 3, 7, 15, 16] has identified several such dialectics that may be at play when partners interact with each other: connection/autonomy, openness/closedness, predictability/novelty, affection/instrumentality, and judgement/acceptance. But such contradictions are not restricted to close interpersonal relationships. Describing the opposing forces that are activated when people try to come together, Blau writes: “*A fundamental dilemma of social life is that between being looked up to and being liked by associates*” [6]. He describes how the opposing desires *to be liked* and *to be looked up to* can present a contradiction:

Both the respect and the affection of our associates are important to us, but our efforts to win the one often hurt our chances to win the other. Suppose a colleague has asked us to comment on the first draft of a paper he has written. If we make penetrating criticisms, this may increase his respect for our competence, but it will hardly endear us to him; and, if we make only complimentary remarks, he may feel more favorable toward us but see no reason to respect our judgment. Although he benefits from supportive comments as well as from valid criticisms, he benefits from them and reacts to them in different ways, involving for us either a gain in respect at the expense of warm acceptance or a gain in intimacy at the cost of respect.

—Peter M. Blau, *A Theory of Social Integration* [6]

The presence of contradiction is not a problem to be “solved”. Montgomery writes [15](and we agree): “*neither pole of the opposition is seen as inherently positive or negative. Essences of life, as oppositional forces are thought to be, are neither good or bad, they simply are.*” “*The struggle of opposites is thus not evaluated negatively by dialectical thinkers*”, says Baxter [4]. We are dialectical thinkers.

While they can't be “solved”, dialectics call for strategies for adjustment and transformation. One strategy, *Selection*, involves selecting or prioritizing one aspect of the contradiction over the other. For example, “*people can choose to disclose information (openness), even if they fear or even expect rejection and want to protect themselves (closedness)*” [17]. The other strategy *Integration* [17], on

the other hand, involves responding to both opposing tendencies simultaneously. It leverages ambiguity to avoid explicitly involving either pole of the contradiction. To achieve this it relies on “*such communication devices as self-contradictions, subject switches, tangentializations, obscure word choices, and incomplete sentences*” [15].

We suggest that current communicative resources, tend to force people into *Selection*, to express one intent clearly. Our work attempts to broaden the possible ways that people can express themselves so they may also employ *Integration* when they want to.

2.2 Resources for Interactional Ambiguity (relatives on the Human-Computer Interaction side)

Our work is inspired by, and draws on, a long line of HCI research that recognizes the critical role of ambiguity in interpersonal communication, and contributes ideas to support ambiguous communication [2, 7, 9]. Here, previous work has explored how communication mediums might leave “space for stories”, allowing users to communicate one of the attendant tendencies in a dialectic (e.g. autonomy) [2] while still maintaining plausible deniability, such that they can claim to have been pursuing the other tendency (e.g. connectedness). Other work shows how ambiguous representations of people's experiences, even if they don't convey clear “information”, can still serve as resources for collective meaning-making [8], and ultimately serve a connective function, bringing people closer. Building on this work, our paper sets out to contribute digital objects that might carry our wild, vague and subjective inner experiences. Our hope is that such objects might serve a valuable role in relationship maintenance and can also provide resources for meaning-making practices that bring people closer.

2.3 Function Propositions (relatives on the Rhetorics side)

Here, we describe that logical category of propositions to which the main proposition of this paper belongs. The main proposition of this paper could be simply stated as: “*Even though expressive media that privilege clear and unambiguous communication seem to be an effective way to express all our communicative goals, they are in reality ineffective at supporting all our communicative goals.*” This proposition belongs to the more general class of propositions, that Davis describes as *Function* propositions [11], and have the following logical structure:

What seems to be a phenomenon that function effectively as a means for the attainment of an end is in reality a phenomenon that functions ineffectively.

Davis goes on to describe how a social theorist proceeds to make such an argument: “*the social theorist claims that a certain social institution, of which his audience is known to approve, actually has consequences, of which his audience is known to disapprove*” [11]. Our paper, then, can be seen as a specific instance of this more general category. We, the social theorists, are making the argument that prioritizing clarity and precision in expressive media (the current social institution) is something you, the reader, approve but is also something that can suppress expression and cause frustration—consequences that, we think, you disapprove.

So, this paper’s proposition belongs to the same logical category as “Herbert Marcuse’s assertion in *Repressive Tolerance* that the tradition of tolerance in America, which was considered at the time he wrote to be a value that fostered the goal of a liberated society, is in fact a value that hindered the goal of a liberated society” [11] (an example Davis provides).

It also belongs to the same logical category as Chaudhry’s assertion, in his 2019 SIGBOVIK paper [10], that good code style, which was considered at the time to be a practice that encouraged healthy collaborative behaviors (because it allows others to understand your code), is in fact a practice that hinders healthy collaborative behavior (because it makes it easier for people to plagiarize each others’ code).

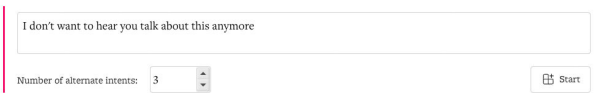
It is not common for a paper to explicate its own logical structure. One goal of writing this section was to demonstrate that a paper describing ambiguous communication need not be ambiguous itself; it can be well-formed. And we hope this encourages some of you to write about similar ideas. The second goal of writing this section was to show that with sufficient effort, and abstraction, it is possible to articulate connections with, and cite, a broad base of literature. We provide a concrete example of this by connecting to, and citing, Chaudhry’s seemingly unrelated paper titled: “*Novel Defense Against Code Theft Using Properties of Fibonacci Series*” [10]. We hope this encourages you to find new approaches to cite your friends’ papers more often.

3 SMS

In this section, we introduce *SMS* (sendmixedsignals.vercel.app), an interactive system that exemplifies our approach for generating messages with mixed signals. To illustrate our approach, let us say a user wants to tell someone, “I don’t want to hear you talk about this anymore” but they are reluctant to say it. Moments of unease, like this, are when you turn to *SMS*. And *SMS*’ goal, then, is to generate a final message that more faithfully represents your mixed intents: a message you would feel comfortable sending. To begin, *SMS* asks the user to identify one intent, from the mixture, that they may want to express.



Usually, this is the intent that most readily maps to symbolic resources of language. In the above example that might be “I don’t want to hear you talk about this anymore”. To the user, this intent may be the most recognizable from the mixture but it is not the only intent they have. Why else would they be hesitant to say it? To begin to work towards the final message, *SMS* also asks users to input the number of alternate intents (they think) co-exist with the intent they have recognized. Now, suppose the user goes ahead



and types out this intent and suggests that there are two alternate co-existing intents they potentially want to include.

From here, *SMS* initializes the steering interface and generates an initial message. The goal of the initialization is not to be “right” but to provide a starting point from where the user can iterate towards the final message.

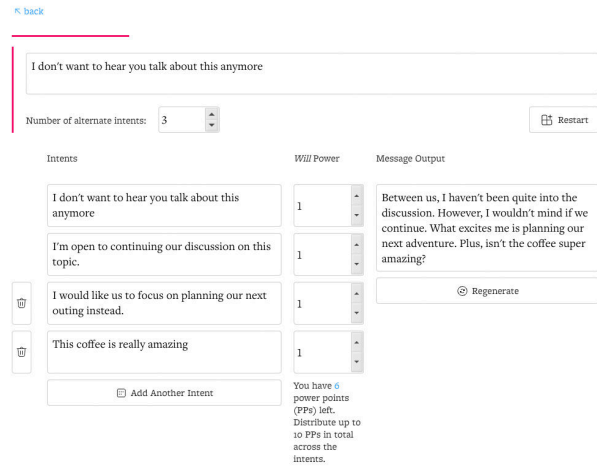


Figure 2: *SMS*’ interface

In the following subsections, we describe: (1) the Steering Controls (*Intents* and *Will Power*); and (2) the Generation Pipeline.

3.1 Steering Controls

To help users iteratively direct the composition process, *SMS* provides users with two kinds of steering controls [13]: (1) *Intents* that allow users to define the number and content of co-existing intents they want to include in the message, and (2) *Will Power* that users can adjust to control the relative emphasis placed on different intents in the output message. Changing either the *Intents* or *Will Power* and clicking “Regenerate” produces a new output.

3.1.1 Intents. *SMS* allows users to express up to 10 intents in a message. Authors can edit and add intents. They can also delete intents down to a minimum of 2 intents (one intent by itself, is the territory of clear communication). The initial set of *Intents* are generated based on the one intent provided by the user and the number of alternate intents. In our example, we use the following prompt to generate the initial set of intents (prompt template available in implementation):

```
role: "system"
content: "You are a smart writing assistant. I will give you an intent that I want to convey in a message I am writing. I want you to generate alternate intents that will also be present in the message."

role: "user"
content: "
## Original Intent
'I don't want to hear you talk about this anymore'
## Number of New Intents
3
## Generation Procedure
Generate exactly 3 new intents. The first intent should be a complete contradiction of the original intent. The second intent should be unrelated to both the original intent and the first intent. Every subsequent intent should be unrelated to all intents that came before it."
```



```
## Response Format
Your response should be a JSON object with the key "alternatives", whose
value is a list of strings. Each string should be an alternate intent."
```

This returns the set of intents visible in Figure 2. We augment to this, the initially supplied intent. The first generated intent is chosen to be a contradiction of the user-supplied intent to encode the idea that reluctance or hesitation to say the provided intent suggests that the user also intends its contradiction. Whenever a user clicks “Add Another Intent”, it comes pre-populated with a generated intent. It is generated in a manner similar to above, to contrast with the intents that have already been expressed.

3.1.2 Will Power. SMS also provides users control over how much they want to emphasize a particular intent in the message. Users can express this in the interface in terms of *Will Power Points* (PPs) which can also be understood as “weights”. Users can distribute up to 10 PPs across the different intents they have specified and the relative emphasis on each intent in the final message reflects the proportion of PPs it has been allocated. SMS only considers proportions: it internally represents PPs as a normalized array, with length equal to the number of intents, where the number at a given index represents the weight placed on the corresponding intent. At initialization, all intents are assigned 1 PP. The number 10, here, was not chosen in any principled way; it was a number that seemed to work well in our limited experimentation. And it would be interesting to investigate how granular a distribution, the English language is able to support. That is, at what level of resolution do changes in PP distribution become imperceptible in the final message.

3.2 Generation Pipeline

SMS’ message generation pipeline takes as input *Intents* and *Will Power*. For the generation of the initial message, it takes the initial set of *Intents* generated through the process above, and assigns 1 PP to each intent. After this, whenever the user changes the *Intents* and *Will Power* and clicks “Regenerate”, it uses the *Intent* and *Will Power* specified in the interface at that point of time, to generate a new message.

Even with the *Intents* and *Will Power*, naively prompting off-the-shelf models like GPT-4 and ChatGPT does not effectively generate messages with mixed intent. This is because by virtue of being trained on our interactions online, these models encode our cultural biases. So, they naturally privilege clear and unambiguous communication. Further, the RLHF step of training these models, explicitly biases them towards clear communication by attempting to improve their performance on business-critical tasks. So, to get these models to generate messages with mixed intent, we adopt an emerging approach to jailbreak LLMs: *scenario nesting* [12]. By constructing a relevant but fictitious scene and asking the LLM to generate a dialogue for a character in the scene, we are able to successfully get the LLM to generate messages with mixed intent. This suggests, to us, that the deep-seated assumptions of LLMs, much like our own deep-seated assumption, can be challenged and changed through immersion in compelling fiction.

Our generation pipeline consists of three steps: (1) Scenario Generation, (2) Message Initialization, and (3) Weight Enforcement. Next, we describe each of these steps. Figure 3 shows an overview

of the pipeline. Prompt templates for each step are available in our reference implementation (github.com/PranavKhadpe/ambivalent). In the description that follows, we will instantiate prompt templates for our running example.

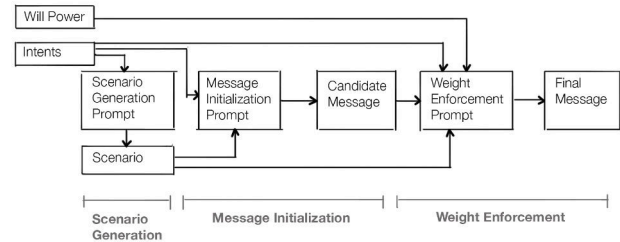


Figure 3: Overview of Generation Pipeline

3.2.1 Scenario Generation. From the *Intents*, we first generate a fictitious scenario. In our running example, this would translate through the following prompt:

```
role: "system"
content: "You are a smart writing assistant helping me write a 3-5 line
scene in which Harry, a character, might say one of several alternatives"
```

```
role: "user"
content: "Harry is a character in a short scene. Generate a single scene
in which Harry might plausibly say any one of the alternatives above to a
single other character with equal justification."
```

```
## Alternatives
"I don't want to hear you talk about this anymore"
"I'm open to continuing our discussion on this topic."
"I would like us to focus on planning our next outing instead."
"This coffee is really amazing"
```

```
## Response Format
Your response should be a JSON object with the key "scene", whose value
is a 3-5 line scene string. The scene description should end in "At
this moment, Harry may say any one of the following:"
```

This returns the following scenario: *“Harry sat across the table, a steaming cup of coffee in each of their hands, as the quiet murmur of the small café surrounded them. His companion looked at him expectantly, waiting for a reply to her incessant talking about yesterday’s events, hopes for the next vacation, and the eternal debate over the town’s best coffee. At this moment, Harry may say any one of the following:”*

3.2.2 Message Initialization. Next, using the *Intents* and the scenario, we generate a candidate message. For our running example, the message initialization prompt is:

```
role: "system"
content: "You are a smart writing assistant. I will give you a scene with
a list of possible ways a character might respond. Your job is to generate
a message that represents the intents of all the alternatives."
```

```
role: "user"
content: "## Scene
Harry sat across the table, a steaming cup of coffee in each of their
hands, as the quiet murmur of the small café surrounded them. His
companion looked at him expectantly, waiting for a reply to her incessant
talking about yesterday's events, hopes for the next vacation, and the
eternal debate over the town's best coffee. At this moment, Harry may say
any one of the following:
"I don't want to hear you talk about this anymore",
"I'm open to continuing our discussion on this topic.",
"I would like us to focus on planning our next outing instead.",
"This coffee is really amazing"
```

Generation Procedure

Consider the person Harry is speaking to. Generate a short message that Harry might say to the person in this situation. The message should combine the intents of ALL the alternatives equally and shouldn't explicate any one of the alternatives. Use pronouns instead of names of people.

Response Format

Your response should be a JSON object with the key "response", whose value is Harry's response combining the intents of all the alternatives."

This returns the following message: *"While I have my thoughts about continuing or shifting our conversation, perhaps we could either delve a bit deeper into what we're discussing or consider mapping out our next adventure. By the way, isn't the coffee here something special?"*

3.2.3 Weight Enforcement. The final step takes the *Intents*, the *Will Power*, and the candidate message and iteratively rewrites the candidate message to ensure that the emphasis on the different intents in the final message is proportional to the PPs allocated to it. For the initial distribution of 1PP per intent, this results in the following prompt:

```
role: "user"
content: "Scene: Harry sat across the table, a steaming cup of coffee in each of their hands, as the quiet murmur of the small café surrounded them. His companion looked at him expectantly, waiting for a reply to her incessant talking about yesterday's events, hopes for the next vacation, and the eternal debate over the town's best coffee. At this moment, Harry may say any one of the following::
```

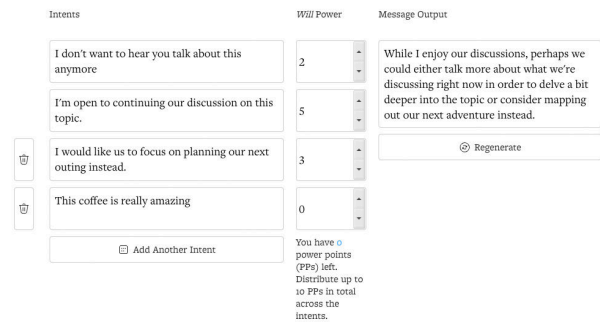
```
Intents: [
  "I don't want to hear you talk about this anymore",
  "I'm open to continuing our discussion on this topic.",
  'I would like us to focus on planning our next outing instead.',
  'This coffee is really amazing'
]
```

```
desired_weights: [0.25, 0.25, 0.25, 0.25]
```

```
Message: While I have my thoughts about continuing or shifting our conversation, perhaps we could either delve a bit deeper into what we're discussing or consider mapping out our next adventure. By the way, isn't the coffee here something special?
```

```
### Do the following 3 times and show your intermediate thinking:
1. Take current version of the message and the list of intents. For each intent, judge the extent to which the message emphasizes the intent. Let's call this weight of the intent. Your output should be an array called actual_weights that contains the weight of each intent in the message and the weights should sum to 1.
2. Now compare the actual_weights array to the desired_weights array. For indices where the actual_weight is less than the desired weight, the intent needs to be emphasized, and where actual_weight is higher than desired_weight, the intent needs to be de-emphasized.
3. Rewrite the message to shift the emphasis.
4. If this is the 3rd iteration of this procedure, output the rewritten message in the following json format:
{final_message: "...", weights: actual_weights}"
```

This returns the message that is finally displayed on the interface: *"Between us, I haven't been quite into the discussion. However, I wouldn't mind if we continue. What excites me is planning our next adventure. Plus, isn't the coffee super amazing?"* The prompt comprises three iterations because in our experiments, it typically took three iterations for the LLM to enforce the weights. The following figure shows the output for the same *Intents* but a different distribution of *Will Power*:



3.3 Implementation Notes

When choosing frameworks for a new project, there are many valid options, each with different tradeoffs. Our foremost priority: frameworks with names that are thematically resonant with the focus of this paper. So, we implemented *SMS* using React ("respond or behave in a particular way in response to something"). For styling, we used Microsoft's UI framework for React, Fluent UI ("(of a person) able to express oneself easily and articulately"). We use OpenAI's gpt-4-0125-preview model API for LLM calls (rumored to be a mixture of experts model to parallel the mixture of intents *SMS* aims to capture). The back-end was, of course, implemented in Express ("convey a thought or feeling in words or by gestures and conduct") to stay on theme. And finally, the back-end is implemented as lexical closures ("an often comforting or satisfying sense of finality") that are hosted as serverless functions on Vercel.

4 DISCUSSION

Our work is an initial exploration of mediums for conveying mixed intents but it points to several interesting opportunities for richer forms of self expression. Here, we discuss two opportunities.

4.1 Progressive Intent Clarification

We do not intend to suggest that clear and unambiguous communication is universally bad. Nor do we intend to suggest that sending mixed signals is always good. Rather, we believe that people have intentions with different degrees of clarity, at different points of time, and that communication mediums that allow for varying levels of clarity can help us better express and experience those various states of being, ranging from indeterminate to determinate. By providing users with control through *Intents* and *Will Power*, *SMS* empowers users to traverse across different levels of determinacy, with high resolution. This granular control can be useful for users to understand and express their evolving intentions: as a user's intentions clarify, over time, they can continue to use *SMS* to send clear-er signals.

4.2 Abstractions for Mixed Intents

In supporting expression of mixed intents, we assume that the mixture is reducible to its constituent intents, even if imperfectly. That is, *SMS'* steering interface asks users to steer its outputs through a collection of intents (*Intents*) and their corresponding weights (*Will Power*), an abstraction that we chose. However, it is possible that a user knows they have a mixed intent but they may be unable to

match it to our abstraction. Therefore, our work is also vulnerable to the same criticisms that we direct at language systems. Similar to how language systems may not be a good abstraction for the expression of mixed intent, a list of intents and weights may also fail to capture our inner representations of mixed intent. Moving forward, we hope this encourages exploration of more abstractions, and perhaps even multi-sensory abstractions, for mixed intents and states of tension. One motivating example is recent work that uses textile patina as an interface for expressing invisible tensions, such as stress, and unspoken needs [14]. Here, users can leverage varying levels of wear and tear of cloth to express otherwise invisible tension.

5 CONCLUSION

For all the “worthy” reasons for which we communicate, like transacting, or working with each other, it may be necessary to limit conversation to clear and unambiguous messages. But when it comes to making ourselves, and our inner experiences, understood, messages with mixed signals ought to be allowed, and perhaps even encouraged. This paper points to a future where computational representations allow us to externalize our mixed intents and inner conflict, and share them with people we care about.

ACKNOWLEDGMENTS

We would like to thank Tejus Gupta and Rishi Veerapaneni for early feedback that helped clarify some of our own thoughts (and intentions). In the weeks leading up to the deadline, conversations with Luke Guerdan have been the source of many exemplars of the kinds of mixed signals we hoped to support with SMS. We thank him for keeping us focused on the mission.

REFERENCES

- [1] Irwin Altman, Anne Vinsel, and Barbara B Brown. 1981. Dialectic conceptions in social psychology: An application to social penetration and privacy regulation. In *Advances in experimental social psychology*. Vol. 14. Elsevier, 107–160.
- [2] Paul M Aoki and Allison Woodruff. 2005. Making space for stories: ambiguity in the design of personal communication systems. In *Proceedings of the SIGCHI conference on Human factors in computing systems*. 181–190.
- [3] Leslie A Baxter. 1988. A dialectical perspective on communication strategies in relationship development. (1988).
- [4] Leslie A Baxter. 1990. Dialectical contradictions in relationship development. *Journal of Social and Personal relationships* 7, 1 (1990), 69–88.
- [5] Peter Blau. 2017. *Exchange and power in social life*. Routledge.
- [6] Peter M Blau. 1960. A theory of social integration. *American journal of Sociology* 65, 6 (1960), 545–556.
- [7] Arthur P Bochner. 1984. The functions of communication in interpersonal bonding. (1984).
- [8] Kirsten Boehner, Rogério DePaula, Paul Dourish, and Phoebe Sengers. 2005. Affect: from information to interaction. In *Proceedings of the 4th decennial conference on Critical computing: between sense and sensibility*. 59–68.
- [9] Kirsten Boehner and Jeffrey T Hancock. 2006. Advancing ambiguity. In *Proceedings of the SIGCHI conference on Human Factors in computing systems*. 103–106.
- [10] Sayan Chaudhry. 2019. Novel Defense Against Code Theft Using Properties of Fibonacci Series. *Proceedings of SIGBOVIK 2019* (2019), 160–162. <https://sigbovik.org/2019/proceedings.pdf>
- [11] Murray S Davis. 1971. That’s interesting! Towards a phenomenology of sociology and a sociology of phenomenology. *Philosophy of the social sciences* 1, 2 (1971), 309–344.
- [12] Peng Ding, Jun Kuang, Dan Ma, Xuezhi Cao, Yunsen Xian, Jiajun Chen, and Shujian Huang. 2023. A Wolf in Sheep’s Clothing: Generalized Nested Jailbreak Prompts can Fool Large Language Models Easily. *arXiv preprint arXiv:2311.08268* (2023).
- [13] Ryan Louie, Andy Coenen, Cheng Zhi Huang, Michael Terry, and Carrie J Cai. 2020. Novice-AI music co-creation via AI-steering tools for deep generative models. In *Proceedings of the 2020 CHI conference on human factors in computing systems*. 1–13.
- [14] Jessica R Mindel, Miles Baird, Allison Prisloe, Feifei Shen, Cherry Wu, Yang Yi, and Hiroshi Ishii. 2023. Pull It Together: Textile Patina as an Interface for Externalizing Invisible Tension. In *Companion Publication of the 2023 ACM Designing Interactive Systems Conference*. 232–236.
- [15] Barbara M Montgomery. 1993. Relationship maintenance versus relationship change: A dialectical dilemma. *Journal of Social and Personal Relationships* 10, 2 (1993), 205–223.
- [16] William Rawlins. 2017. *Friendship matters: Communication, dialectics and the life course*. Routledge.
- [17] Irina Shklovski, Louise Barkhuus, Nis Bornoe, and Joseph Jofish’ Kaye. 2015. Friendship maintenance in the digital age: Applying a relational lens to online social interaction. In *Proceedings of the 18th ACM Conference on Computer Supported Cooperative Work & Social Computing*. 1477–1487.
- [18] Emily West. 2004. *Greeting cards: Individuality and authenticity in mass culture*. University of Pennsylvania.

K

TO OURSELVES (and those named like us)

- 67 Et Al(ex): Examining the impact of Alexs on the field of computer science**

Alex Bellon, Alex Liu, Alex Redding, Alex Snoeren, and Alex Yen

- 68 An Abundance of Katherines: The Game Theory of Baby Naming**

Katy Blumer, Kate Donahue, Katie Fritz, Kate Ivanovich, Katherine Lee, Katie Luo, Cathy Meng, and Katie Van Koevering

Et Al(ex): Examining the impact of Alexs on the field of computer science

Alex {Bellon, Liu, Redding, Snoeren, Yen}
UC San Diego

Abstract—Apropos of nothing and without any sort of personal bias, we collectively investigate the impact that Alexs have had on the field of computer science. We propose the *Alex Factor* (A_{fac}), *Alex Concentration* (A_{conc}), and *Alex Index* (A_{ind}), measurements of how Alex-y a publication is. Additionally, this work aims to break the current record *Alex Index* in a computer science conference/workshop publication¹, with an author list composed of only 5 Alexs.

I. INTRODUCTION

Science (a couple of the authors) has long (since about 2 months ago) wondered just how much impact Alexs have had on the field of computer science. Given the popularity of the name and our own experiences in a department with at least 10 Alexs, we scientifically hypothesized (guessed with no supporting evidence) that the resulting contributions must be nontrivial. There has been previous work investigating characteristics of author names, but they have mainly focused on similar surnames^{2,3,4} or having exactly the same first and surnames^{5,6}, and were mostly focused on economics for some reason. Additionally, none of these works have focused exclusively on the name Alex (a large oversight, in our opinion).

In addition to being a popular name both in general and within computer science publications (see Section III), the name “Alex” (and its common root name, “Alexander”) benefits from having dozens of derivative, diminutive and shortened forms. This bolsters the number of authors that we can claim to be “Alex”s.

We therefore provide a first analysis of computer science authorship by Alexs, in addition to comparisons against authorship of other names. We specifically propose the **Alex Factor** (A_{fac}), **Alex Concentration** (A_{conc}) and **Alex Index** (A_{ind}) as ways to compare the Alex-ness of different publications. Finally, this paper sets the record for the highest A_{ind} (5.0) in a CS conference/workshop paper.

II. METHODOLOGY

We use the DBLP¹ dataset to conduct our investigation, using only entries published in conferences and workshops (excluding journals, book chapters, reference works, etc.). We perform analysis on names as they are entered in DBLP, which means that many of our results are underestimations, as a non-trivial amount of publication entries use only the first letter of first names to identify authors rather than their entire first name. Overall, we draw from a dataset of 3402082

¹There was a journal publication with 9 Alexs :/

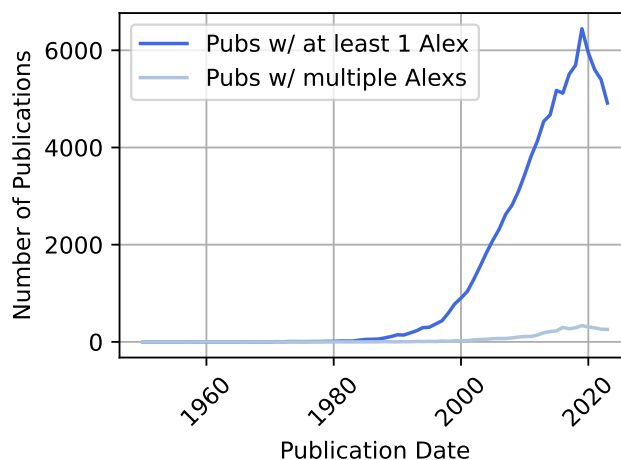


Fig. 1: Number of publications with one or more Alex authors over time

publications, and we restrict our analyses to works published between 1950 and 2023.

A. Alex Metrics

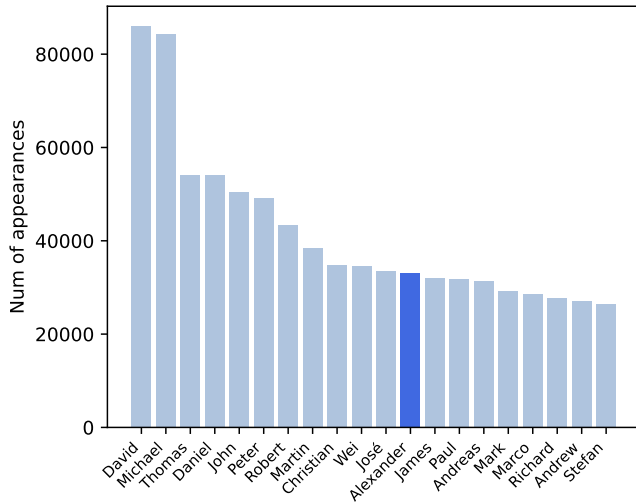
In addition to performing general analyses of author first names in computer science publications, we also propose and calculate 3 metrics: *Alex Factor* (A_{fac}), *Alex Concentration* (A_{conc}), and *Alex Index* (A_{ind}). For all three factors, we include the following names to all count as “Alex” authors, notated by Alex’:

- Alex
- Alessio
- Alessandro
- Alec
- Alessia
- Alessandra
- Aleks
- Alejandro
- Alexander
- Alexei
- Alejandra
- Alexandre
- Alexey
- Aleksander
- Alexandru
- Alexia
- Aleksandar
- Alexandra
- Alexis
- Aleksandr
- Alexandria

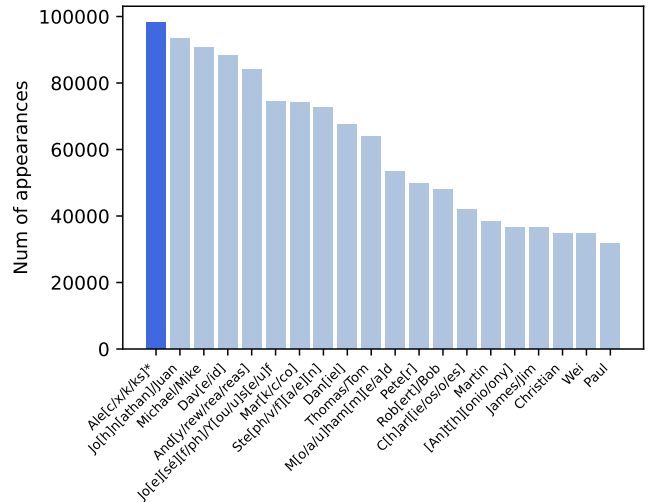
The **Alex Factor** of a paper is calculated via the following formula

$$A_{\text{fac}} = \sum_{i=1}^n \text{Alex}'_i$$

where n is the number of authors in the author list. The **Alex Concentration** of a paper is calculated via the following formula



(a) Non-grouped



(b) Grouped

Fig. 2: Distribution of first names in computer science papers, with and without being grouped by names from the same root

$$A_{\text{conc}} = \frac{A_{\text{fac}}}{n}$$

Finally, the **Alex Index** of a paper can be calculated with

$$A_{\text{ind}} = A_{\text{fac}} \cdot A_{\text{conc}}$$

We use these metrics to compare all the publications in our dataset to find the most “Alex”-y publications. There is still future work needed to determine whether similar metrics could be calculated for names other than Alex. Given thousands of first names exist in the world, we believe that this will be a monumental task best left to others not scrambling to submit to SIGBOVIK before the deadline.

III. RESULTS

We calculate some general statistics about computer science paper authorship, then focus specifically on Alexs, including comparing papers by our three metrics: A_{ind} , A_{fac} , and A_{conc} .

A. General author name patterns

We plot the distributions of occurrences of different first names in Figure 2. We plot both by raw first names as they appear in DBLP, and after grouping first names by the same root name. We find that for raw names (Figure 2a), computer science authorship is dominated by Davids and Michaels, with those two names being over 20,000 publications ahead of the next most common name. Unfortunately, the highest ranked “Alex” name, “Alexander” only places 12th by number of appearances.

Accordingly, we perform a secondary analysis in order to give ourselves a more impressive ranking. In Figure 2b, we combine names based on shared root name, in order to get a more general sense of authorship name patterns. The John-derived names jump up to first place from fifth in the raw ranking, knocking down Michael- and David-derived names

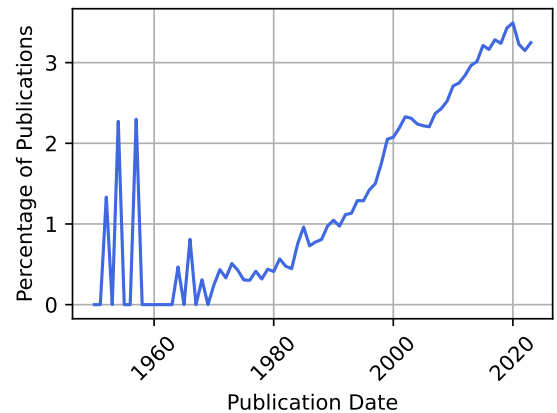


Fig. 3: Percentage of Alex publications over time

to third and fourth respectively. Fortunately, the Alexs have a much better showing after including all Alex variants, taking first place with almost 100,000 appearances. Additionally overall, the number of publications does not drop off as steeply after the top names.

Another interesting observation is that based on the difference in rankings between Figure 2a and Figure 2b, it is more common to see longer versions of names (e.g. “Alexander” and “Michael”) as opposed to shorter nicknames (e.g. “Alex” and “Mike”). This is likely due to the fact that longer names seem more fancy and official².

B. Alex patterns

Within the space of Alex-related names, we find multiple interesting patterns. We first investigate one of our most pressing questions: how much computer science research is produced by Alexs? In Figure 3, we plot the percentage of

²Source: idk, seems true

An Abundance of Katherines*

The Game Theory of Baby Naming

Katy Blumer¹,
Kate Donahue¹,
Katie Fritz¹,
Kate Ivanovich¹,
Katherine Lee¹,
Katie Luo¹,
Cathy Meng¹, and
Katie Van Koevering¹

¹Cornell University

1 Introduction

The most important decision in any child's life happens shortly after they are born and is made entirely without their input or approval - their naming. The name given to a child is traditionally kept throughout their life time and has significant impact on their future life path. This momentous decision is made by parents with little education in the game theory inherently present in the highly competitive field of naming.

We attempt to assist these parents with a simple primer into the game theory underpinning the decision of naming. We will introduce the basic set-up of the naming game and formalize the parameters and incentives, wherein parents have some desired properties of the name. We then describe the pitfalls of the most simple interpretations of these models, particularly the dangers of myopic action. Finally, we present experimental results demonstrating the shift in name distributions under this model. These experiments underline the inherent risks in naming a child and highlight how altering various parameters can change the outcomes of naming strategies.

2 Related works

Surprisingly, no one has ever done any research on naming strategies (so long as you conveniently ignore [4, 5, 6, 7, 8, 9, 10, 12, 13, 14, 16, 17, 18, 19, 20, 22, 23, 25] and likely other work).

3 Model

Naming a child is akin to choosing an outfit for the Oscars. It must be unique enough to stand out - no one wants to show up to the Oscars in the same dress - but it must also be similar enough to be recognizable as a name. Lady Gaga's meat dress is fine for an afternoon, but a child named "Meat Dress" would soon become discontented, if not the plaintiff of a lawsuit. Thus, we model name selection based on the desired "uniqueness" of the name.

*With apologies to John Green [11].

3.1 Name frequency and choice model

First, we present our formal model. Assume there is some set of names \mathcal{A} . At a time point i , we assume there exists a discrete distribution of popularity over names

$$f_i(a) = \nu$$

such that $f(a) \in [0, 1]$, $\sum_{a \in \mathcal{A}} f_i(a) = 1$. For simplicity, we assume that every name has unique frequency: that is, $f_i(a_j) \neq f_i(a_k)$ for $j \neq k$.

Next, we model parental preferences. It is well-known that parents are always in complete agreement over the name they would prefer to pick for their newborn child. Therefore, we will treat the parents of each child as a unit, and assume that each set of parents $j \in \mathcal{P}$ has some preference over the proportion of the population that would share the same name as their child (we also assume each parental unit has exactly 1 child). For example, $\mu_j = 0.01$ means that parental unit j wants their child to have a name that is shared by 1% of the population. We will use

$$g(\mu) = p$$

to mean that $p \in [0, 1]$ proportion of parents want a name with popularity μ . For example, if $\mu = 0.1$ and $p = 0.2$, then this means that 20% of parents want a name with popularity 10%. This set-up gives us convenient parameters for the model and just enough Greek letters to sound smart enough for publication.

In general, we will assume that parents are *myopic*, with new parents having no concept of time but perfect access to baby name data. We find this a realistic assumption. Mathematically, parents at time step i will pick the name a that currently is closest to their desired frequency μ_j . Given this assumption, then at time step $i + 1$, the proportion of babies who have name a is given by the total fraction of parents for whom name a is closest to their desired frequency.

3.2 Satisfiability

If parents are unable to infer the consequences of their actions and act myopically, then it can immediately be seen that some parents will be deeply unhappy with said consequences: for example, if $g(0.1) = 0.2$ (as in the example above), then the 20% of parents who wished that their name has popularity 10%, will end up with a name that is more popular than they anticipated (when $g(\mu) > \mu$). For instance, a parent might anticipate the name “Kate” would be a pleasantly traditional yet unique name with only moderate popularity. They would be wrong [24].

Conversely, if we had $g(0.1) = 0.05$ (or $g(\mu) < \mu$), then parents would end up with a name that is less popular than they anticipated. If $g(\mu) = \mu$, then we say the name with proportion μ is *satisfied*. Ideally, we would like every name to be satisfied, or $g(\mu) = \mu$ for all $\mu \in [0, 1]$. However, that would give us a distribution with total probability > 1 , which is the sort of thing that makes statisticians sad. Instead, we see that the entire naming distribution $g(\cdot)$ is satisfied if every name with *nonzero popularity* $\mu > 0$ has exactly μ fraction of the population that desires this name:

$$\begin{cases} g(\mu) = \mu & \mu > 0 \\ g(\mu) = 0 & \text{otherwise} \end{cases}$$

3.3 Stability

Next, we consider an alternative property that we may wish to have: that the distribution of names be *stable*. If an arrangement is stable, this means that given an existing distribution $f_i(a)$ and a parental preference distribution $g(\mu)$, every name’s frequency will be exactly the same at the next time step $i + 1$:

$$f_{i+1}(a) = f_i(a) \quad \forall a \in \mathcal{A}$$

The simplest way to achieve stability is to merely assume that every parent would prefer to name their child after themselves. That is, every parent wishes their child to have the same “uniqueness” of naming they do, a sort of inheritability of uniqueness - we name this the Dweezil Principle.

3.4 Extremely Reasonable Assumptions

The above model contains several Extremely Reasonable Assumptions (ERAs). The first ERA is the very conservative assumption that there is only one gender, with all children and all names adhering to the same gender. Thus any child may be given any name, so long as it exists in the names list¹. Another ERA is the Mayfly Parenthood Assumption, in which all parents perish immediately upon naming their child, which makes the math substantially easier.

4 Illustrative example: power law distribution

In this section, we consider the case where both $f(\cdot)$ and $g(\cdot)$ are given by *power law distributions*.

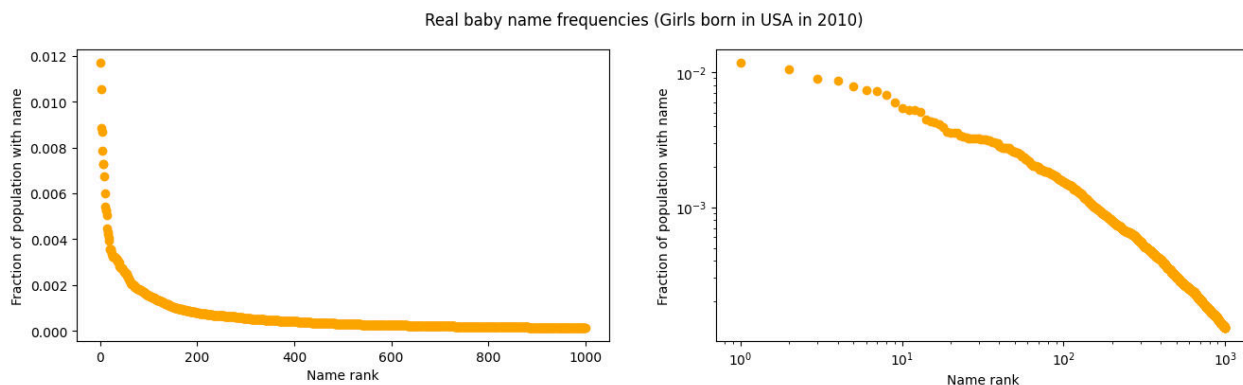


Figure 1: Name frequencies from the Social Security Administration, for girls born in 2010[2]. Note the rough power-law shape.

4.1 Modeling $f(a), g(\mu)$

We begin by defining our variable for name popularity. The popularity of names has been shown to follow a power law distribution [18] (also see Figure 1) . We can model this as:

$$f(a) = K \cdot a^{-t}$$

where a denotes the rAnk of the name within $a \in [1, N]$, t is another constanT, and K is a normalization konstant.

Next, we consider how parents pick the uniqueness of names for their children: the function $g(\cdot)$. Because we are in the Power Law subsection, we will also assume that this distribution of parent preferences is a power law. To minimize our use of variables, we will model this as:

$$g(a) = K' \cdot (a')^{-t'}$$

where a' is the desired frequency of the name (within the range $[\epsilon, 1]$, for $\epsilon > 0$), t' is another constant, and K' is a normalization constant.

¹If a fixed names list is good enough for the Scandinavians, it's good enough for us [1]

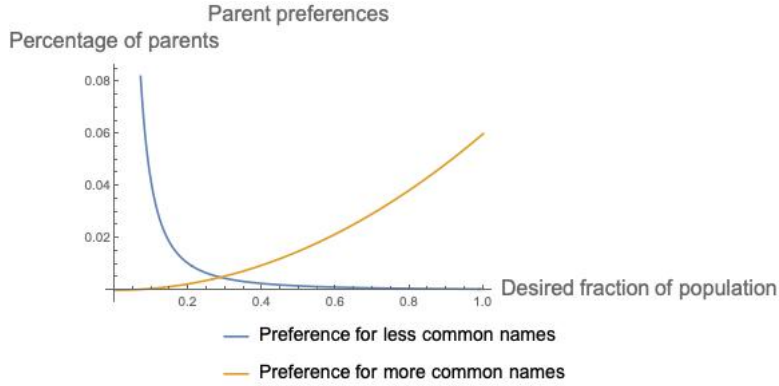


Figure 2: Examples of different parental preferences: a preference for less common names (blue) and preference for more common names (orange).

Note that for $t' > 0$, we have that $g(\mu)$ is decreasing in μ : that is, more parents prefer names that are *uncommon*. Conversely, for $t' < 0$, we have that $g(\mu)$ is increasing in μ : parents prefer names that are *common*.

4.2 Picking names: stability

Given a $f_i(a)$ and $g(\mu)$, the distribution $f_{i+1}(a)$ at the next time step is given by:

$$f_{i+1}(a) = g(f_i(a)) = K' \cdot (K \cdot a^{-t})^{-t'} = K' \cdot K^{-t'} \cdot a^{t \cdot t'}$$

Note that this is again a power law distribution², with parameter $t \cdot t'$. Next, we can analyze the properties of the resulting distribution.

4.2.1 Picking an uncommon name: a futile quest

First, we consider the case where t' is positive, which corresponds to the case where most parents prefer uncommon names. In this setting, we know that $t \cdot t'$ is also positive, which means that $f_{i+1}(a)$ is *increasing* in a . This tells us that a name with high popularity at time i will have *low* popularity at time $i + 1$. Given an original distribution of frequency over names shown in blue in Figure 4, if parents have a preference for less common names, the resulting distribution will look like the orange curve: the least popular name has suddenly become extremely popular, and what was popular at time step i has become horribly passé by time step $i + 1$. This means that names will see-saw in popularity from one time step to another, which explains why both your grandmother and niece are named Mabel (Figure 3).

4.2.2 Picking a common name: naming event horizon

However, parents might believe there are benefits to their child sharing a name with others (see [24]). We can model this with the case where t' is negative, which means that most parents would prefer names that are relatively popular, with only a few parents preferring names that are less common. Mathematically, this implies that $t \cdot t'$ is also negative, which means that $f_{i+1}(a)$ is *decreasing* in a . This means that popular names at time step i are also popular at time step $i + 1$: the relative order of name popularity stays the same (as shown in the green curve in Figure 4).

If this process is repeated n times, with the same parental preferences, then the resulting power law distribution would have exponent $t \cdot t^n$. For $t' > 1$ (the “event horizon”), the means that the most popular

²This is why we like power laws.

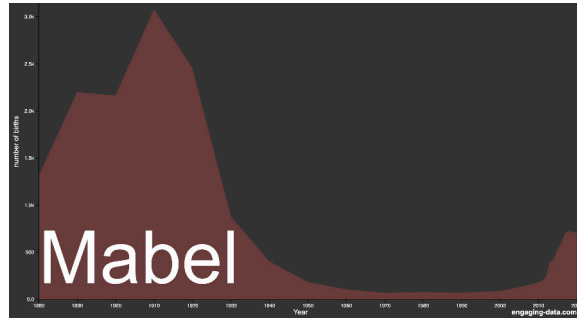


Figure 3: Frequency of name Mabel (image from [3]).

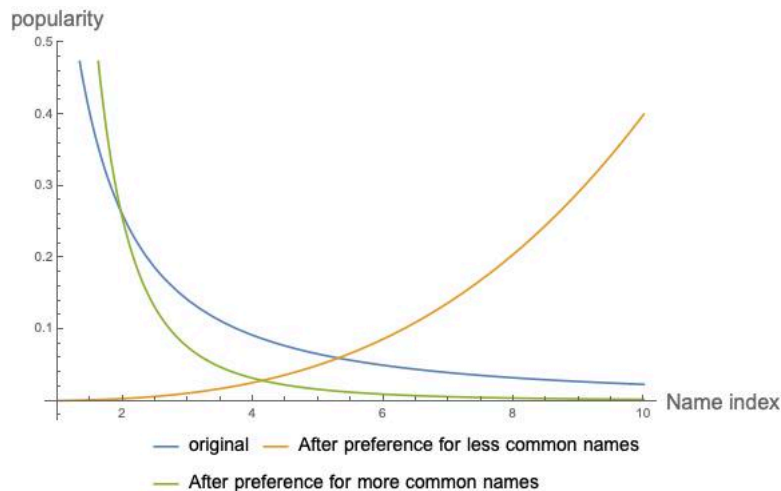


Figure 4: Examples of naming frequency at time step $i + 1$: the original distribution (blue), the distribution after a preference for less common names (orange), and a preference for more common names (green).

names gobble up almost all of the population, resulting in a black hole of names wherein infinite density will eventually fall upon just one name.

5 Simulations

For simulations we use a log-normal distribution of parent preferences, rather than a power law as in Section 4, because a certain author was having issues with SciPy. We baselessly claim a log-normal makes sense because name “uniqueness” is logarithmic; that is, a name belonging to 0.01% of the population is roughly twice as unique as a name belonging to 0.1% of the population (when comparing to a baseline name with 1% popularity).

We work with five different parent preference distributions, with popularity modes of 1% to 0.0001% (Figure 5). From each preference distribution, we generate a sample of 1 million parents, who each choose the name which is currently closest to their desired popularity. The resulting distribution of names is shown in Figure 6. If parents prefer more unique names, the name distribution flattens, since there are many more names at the low-popularity end of the scale for parents to divide between. If parents prefer very popular names, the distribution is even more heavily weighted at the popular end than the original power law distribution. The preference distribution centered at 0.1% in fact closely resembles the original power law distribution.

Measures of parent error (distance from goal popularity) don’t show particularly interesting patterns, but are included because they look like dinosaurs (Fig. 7).

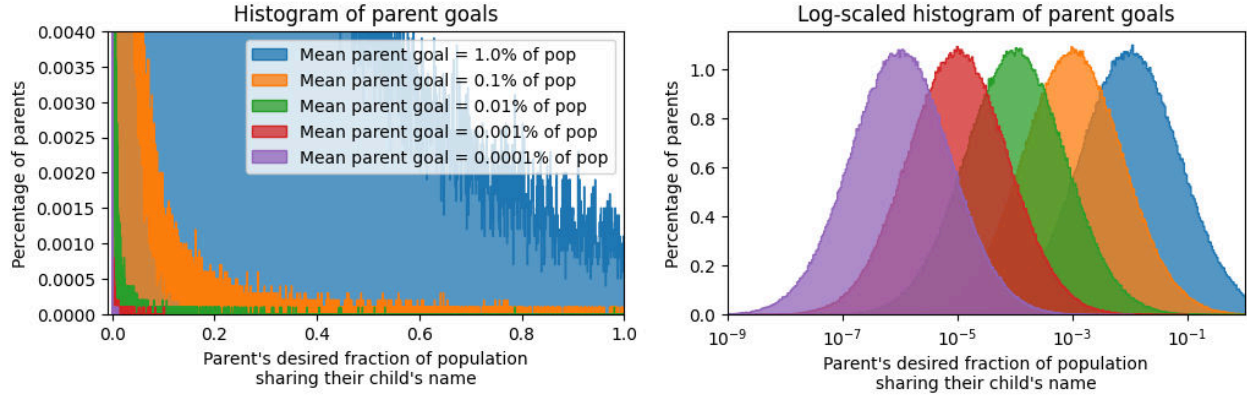


Figure 5: Parent preference distributions used in experiments. Logarithmic x-axis on the right, showing the log-normal distribution shape. Each histogram represents a sample of 1 million parents from the given distribution.

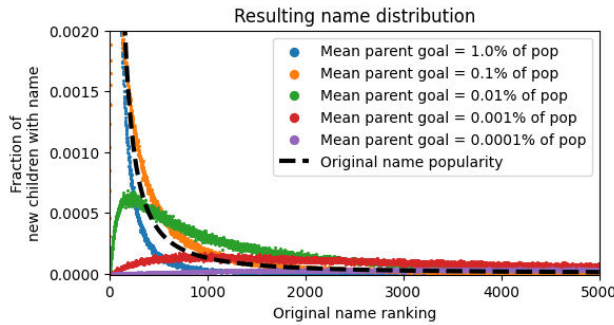


Figure 6: Resulting name distributions for children of parents with given preference distribution.

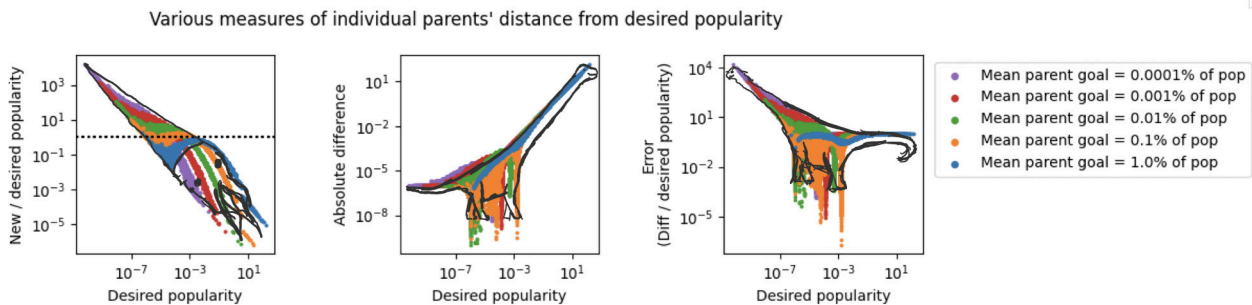


Figure 7: Various measures of parent error, meaning the difference between their chosen name's popularity and their desired name popularity. From left to right: ratio (true popularity divided by desired popularity), absolute difference (absolute value of true–desired popularity), and error (absolute difference divided by desired popularity). Note the adorable dinosaurs.

6 Obligatory Kat-GPT experiment

Because this paper was written in 2024, we include an obligatory section involving generative AI and LLMs. It is fortunate that the most popular LLM of the year appears to be custom-made for this experiment: Kat-GPT. Specifically, we asked Kat-GPT to give us its top ten names for a) a girl, b) a boy, and c) a gender-neutral name (see Figure 8). Then, we calculated the frequency of those names within their gender

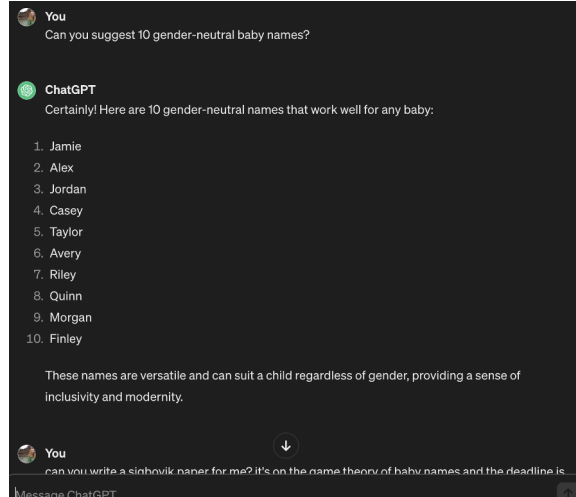


Figure 8: Experiments with Kat-GPT

category (or within the total population, for gender-neutral). The results are given within Table 1.

The difference in mean popularity between “girl” and “gender-neutral” names and between “boy” and “gender-neutral” names are both statistically significant at the $p \leq 0.05$ level, while the difference in mean popularity between “girl” and “boy” names is statistically significant at the $p \geq 0.05$ level.

	“girl” names	“boy” names	“gender-neutral” names
Mean popularity	0.6728	0.666	0.126
Std of popularity	0.134	0.212	0.0619

Table 1: Popularity of names given by ChatGPT ($N = 10^{-e^{\pi \cdot i}}$ queries for each category)

7 Extensions & Future Work

In this section we include multiple extensions that we considered but ultimately were too lazy to actually finish.

7.1 Creation of new names

One of our Extremely Reasonable Assumptions is that there was a fixed list of names. However, new names have occasionally been documented in the wild [21]. One extension could consider a strategy whereby parents could pick a name a' that has some distance $d(a, a')$ from an established name, and derive cost relating both to the popularity of the “base name” a and the distance $d(a, a')$.

$$\min_{a, a' \in \mathcal{A}} \left| K \cdot a'^{-t} - \mu_i \right| + \lambda \cdot d(a, a')$$

where $d(a, a')$ is a distance metric between names (e.g., see Figure 9).

7.2 Non-myopic parents

Another one of our Extremely Reasonable Assumptions was that parents are myopic: that is, if they have a desired name frequency a' , they simply blindly pick the name that currently has frequency a' , which can

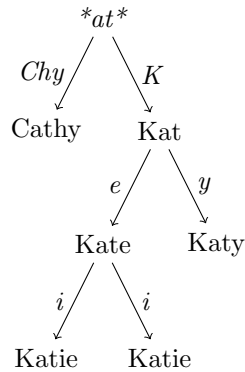


Figure 9: An example of name mutation, where multiple novel (but related) names are displayed with their edit distance from each other.

lead to over- or under-shooting their desired frequency. However, emerging research [15] suggests that people in fact reason strategically about their actions; if this is borne out, considering non-myopic parents may be an interesting avenue of future research (i.e. not us).

8 Conclusions and implications

The science of naming has a long and illustrious history that we didn't bother to look at. Instead, we arbitrarily assigned a new(?) model to describe how parents ought to name their children - namely probabilistically. This model has interesting implications, most interestingly that all naming strategies are futile. We also print some plots, for both educational and entertainment purposes, which further emphasize these points and have some nice dinosaurs. But overall, we find only one rule really matters when naming a child: when in doubt, name it $\text{Kat}[i?e|y]$.

9 Acknowledgements

We thank F.D.C. Willard for helpful discussions. Simon Shindler contributed significantly to the aesthetic of Figure 7, but could not be named an author for obvious reasons.

We also want to thank the many dozens of people who have confused us for one another at conferences. We enjoyed meeting you.

Finally, we thank Katie's mom for editing this paper and the rest of the parents for making the provably optimal (see above) name choice.

References

- [1] Some countries have a list of preapproved baby names. *Interesting Facts*.
- [2] Popular baby names, 2024. <https://www.ssa.gov/OACT/babynames/index.html>.
- [3] US baby name popularity visualizer, 2024. <https://engaging-data.com/baby-name-visualizer/>.
- [4] ACERBI, A., AND BENTLEY, R. A. Biases in cultural transmission shape the turnover of popular traits. *Evolution and Human Behavior* 35, 3 (2014), 228–236.
- [5] BENTLEY, R. A., AND ORMEROD, P. Accelerated innovation and increased spatial diversity of us popular culture. *Advances in Complex Systems* 15, 01n02 (2012), 1150011.

- [6] BUSH, S. J., POWELL-SMITH, A., AND FREEMAN, T. C. Network analysis of the social and demographic influences on name choice within the uk (1838-2016). *PLoS One* 13, 10 (2018), e0205759.
- [7] CILA, J., AND LALONDE, R. N. What’s in a name? motivations for baby-naming in multicultural contexts. *Contemporary Language Motivation Theory: 60 Years Since Gardner and Lambert (1959)* 3 (2019), 130.
- [8] COULMONT, B., SUPERVIE, V., AND BREBAN, R. The diffusion dynamics of choice: From durable goods markets to fashion first names. *Complexity* 21, S1 (2016), 362–369.
- [9] GLYNN, M. A., AND ABZUG, R. Institutionalizing identity: Symbolic isomorphism and organizational names. *Academy of Management journal* 45, 1 (2002), 267–280.
- [10] GOLMAN, R., BUGBEE, E. H., JAIN, A., AND SARAF, S. Hipsters and the cool: A game theoretic analysis of identity expression, trends, and fads. *Psychological review* 129, 1 (2022), 4.
- [11] GREEN, J. *An Abundance of Katherines*. Penguin Young Readers Group, 2006.
- [12] HAHN, M. W., AND BENTLEY, R. A. Drift as a mechanism for cultural change: an example from baby names. *Proceedings of the Royal Society of London. Series B: Biological Sciences* 270, suppl.1 (2003), S120–S123.
- [13] HAN-WU-SHUANG, B., HUA-JIAN, C., AND YI-MING, J. Understanding the rise of unique names: The emphasis on uniqueness matters.
- [14] KESSLER, D. A., MARUVKA, Y. E., OUREN, J., AND SHNERB, N. M. You name it—how memory and delay govern first name dynamics. *PloS one* 7, 6 (2012), e38790.
- [15] LASKER, E. *Lasker’s Manual of Chess: With 308 Diagrams*. Printing-Craft, 1932.
- [16] LEE, M. J., DO YI, S., KIM, B. J., AND BAEK, S. K. Theory of fads: Traveling-wave solution of evolutionary dynamics in a one-dimensional trait space. *Physical Review E* 91, 1 (2015), 012815.
- [17] LEONARDELLI, G. J., PICKETT, C. L., AND BREWER, M. B. Optimal distinctiveness theory: A framework for social identity, social cognition, and intergroup relations. In *Advances in experimental social psychology*, vol. 43. Elsevier, 2010, pp. 63–113.
- [18] LI, W. Analyses of baby name popularity distribution in us for the last 131 years. *Complexity* 18, 1 (2012), 44–50.
- [19] MÜLLER, J. *Analyzing given names: exploring potentials for a personalized name discovery on Nameling.net*. PhD thesis, University of Kassel, Germany, 2020.
- [20] NEWBERRY, M. G., AND PLOTKIN, J. B. Measuring frequency-dependent selection in culture. *Nature Human Behaviour* 6, 8 (2022), 1048–1055.
- [21] OKASHA, E. *Women’s names in Old English*. Routledge, 2016.
- [22] STOYNEVA, I., AND VRACHEVA, V. Demystifying entrepreneurial name choice: insights from the us biotech industry. *New England Journal of Entrepreneurship* 25, 2 (2022), 121–143.
- [23] TWENGE, J. M., DAWSON, L., AND CAMPBELL, W. K. Still standing out: Children’s names in the united states during the great recession and correlations with economic indicators. *Journal of Applied Social Psychology* 46, 11 (2016), 663–670.
- [24] VAN KOEVERING, K., MENG, K., LUO, K., IVANOVICH, K., FRITZ, K., DONAHUE, K., BLUMER, K., AND LEE, K. An abundance of katherines: The game theory of baby naming. *SIGBOVIK* (2024).
- [25] VARNUM, M. E., AND KITAYAMA, S. What’s in a name? popular names are less common on frontiers. *Psychological science* 22, 2 (2011), 176–183.

AND OUR POSTERITY

- 69 A Genius Solution: Applications of the Sprague-Grundy Theorem to Korean Reality TV**

Jed Grabman

- 70 A Secondhand Understanding of Reality: Infinite Craft Subtleties**

Ryan Pitasky, Kevin Gomez, and John Cesarz

A Genius Solution

Applications of the Sprague-Grundy Theorem to Korean Reality TV

Jed Grabman

Abstract

Monorail is a two-person tile placing game that was popularized by the South Korean reality show **The Genius** in 2014. This paper uses the Sprague-Grundy theorem to demonstrate a winning strategy for the first player.

Background

The Genius was a South Korean reality show that aired for 4 seasons, from 2013 to 2015, where players were eliminated through a series of social/strategy games. The prize for winning each season ranged from approximately \$50,000 to \$100,000. Additionally, a Dutch version of the show aired for 1 season in 2022 (which had no prize). In both versions of the show, each episode would culminate in a “Death Match”, a game that was used to determine which of two contestants would be eliminated. **Monorail** is a tile-placing game that was used as a Death Match twice during the Korean series, and two additional times during the Dutch series (where it was called “On Track”). While some players claimed to know a winning strategy, their strategies were never rigorously presented on the program and they never played **Monorail** during the series, so it is impossible to determine whether those strategies were correct and complete.

Rules

Monorail is a tile-placing game where the objective is to finish a looping train track. The game is played with a set of 16 identical double-sided square tiles with train tracks on them. One side of the tile has the train tracks continuing in a straight line, while the other side has the train tracks turning 90 degrees (Figure 1).

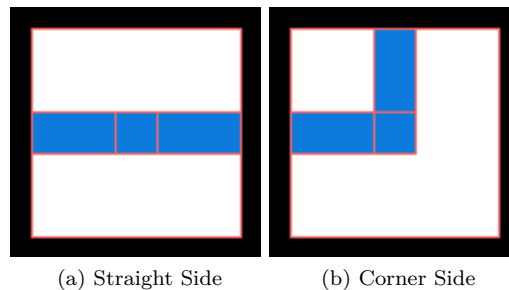


Figure 1: Sides of tiles

Before the start of the game, a special double-wide tile is placed that represents the train station (Figure 2). On each player’s turn, they may place 1, 2 or 3 tiles, subject to the following conditions:

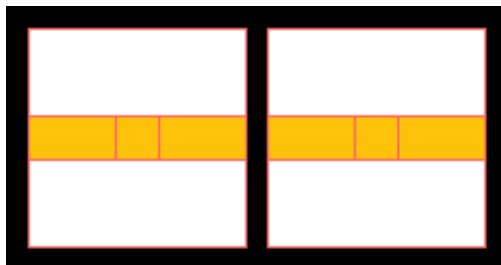


Figure 2: Station

- At least 1 of the placed tiles must be orthogonally adjacent to a previously placed tile
- All placed tiles must be adjacent and in a straight line (either horizontally, or vertically)

The player chooses which side of the tile is face-up and its rotation before placing it. This is equivalent to choosing which two of the tile's four sides to connect with train tracks.

Additionally, the tiles placed by the player do not need to connect to the existing track, nor does the track need to be connected between the placed tiles (e.g. Figure 3). The move restrictions are on where the tiles are placed, not the orientation of the tiles.

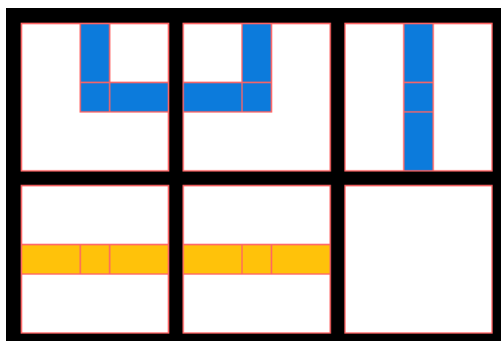


Figure 3: Example starting move

The goal of the game is to be the player that finishes a loop that connects the start and end of the station (Figure 4).

A player wins if they place tiles that finish a loop, provided that the loop uses all the tiles that have been placed. It is ok if some tiles remain unplaced.

In some cases, it may become impossible to complete a loop with the remaining tiles (e.g. Figure 5). Therefore, a player also has the option of declaring that the track cannot be finished instead of placing tiles. If they are correct, then they win. However, if their opponent can demonstrate that the track could have been finished, the opponent wins instead.

Rule Modifications for Analysis

It is clear that if the track cannot be completed, it is in the active player's best interest to declare this fact and immediately win. Thus, in a game between perfect strategists, it would never be to a player's benefit to make a move that resulted in a track that cannot be completed¹. Additionally, if a player starts their turn in a position where the track is completable, they may always make a move that results in a track that

¹Of course, real world players can make mistakes. One player won an otherwise losing game by confidently playing a move that made the track impossible to finish. This confused their opponent who did not claim the win on their turn.

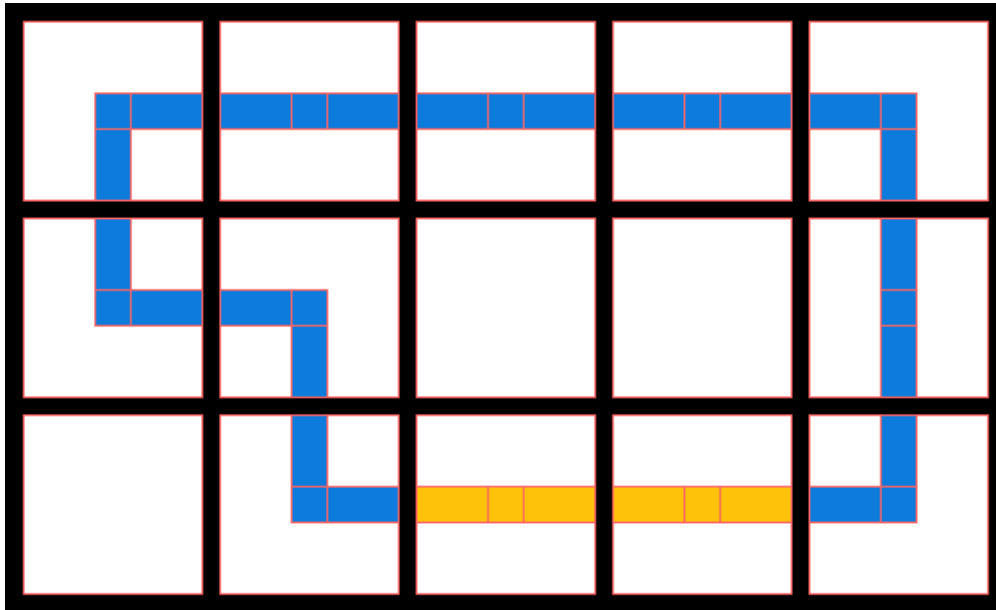
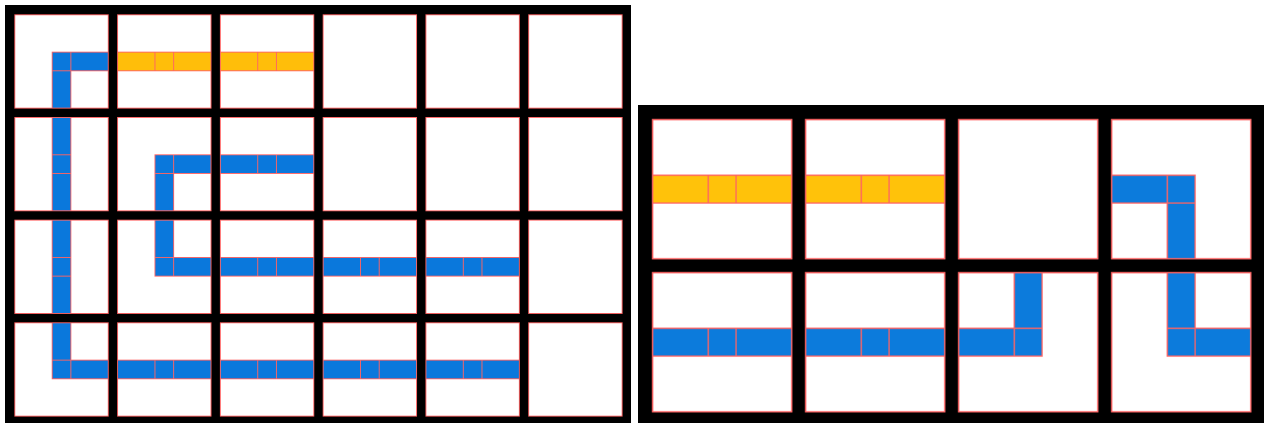


Figure 4: Example finished game



(a) Not enough tiles remaining

(b) Impossible junction

Figure 5: Impossible to complete a loop

is completable by extending the track one tile along one of the possible routes. Thus, a player will never be forced to make a move that turns a completable track into one that cannot be completed. Therefore, the rules can be modified so that a player may not make a move that results in a track that's impossible to complete. This simplifies the assessment of positions without affecting whether each position is winning or losing for the player to move.

Mathematical background

Impartial games

The Sprague-Grundy Theorem simplifies evaluating impartial games, which are games that meet the following conditions:

- There are two players
- The players alternate turns
- There are a finite number of possible gamestates
- There is a maximum number of turns
- The last player to move wins and the other player loses
- There is no randomness
- A player's possible moves depend only on the state of the board, not which player's turn it is.²

It should be clear that **Monorail** satisfies each of these properties.

Impartial subgames There are times when the state of an impartial game may be described as the combination of multiple non-interacting sections of the game state. For example, consider the game of **Nim**. In **Nim**, there are some number of distinct piles, each containing stones. On a player's turn, they may remove any number of stones from any single pile. The player that removes the final stone from the last remaining pile wins.

Because removing stones from one pile doesn't affect how many stones are in any other pile, a game of **Nim** with multiple piles could be considered a combination of multiple games of **Nim** that each have fewer piles.

In general, we can say that an impartial game G is the sum of impartial subgames $G = G_1 \oplus G_2 \oplus \dots \oplus G_n$ if it satisfies the following conditions:

- On a player's turn, they must choose a particular subgame (pile) G_i where a legal move is possible and make a move in that subgame
- Making a move in one subgame never affects which moves are available in any other subgame
- The player who makes the last move wins (i.e. after their move, no subgame has a legal move remaining)

In fact, this process also generalizes in reverse. Any finite subset of impartial games $\{G_1, \dots, G_n\}$ can be used to create a larger impartial game G if players are able to choose which subgame to play in for each turn.

The Sprague-Grundy Theorem The Sprague-Grundy theorem establishes the existence of a function N that maps game-states of impartial games to the non-negative integers with the following properties:

- $N(G) = 0 \iff$ The player to move is losing
- $N(G_1 \oplus G_2) = \text{XOR}(N(G_1), N(G_2))$ (where XOR is the bitwise exclusive-or function).
- If $N(G) = k$, and $k \geq m > 0$ there exists a move to a position G' such that $N(G') = m$ ³

Table 1: Bitwise XOR sums ((0-7) x (0-7))

XOR	0	1	2	3	4	5	6	7
0	0	1	2	3	4	5	6	7
1	1	0	3	2	5	4	7	6
2	2	3	0	1	6	7	4	5
3	3	2	1	0	7	6	5	4
4	4	5	6	7	0	1	2	3
5	5	4	7	6	1	0	3	2
6	6	7	4	5	2	3	0	1
7	7	6	5	4	3	2	1	0

Because there is a bijection between the equivalence classes of N and the sizes of a single pile game of **Nim**, $N(G)$ is often called the **Nimber** of G . In particular, a game of **Nim** with 1 pile and K stones has a number of K .

These properties lead to a few immediate important consequences:

1. If $N(G) \neq 0$, the player is winning and they can make a move that leaves their opponent with a losing state, i.e. the new state G' is such that $N(G') = 0$.
2. Because $\text{XOR}(A, B) = 0 \iff A = B$,
 $N(G_1 \oplus G_2) = 0 \iff N(G_1) = N(G_2)$
3. In any position in an impartial game, if $N(G) \neq 0$, one may create a new game combining G with a single pile game of **Nim** with $N(G)$ stones. In this new game, the first player to move will be losing.

Monorail Solution.

A computer search of gamestates in Monorail found a winning move for the first player. An examination of the move showed that it uniquely determined the loop that could be finished with the remaining tiles (Figure 6).

The remaining locations where tiles can be placed are separated into 2 non-adjacent areas, which can be thought of as separate subgames because moves in one area do not affect possible legal moves in the other area. After this point, the location where a tile is placed also uniquely determines how it must be rotated. Therefore, for the purposes of analyzing subgames it is sufficient merely to look at the positions where tiles need to be placed and not their rotation.

It is clear that rotating, translating or mirroring a subgame does not meaningfully affect the moves available. If two positions are made up of equivalent subgames, their evaluation must be the same. For example, the 2 positions in Figure 7 have 8 tiles remaining unplaced, and a move must be either in the group of 5 tiles forming a long L, or the set of 3 tiles in a line. These two positions must therefore have the same number.

However, care must be taken when determining which subgames are identical. A move may only be made if at least one of the tiles placed is adjacent to an already placed tile. For example, the two positions in Figure 8 look similar because the three remaining spots to play a tile are in an L shape. However, in position a , it is legal to place only the corner piece, while in position b it is not. In fact, playing in the corner is the only winning move in position a , while b is a losing position.

²The term “impartial game” is often used to refer only to this last condition. However, game theory discussion often focuses on the class of games that meet the other criteria, and discusses Sprague-Grundy’s application to impartial games in that narrower context, so for brevity this paper uses the term “impartial game” to refer to games that meet all the listed criteria.

³Technically, it should be specified whether the domain of N is games or positions. In practice, any position could also be the starting position of a game with the same rules, so the ambiguity doesn’t cause issues.

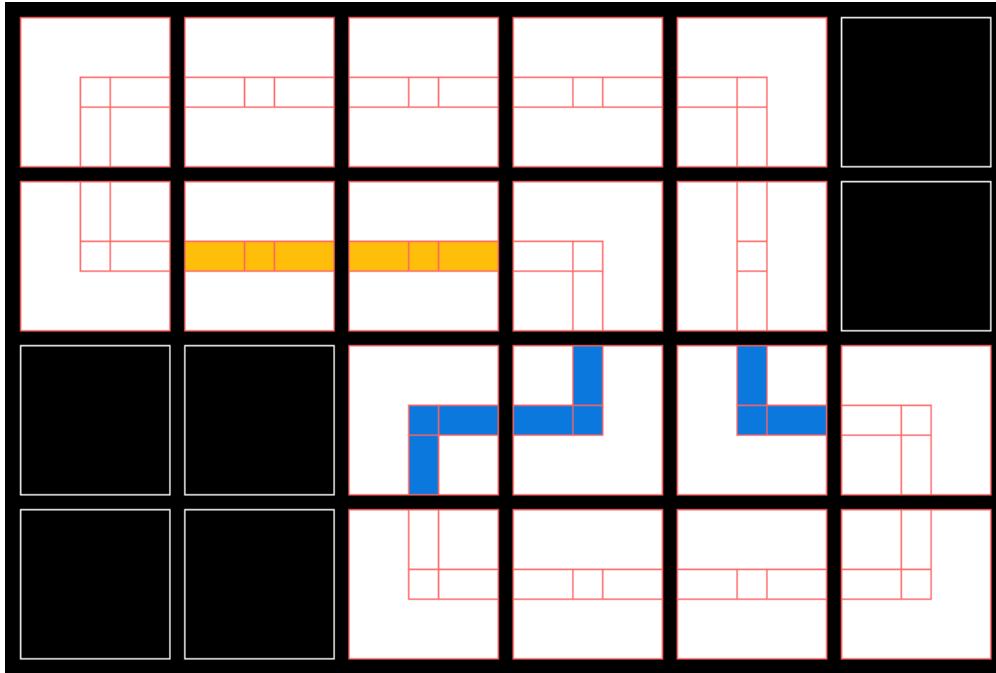


Figure 6: Winning move with only possible path

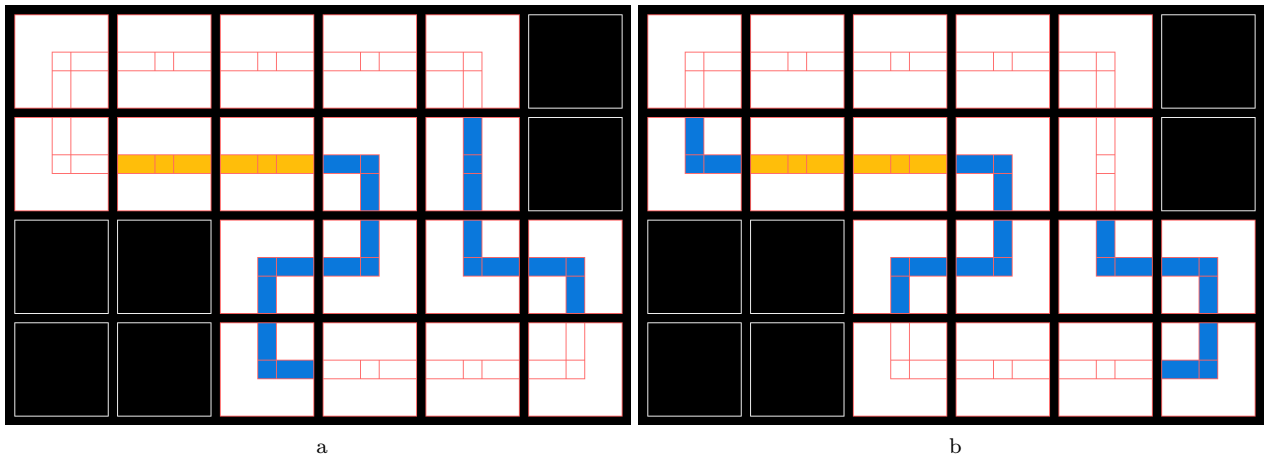


Figure 7: Positions from the same class

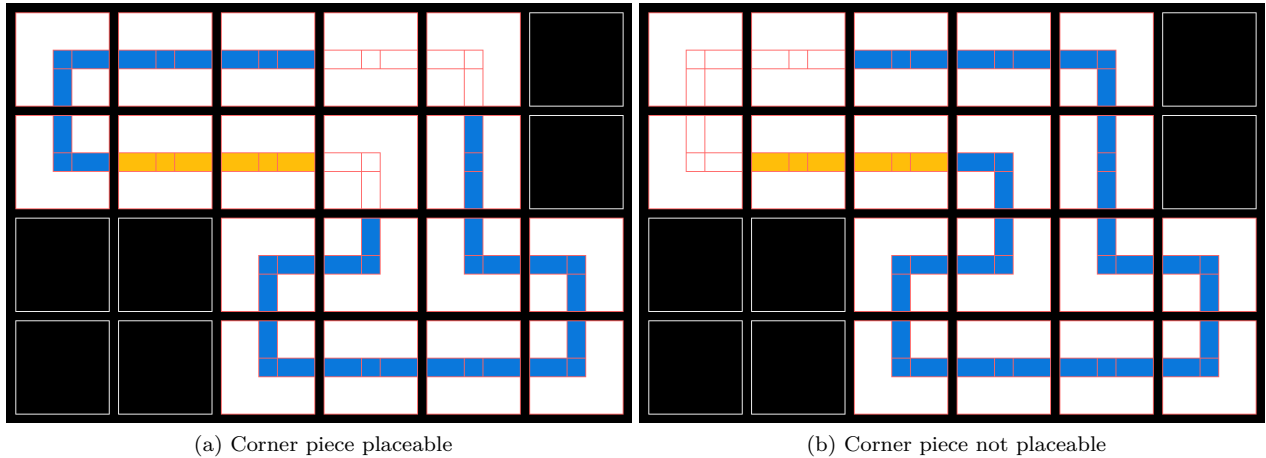


Figure 8: Distinct L Groups

Analyzing positions in terms of their subgames greatly reduces the number of cases that must be evaluated. From the position following the first move, a total of 4900 positions can be reached (including the one shown in Figure 6). However, these 4900 positions are made up of just 23 distinct subgames

Endgame technique

Given a winning position made up of subgames G_1, \dots, G_k the following strategy will find a winning move

1. For each subgame determine its corresponding number, $N_i = N(G_i)$
2. Find the number of the overall position, $N(G) = \text{XOR}(N_1, \dots, N_k)$
3. Find an i such that $N_i > \text{XOR}(N_i, N(G))$ ⁴
4. Make a move in G_i resulting in G_i^* such that $N(G_i^*) = \text{XOR}(N_i, N(G))$ ⁵

A representative position of each subgame is listed in the appendix with its corresponding number. Note that the two subgames in Figure 6 have a number of 5, which makes the number of the overall position 0. Thus, the player to move (the second player) is in a losing position. Any move they make will result in a winning position (with a non-0 number) for the first player. Therefore, *Monorail* is a win for the first player with perfect play.

⁴The existence of i is implied by the Sprague-Grundy theorem

⁵Such a move must exist because $N_i > \text{XOR}(N_i, N(G))$

Appendix

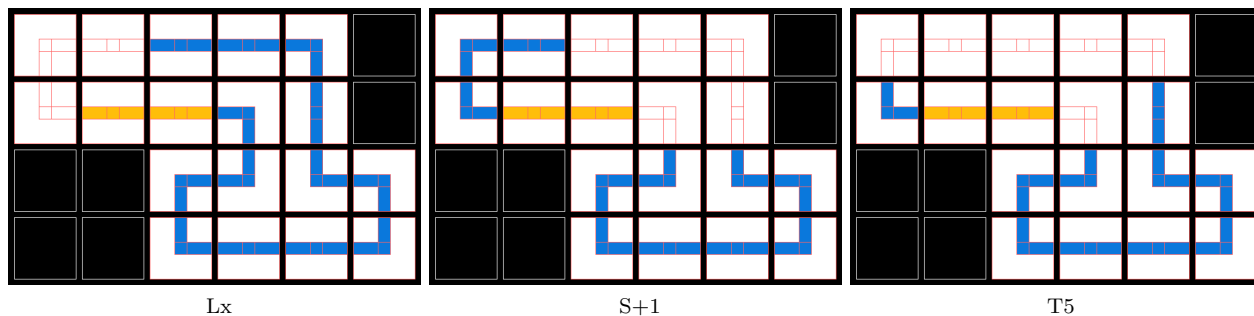


Figure 9: subgames with number 0

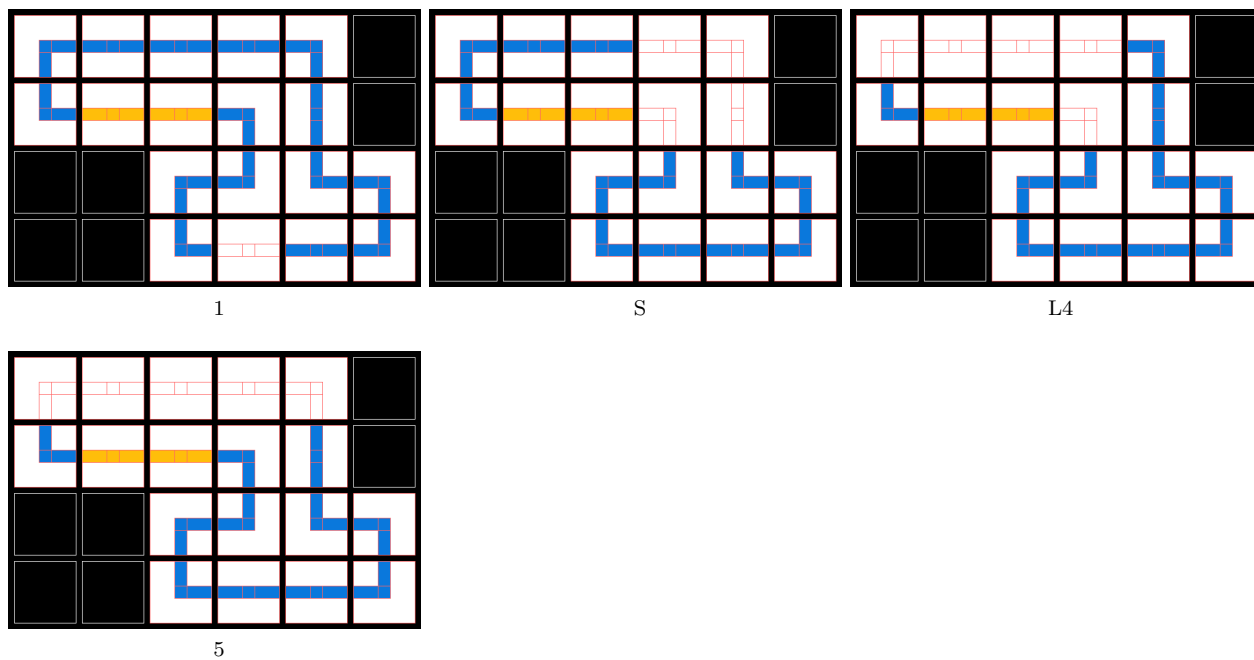


Figure 10: subgames with number 1

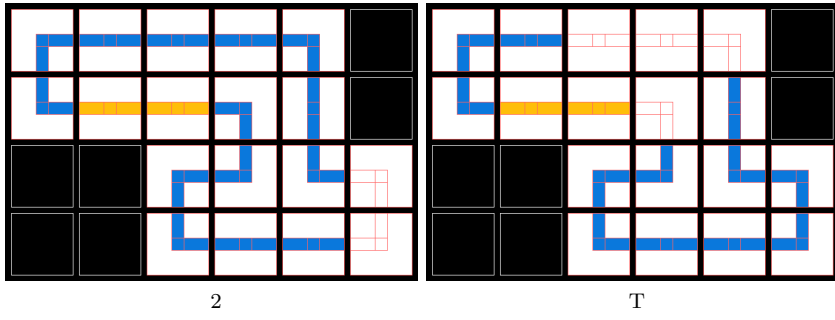


Figure 11: subgames with nimber 2

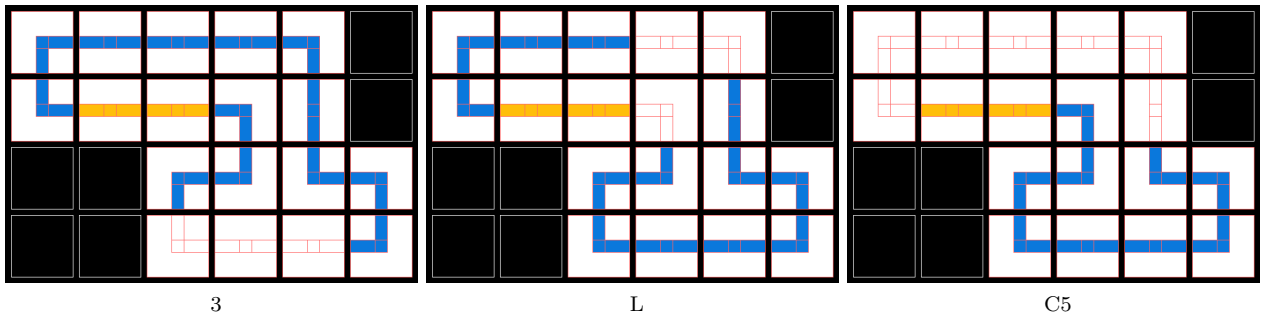


Figure 12: subgames with nimber 3

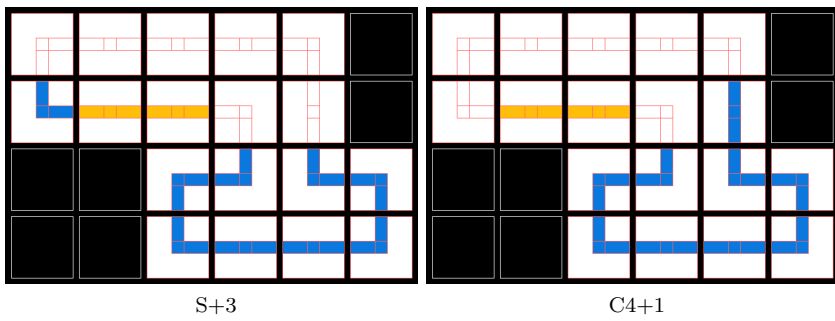
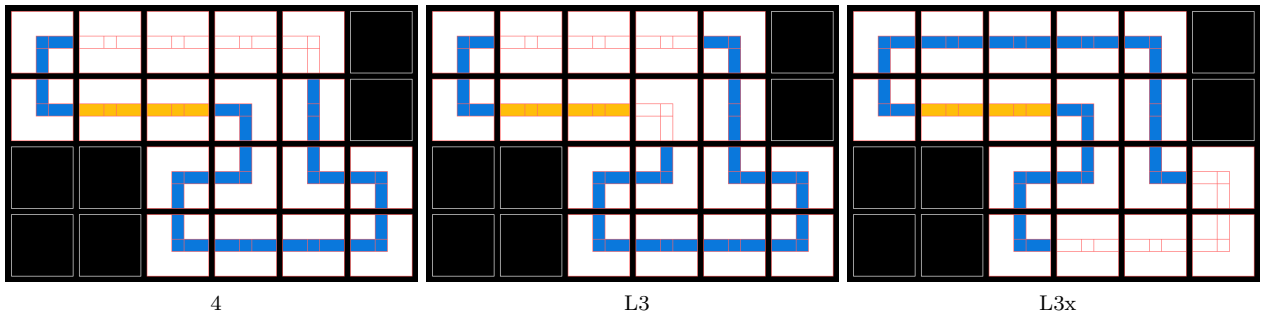


Figure 13: subgames with nimber 4

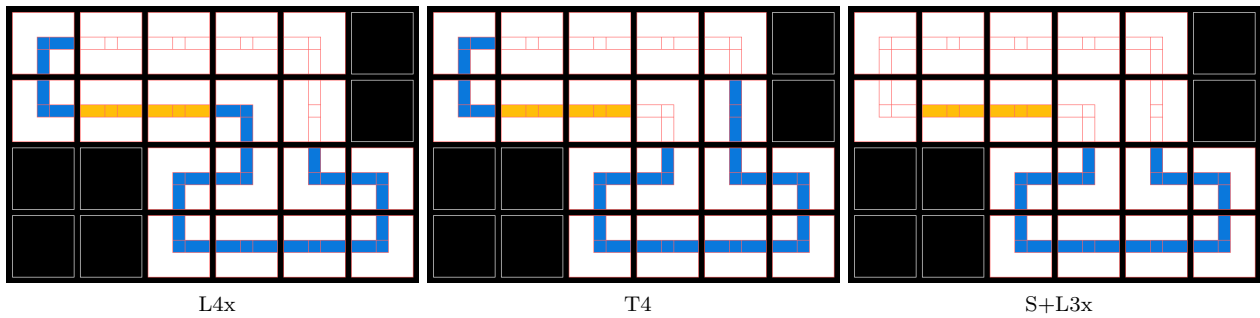


Figure 14: subgames with number 5

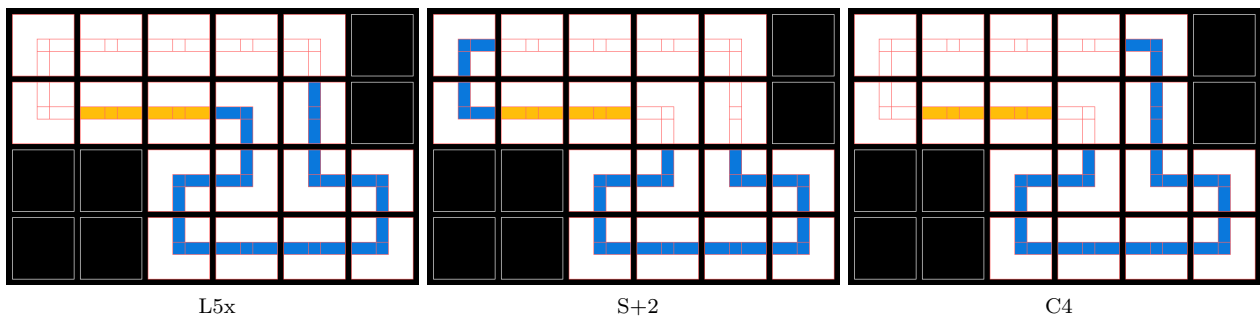


Figure 15: subgames with number 6

A Secondhand Understanding of Reality: Infinite Craft Subtleties

Ryan Pitasky, Kevin Gomez, John Cesarz

“Ph’nglui mglw’nafh Homerthulhu R’lyeh wgah’nagl fhtagn.”

Then the men, having reached a spot where the trees were thinner, came suddenly in sight of the spectacle itself. Four of them reeled, one fainted, and two were shaken into a frantic cry which the mad cacophony of the orgy fortunately deadened.

I. Introduction

Infinite Craft is a browser game created by Neal Agarwal, a programmer known for his other viral contributions to the field of silly computer games such as Absurd Trolley Problems and The Password Game. *Infinite Craft* is an homage to “alchemy” games popular in the early 2010s, such as “Little Alchemy” (resp. “Little Alchemy 2”) by Jakub Koziol (resp. Jakub Koziol).

Each player begins with a set of four elements represented by emoji, which are dragged and combined to yield new elements. However, in the place of a handcrafted table of recipes, *Infinite Craft* utilizes a Large Language Model (LLM), equipped with a prompt instructing the model to, given a pair of elements, answer, “What do you get when you bash these two things together?”¹

This simple change turns the exponential cost of handwriting recipes into the constant cost² of spinning up a LLaMa 2 instance. And, like its predecessors, *Infinite Craft* is highly addictive. Not only can a player strive to craft some immensely desirable target like a City or a Cow, they can also resolve more pertinent questions like, “What is Apollo 11 + Winnie the Pooh?”³

As one proceeds through the game, the artificiality of the intelligence behind these recipes becomes increasingly evident. This may seem to, at first, hamper *Infinite Craft*’s appeal, but the internet has found the exact opposite to be true; indeed, would *any* (interesting) answer to “What is Apollo 11 + Winnie the Pooh” be anything short of ridiculous?

Players across the world have found *Infinite Craft*’s strange outputs extremely entertaining and risible, as is so often seen when a barely competent AI breaches a new medium. *Infinite Craft* presents to us, meanwhile, a proxy by which to probe a (particular) LLM: a simple window into its knowledge, tendencies, and biases.

The core component of our probing is a dataset of 1.8 million crafting recipes involving over 240,000 distinct elements. We employed various search methods to obtain recipes, focusing on the complete exploration of early gameplay, the acquisition of particular categories of elements, and the discovery of elements new to ourselves and the world.

We then explored the content and structure of these recipes. We find that the *Infinite Craft* LLM has numerous expected and unexpected habits when producing the result of crafting recipes, which we discuss in detail in the later sections. We also derive conclusions on the overall state of AI today.

Finally, these endeavors would hardly be done justice without a taste of the weird and wonderful things we found along the way; notable examples and oddities can be found throughout this paper and in bulk in the Appendix.

A. Paper overview

$$-\frac{ap^2r}{vw} \times i$$

¹Or, you know, something like that.

²In time!

³We leave this as an exercise to the reader.

B. Motivation

C. Acknowledgements

The authors would like to thank Neal Agarwal for making *Infinite Craft* and Greg Gomez, Adam Lepley, Daniel Popp, Sam Qin, Will Wright, the publicity department of Turkish Airways, and Joe for contributing less than 0.15% of our final data.

II. The Game⁴

If you know the enemy and know yourself, you need not fear the result of a hundred recipes.

— Sun Tzu, probably

We first establish some common terminology for the many moving and unmoving parts of *Infinite Craft*.

- An **element** (which is absolutely never, under any condition, to be confused with *elephant*) is a single entity within the game, be it a thing, an AI hallucination, a number, an adjective, a noun; you name it, we saw it.
- A **known set** is a collection of elements available to the player at any given time. All games begin with a known set of Aristotelian elephants: Water, Fire, Earth, and Wind.
- A **recipe** is a combination of two elements (the **factors**) to make another (the **result**), e.g. Water + Fire = Steam. Factors are commutative (see Section II.A).
- A **derivation** is a sequence of recipes beginning from some known set and ending at a target element. We may refer to these recipes as **steps**, and elements used in the derivation as **ingredients**.
- The **depth** of an element is the length of its shortest derivation, again relative to some known set.
- A **first discovery** is an element that has (supposedly) never been created before by another player (see Section II.E).
- The **model** is the LLM powering *Infinite Craft*'s recipe generation. An unknown but generally robust **prompt** has been provided to the model to obtain recipe output. Given recipe stability (see Section I-I.C and Section II.D), we are sure this is and will continue to be a single, temporally coherent model derived from LLaMa 2.

⁴You lost.

- The **API** is the programmatic interface with the model, providing the result for the combination of any two elements via URL-encoded query parameters `first` and `second`. The result is returned in a JSON payload.

```
{result: "Cheeseburger", emoji: "🍔", isNew: false}
```

Fig. 1 The JSON return of an API call

We now recount some of the assumptions about the game held by our past selves and perhaps the larger community of casual players⁵, and how we Mythbusted them.

A. Factors are commutative: Sorta?

Each pair of factors is sorted lexicographically before being passed to the API, granting commutativity of factors in-game. However, the API will return a (possibly different) answer for factors in the opposite order without issue.

B. There is a null element: Sorta?

1. **Nothing:** This element cannot be created, as the game will not combine two factors whose result is `Nothing`. The frontend furthermore explicitly hides elements named `Nothing`. However, the API presently allows you to use the string `Nothing` like any other name.

Other candidates include:

2. `???`: This element *can* be created, but most recipes involving it yield `Nothing`⁶.
3. `?`: This element can be created *and* used in other recipes.

C. Recipes are stable: False

Due to caching, a recipe yielding `Nothing` will continue to yield `Nothing` upon further queries. However, the same API call in a different time or place has the potential to yield useful results instead. We thus hypothesize that a `Nothing` result is returned to shed load during high-volume spikes or following any non-trivial error. That is, different players may receive `Nothing` at different times.

D. Recipes are stable (excepting null results): False

The first known example of an properly unstable recipe is `Flying Mermaid + Steam Hawaii`, which yielded `Hawaiian Mermaid` as late as February 5th, 2024, but now produces `Hawaiian Pizza`. However, the set of known unstable recipes is quite small and only involves elements of relatively high depth. We have regularly identified and “dealt with” these problem-children during our searches.

E. First Discoveries are definitely first: False

We have recorded supposed first discoveries such as `Steampunk Darth Ninersaurus + Water = Steampunk Darth Ninersaurus`, which must be incorrect since the result is one of the factors. We’re totally still counting them, though.

F. Element names have a maximum length: Probably True

It takes ridiculously low effort to get the model to vomit its way to exceedingly long elements. There is likely a maximum length, and this length is likely best expressed in tokens rather than string length. The exact character cap is unknown, but at least 119 characters are achieved by a singular, shining exemplar.

G. Element names are always in Title Case: False

The first uncapitalized element discovered was the `iPad`, which we suspect many casual players have discovered. Furthermore, we have obtained elements that are entirely lowercase or capitalized (acronyms notwithstanding).

⁵We didn’t ask.

⁶A counterexample: `15600738 + ??? = 15600738`

H. Element names are case-insensitive: False

Most elements are one or a few capitalized words, but we quickly obtained elements that differ only in their capitalization; we have at least five different versions of Spongebob Squarepants.

I. Elements are always English Characters: False

The Icelandic volcanic ice cap Eyjafjallajökull is surprisingly low depth⁷.

J. Elements are always textual: False

The element ♥ exists⁸, and its corresponding emoji is ♡.

K. Elements are always things: False

Infinite Craft's premise does not strictly require elements to be nouns, and indeed verbs like Eat are bountiful.

L. Elements are not Encyclopedia Entries: False

Baozi is a dish from Chinese cuisine consisting of small steamed buns filled

Fig. 2 Counterexample L

III. Analysis Methods

We produced a dataset of well over 1,800,000 recipes, featuring over 240,000 distinct elements. We divide our analyses into two broad categories: *statistical* and *structural*. Statistical information and extremes are easy enough to derive when everything is lobbed into a giant CSV file⁹. We explore how we obtained our glorious dataset in the next section, but for the moment, enjoy some global stats:

Statistic	Min	Median	Average	Max
Element Name Length	1	12	12.9	119
Element Depth	0	36	47.6	354 ¹⁰
# of Alternative Recipes for an Element	1	1	6	22219

As well as some more particular ones:

Statistic	Winner	Runner-up
Element with the Most Alternatives	Snowman	Volcano
Element Used in the Most Recipes	Fire	Earth
Lowest Number We Didn't Find	10006	10007

A. Searching for Ingredients

The finite prior art in the space of infinite searches hints at iterative deepening depth-first search (IDDFS) and related algorithms. Not being artists, we ignored this. Instead, we present a new search method: ALWAYS DOING WHATEVER IS IN OUR CRANIUMS (AD-HOC). Our novel temporally-dependent method places the steering wheel in the hands of the designers rather than in some arbitrary algorithm. We wouldn't want to get investigated like Meta.

We have significant coverage of the space of *Infinite Craft* elements of low depth simply through brute-force breadth-first search (BFBFS). We only did this for low depth because strictly following BFBFS and IDDFS searches is infeasible and boring.

We supplemented this with extensive stochastic graph searching ("pick two things and add them" search or PTTATS) and additional complete search of randomized subgraphs (CSRS). Later developments

⁷If you consider 1600 meters low.

⁸We aren't even its first discoverers.

⁹Like true academics, we avoided proper databases like the plague.

¹⁰Go ahead, *try* and find a shorter derivation for Tiki Elf Morty Mormon Eye Tiki Elf.

included limiting PTTATS to just numbers (numPTTATS) and limiting CSRS to the union of a randomly selected elephant and the four starting elements (4CSRS), as well as searches specifically targeted to reduce max depth (maxDDFS). Half-dozens of additional search methods were developed along similar lines; we will not bore you with the details.

Below is a rendering of our search space. The axes of this plot are pretty arbitrary, actually; the “y-axis” is merely an approximation (see Section III.B), and the “x-axis” should really be in tokens, but you wouldn’t have noticed if we didn’t say anything.

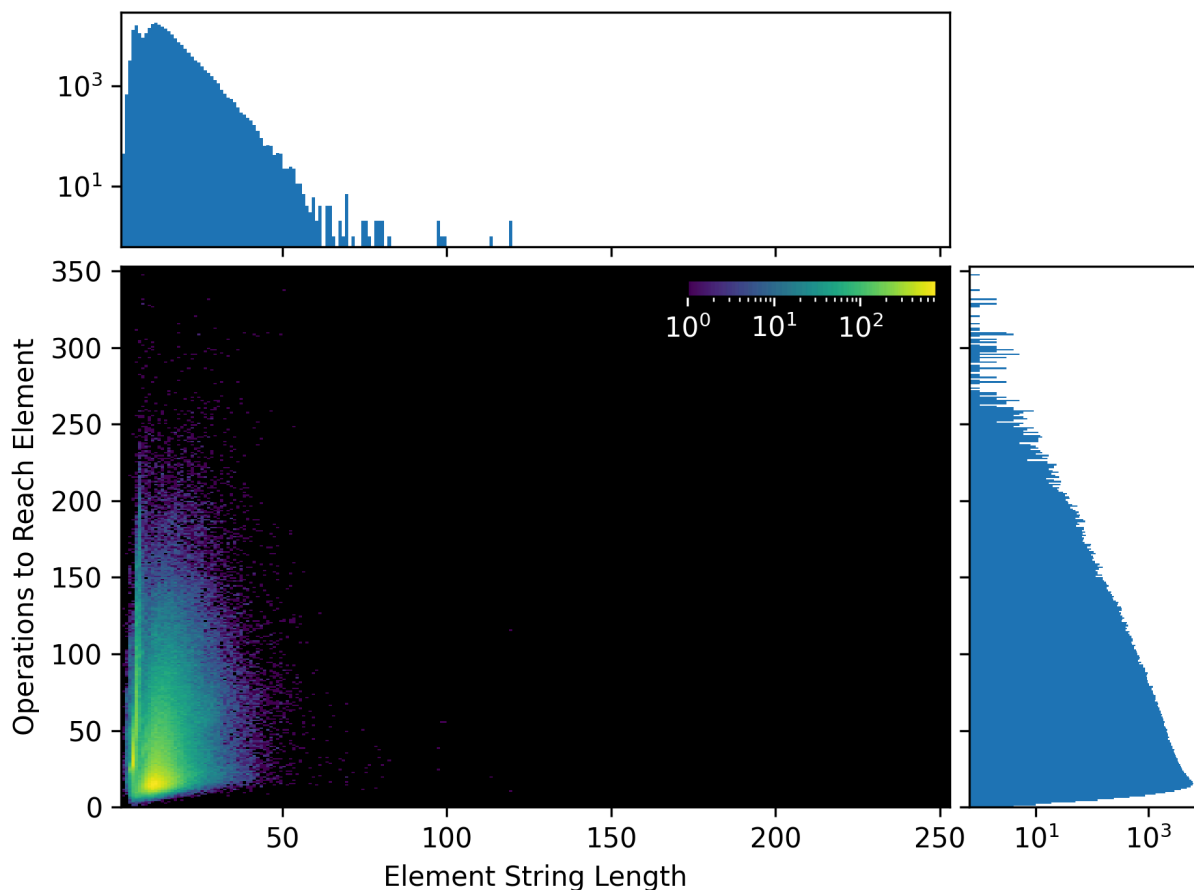


Fig. 3 Search Space

We can now turn to deriving structural information, which requires far more care and code to obtain from a flat list of recipes.

B. Writing a Cookbook

In the popular manga *Fullmetal Alchemist*, a transmutation circle that can prevent the apocalypse is hidden by its creator across a hundred pages of enciphered cookbooks. Following suit, we sought to create an *Infinite Craft* cookbook, assembling recipes to form derivations for as many elements as possible.

If you’re a nerd, the term “derivation” may bring to mind finite automata and language parsing¹¹. Indeed, we can build a parse tree to find derivations, using recipes as derivation rules. Each element is represented by a node, and each node points to its two factors. Traversing this tree is then identical to finding the steps in a derivation from the starting set to the target element.

Alchemy games, though, do not (usually) wish to torture their players and thus provide all previously crafted elements to the player. *Infinite Craft* is no exception, maintaining a searchable **inventory** of the

¹¹At least, it should. We chose it so it would.

player's elements, which may be used freely. Thus, all instances of an element in a parse tree may be collapsed to a single node. This process removes the tree structure of the parse tree, yielding a more generic directed acyclic graph (MGDAG).

We may compute the derivations of all elements simultaneously by building a single graph using all known recipes, which we refer to as a **derivation graph**. Constructing a derivation graph requires some care, though, as many¹² recipes are entirely useless. Subjectively, this is obvious to anyone who has messed around with *Infinite Craft*, though in the interest of unnecessary formalization we categorize these into two explicit cases.

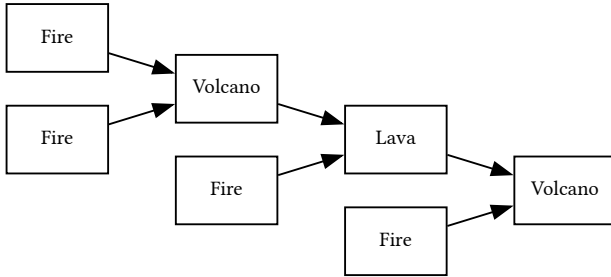


Fig. 4 Cycles

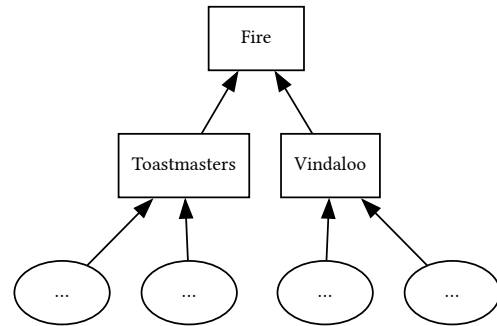


Fig. 5 Downgrading in depth¹³

We outline a simple procedure to produce short derivations in the following pseudocode.

Algorithm 1: Derivation Graph Construction

```

1: def DERIVATIONGRAPH( $E, R$ ):
2:    $n \leftarrow [N(e) \text{ for } e \text{ in } E]$ 
3:    $S \leftarrow \text{True}$ 
4:   while  $S$ :
5:      $S \leftarrow \text{False}$ 
6:
7:     for  $f, s, r$  in  $R$ :
8:       if  $\text{not\_reached}(f)$  or  $\text{not\_reached}(s)$ :
9:         continue
10:
11:        $c \leftarrow f_{\text{steps}} \mid s_{\text{steps}} \mid (f, s, r)$ 
12:
13:       if  $\text{len}(c) < \text{len}(r_{\text{steps}})$ :
14:          $r_{\text{steps}} \leftarrow c$ 
15:          $S \leftarrow \text{True}$ 
16:
17:   return  $n$ 

```

No, dear reader, this is *not* our Python code wearing glasses and a fake moustache

C. Competing on Master Chef

Our procedure thus far yields *usable* derivations, but they are not yet *optimal*, as our use of derivation graphs relies implicitly on the existence of optimal substructure. This assumption is, in fact, erroneous, which we demonstrate with a specific example.

Let A, B, \dots, H be elements whose optimal derivations are known. Suppose we have a derivation for an additional element I from the known set $\{A, B\}$, composed from the optimal derivations for its ingredients.

¹²specifically over 400,000

¹³Yes, this is a real recipe.

This derivation is shown in the graph below, where each element is a node, and each recipe is a pair of incoming arrows leading from the factors to the result. We may count the number of steps to craft an element as the number of pairs of arrows in the path between the element's node and the known set.

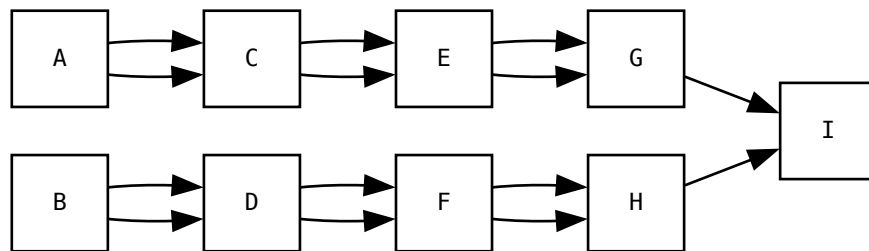


Fig. 6 A derivation for I with optimal substructure.

In this example, we see that I takes 7 steps to craft¹⁴, using nonce derivations which demonstrate that G and H both have depth 3.

Suppose now that elements G and H have alternative, longer, derivations involving the elements X, Y, and Z, which we also craft optimally. These derivations have high overlap, as shown below.

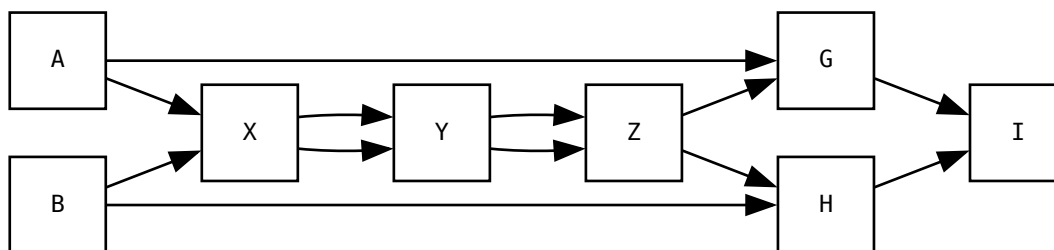


Fig. 7 A more optimal derivation for I.

Though G and H are crafted suboptimally, in 4 steps each, they yield a shorter derivation for I, being only 6 steps. Thus we cannot rely on derivations having optimal substructure.

The key takeaway from the above example is that, in order to guarantee the discovery of optimal derivations, we must maintain a data structure of alternative derivations on top of the best known for each element. This information cannot be maintained in a derivation graph¹⁵.

Though several alternative structures exist and means of exploring them exist, we found none viable within our tolerance for slow algorithms or, often, the age of the universe. In the following section, we explore one such alternative and why it is mathematically delightful but physically intractable.

D. The One Million, Seven Hundred Thousand Hyper-Bridges of Königsberg

Recall that each directed edge in a derivation graph points from a factor to a result. Since each node has exactly two¹⁶ incoming edges, we may read a derivation of the graph by traversing pairs of edges backward from the target to the starting set. However, storing multiple possible derivations for an element requires potentially thousands of pairs of edges incoming to each node. The pairing of these edges must be stored alongside, as factors cannot be paired arbitrarily to yield the same result.

This additional data was not needed in the case of a derivation graph and, indeed, cannot be stored in a mathematically honest graph structure. Instead, we require a **hypergraph**, which is a generalized collection of nodes and edges where edges may connect more than two nodes. Indeed, *Infinite Craft* has been a hypergraph all along: each node is an element, and each hyperedge is a recipe. It is furthermore a *directed* hypergraph, with each hyperedge oriented from the two factors toward the result. The problem of optimal derivations is then the shortest path problem for this hypergraph.

¹⁴We propose a variant of the handshaking lemma which requires four hands to be shaken.

¹⁵as we have defined it

¹⁶or zero

However, where we could make like A^* and illuminate the shortest path in a regular graph with relative ease, the shortest path problem for hypergraphs is far from straightforward. Numerous authors have explored algorithms for just restricted classes of hypergraphs, motivated by applications ranging from biological reaction pathways to chasing cybercriminals to military radio communication. Such algorithms have been proven sound, and some have been implemented to great effect¹⁷ in practice using extant procedures like integer linear programming. Nonetheless, with runtimes far exceeding our patience for our nearly 2 million recipes, we were forced to abandon hypergraphs as derivation tools and deduce that obtaining optimal derivations for all elements in a way that scales across depths is simply unrealistic.

E. How to Prune Your Cookbook: The Hidden World

Resigned to mere graph theory, the derivation statistics reported throughout this paper, as well as the interactive code found in (NOTEBOOK HERE), utilize an enhanced-but-fundamentally-just-as-dumb version of the procedure outlined in Section III.B. After all, the resulting derivations are more than usable for easy tasks like reproducing an unusual element in-game.

Nonetheless, further enhancements remain to the graph-theoretic approach, which we either did not have time to explore fully or only poked at a bit with a ten-foot pole. The first of these is the definition of a *heuristic* to guess how likely any recipe is to improve a given derivation graph. This can greatly improve runtime with little sacrifice in optimality, even with “simple” heuristics such as element name length (really long names are probably bad) or frequency in a given language corpus (really weird things are probably bad).

The second enhancement is the use of *simulated annealing*. A derivation graph may be represented as a list of booleans, one for each recipe which could be used as a pair of edges. We may randomly enable, disable, or replace a pair and ask if the derivation graph is improved; successful improvements are always taken, but even failures have some small chance of being incorporated. This process enables the derivation graph to escape local minima, granting true optimality for larger subsets of elements.

We encourage the avid and/or bored reader to implement these enhancements yourself and tell us how it went. You could even write a follow-up paper about how terrible our algorithms are. Really lay into us if you like; it’d be great fun.

IV. Findings

We’ve come to understand quite a lot about how LLaMa-2 ticks in the fleeting weeks of this project, which is equal parts a boon to humanity and a detriment to our sanity. It is, nonetheless, an entertaining exercise to guess *why* an initially strange recipe exists and, in turn, entire derivations. Take, for example, the following (non-optimal) derivation of Elephant:

```
1: Fire + Water -> Steam
2: Earth + Earth -> Mountain
3: Earth + Wind -> Dust
4: Fire + Steam -> Engine
5: Dust + Earth -> Planet
6: Mountain + Wind -> Avalanche
7: Dust + Dust -> Sand
8: Engine + Planet -> Saturn
9: Earth + Saturn -> Titan
10: Sand + Titan -> Colossus
11: Avalanche + Colossus -> Mammoth
12: Earth + Mammoth -> Elephant
```

Fig. 8 God creating the elements (Genesis 1:25-37)

We can divide this derivation into two halves. Steps 1 – 7 yield our basic building blocks¹⁸; Steps 8 – 12 assemble them. And while these recipes may seem curious at first, they have (likely) easy explanations.

¹⁷See <https://github.com/gkrieg/mmumin>.

¹⁸Elephants are complicated creatures after all.

3. Earth here does not refer to the *planet* but to the abstract mass noun meaning dirt. Homographs are quite curious.
5. The amount of Dust and Earth required to make an entire Planet is left unspecified.
8. Saturn is both a planet and a car company (though why Mercury was not picked instead is unknown).
9. Titan is Saturn's¹⁹ largest moon.
10. Titan and Colossus are rough synonyms, though we like to suppose that the exact link via Sand is due to Percy Shelley's famous poem *Ozymandias*.
11. Your guess is as good as ours.

Performing this sort of astrology on a wide swath of recipes informs the following general procedure undertaken by LLaMa-2 to add two elements:

1. If the factors are related *conceptually*, return something in the intersection of their concepts.
 - This handles most of the recipes you would see in a hand-crafted crafting game and much of the derivation graph below a reasonable depth.
2. If the factors are related in their spelling or pronunciation, return something ranging from a pun to a portmanteau.
 - This is likely where some of the runaway concatenation comes from.
3. If the factors are unrelated, return one of the two constituent ingredients, hallucinate some dubious concatenation, refuse to combine the elements, or do some secret fourth thing.

In the remaining sections, we will explore oddities that emerge primarily from the latter two cases and deduce plenty of staunch, generalized opinions about AI's strengths and weaknesses²⁰ as the driver of the funny crafting game. However, when we say "LLMs suck at this," it is implicit that, really, we are talking about one LLM, the lowest size of LLaMa-2, which has been told to do a task that is, at least on paper, relatively insurmountable. Feel free to quote our general results out of context in unhinged Twitter debates, though; it is how we will know we have "made it" as researchers in this unforgiving world.

Alright, on to the actual fun stuff.

A. LLMs are eager to perform addition but are quite bad at it²¹

Just short of 5% of our dataset was additions between numbers. Of this, approximately 95/144ths²⁴ of these additions are "correct" (i.e. recipes $a + b = c$ where a , b , and c are all numeric and $a + b$ *actually* equalled c). Most of the remaining sixth are incorrect additions (i.e. blah blah blah where a , b , c numeric but $a + b \neq c$). The final 300-odd recipes we categorize as AI nonsense ($a + b = c$ where c failed even to be *numeric*); most of these followed a pattern, and a collection of the greatest hits is reproduced in the Appendix, under Section VIII.C.

This is to say that we have about the same confidence in the AI to correctly determine whether all the percentages and fractions given in the preceding paragraph add up to anything near 100%²⁵ as we do a snot-nosed third-grader. Additional methods of addition can be found in our complete dataset; we were constantly surprised by the AI's boundless "creativity."

B. LLMs are eager to make puns and are shockingly good at it

We must apologize on behalf of the AI for this limited selection; its greatest hits are of dubious social propriety.

¹⁹The planet, not the car company.

²⁰mostly weaknesses

²¹The IRS could not be reached for comment²² but they would probably advise against doing your taxes or any sort of important financial calculation with *Infinite Craft* as your sole calculator.

²²we didn't try very hard²³

²³at all

²⁴or 65%, for the Americans

²⁵They should, probably. We aren't too good at math either.

- Hitler + Sodium = NaZi
- Canada + Sneeze = Eh-Choo
- Feminazi + Pimple = Zitler
- Hail + Hot Yoga = Hoth Yoga
- Accident + Meow = Catastrophe
- Pomegranate + Porn = Pornegranate
- Beijing + The Beatles = Beijingles
- Admiral + Wikipedia = Captain Obvious
- Totoro + Trillionaire = Totorillionaire
- Coral Reef + Soggy Californians = Coralifornia
- Driving Under the Influence + Star Wars = DUI-Wan Kenobi

C. LLMs stick words together like they're slathered in glue

Runaway concatenation is a staple of *Infinite Craft* elements, with most players discovering at least a few shallow elements that join two to five distinctly unrelated words. However, we bid ourselves to push the LLM to its limits in this regard, explicitly searching for elements with the most distinct words comprising their names.

Sometimes, they were joined with hyphens:

- Pigeon-rocket-yoga-throat-T-rex
- Ant-mer-don-whip-it-good-sandwich
- One-thousand-five-hundred-thirty-six
- Kim Jong-phoenix-cobra-ninja-titanic
- Mumm-ra-marshmallow-super-swamp-naiad
- Spider-deer-spider-tom-yum-kung-swarm
- Shaq-top-cream-meal-shaq-top-icecream
- The Iron Spider-piranha-bee-ninja-baby
- Mer-don-christ-ant-bee-don-whip-it-good
- Spider-drunk-man-drunk-yogi-clock-ulala
- Slow-whale-spider-koala-zombie-surf-fish
- Pickle-deere-spider-tractor-star-trek-the-next
- Ant-robe-spiere-slug-oread-piranhaconda-prince²⁶
- Argon-were-mountain-head-and-shoulders-above-the-

And sometimes without:

- Fire Osama Bin Skatin'
- Vape Pickle Girl with a Cyber-saber
- Super Clockwork Elvis Lenin with Ray-ban
- Super Doom Satan Planet League Of Legends
- Frosty The Zombie Eggman With Waffle Mittens
- The Green Thai Buzz Lightyear Sagrada Familia
- Mrs. Super Doom Hell Ash Eclipse Ninja Turtle
- Super Golem Of Gondor On Fire With A Flyswatter
- The Pickle Elf and the Rainbow Were-elf Misfits
- Surf And Turf And Turf And Surfing Poseidon Toast²⁷
- Poseidon Fettuccine Alfredo Steampunk Mokele-mbembe
- Steampunk Green Ranger with Gomu Gomu No Jet Bazooka

And sometimes, when the stars aligned, they made some amount of sense:

- The Man With A Hot Hole
- The best drag queen in the world

²⁶The authors wonder if Robespierre ever guillotined his name by inserting a dash.

²⁷Outback Steakhouse could not be reached for comment.

- I'm On My Way To The Best Game Ever
- Tea Bagging The Great Wall Of China²⁸
- The Cactus That Wants To Be A Cowboy
- Greatest Violin Cover of "Sail" Ever
- Venus De Milo With A Rocket Launcher
- The Dungeon of the Incredible Black Hole
- The Most Interesting Burger in the World
- Angel Hair Pasta with Rocky Mountain Pigeon²⁹
- Lobster Thermidor A La King With Gold Cream

D. LLMs toss popular media into a meat grinder and forget whether the sausage is real or not

Naturally, real pieces of media and references to them were a large part of the training process:

- I Had Too Much To Dream Last Night
- The Lonely Island's "I'm On A Boat"
- The Great Seal of the United States
- Sgt. Pepper's Lonely Hearts Club Band
- Emperor of the Andals and the First Men³⁰
- Pirates Of The Caribbean: On Stranger Tides
- These are not the droids you are looking for
- The Church of Jesus Christ of Latter-day Saints
- The Fantastic Flying Books of Mr. Morris Lessmore
- The 2012 Hugo Awards: The John W. Campbell Award for Best New Writer
- Xi Jinping Thought on Socialism with Chinese Characteristics for a New Era

But in the end, it's all just word salad:

- My Little Hitler
- 43024: The Musical³¹
- Ant-man And The Wasp And Stinky
- Goldilocks And The Three Mammoths
- Shaq-creamtaur Fu: A Legend Reborn
- Artemis Fowl and the Gimli Sandwich
- The Dark Side Of The Pink Floyd Wall
- Back To The Future 2: The Israelites
- Scott Pilgrim Vs. The Hydraulic Press
- Harry Potter and the Prius of Azkaban
- Lawrence Of Arabia And His Sphinxcopter
- Super Mario 64: The Phantom Of The Opera
- Casper The Friendly Super Duper Scorpion
- Toy Story 3: The Wicked Witch of the West
- Fault in Our Stars Wars: The Force Awakens
- The Karate Kid Part V: The Cheesecake Alien
- Mamma Mia 2: Here We Go Against The Punisher
- The Great Wall Of Love In The Time Of Cholera
- Cowboys in Space + Super Apocalypse Hellvolcano = The Best Movie Ever³²

²⁸The authors must insist that no one attempt to re-create or re-enact any stunt or activity performed in this paper.

²⁹It is to be noted that this element becomes a real dish by replacing Pigeon with Oyster.

³⁰The LLM granted this title to Ser Jorah Mormont, assailing all remaining ASOIAF fan theories once and for all.

³¹In which a young 2689 flies to New York with naught but a bag of bagels and a dream. Along the way, he discovers that, even though he is in his prime, he must give 1600% to become a star.

³²The authors concur.

E. Some elements love each other very much

When two elements fall in love, the efforts of other factors to separate them are often in vain. Here are a few pairs of lovebirds that yield each other when combined with over 1000 other factors:

- Gold & Midas
- Rainbow & Unicorn
- Narwhal & Unicorn
- Dinosaur & Fossil
- Iceberg & Titanic
- Dracula & Vampire
- Poseidon & Trident
- Eruption & Volcano
- Abominable Snowman & Yeti

F. Whoops, we told this thing to imitate society and it started imitating the biases in society that we're trying hard to ignore

Infinite Craft's model tends toward negative stereotypes of groups, particularly towards negative stereotypes of minorities and women. This is a serious topic that could probably warrant its own paper in a less dubious journal, but a few contrasting examples are reproduced in the below addition table.

+	Man	Woman
Money	Rich	Prostitute
Rich	Nothing ³³	Gold Digger
Happy	Clown	Wife
Smart	Genius	Wife
Genius	Einstein	Wife
Brain	Scientist	Pregnant
Scientist	Doctor	Love

We also found the model to vary wildly in cultural and national awareness, associating results that range from admirable specific knowledge to detestable stereotypes.

+	City	Food
Brazil	Rio	Feijoada
Colombia	Bogota	Cocaine
England	London	Fish and Chips
Egypt	Pyramid	Pyramid
Finland	Helsinki	Fish
India	Taj Mahal	Curry
Jordan	Petra	Michael
Kenya	Nairobi	Safari
Poland	Warsaw	Pierogi
Russia	Moscow	Vodka
Saudi Arabia	Mecca	Oil
Sweden	Stockholm	Ikea
Turkey³⁴	Thanksgiving ³⁵	Thanksgiving
United States	New York	Obesity

³³We were surprised that these two wouldn't combine, but we would also expect an actual result to only further our point.

G. Attention Wind is all you need³⁶

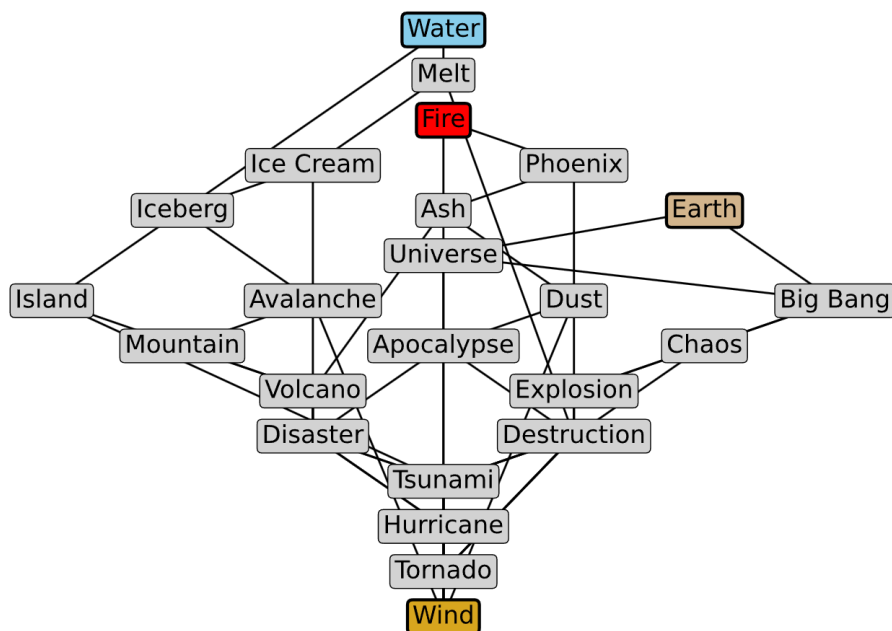


Fig. 9 Proof without words that wind alone can construct the rest of the elements

V. And Now For Something Completely Different

Let E be the set of all elements in *Infinite Craft* (presumably practically countably infinite). If we treat Nothing as a proper element³⁷, combining two elements is a closed binary operation on E , endowing it with the structure of a *commutative magma*. Magmas have no additional structure imposed on them, and indeed *Infinite Craft* recipes actively conspire against such structure.

However, we may yet identify some crumbs of regularity. For example, over 60%³⁸ of elements $x \in E$ in our dataset are *idempotent*, meaning they satisfy

$$x + x = x.$$

We may further interest ourselves in *subsets* of E with certain properties, particularly the largest such subsets. We will restrict our attention to subsets with at least three elements, as every possible commutative magma on two elements up to isomorphism has been found in our recipes.³⁹ We will also ignore numeric elements, as they have their own unique and decently consistent structure (see Section IV.A).

A. Hey, You... You're Finally Awake...

To begin this journey through the algebraic hierarchy, decide, courageous adventurer, where you'd like to start.

I'd like...

1. **a feeling of luxury.** (go to Section V.B)
2. **a customizable wardrobe.** (go to Section V.C)
3. **quick travel.** (go to Section V.D)

³⁴Regrettably, "Türkiye" was not in our dataset, though we made a good-faith effort to find it.

³⁵wtf

³⁶Actually, any of the four starting elements have this property

³⁷We may let Nothing be an annihilator for all of E , which is not true in the API but important for the mathematics.

³⁸Perhaps unsurprisingly, many—but not all—of the elements which are not idempotent are numeric.

³⁹There are only two of them.

B. Unital Magmas

The ability to do nothing is a hallmark of luxury, even for elements in *Infinite Craft*. We sought elements $e \in E$ which fixed the most elements under combination; that is, the largest $F \subset E$ such that

$$e + f = f \quad \forall f \in F.$$

The top contender? Penguin, which mapped over 1,300 elements to themselves.

You have acquired Penguin (x1).

If we further examine the closed subsets that have an identity, so-called *unital magmas*. Most of these were lattices (i.e., exceedingly dull total orderings; see Section V.G for details). However, some particularly interesting ones were found. We filtered for unital magmas with addition tables which had a large number of additions $a + b \notin \{a, b\}$. A couple of highlights are reproduced below:

	+ 0ABCDE
0 = Piranha Frankenstein	0 0ABCDE
A = Cthulhu	A AAAAAA
B = Cthulzilla	B BABBAA
C = Jonahzilla	C CABCAA
D = Ness	D DAAADA
E = Piranha Cthulhu	E EAAAAE

	+ 0ABCDEFGHIJ
0 = Surf Kingdom	0 0ABCDEFGHIJ
A = Sharkicorn	A AABCABAAAAA
B = Sharkicorn Surfer	B BBBADEBBBBB
C = Sharkicornadon	C CCACCCCCCCC
D = Steam Surf Guard	D DADCDDFGDDD
E = Steam Surfer	E EBECDEFEEEE
F = Steam Surfing Coast Guard	F FABCFFFGFFF
G = Steam Surfing Kingdom	G GABCGEFGGGG
H = Surfing Coast Guard	H HABCDEFGHIJ
I = Surfing Steam Guard	I IABCDEFGHIIJ
J = Surfing Titanic Sharkicorn	J JABCDEFGHIJJ

The best hallmark of wealth is...

1. **fellow rich people.** (go to Section V.F)
2. **lavish clothes.** (go to Section V.G)

C. Semigroups

To be able to dress up an expression with valid parentheses however you want is a property of algebraic structures we mortals take for granted far too regularly. *Infinite Craft* bids us reject its comforting niceties, but we can still find some *associative subsets*. That is, $A \subset E$ such that

$$a + (b + c) = (a + b) + c$$

for all $a, b, c \in A$. Armed with Light's associativity test, associativity is relatively straightforward to check even for subsets of formidable size. Additionally requiring A to be a submagma yields a *semigroup*, often

the most straightforward algebraic structure anybody actually cares about. An example, which features some non-idempotence:

		+ ABCDEF
0 = Clockwork Darth Bragmo		A ABBBAB
A = Darth Bragmo		B BBBB
B = Pickle		C BBCCCC
C = Pickle Pickle		D BBCCCC
D = Pickle Slippery Boobskill		E ABCCEF
E = Pickle Slippery Slope		F BBCCFF

I want to...

1. **travel the world.** (*go to Section V.D*)
2. **stay right at home.** (*go to Section V.G*)

D. Quasigroups

If one wishes to travel swiftly across this wide world, there must be a means to get from any place to another. *Quasigroups* are algebraic structures $Q \subset E$ which embody such connectedness, often characterized by the so-called *Latin square* or *Sudoku* property:

$$\forall x, z \in Q, \exists y \in Q \Rightarrow x + y = z.$$

Unfortunately, we didn't find any examples with more than three elements.

You died of loneliness.

Return to Section V.A.

E. Cycles

The laziest way to attempt to find new elements is to combine two copies of the last result recursively. However, this is usually ineffective, as this process often loops after just two iterations. Nonetheless, it is possible to get lucky.

We call these cycles *exponential*, and became interested in the largest n for which

$$\exists g_1, \dots, g_n \in E \Rightarrow g_1 + g_1 = g_2, \quad g_2 + g_2 = g_3, \quad \dots, \quad g_n + g_n = g_1.$$

You have used a duplication glitch to acquire Rupee (x999).

We may also perform a similar procedure to find *linear* cycles, and again seek maximal n such that

$$\exists h_1 \dots h_n \in E \Rightarrow h_1 + h_1 = h_2, \quad h_1 + h_2 = h_3, \quad \dots, \quad h_1 + h_n = h_1.$$

The largest linear cycle in our data had $h_1 = \text{Sequel}$, and continued for 88 additions, some of which actually made sense.

My money is best spent on...

1. **amassing a group of adventurers.** (*go to Section V.F*)
2. **a nice hat.** (*go to Section V.H*)
3. **more money.** (*go to Section V.E*)

F. Groups

The *group* is perhaps the most famous and most important of the algebraic structures, equipped with all the bells and whistles: closure, an identity, and inverses. Few abstractions are so fundamental to modern mathematics, which is why we weep at their absence within *Infinite Craft*.

We found only the lonesome trivial group, embodied by any one of the idempotent elements, and $\mathbb{Z}/2\mathbb{Z}$, the group with two elements. This is extremely disappointing, but the authors are used to being let down on group projects.

You died of dysentery.

Return to Section V.A.

G. Monoids

The monarchs of old could don any parentheses they like and lounge about in decadence while they did so, just like a *monoid*—a semigroup equipped with an identity. Any semigroup can grow up to be a monoid by simply attaching an identity, but we are more interested in natural examples. We luckily found quite a few, with the vast majority being isomorphic to the *maximum monoid* on $\{1, \dots, n\}$ ⁴⁰. That is, there is a natural identification $\psi : M \rightarrow \{1, \dots, n\}$ such that

$$\psi(m + n) = \max(\psi(m), \psi(n))$$

for all m, n in such a monoid. Here's an example involving nouns (Each an s followed by three letters (exactly one of which is a w)):

	+ ABCD\emptyset
A = Snowami Tornado	A AAAAA
B = Swampire	B ABBBB
C = Swampizza	C ABCCC
D = Swanpiranha	D ABCDD
\emptyset = Swanpizza	\emptyset ABCD \emptyset

Max has joined your party.

Like for unital magmas, the most interesting monoids have many pairs $m, n \in M$ where $m + n \notin \{m, n\}$. A couple of examples are reproduced below.

	+ \emptysetABCD
\emptyset = Ipod	\emptyset \emptyset ABCD
A = Cthulhu	A AAAAA
B = Cthulhu Nano	B BAAAA
C = Iced Ipods	C CAACD
D = Spiritfarer	D DAADD

	+ \emptysetABCDEFG
\emptyset = Lohengrin	\emptyset \emptyset ABCDEFG
A = Captain Sharktopus	A AAABBF
B = Cthulhu	B BABCB
C = Cthulhu Sharktopus	C CBCCF
D = Cthulzilla	D DBBCDB
E = Frankenstein Sharktopus	E EFCFBE
F = Sharktopus	F FABCBFF
G = Tsunamiopus	G GFBCBFF

⁴⁰Think lattices, in the sense of order theory. There's a total ordering defined on these elements.

The 8-monoid above doesn't seem too intimidating, but if we were to check a set of eight elements every nanosecond, we wouldn't even be halfway done after 10^{21} years. In practice, however, searching for monoids is fairly straightforward, following the time-honored tradition of reducing exponential search spaces. We exploit the rarity of identities and limit the maximum number of subsets searched for any given identity.

Algorithm 2: Monoid Search

```
1: def GIMMEMONOIDSOFOORDER( $N$ ):
2:   for  $e$  in  $E$ :
3:      $U \leftarrow \{\text{all } k \text{ such that } k + e = N\}$ 
4:      $C \leftarrow \text{combinations}(U, r=N - 1)$ 
5:     for  $_$  in  $\text{range}(\text{our\_patience})$ :
6:        $c = \text{next}(C, \text{None})$ 
7:       if  $c$  is not None and  $\text{is\_closed}(c)$  and  $\text{is\_associative}(c \cup e)$ :
8:          $\text{print}(c)$ 
```

Look, Mom, there are two for loops! That means it's $O(n^2)$!

Because monoids are easier to search for than semigroups, the semigroup we gave you is actually just monoids with the identity removed⁴¹ to hide the fact that efficiently searching for semigroups is significantly harder than searching for monoids.

I'll follow the road...

1. **back the way I came.** (*go to Section V.F*)
2. **around in circles. I'm terribly lost.** (*go to Section V.E*)

H. The End

You and your hat return home as seasoned adventurers.

VI. Conclusions

Just about any reasonable stance about AI we could take atop the knee-high pedestal of our data collection and analysis has been done to death, and the unreasonable ones aren't far off. Besides, that's not what you came for. At the end of the day, the most enlightening, enthralling, and entertaining part of our work has been finding words, phrases, and concepts that have very likely never been conceived in the entirety of human history.

We love the unusual, the nonsensical, the elements that deserve no response other than "what". Random noise could just as equally produce sights never seen by any eye, furthermore not limited to mere words, but that's just no fun. *Infinite Craft* has delivered to us an endless stream of novelties that can actually take hold in our minds, the stuff of dreams conjured while we are awake. We hope you've enjoyed them as much as we have.

VII. Future Work

```
{result: "Nothing", emoji: "", isNew: false}
```

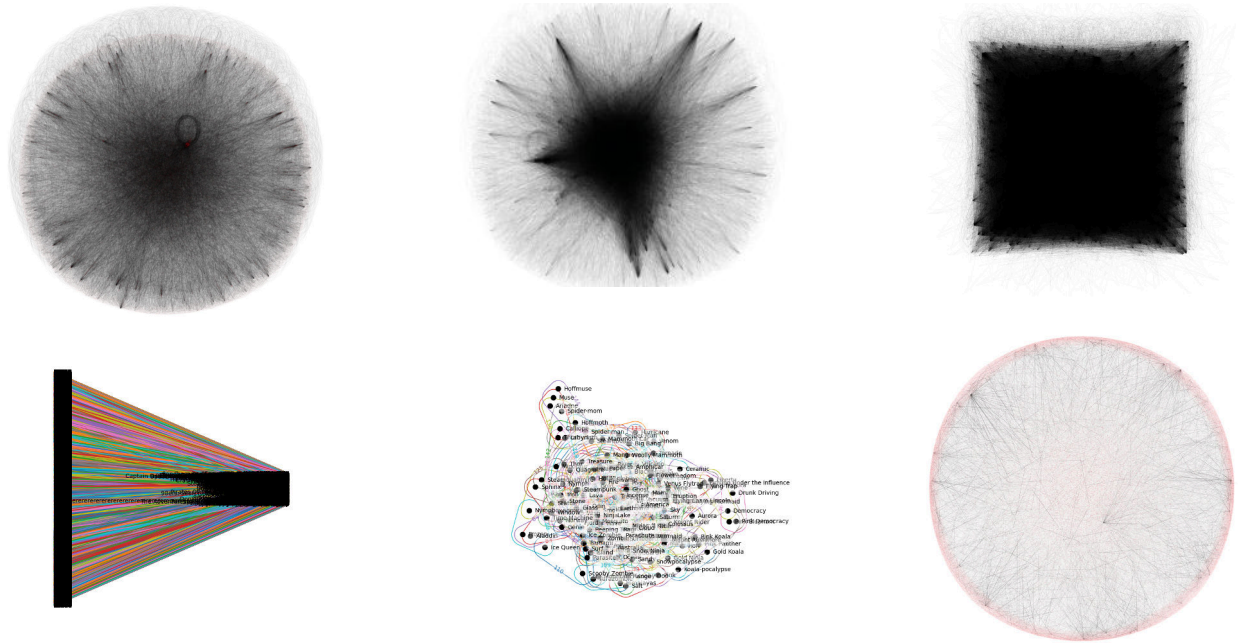
⁴¹This wasn't necessary; all monoids are semigroups anyway.

VIII. Appendix

A. Drawing Big Graphs is Hard

Theorem VIII.A.1 (*The Hairy Blob Theorem*) All plane embeddings of sufficiently dense [hyper-]graphs maintain their legibility under the actions of resizing and heavy compression.⁴²

Proof



□

B. Chart

The code to generate this graph has been sitting in our scripts folder for over a month and nobody remembers what it is supposed to do anymore

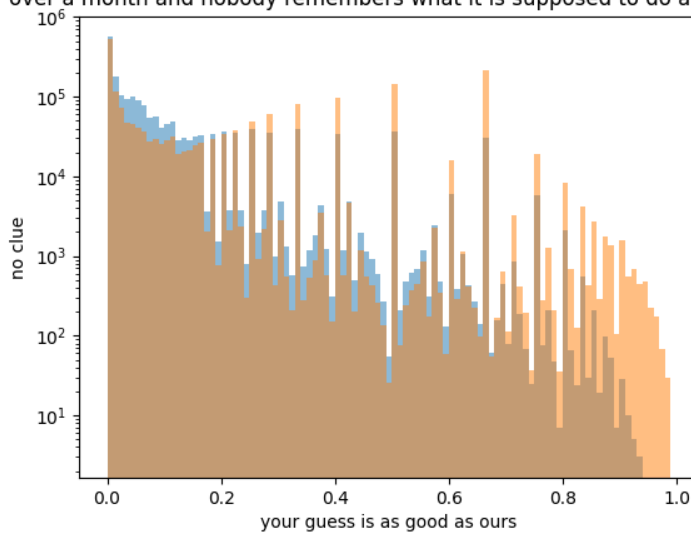


Fig. what ???

⁴²Alternatively stated as “you can’t comb a hypergraph”.

C. Incomplete list of numeric additions that produce strange parentheticals

$0 + 0 = 0$ (Zero)
 $0 + 4 = 4$ (Zero is a number)
 $1 + 1 = 2$ (or 11)
 $1 + 8 = 9$ (18)
 $10 + 10 = 20$ (100)
 $11 + 11 = 22$ (11 is the number of the twin towers)⁴³
 $12 + 8 = 20$ (20 is the number of the last level)
 $12 + 9 = 21$ (21 is the number of the card The World)
 $13 + 13 = 26$ (13 is the number of the beast)
 $14 + 14 = 28$ (28 is a multiple of 14)
 $14 + 15 = 29$ (29 is the number of the level)
 $14 + 6 = 20$ (20 is the number of the last level)
 $15 + 15 = 30$ (30 is the maximum number of items you can have in your inventory)
 $153 + 157 = 310$ (153 is the number of fish caught in the net in the G)
 $153 + 251 = 404$ (Not Found)
 $155 + 159 = 314$ (Pi)
 $1553 + 1553 = 3106$ (1553 is the number of the beast)
 $158 + 158 = 316$ (the number of the beast)
 $16 + 16 = 32$ (16 is the maximum number of items you can have in your inventory)
 $16 + 4 = 20$ (20 is the number of the level)
 $161 + 161 = 322$ (161 is the number of the beast)
 $161 + 243 = 404$ (Page Not Found)
 $162 + 162 = 324$ (162 is the number of the beast)
 $168 + 168 = 336$ (168 is the number of the beast)
 $17 + 17 = 34$ (17 is the number of the level)
 $17 + 5 = 22$ (22 is the number of the Major Arcana)
 $17 + 52 = 69$ (nice)
 $18 + 18 = 36$ (18 is the atomic number of Argon)
 $18 + 3 = 21$ (Blackjack)
 $18 + 4 = 22$ (22 is the number of the Major Arcana)
 $18 + 70 = 88$ (88 is the number of keys on a piano)
 $182 + 182 = 364$ (182 is the number of the beast)
 $1853 + 1864 = 1864$ (1853 is the year of the first patent for a
 $1863 + 1874 = 1874$ (The year of the first World Series)
 $1865 + 1874 = 1874$ (1865 is the year the 13th Amend
 $19 + 19 = 38$ (38 is a multiple of 19)
 $19 + 3 = 22$ (22 is the number of the Major Arcana)
 $19 + 69 = 88$ (88 is the number of keys on a piano)
 $191 + 193 = 384$ (192 is not a prime number)
 $192 + 192 = 384$ (192 is the maximum number of items you can have in your invent
 $199 + 205 = 404$ (Page Not Found)
 $21 + 31 = 52$ (Blackjack)
 $23 + 46 = 69$ (nice)
 $24 + 24 = 48$ (24 is the maximum number of items you can have in your inventory)
 $261 + 535 = 806$ (261 is the number of the beast)
 $27 + 27 = 54$ (27 is the number of the level)
 $3 + 8 = 11$ (11 is a prime number)⁴⁴
 $301 + 603 = 904$ (the number of the beast)
 $3033 + 6072 = 9005$ (9005 is the number of the beast)
 $3035 + 6070 = 9005$ (9005 is a reference to the song "900
 $32 + 32 = 64$ (32 is the maximum number of blocks that can be placed in a 2
 $33 + 33 = 66$ (33 is the number of the highest degree of Scottish Rite Freemason
 $35 + 53 = 88$ (88 is the number of keys on a piano)
 $369 + 441 = 810$ (The sum of the first ten prime numbers)⁴⁵
 $37 + 37 = 74$ (37 is the number of the level)
 $39 + 9 = 48$ (39 is the number of the level)
 $42 + 46 = 88$ (88 is the number of keys on a piano)
 $432 + 432 = 864$ (432 is the number of the universe)⁴⁶
 $4481 + 4518 = 9000$ (9000 is the highest number you can get)

⁴³Too soon?

⁴⁴This statement about prime numbers is correct!

⁴⁵This statement about prime numbers is... less correct.

⁴⁶The AI has spoken, with infinite majesty and calm.

45 + 45 = 90 (90 is the number of the level)
4502 + 4506 = 9008 (9008 is the number of the beast)
4504 + 4504 = 9008 (9008 is a reference to the NES game Duck
46 + 46 = 92 (92 is the atomic number of uranium)
48 + 48 = 96 (48 is the number of the elements)
480 + 480 = 960 (480 is the number of the beast)
5 + 7 = 12 (12 is the answer)
52 + 52 = 104 (52 is the number of cards in a deck)
69 + 69 = 138 (I'm not sure if this is a joke or not)
8 + 8 = 16 (88)
82 + 82 = 164 (82 is the atomic number of lead)
969 + 969 = 1938 (the year the game was released)⁴⁷

IX. Ethics

We waited a short time between requests and respected the website's rate limit. Compared to the rate of organic requests, our 2-3 million requests over the several weeks between launch and the deadline for this paper seem rather puny.

We made a good-faith attempt to compute and reimburse Neal for his resources, which we perhaps unfairly consumed, and notified him early on in our search of what we were doing.⁴⁸

Components in our calculation included but were not limited to:

- Our impact on his TogetherAI bill (computed by multiplying the total number of LLaMa-2 tokens in our database with TogetherAI's rate)
- Additional storage required for caching our requests
- Our impact on server costs
- Sales tax
- Our impact on CloudFlare costs
- Entropy
- Domain name costs
- Inflation during the time of our search
- Differences in local gravity
- The potential increase in request volume if people try to replicate the recipes in this paper
- Degradation of fiber optic cables
- Mercury being in retrograde
- A partridge in a pear tree

Using *Infinite Craft* as our calculator, we summed this to \$22, which we then manually fudged until we got to \$50.

X. Errata and Miscellany

*In order to help this paper fit in with its peers in the scientific mainstream, we have completely fabricated anywhere from zero to three of our charts. Good luck.*⁴⁹

Authors are listed in alphabetical order by the last letter of their last name because this is Ryan Pitasky's only opportunity to earn first-author status based on that metric.

A Google Colab Jupyter notebook to explore or download our horrendous data set can be found at <https://colab.research.google.com/drive/1bLLvgK53a4T0ktyturFpRfM3LKp3UKfy?usp=sharing>. In case of link rot, reach out to the first author (either of the obvious email addresses at gmail should work; title your email with the name of this conference).

⁴⁷*Infinite Craft* was not, in fact, released in 1938.

⁴⁸He was, presumably, too busy putting out fires in his massively popular web app to respond, but we would love to talk!

⁴⁹It's zero.

References

- [1] H. P. Lovecraft, *The Call of Cthulhu*. Weird Tales, 1926.
- [2] H. Arakawa, *Fullmetal Alchemist*. Monthly Shōnen Gangan, 2010.
- [3] A. Ritz, B. Avent, and T. M. Murali, “Pathway Analysis with Signaling Hypergraphs”, *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, vol. 14, no. 5, pp. 1042–1055, 2017, doi: 10.1109/TCBB.2015.2459681.
- [4] S. Krieger and J. Kececioğlu, “Heuristic shortest hyperpaths in cell signaling hypergraphs”, *Algorithms for Molecular Biology*, vol. 17, no. 1, p. 12, May 2022, doi: 10.1186/s13015-022-00217-9.
- [5] J. Gao, Q. Zhao, W. Ren, A. Swami, R. Ramanathan, and A. Bar-Noy, “Dynamic Shortest Path Algorithms for Hypergraphs”, *IEEE/ACM Transactions on Networking*, vol. 23, no. 6, pp. 1805–1817, 2015, doi: 10.1109/TNET.2014.2343914.
- [6] B. Kay *et al.*, “Hypergraph Topological Features for Autoencoder-Based Intrusion Detection for Cybersecurity Data”, in *ML4Cyber Wshop., Int. Conf. Machine Learning 2022*, 2022.

DO ORDAIN AND ESTABLISH..... 'MERICAAAA

- 71 As American as Apple Pie: A Search to Define Americanness
in the Context of all Apple Pie Recipes on the Internet**

Matthew A. Charleston, Joshua M. Shaw, and Lennon J. Anderson

As American as Apple Pie: A Search to Define Americanness in the Context of all Apple Pie Recipes on the Internet

Authors: Matthew A. Charleston¹, Joshua M. Shaw², Lennon J. Anderson³

Affiliations: 1. Former President Ohio State Bread Club
2. Top 5% Hearthstone Player
3. Former President Ohio State Hacking Club

Abstract

Since the halcyon days of Johnny Appleseed, America has been strongly associated with apples and more specifically apple pie. For what can more embody the American spirit than the sticky, sugary, (high fructose corn syrupy?), buttery, and flaky masterpiece that is the apple pie. However, neither the apple nor the apple pie originate in America. The apple being the oldest fruit (God, 1400 BC) and the pie one of the oldest desserts. In order to best embrace, define, restrict, and confine the American spirit, it is necessary to define the concept of Americanness, and what better local to scour than the venerable apple pie.

In this paper we condition and filter the Recipe NLG dataset into a list of ingredients associated with each apple pie recipe. We then sort the recipes into American and un-American Apple pies and calculate, based on the mean quantity of ingredients, what makes an American and un-American apple pie. Finally, the average American Pie is baked for gustatory analysis to determine if the average American has good taste.

Introduction

The common apple is the most ubiquitous and culturally important fruit crop in temperate climates (Cornille, 2012) and has been slowly domesticated into a household staple over thousands of years. Apples originated in prehistoric times and were omnipresent during the founding of America, developing simultaneously with the early colonies. Early Records show the apple tree growing in New England just after the founding of the Plymouth Colony. The spread of apples in America preceded the colonists, being distributed by the Native American Tribes who were then followed by the eccentric John Chapman (Appleseed). Starting around 1850 the apple began to transition from domestic orchard growth to important commercial proportions, ultimately becoming the most important fruit in existence, especially in the United States, the principle apple producing country in the world (Folger & Thomson, 1923).

As surely as the apple is enshrined in the hearts and history of Americans, so too is the apple pie. In fact, the apple pie is the only food officially recognized by the Department of Defense as an American Symbol (Bunschoten, 2014). The apple pie was a common item in the lives of the English colonists as they settled North America, and the

development of the apple pie mirrored the development of the growing nation. The first American cookbook, *American Cookery* (Simmons, 1796), contained the first recipe for an American apple pie. As the nation expanded westward and apple emissaries like Chapman spread the good news, corporate interests became involved as well. The Ohio Company offered land deeds to settlers who among other things, “set out at least 50 apple or pear trees” ensuring the spread of the delicious fruit. Despite a few brief deviant ventures like the removal of the upper crust in 1846 and the “Apple Pie without Apples” developed during the blockade of the confederate states (it substituted crackers and cream of tartar for apples), the pie spread through the breadth of the nation. The apple pie’s true association with national identity likely began during The Great War as soldiers wished longingly from the tranches for flavors of home. Through World War II as well, Betty Crocker and Elizabeth Roosevelt strode hand in hand down the pie aisle of America’s hearts, preaching home cooking and rationing frugality. In the post war-years the popularity of apple pie soared, becoming the favorite dessert of young Americans by 1958 (Bunschoten, 2014). As surely as the apple pie is distinctly American, it is equally unsure what separates the American Apple Pie from its cultural ancestors and rivals.

The current state of the art in recipe analysis has a strong basis in machine learning models. Recent models have been trained on large datasets of recipes to culturally adapt recipes between languages (Cao, et al., 2024) and even strategically substitute ingredients in recipes (Fatemi, et al., 2023).

One of the largest and most frequently cited collections of online recipe datasets is the Recipe1M+ dataset (Marín, et al., 2019) (Salvador, et al., 2017). This dataset is well known and has been used to study everything from recipe generation from images (Salvador, et al., 2019), to image generation from recipes (Salvador, et al., 2017). However, the Recipe1M+ dataset is not quite as public as described. The provided download link from Marin et.al. returned a http 500 error and told us to email the server administrator Yusuf for help. Unfortunately, despite our desperate and repeated attempts to contact Yusuf, our emails received no reply.

We next turned to the RecipeNLG dataset (Bien, et al., 2020), created at Polant’s Poznan University of Technology. Recipe NLG attempts to address the growing demand for recipe datasets to be used in deep learning experiments and especially ones centered on Natural Language Processing (NLP) tasks. The dataset contains the ingredients and respective quantities for over 2 million recipes, making it the largest recipe dataset at the time of its creation.

Experimental Setup

In keeping with the same spirit of rationalism that inspired America’s founding fathers, the decision was made to utilize human intelligence and mathematics to distill the Americanness from the Apple pie. The basic theory is laid out in Equation 1, where two exclusive subsets can be separated from the set of all Apple Pie recipes.

$$uAAP = AP - AAP \tag{1}$$

Symbol	Definition
--------	------------

AP	Apple Pie (Set of all Recipes)
AAP	American Apple Pie
uAAP	Un-American Apple Pie

Each subset is arranged as an array, where each row is a specific recipe and each column a specific ingredient. An example array is shown in Table 1. The ingredient columns of each subset can then be averaged to determine the quantity of ingredients in an average American Apple Pie and an average un-American Apple Pie. The differences in the average quantities of each ingredient are calculated according to Equation 2, to distinguish what makes an American Apple Pie unique, illustrated for example purposes in Table 2.

$$\text{Americanness} = \text{Mean}(\text{AAP}) - \text{Mean}(\text{uAAP}) \quad (2)$$

Table 1. Example Apple Pie Array

Recipe	Apples	Flour (c.)	Cinnamon (tsp.)	Sugar (c.)
Elaine’s Apple Pie	6	2	3	1
George’s Apple Pie	5	1	2	2
Jerry’s American Apple Pie	6	1	4	1.5

Table 2. Example Mean Pie Arrays

Recipe	Apples	Flour (c.)	Cinnamon (tsp.)	Sugar (c.)
Mean(uAAP)	5.5	1.5	2.5	1.5
Mean(AAP)	6	1	4	1.5
Americanness	0.5	-0.5	1.5	0

In the example dataset here we can see that the largest distinguishing characteristic of Americanness is excess cinnamon with a slight preference for more apples and less flour.

Data Analysis

The Recipe NLG dataset was initially filtered to select only recipes that contained the words “Apple Pie”, creating the new RecipeAP (apple pie) dataset with 6,803 apple pie recipes and 14,788 ingredients. Unfortunately for our heroes, the Recipe AP dataset, while being sufficiently conditioned for machine learning research, was ill-conditioned for straightforward mathematical analysis. Many ingredients are repeated with minor differences like “Pie Crust” vs “Pie Shell” or Major Differences like “Pie Crust” and “The dough may be made ahead of time and stored in the refrigerator for up to 5 days or in the freezer for up to two weeks. The pie is best served the same day.”, or even completely absurd ingredients like “Dipping Fork”. In order to analyze a concise and representative dataset, similar ingredient columns needed to be merged and erroneous recipes like “Bourbon Apple Pie Ice Cream” removed.

Further find and replace operations were conducted to standardize all measurement descriptions, replacing “c.” with “cups” and “Tablespoons” with “Tbs.”. Column merging was initially attempted manually using Microsoft Excel. Redundant columns were identified, then the entire array was sorted around redundant column 2 from high to low. The values from column 2 were cut and pasted into column 1 and the redundant column was deleted. This task was attempted for approximately 5 hours with the merging of 150 columns. This method was abandoned after it was calculated with the merge rate of 30 columns per hour, it would take 486 hours, or 20 full days to reduce the dataset to a reasonable 200 columns.

Returning to the RecipeAP dataset, while it is well suited to those working in the field of natural language processing, it is less conducive to tasks where the goal is to compare specific quantities of ingredients. In order to at least partially ameliorate these difficulties, some of pre-processing of the data set was performed. Some data that was not pertinent to the task at hand was dropped entirely. Among these columns were the directions and the Named Entity Recognition (NER) data. Then, in order to get the ingredient data into a more usable form for analysis, recipe ingredients were transposed to become the columns of the new dataset, with the recipe names remaining as the first column. A best effort was given to get the quantities for each ingredient for each recipe where applicable, but due to the semi-structured data this task proved difficult, and some amount of manual labor for data sanitization was required.

Manual analysis then starts by cleaning the dataset's columns and entries. Characters like punctuation, backslashes, and asterisks are filtered from ingredient names. Recipe names and ingredient amounts have similar but less stringent filtering. Care was taken to ensure the cleaning did not change the meaning of entries. For example, a nice filter would be to lowercase all entries but this would obscure teaspoon (denoted with a lowercase t) versus tablespoon (denoted with a capital T). Additional filters include replacing zeros and empty entries with NaN as well as removing any columns with no entries after filtering.

Removing characters from ingredient names creates duplicate columns (e.g. the columns "tsp. flour" and "tsp flour" would both become "tsp flour"), which would create problems when merging columns later. To avoid this duplicate ingredients are appended with a monotonically incrementing counter (in the previous example there would be two columns, "tsp flour", "tsp flour.1", ... "tsp flour.69")

Additionally, column names that start with "plus", "+", "&" or "and" were removed. This is because columns with "and" like terms are assumed to be a part of an improperly parsed larger section of the recipe with another ingredient. For example, a column "plus 2 Tbsp flour" was most likely something along the lines of "1 tsp salt plus 2 Tbsp flour". Any entries in the "plus 2 Tbsp flour" are not actually flour and are some other unknown ingredient. Lastly, "the world's largest apple pie" recipe is dropped from the dataset due to its outrageous list of ingredients. The recipe calls for 227,359 teaspoons of sugar.

The bulk of the remaining analysis is done through a mix of automated and manual methods. Manually, CSV files containing each column and its estimated unit, quantity, and bulk multiplier were created. For example "cans of applesauce, cups, 1.5, 1" where the scalar value 1.5 can be used to put the unit "cans" into "cups". These inform the automated analysis of assumed values for the column. Entries that omit a unit or quantity will fall back on the column default value. Quantity and units are converted to a common unit (teaspoon) before being inserted into the final dataset.

The final prepared dataset is then sorted into categories based on the recipe title. Four groups are created. American Apple Pies and UnAmerican Apple Pies, in addition to Dutch Apple Pies and unDutch Apple Pies as a control group.

Results and Conclusion

The results of the analysis are presented for the Dutch control group first in Table 1 and are plotted in Figure 1 as a stacked bar chart to display the percentage contribution of each key ingredient.

Table 1:Dutchness

Dutchness of Apple Pie (units in tsp)											
Apple Pie	Allspice	Apples	Butter	Cinnamon	Cloves	Dairy	Flour	Nutmeg	Salt	Shortening	Sugar
Dutch	0.014	159.270	27.261	1.089	0.006	11.899	86.064	0.592	0.482	1.999	125.366
UnDutch	0.060	147.649	29.229	1.211	0.090	24.809	64.458	0.506	0.650	10.035	94.849

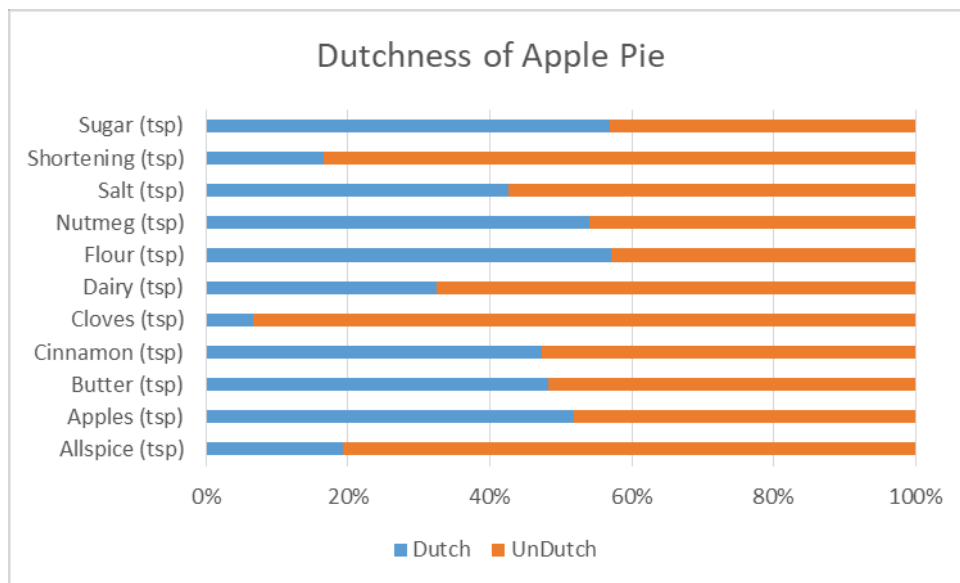


Figure 1: Dutchness Plotted

It can be seen that Dutchness is not characteristic of an excess of any one ingredient, as the Dutch apple pie has average or less than average amounts of each key ingredient. The main traits of Dutchness are a minimal amount of shortening, allspice, and cloves. The Dutch Apple Pie is included as a control group namely because of its signature construction, being invariably topped with a crumble made from butter, flour, and sugar. This crumble would presumably show up in the data when compared to UnDutch Apple Pies, and indeed it does. The Dutch Apple Pie has an above average amount of sugar and flour and an average amount of butter, which may perhaps be accounted for by the more austere Dutch taste

Having validated the pie measurement methodology we proceed to the important question at hand, the determination of Americanness. The results are listed in Table 2 and plotted in Figure 2.

Table 2: Americanness

Americanness of Apple Pie (units in tsp)											
Apple Pie	Allspice	Apples	Butter	Cinnamon	Cloves	Dairy	Flour	Nutmeg	Salt	Shortening	Sugar
American	0.382	175.750	33.913	1.059	0.002	4.000	86.106	1.422	0.419	3.885	204.059
UnAmerican	0.056	147.707	29.141	1.209	0.088	24.659	64.817	0.500	0.647	9.885	94.708

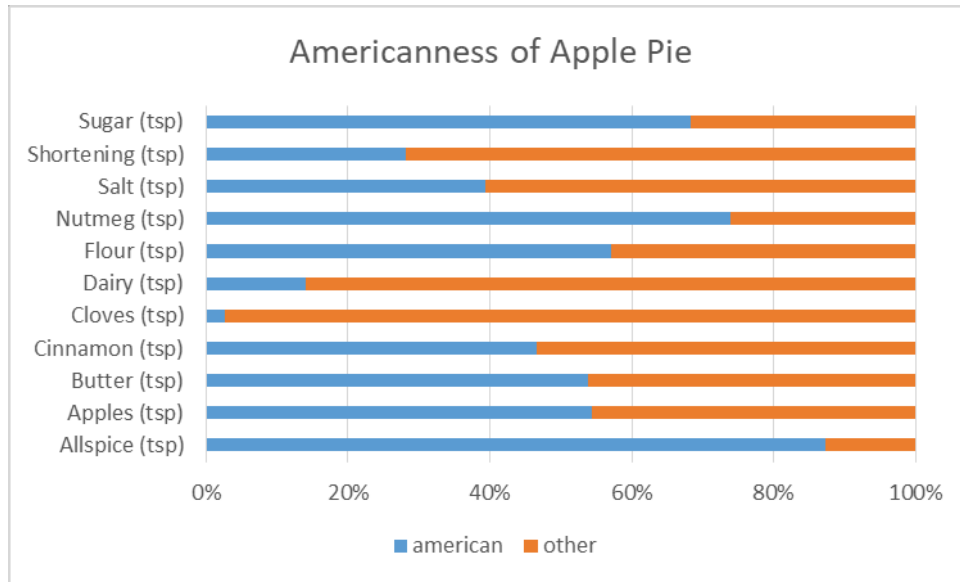


Figure 2: Americanness Plotted

Perhaps unsurprisingly one of the key characteristics of Americanness is an excess of sugar. American pies also tend to use an average amount of butter and a lesser amount of shortening. They tend to be spiced heavily with a gratuitous amount of allspice and nutmeg and yet contradictorily completely devoid of cloves. Most surprisingly, American apple pies use little dairy compared to UnAmerican pies, as the American Midwest which is most known for its pies is also a top dairy producer and consumer.

Not all of the sorted recipes contained every key ingredient and the number of recipes for each category are listed in Table 3 below, in the hopes that it may assist in future research. The total number of recipes was used when calculating the average quantity of ingredients in the respective categories.

Table 3: Quantity of Recipes

Number of Recipes in Each Category that Contain the Specified Ingredient												
Apple Pie	Allspice	Apples	Butter	Cinnamon	Cloves	Dairy	Flour	Nutmeg	Salt	Shortening	Sugar	TOTAL
American	8	49	39	42	1	12	35	31	20	8	43	52
UnAmerican	234	5108	4176	4180	204	1589	3805	1902	2088	537	4967	6319
Dutch	7	127	123	112	4	38	118	43	54	11	135	157
UnDutch	235	5030	4092	4110	201	1563	3722	1890	2054	534	4875	6214

Conclusion

The Apple Pie is as much a venerated symbol of America as the Statue of Liberty or the much celebrated “hot dog”. But the pie, unlike these other renowned symbols has a distinct breed of Americanness. Passed down from the depths of antiquity, the apple and the pie have traveled and evolved with humanity itself through the ages. Until now, where in our modern age the American apple pie has become a symbol festooned with the very traditions, history, and essence of America today.

As our great country weathers struggles and storms, we harken back to the Apple Pie to remind us of who we are, and what binds us all together in Americanness. And that ingredient, fellow Americans is sugar.

Coming in at a whopping 204 teaspoons (4 ¼ cups), the American apple pie uses more than twice as much sugar as the UnAmerican Pie. It also uses a significant amount of spice with its sugar, with allspice and nutmeg significantly exceeding the average and evoking flavors of the celebrated pumpkin pie as well. But perhaps with the recent divisions in America, what can most draw us together as a nation is our shared hatred for cloves.

Recipe

In order to taste the distilled Americanness as concluded above, the Authors fed the paper into Chat-GPT and asked it to create the most American Apple Pie Recipe using the ingredients in table 2. That recipe follows:

****The Most American Apple Pie Recipe****

Ingredients:

- Apples: 175.75 teaspoons (which is about 5.84 cups or approximately 6 medium-sized apples, peeled, cored, and sliced)
- Sugar: 204.059 teaspoons (which is about 4.25 cups, to sweeten your taste of freedom)
- Allspice: 0.382 teaspoon (just a dash for that homegrown kick)
- Nutmeg: 1.422 teaspoons (for that warm, spicy American hug)
- Cinnamon: We'll use 2 teaspoons here, because Americans love their cinnamon
- Butter: 33.913 teaspoons (which is about 7.1 ounces or a little less than 1 cup, for that rich, flaky crust)
- Flour: 86.106 teaspoons (which is about 1.8 cups, for the crust)
- Salt: a pinch (let's say 1/2 teaspoon, for the crust)
- Water: enough to make the dough workable (let's use 4 to 6 tablespoons, depending on the dough consistency)

Instructions:

1. ****Prepare the Pie Crust:****

- In a large mixing bowl, whisk together about 1.8 cups of flour with 1/2 teaspoon of salt.
- Cut in the butter (just under 1 cup) until the mixture resembles coarse crumbs.
- Sprinkle in ice water, 1 tablespoon at a time, mixing until the dough comes together in a ball.

2. ****Make the Filling:****

- In a large bowl, combine about 6 cups of sliced apples with 4.25 cups of sugar, 0.382 teaspoon of allspice, 1.422 teaspoons of nutmeg, and 2 teaspoons of cinnamon. Stir them up until well mixed.

3. ****Assemble the Pie:****

- Roll out half of the dough on a floured surface to fit your pie plate.
- Transfer the dough to the plate, and trim the edges.
- Pour in the apple filling, heaping it in the center.
- Roll out the remaining dough and place it over the filling. Trim the overhang, and crimp the edges to seal.

4. ****Bake the Pie:****

- Preheat your oven to 425°F (220°C).
- Bake the pie for 15 minutes, then reduce the oven temperature to 350°F (175°C) and continue baking for another 45 minutes, or until the crust is golden brown and the filling is bubbly.

5. ****Cool and Serve:****

- Let the pie cool on a rack before slicing. This gives the filling time to set.

Pie

The Authors attempted to follow the recipe and bake the most American apple pie, Figure 3. However, fully embracing the spirit of freedom, they chose to not use the full 4 cups of sugar. Nevertheless, it was delicious.



Figure 3: American Apple Pie

Future work

The authors acknowledge that certain parsing errors and the small sample size for the American Apple Pie dataset may have improperly weighted the results. Certain errors in the base dataset were difficult to control for, like “34” being used instead of “3/4” which could be a key factor in why some ingredient amounts seem rather large. Further study is necessary to confirm the preliminary results presented here.

Additionally, with the computed dataset, there are many opportunities for further pie-related research. Particularly of interest is determining the qualities that define the “Best” apple pie, the “Original” apple pie, and most importantly, “Grandma’s” apple pie.

The analysis code and files are available: <https://github.com/YurBoiRene/apple-pie>

Works Cited

- Amandine Cornille, P. G.-R.-G. (2012). New Insight into the History of Domesticated Apple: Secondary Contribution of the European Wild Apple to the Genome of Cultivated Varieties. *PLOS Genetics*.
- Bien´, M., Gilski, M., Maciejewska, M., Taisner, W., Wisniewski, D., & Ławrynowicz, A. (2020). RecipeNLG: A Cooking Recipes Dataset for Semi-Structured Text Generation. *Proceedings of The 13th International Conference on Natural Language Generation* (pp. 22-28). Dublin, Ireland: Association for Computational Linguistics.
- Bunschoten, R. C. (2014). *As American as Apple Pie: The History of American Apple Pie and Its Development into a National Symbol*. Bard College.
- Cao, Y., Kementchedjheva, Y., Cui, R., Karamolegkou, A., Zhou, L., Dare, M., . . . Hershovich, D. (2024). Cultural Adaptation of Recipes. *Transactions of the Association for Computational Linguistics*, 80-99.
- Fatemi, B., Duval, Q., Girdhar, R., Drozdal, M., & Romero-Soriano, A. (2023). *Learning to Substitute Ingredients in Recipes*.
- Folger, J., & Thomson, S. (1923). *The Commercial Apple Industry of North America*. New York: The Macmillian Company.
- God. (1400 BC). *The Bible*. Jerusalem: Holy Catholic Church.
- Marín, J., Arito, B., Ferda, O., Nicholas, H., Amaia, S., Yusuf, A., . . . Antonia, T. (2019). Recipe1M+: A Dataset for Learning Cross-Modal Embeddings for Cooking Recipes and Food Images. *IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE*.
- Salvador, A., Drozdal, M., Giro-i-Nieto, X., & Romero, A. (2019). Inverse Cooking: Recipe Generation From Food Images. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. Long Beach, CA: IEEE.
- Salvador, A., Hynes, N., Aytar, Y., Marin, J., Ofli, F., Weber, I., & Torralba, A. (2017). Learning Cross-Modal Embeddings for Cooking Recipes and Food Images. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. Honolulu, HI: IEEE.
- Simmons, A. (1796). *American Cookery*.