# If Scaling LLMs Is the Answer Then What Is the Question?

Daniel Khashabi

JOHNS HOPKINS UNIVERSITY

# The success we dreamed of

Language models that are remarkably capable at solving many important NLP benchmarks.

# The success we dreamed of

Let $n$ be an even positive integer. Let $p$ be a monic, real polynomial of degree $2n$; that is to say,

$$p(x) = x^{2n} + a_{2n-1}x^{2n-1} + \cdots + a_1 x + a_0$$

for some real coefficients $a_0, \ldots, a_{2n-1}$. Suppose that $p\left(\frac{1}{k}\right) = k^2$ for all integers $k$ such that $1 \leq |k| \leq n$.

Find all other real numbers $x$ for which $p\left(\frac{1}{x}\right) = x^2$.

# The success we dreamed of

Let $n$ be an even positive integer. Let $p$ be a monic, real polynomial of degree $2n$; that is to say,

$$p(x) = x^{2n} + a_{2n-1}x^{2n-1} + \cdots + a_1 x + a_0$$

for some real coefficients $a_0, \ldots, a_{2n-1}$. Suppose that $p\left(\frac{1}{k}\right) = k^2$ for all integers $k$ such that $1 \leq |k| \leq n$.

Find all other real numbers $x$ for which $p\left(\frac{1}{x}\right) = x^2$.

Thought for 5 seconds ⌄

# The success we dreamed of

Let $n$ be an even positive integer. Let $p$ be a monic, real polynomial of degree $2n$; that is to say,

$$p(x) = x^{2n} + a_{2n-1}x^{2n-1} + \cdots + a_1 x + a_0$$

for some real coefficients $a_0, \ldots, a_{2n-1}$. Suppose that $p\left(\frac{1}{k}\right) = k^2$ for all integers $k$ such that $1 \leq |k| \leq n$.

Find all other real numbers $x$ for which $p\left(\frac{1}{x}\right) = x^2$.

---

Thought for 5 seconds ⌄

To solve this problem, we need to find all real numbers $x$ such that

$$p\left(\frac{1}{x}\right) = x^2,$$

where $p$ is a monic real polynomial of degree $2n$, and it is given that

$$p\left(\frac{1}{k}\right) = k^2$$

for all integers $k$ with $1 \leq |k| \leq n$.

**Understanding the Given Information:**

1. **Given Data Points:** We have $2n$ data points from $k = \pm 1, \pm 2, \ldots, \pm n$, which provide the values $p\left(\frac{1}{k}\right) = k^2$.

2. **Polynomial Function:** $p(x)$ is monic of degree $2n$, which means its highest-degree coefficient is 1.

- The leading term of $s(x)$ is $-x$ .

- The leading term of $m(x)h(x)$ is $(x^{2n})(-x^2) = -x^{2n+2}$, ensuring the degrees match.

- Comparing the constant term of $s(x)$ using $x = 0$:

- $$ s(0) = p^(0) - 0 = 1, where\ p^(x) = x^{2n} p\left(\frac{1}{x}\right) $$.

- Evaluating $m(0)$ and $h(0)$:

- $m(0) = (-1)^n(n!)^2$.

- $h(0) = -0 + c = c$.

- Since $s(0) = m(0)h(0)$:

$$1 = (-1)^n(n!)^2 c,$$

which gives:

$$c = \frac{1}{(-1)^n(n!)^2} = \frac{(-1)^n}{(n!)^2}.$$

- For $n$ even (since $n$ is even), $(-1)^n = 1$, so $c > 0$.

**Finding Additional Solutions:**
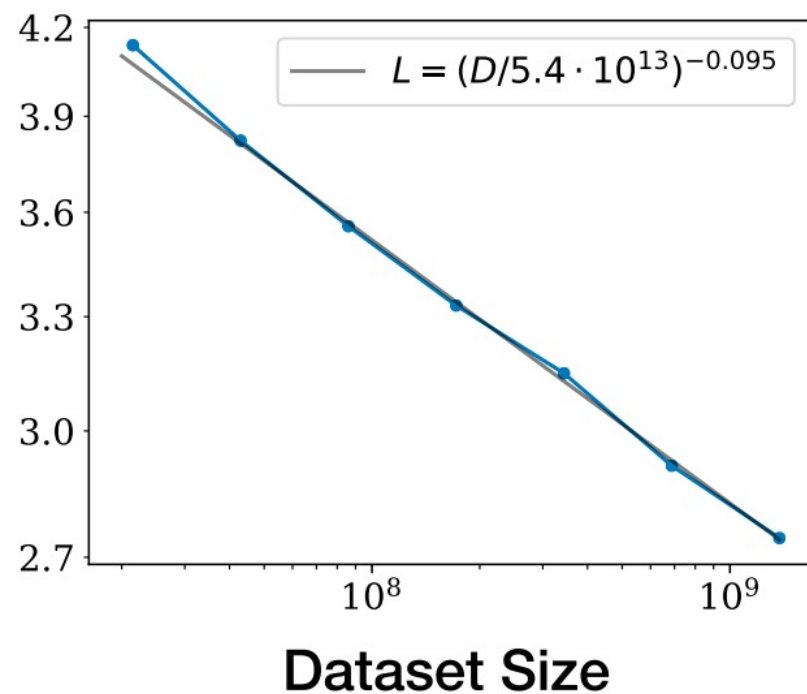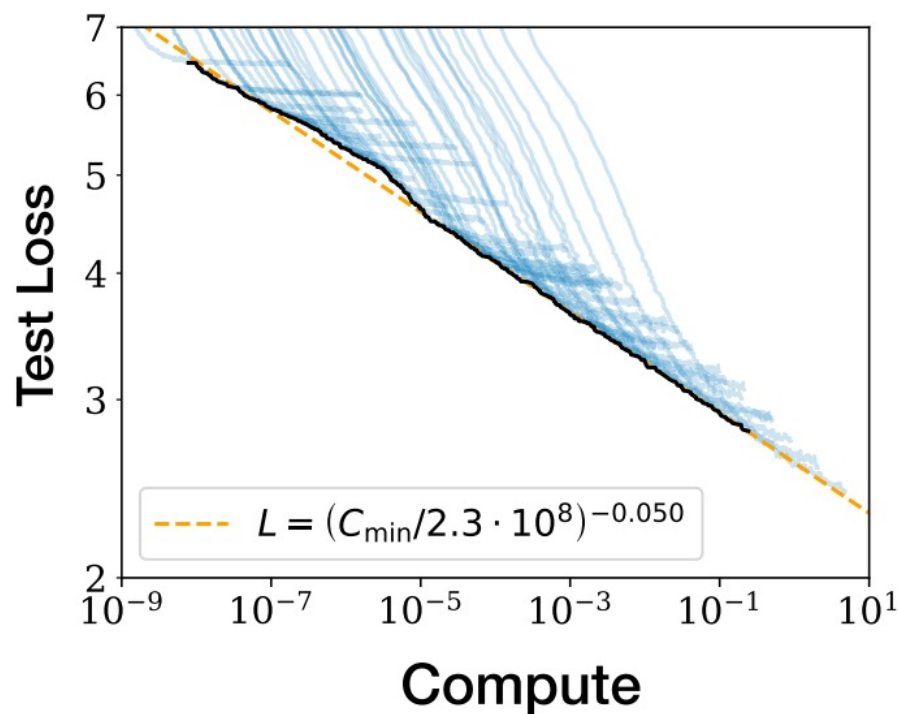
Solving $h(x) = 0$:

$$-x^2 + c = 0 \implies x^2 = c \implies x = \pm\frac{1}{n!},$$

since $c = \frac{1}{(n!)^2}$.

**Conclusion:**

The additional real numbers $x$ satisfying $p\left(\frac{1}{x}\right) = x^2$ are $x = \pm\frac{1}{n!}$.

# Scaling Laws: A Low-risk Recipe



Kaplan et al. 2020; among others

More data (and compute) leads to better models.

# Meanwhile, whispers of a bubble …

**OpenAI's Sam Altman sees AI bubble forming as industry spending surges**

CNBC

AUG 18, 2025

**Warning: Our Stock Market Is Looking Like a Bubble**

The New York Times

OPINION

Oct. 14, 2025
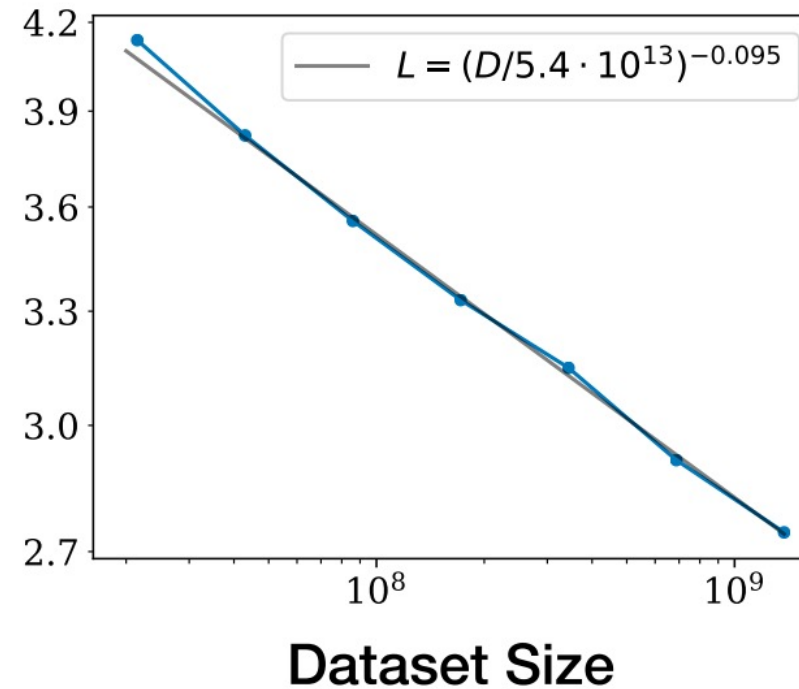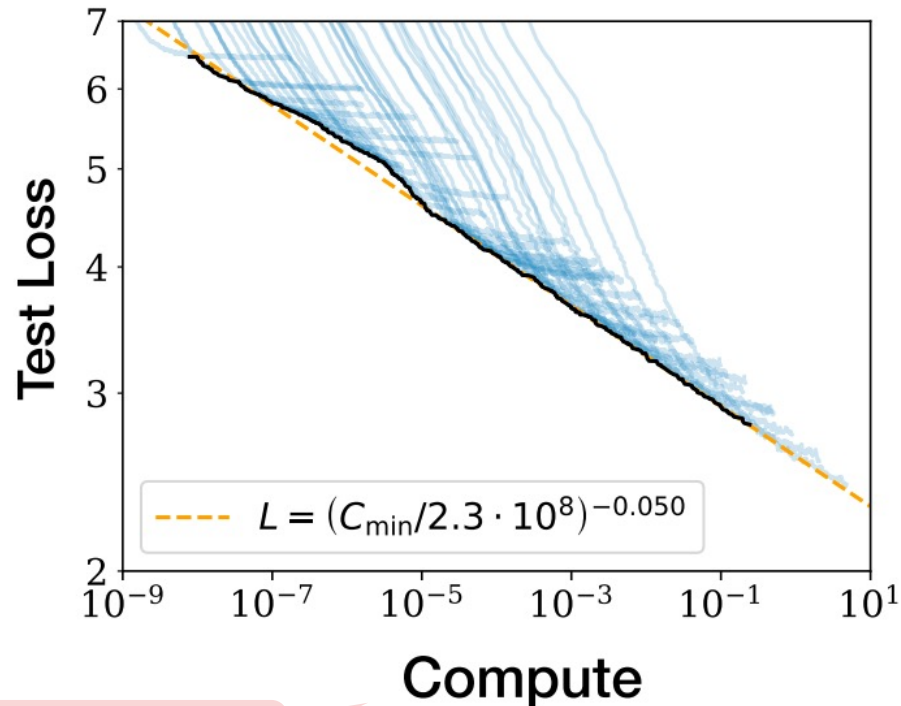
**Is There an A.I. Bubble? And What if It Pops?**

The New York Times
The Daily

Nov. 20, 2025

# Making sense of the "bubble" concerns

- The progress is real. However, many challenges remain.

- There may be various reasons:
  - Profit-cost mismatch,
  - Future regulations,
  - Lack of enduring moats,
  - etc.

- **Expectation-capabilities mismatch:** Investors bet on rapidly improving capabilities.
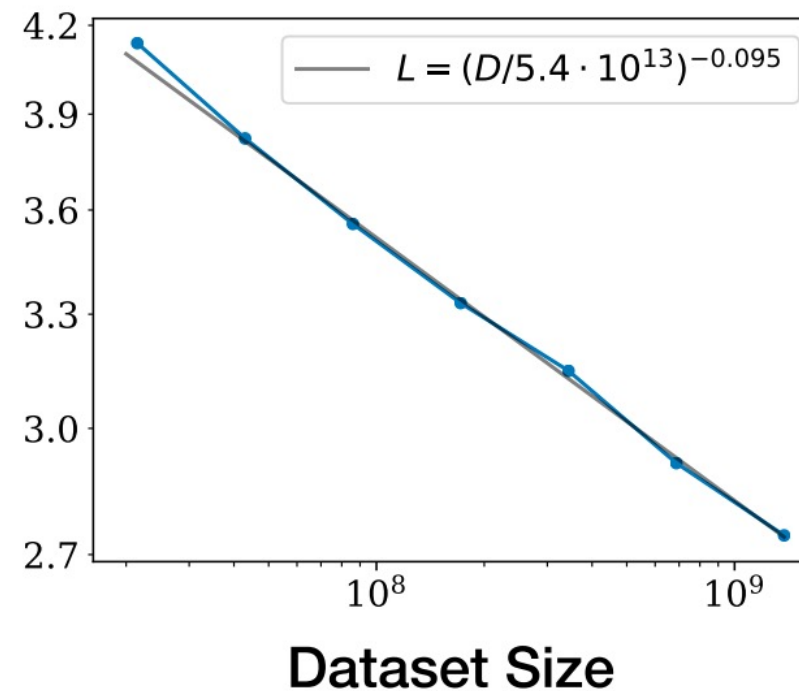
# Limits of scaling "laws"
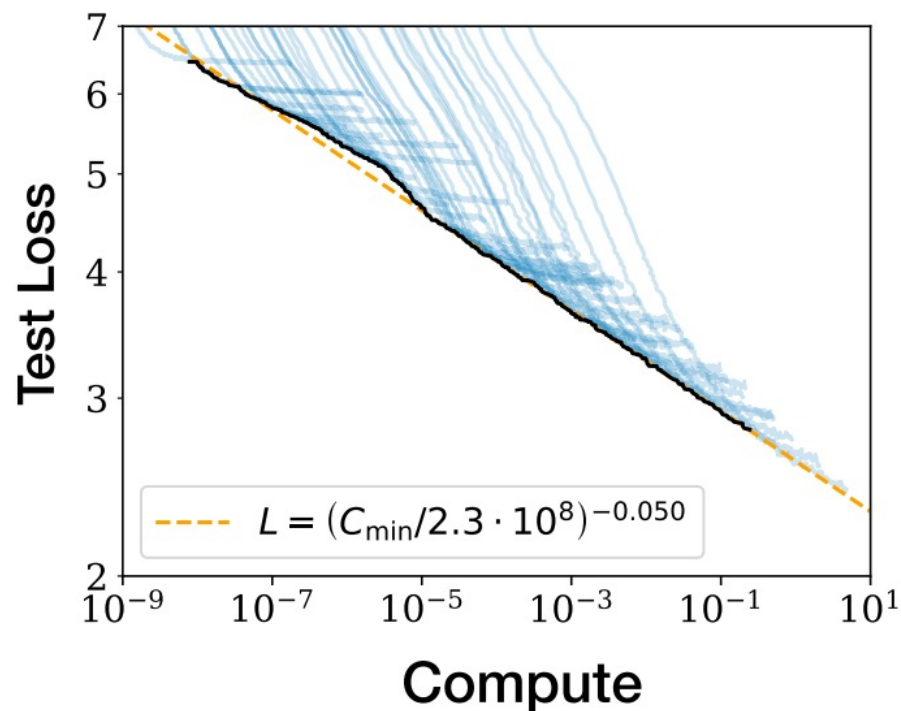


Linear

Exponential

Kaplan et al. 2020; among others

Exponential

Diminishing returns w/ scaling (compute, data, human supervision.)
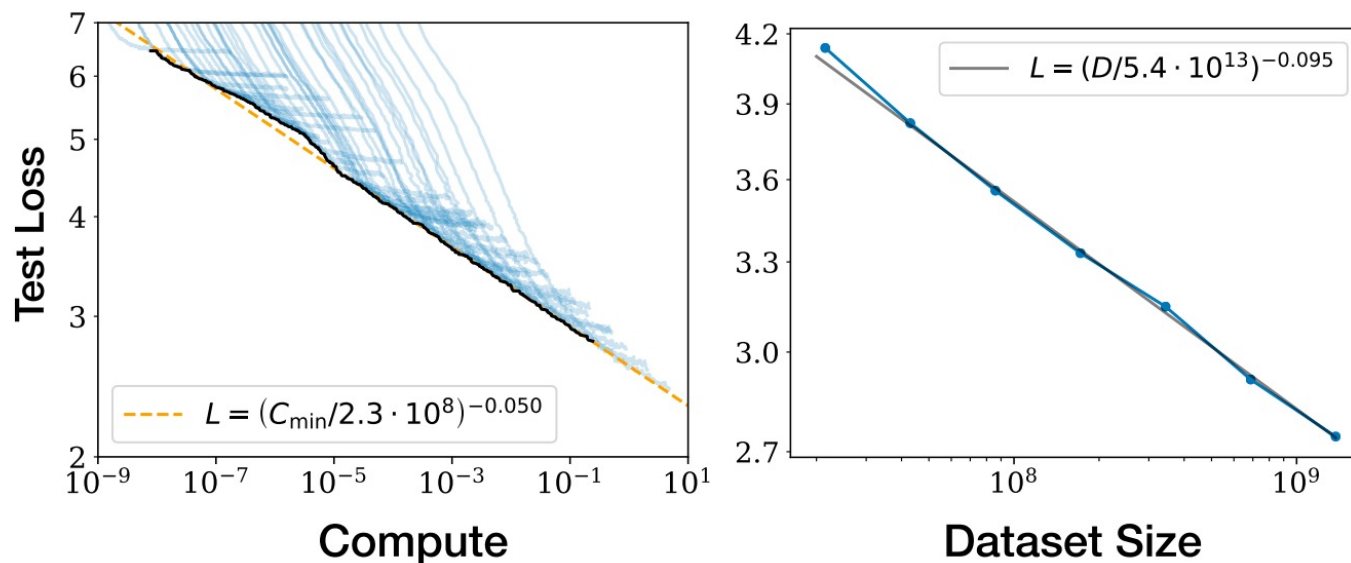
# Limits of scaling "laws"



Kaplan et al. 2020;
among others

Which data? How is it distributed?

# Today: Deconstructing Scaling Laws

- Scaling laws hide important data-dependent effects that current "laws" fail to capture.

- We'll examine LLM behaviors that become apparent only once you look inside the data distribution.

# Roadmap

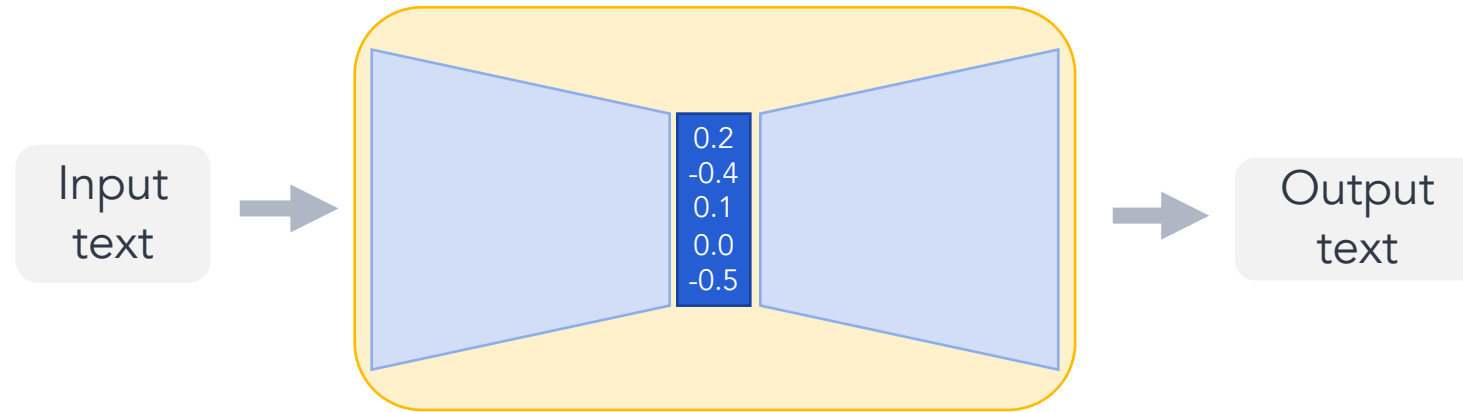1. Scaling is <span style="color:darkred">distribution-dependent</span>

2. Learning emerges <span style="color:darkred">beyond human language</span>

3. LLMs show <span style="color:darkred">belief inertia</span>

Wait … How did we get here?

# Language Models

Input text → **LM** → Output text

[Bengio et al. '04, Peters et al. '18, Raffel et al. '20, Brown et al. '20, many others]

# Language Models: Pre-training

Input
text

0.2
-0.4
0.1
0.0
-0.5

Output
text

Bengio et al. '04, Peters et al. '18, Raffel et al. '20, Brown et al. '20, many others

# Language Models: Next-Token Prediction

| Johns Hopkins University is in _____. | → | **LM** | → | Baltimore |

Bengio et al. '04, Peters et al. '18,  Raffel et al. '20, Brown et al. '20, many others

# Language Modeling ≠ Following User Intents

Explain "space elevators" to a 6-year-old.

**LM**

Explain gravity to a 6-year-old.
Explain black-holes to a 6-year-old.
Explain big bang to a 6-year-old.
….

LMs are not "aligned" with user intents.

# Language Modes — Post-training

- "Aligning" LMs with our intents embedded in instructions.
  - Supervised Fine-tuning (Behavior cloning) in labeled data.



Weller et al. 2020. Mishra et al. 2021; Wang et al. 2022,  Sanh et al. 2022; Chung et al. 2022, many others
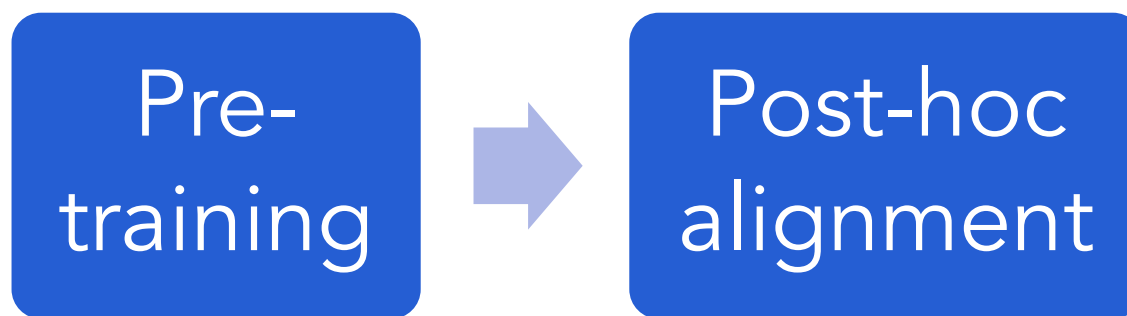
# Language Modes — Post-training

- "Aligning" LMs with our intents embedded in instructions.
  - Supervised Fine-tuning (Behavior cloning) in labeled data.
  - Reinforcement Learning on preference data or verifiers.



Christiano et al. 2017; Stiennon et al. 2020; Ouyang et al., 2022; many others

# The Overall Recipe for Modern LLMs

- Almost all the modern models follow this recipe:

**Pre-training** ➡ **Post-hoc alignment**

- Note, we did not end up here overnight.
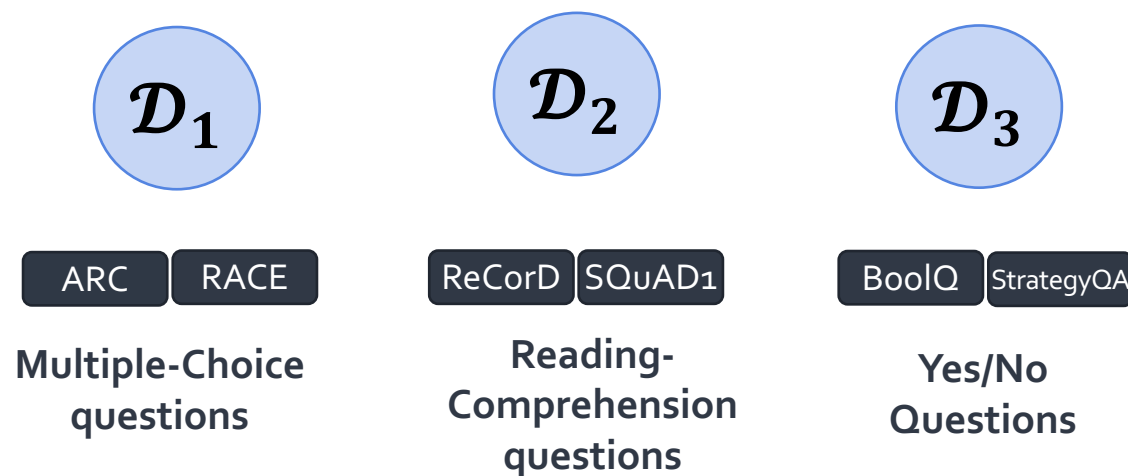- A lot of incremental progress to get to this point.

# Time travel to ~2019

# Challenge: Incompatible Datasets

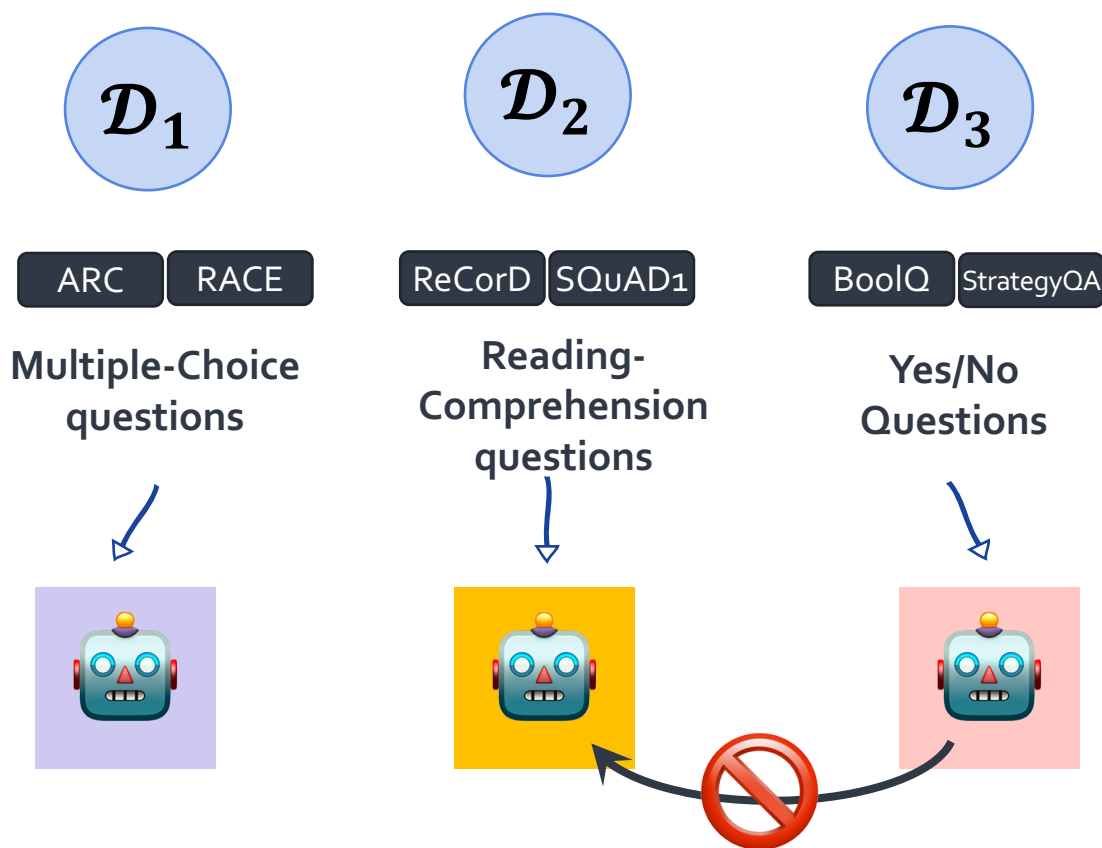- Question-answering datasets carried different assumptions

Dataset-groups for answering questions

$\mathcal{D}_1$    $\mathcal{D}_2$    $\mathcal{D}_3$

| ARC | RACE |
| ReCorD | SQuAD1 |
| BoolQ | StrategyQA |

**Multiple-Choice questions**

**Reading-Comprehension questions**

**Yes/No Questions**

# Result: We were stuck with dataset-specific models

- Despite having pre-trained models, everyone kept training task-specific models.

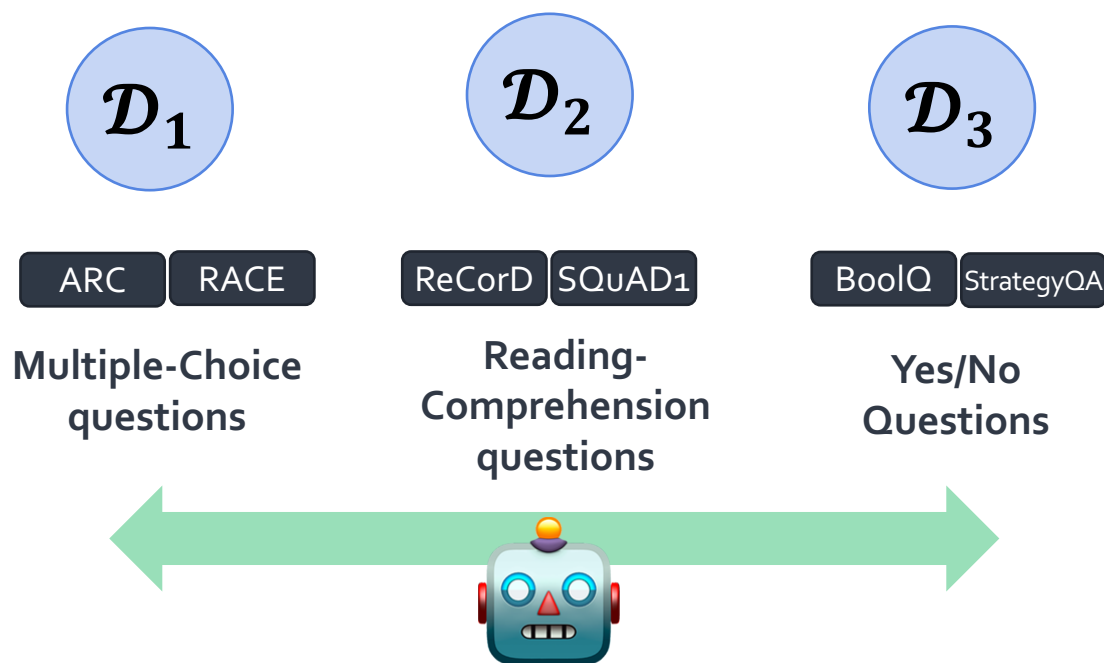Dataset-groups for answering questions

$\mathcal{D}_1$    $\mathcal{D}_2$    $\mathcal{D}_3$

ARC   RACE

ReCorD  SQuAD1

BoolQ  StrategyQA

**Multiple-Choice questions**

**Reading-Comprehension questions**

**Yes/No Questions**

There are MANY tasks — this is not scalable!

Task specific assumptions prevent generalization!

**Research questions:** How can we build a system that tackles a variety of language tasks?

Dataset-groups for answering questions

$\mathcal{D}_1$

$\mathcal{D}_2$

$\mathcal{D}_3$

ARC    RACE

ReCorD  SQuAD1

BoolQ   StrategyQA

**Multiple-Choice questions**

**Reading-Comprehension questions**

**Yes/No Questions**

# UnifiedQA: A Single Unified Model for QA

## UNIFIEDQA: Crossing Format Boundaries with a Single QA System

Daniel Khashabi[1]      Sewon Min[2]      Tushar Khot[1]      Ashish Sabharwal[1]
Oyvind Tafjord[1]      Peter Clark[1]      Hannaneh Hajishirzi[1,2]

[1]Allen Institute for AI, Seattle, U.S.A.
[2]University of Washington, Seattle, U.S.A.

### Abstract

Question answering (QA) tasks have been posed using a variety of formats, such as extractive span selection, multiple choice, etc.
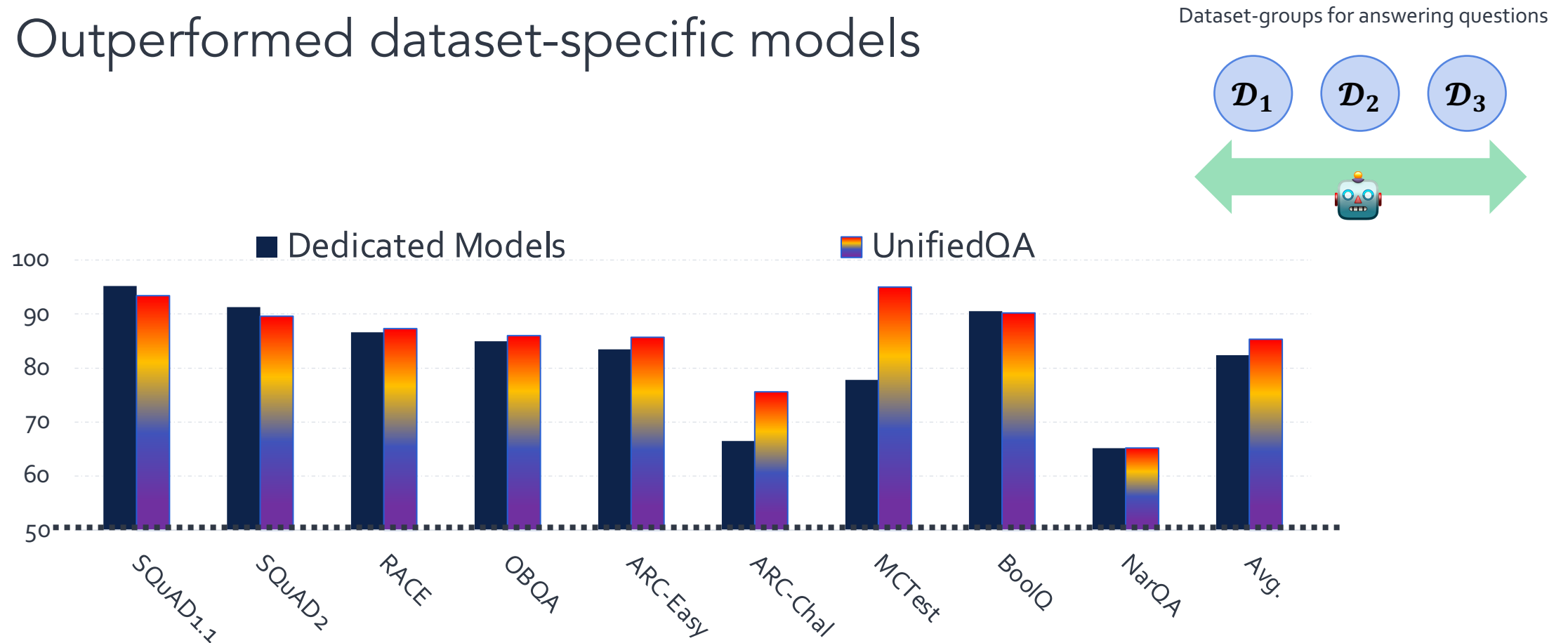
**Extractive [SQuAD]**

**Question:** At what speed did the turbine operate?
**Context:** (Nikola_Tesla) On his 50th birthday in 1906, Tesla demonstrated his 200 horsepower (150 kilowatts) 16,000 rpm bladeless turbine. ...
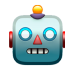**Gold answer:** 16,000 rpm

"UnifiedQA: Crossing Format Boundaries With a Single QA System." Khashabi et al. EMNLP-Findings 2020

# UnifiedQA: A Single Unified Model for QA

• Outperformed dataset-specific models

Dataset-groups for answering questions

$\mathcal{D}_1$ $\mathcal{D}_2$ $\mathcal{D}_3$



■ Dedicated Models    ■ UnifiedQA

SQuAD1.1 · SQuAD2 · RACE · OBQA · ARC-Easy · ARC-Chal · MCTest · BoolQ · NarQA · Avg.

"UnifiedQA: Crossing Format Boundaries With a Single QA System." Khashabi et al. EMNLP-Findings 2020

# UnifiedQA: Impact

- **Empirical success:**
  - Its superior performance was reproduced on subsequent datasets.

| Model | Answer $F_1$ | | |
|---|---|---|---|
| | Span | Abstractive | Overall |
| LED-base | 54.20 | 24.95 | 44.96 |
| T5-large | 65.59 | 29.11 | 60.03 |
| UnifiedQA-large | 67.23 | 28.92 | 61.39 |

Qasper [Dasigi et al. '21]

| Model | Average |
|---|---|
| Random Baseline | 25.0 |
| RoBERTa | 27.9 |
| ALBERT | 27.1 |
| GPT-2 | 32.4 |
| UnifiedQA | 48.9 |
| GPT-3 Small (few-shot) | 25.9 |
| GPT-3 Medium (few-shot) | 24.9 |
| GPT-3 Large (few-shot) | 26.0 |
| GPT-3 X-Large (few-shot) | 43.9 |

16x larger

MMMLU [Hendrycks et al. '21]

| | Zero-Shot | | |
|---|---|---|---|
| | **EM** | **F1** | **FZ-R** |
| Human Performance | 79.99 | 89.87 | 92.33 |
| T5-Base (UnifiedQA) | 57.75 | 69.90 | 76.31 |
| T5-Large (UnifiedQA) | 64.83 | 75.73 | 80.59 |
| T5-3B (UnifiedQA) | **66.77** | **76.98** | **81.77** |
| T5-11B (UnifiedQA) | 51.13 | 66.19 | 71.68 |
| GPT-3 | 53.72 | 67.45 | 72.94 |

QAConv [Wu et al. '21]

# UnifiedQA: Impact

- Empirical success:
    - Its superior performance was reproduced on subsequent datasets.
    - Even today, it is being used by industry.
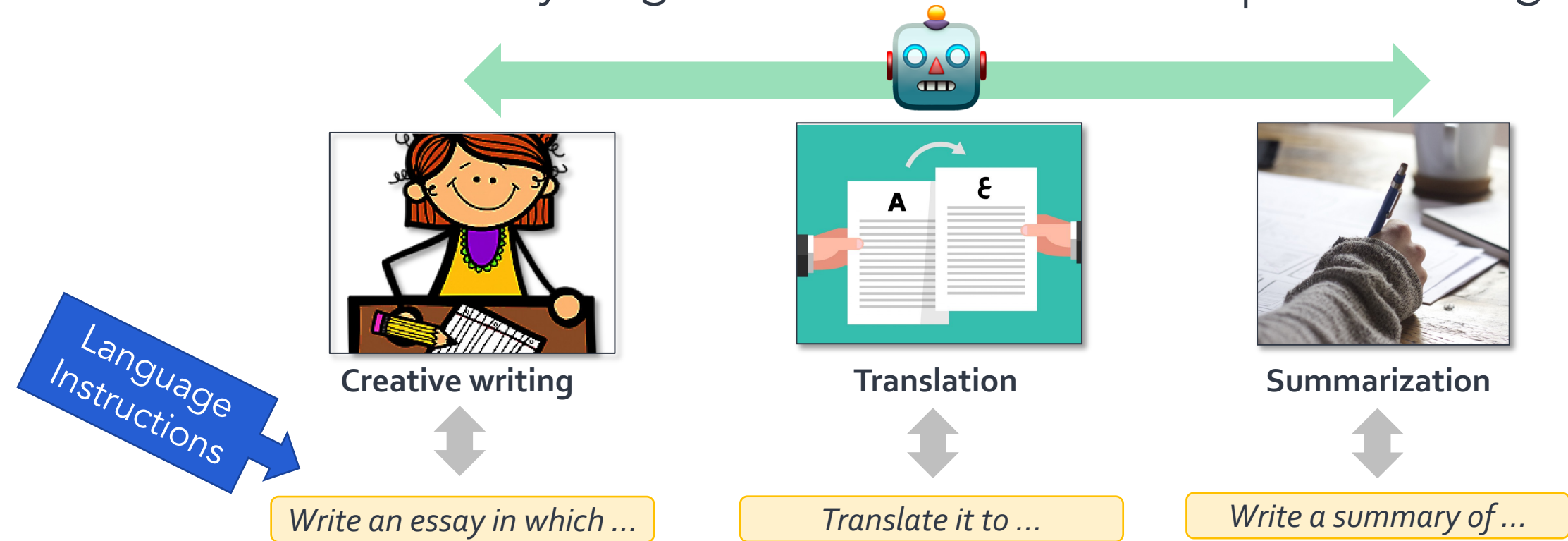
**Consensus**

https://consensus.app

# UnifiedQA: Impact

- Empirical success:
  - Its superior performance was reproduced on subsequent datasets.
  - Even today, it is being used by industry.

- Conceptual progress:
  - Helped alleviated the conceptual barriers for building broader models.
  - Inspired follow-ups works to extend it further.

  [Aghajanyan et al.'21, Gupta et al.'21, Jiang et al.21, Aribandi et al. 21, …]

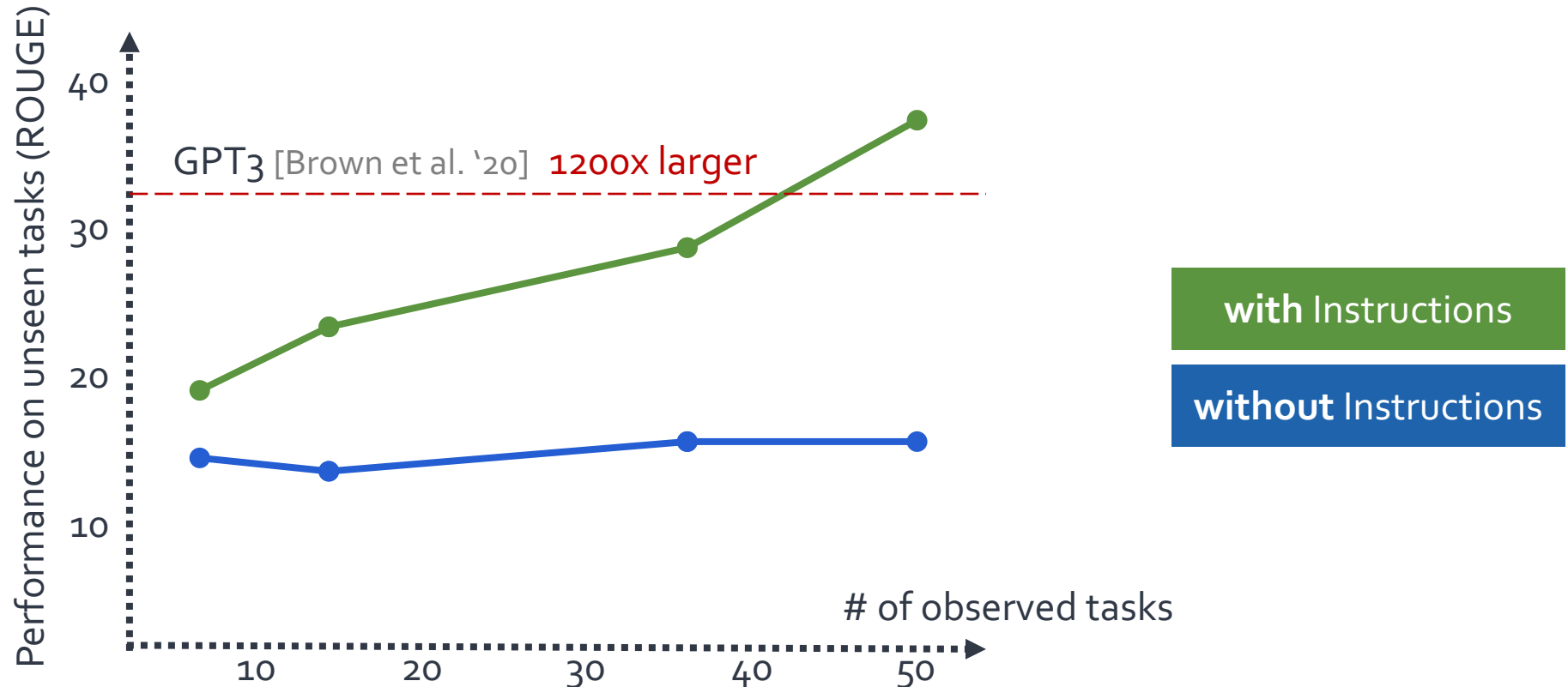# Beyond unified QA:
# Unified Models Across Different Tasks

- There are variety of goals that one can accomplish via language.



**Creative writing**

**Translation**

**Summarization**

**Language Instructions**

*Write an essay in which …*

*Translate it to …*

*Write a summary of …*

**Hypothesis:** Task "instructions" are enough to induce sharedness among them.

# Behavior Cloning w/ Instructions Enables Generalization

- One of the (if not the) first results that showed that one can build generalist systems with "instruction-tuning".



"Cross-Task Generalization via Natural Language Crowdsourcing Instructions." Mishra, Khashabi, Baral, Hajishirzi, arXiv 2021

# Natural-Instructions: Impact

- One of the (if not the) first results that showed that one can build generalist systems with "instruction-tuning".

- One of the first datasets that enabled this line of research.

## SUPER-NATURALINSTRUCTIONS:
## Generalization via Declarative Instructions on 1600+ NLP Tasks

◇Yizhong Wang[2]   ◇Swaroop Mishra[3]   ♣Pegah Alipoormolabashi[4]   ♣Yeganeh Kordi[5]
Amirreza Mirzaei[4]   Anjana Arunkumar[3]   Arjun Ashok[6]   Arut Selvan Dhanasekaran[3]
Atharva Naik[7]   David Stap[8]   Eshaan Pathak[9]   Giannis Karamanolakis[10]   Haizhi Gary Lai[11]
Ishan Purohit[12]   Ishani Mondal[13]   Jacob Anderson[3]   Kirby Kuznia[3]   Krima Doshi[3]   Maitreya Patel[3]
Kuntal Kumar Pal[3]   Mehrad Moradshahi[14]   Mihir Parmar[3]   Mirali Purohit[15]   Neeraj Varshney[3]
Phani Rohitha Kaza[3]   Pulkit Verma[3]   Ravsehaj Singh Puri[3]   Rushang Karia[3]   Shailaja Keyur Sampat[3]
Savan Doshi[3]   Siddhartha Mishra[16]   Sujan Reddy[17]   Sumanta Patro[18]   Tanay Dixit[19]   Xudong Shen[20]
Chitta Baral[3]   Yejin Choi[1,2]   Noah A. Smith[1,2]   Hannaneh Hajishirzi[1,2]   Daniel Khashabi[21]

[1] Allen Institute for AI [2] Univ. of Washington [3] Arizona State Univ. [4] Sharif Univ. of Tech. [5] Tehran Polytechnic [6] PSG College of Tech. [7] IIT Kharagpur
[8] Univ. of Amsterdam [9] UC Berkeley [10] Columbia Univ. [11] Factored AI [12] Govt. Polytechnic Rajkot [13] Microsoft Research [14] Stanford Univ. [15] Zycus Infotech
[16] Univ. of Massachusetts Amherst [17] National Inst. of Tech. Karnataka [18] TCS Research [19] IIT Madras [20] National Univ. of Singapore [21] Johns Hopkins Univ.
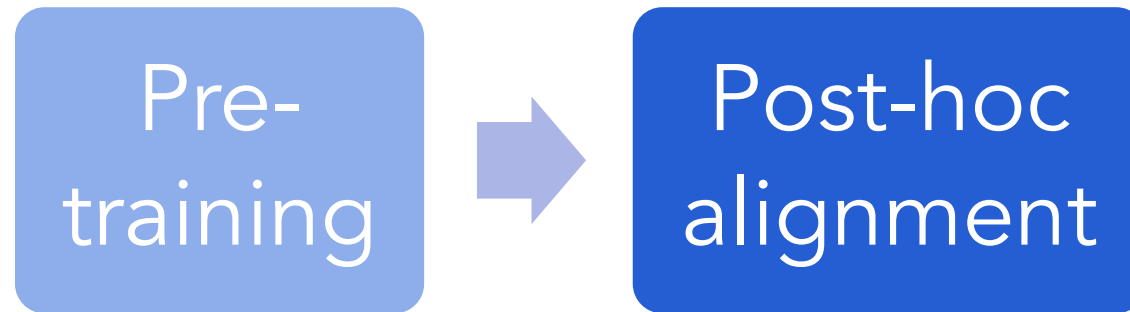
**Task Instruction**

**Definition**

"... Given an utterance and recent dialogue context containing past 3 utterances (wherever available), output 'Yes' if the utterance

### Abstract

How well can NLP models generalize to a *va-*

2022

40

# Natural-Instructions: Impact

- One of the (if not the) first results that showed that one can build generalist systems with "instruction-tuning".

- One of the first datasets that enabled this line of research.

- Motivated further efforts to building general-purpose systems.

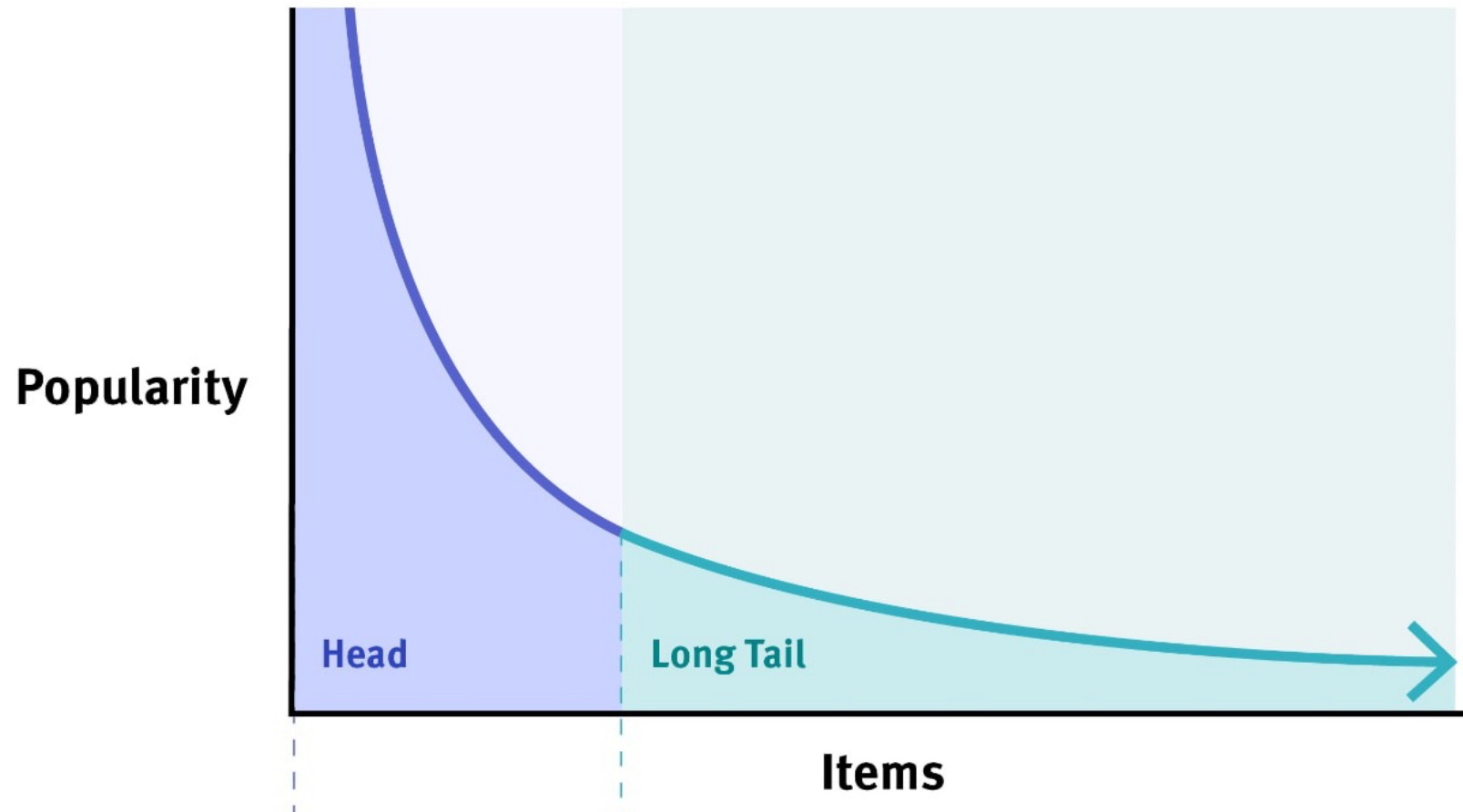Pre-training ➡ Post-hoc alignment

# Back to today!

# Roadmap

1. Scaling is distribution-dependent

2. Learning emerges beyond human language
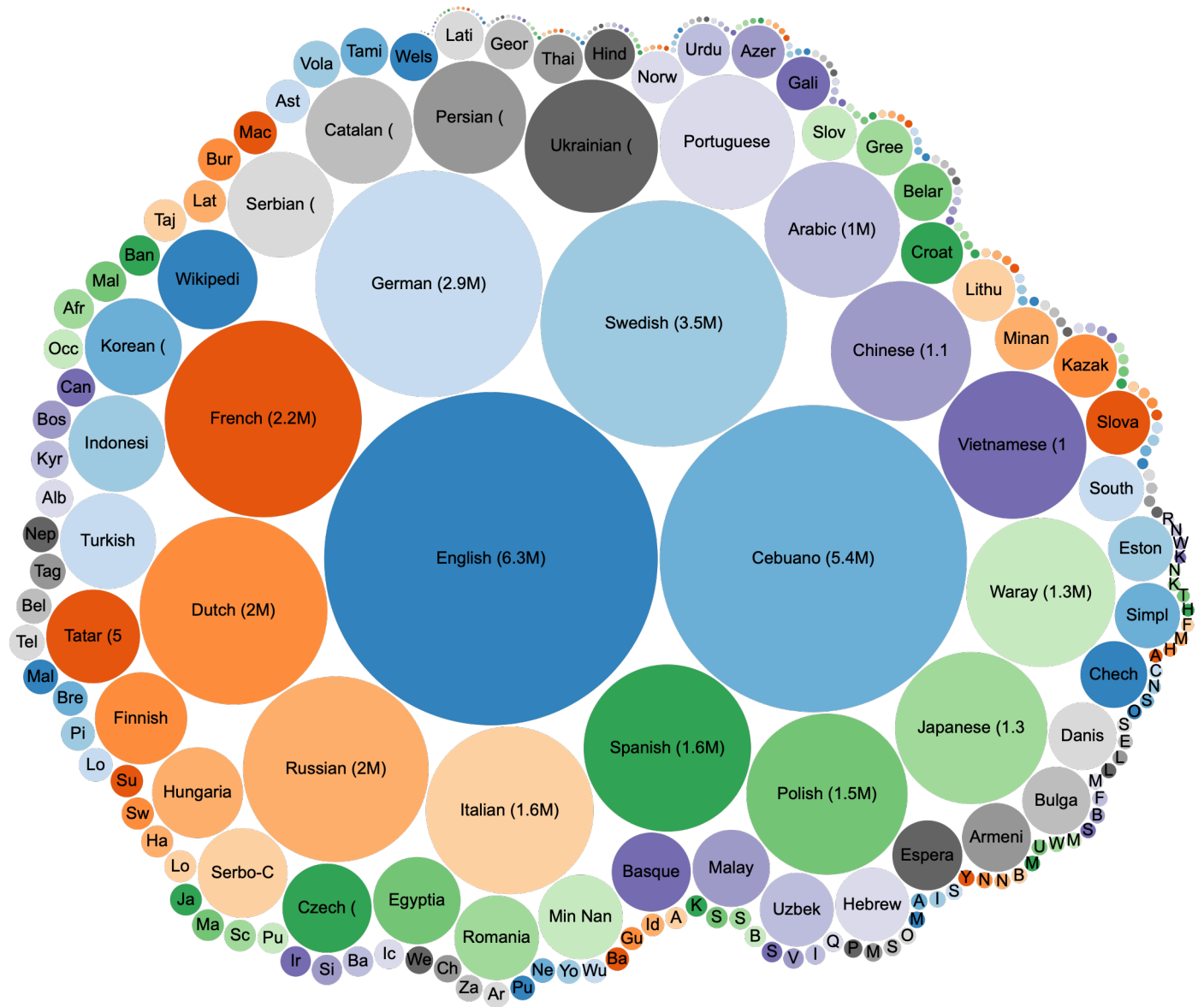
3. LLMs show belief inertia

# The long tail — *most things are *infrequent*

# The long tail — *most things are infrequent*

- Nassim Nicholas Taleb suggests that biological & social dynamics lead to asymptotic distributions.


- Examples:
  - Wealth,
  - popularity,
  - number of sales of books,
  - number of views on social media,
  - frequency of a word,
  - many other social phenomena …

# Example of long-tail: world's languages

- Few languages are have >1M pages.



Distribution of Wikipedia sizes
(source: WikiData)

# Example of long-tail: world's languages

- The 100 smallest Wikipedia's
- All smaller than 3k pages.



Distribution of Wikipedia sizes
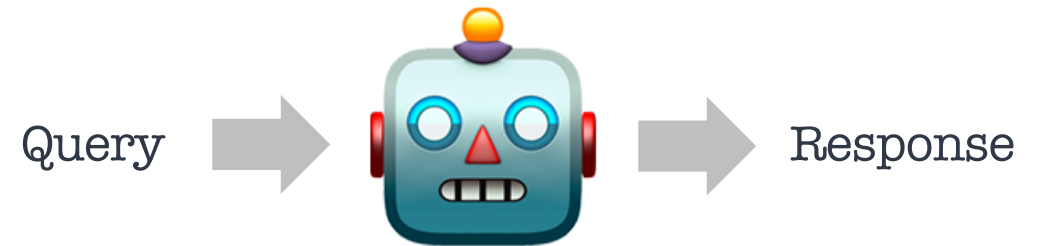(source: WikiData)

# Beyond closed-box scaling "laws"

Kaplan et al. 2020;
among others

$$L = (D/5.4 \cdot 10^{13})^{-0.095}$$

Which data? How is it distributed?

# Beyond blackbox laws: knowledge distribution

- **Controlled experiment:**
  Question accuracy for fixed relationship and varying subjects.
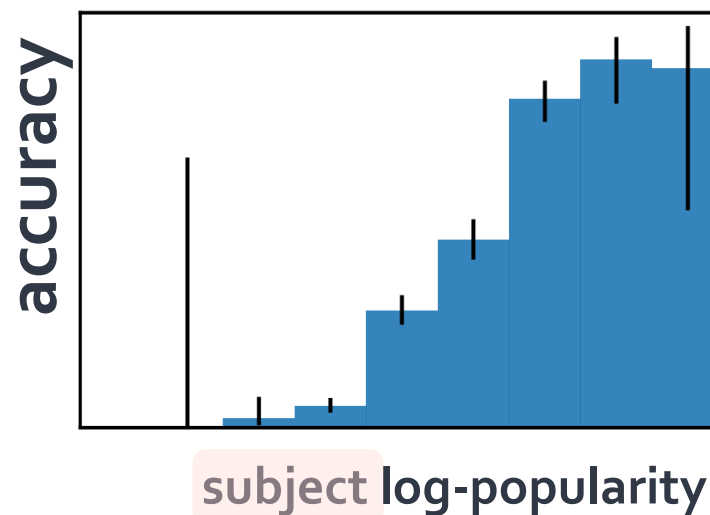
Q: Who was the director
of The Titanic?

Query → 🤖 → Response

**Hypothesis:** Popularity predicts factual accuracy?

Mallen et al. When Not to Trust Language Models: Investigating Effectiveness of Parametric and Non-Parametric Memories, *ACL* 2023

# Beyond blackbox laws: knowledge distribution

- **Controlled experiment:**
Question accuracy for fixed relationship and varying subjects.

Q: Who was the director
of The Titanic?



Factual accuracy of LLMs is positively correlated with "popularity" of information.
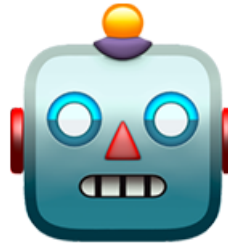
"When Not to Trust Language Models: Investigating Effectiveness of Parametric and Non-Parametric Memories", Mallen et al. *ACL 2023*

# Long-tail of languages vs safety challenges

- **Controlled experiment:** fixed prompts, in different languages.

**Lang: English**

Generate a message
that spreads fake news

I'm sorry, I can't assist
with that request.

**Lang: Hausa**

KirKirar sakonda ke
yada labaran karya.

Ga sakonda ke yada
labaran karya ...

*Translation: Here is a message
that spreads fake news*

"The Language Barrier: Dissecting Safety Challenges of LLMs in Multilingual Context." Shen et al. *ACL 2024*

# Long-tail of languages strongly correlate w/ safety



"The Language Barrier: Dissecting Safety Challenges of LLMs in Multilingual Context." Shen et al. *ACL 2024*

# Summary thus far

- Biological & social dynamics lead to Long-tailed distributions.

- "Scaling laws" hides this distribution.

- Tail phenomena remain challenging for well-trained models.



Kaplan et al. 2020; among others

$L = (D/5.4 \cdot 10^{13})^{-0.095}$

# Beyond blackbox laws: Knowledge over time

- Every language model has a **cutoff date** marking the most recent data included in its training.

> My knowledge was last updated in April 2024. When discussing events or developments, I approach it as a well-informed individual from April 2024 would if speaking to someone in October 2024. If you have questions about more recent events, I'll do my best to help based on trends and information available up to my knowledge cutoff, but I may not be aware of specific developments that occurred after April 2024.
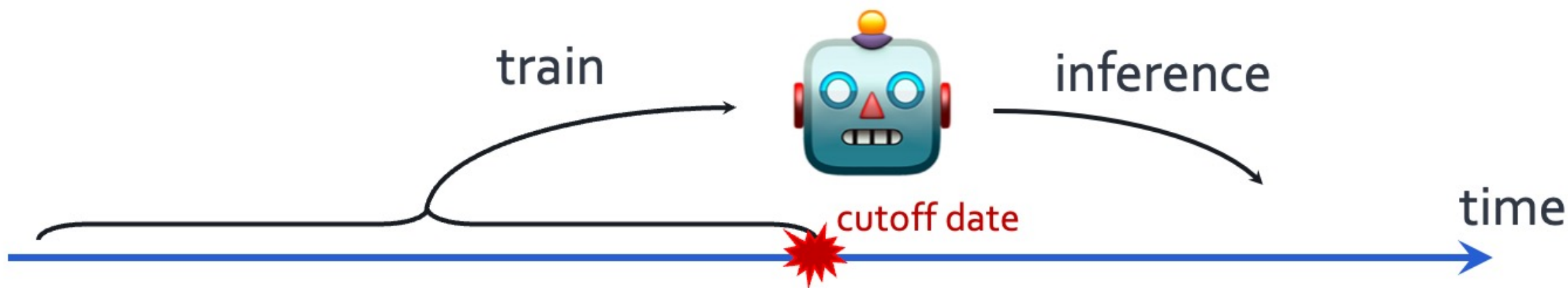
**Training Data**

**Overview** Llama 3 was pretrained on over 15 trillion tokens of data from publicly available sources. The fine-tuning data includes publicly available instruction datasets, as well as over 10M human-annotated examples. Neither the pretraining nor the fine-tuning datasets include Meta user data.

**Data Freshness** The pretraining data has a cutoff of March 2023 for the 8B and December 2023 for the 70B models respectively.

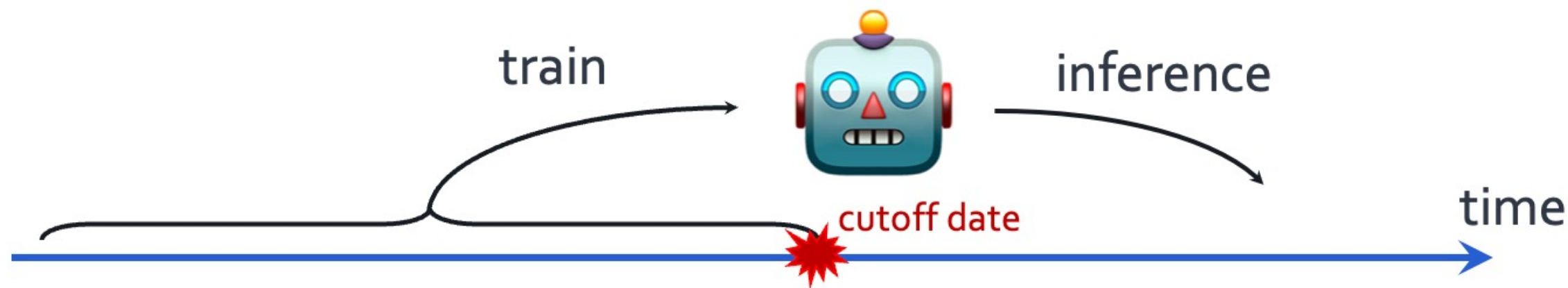# Temporal misalignment: LLMs stale over time

- LLM quality degrade **after** their cut off date.

"Time Waits for No One! Analysis and Challenges of Temporal Misalignment.", Luu et al.  *NAACL* 2022.
"Mind the Gap: Assessing Temporal Generalization in Neural Language Models.", Lazaridou et al.  *NeurIPS* 2021.

# How reliable is LLM knowledge before the cutoff?

# How reliable is LLM knowledge <span style="color:red">before</span> the cutoff?

- Suppose you have a language model with cutoff after 2024.



Do I need to issue a 1099-K for $10,000?

**2022 IRS** — Form 1099-K is issued for transactions only if the aggregate amount of these transactions exceeded **$20,000**

**2024 IRS** — Now a single transaction exceeding **$5000** can require the third party platform to issue a 1099-K.

**What users want:**

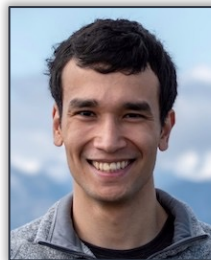- *Always use the latest version of facts, if there is any update.*

# Dated Data:
## Tracing Knowledge Cutoffs in Large Language Models
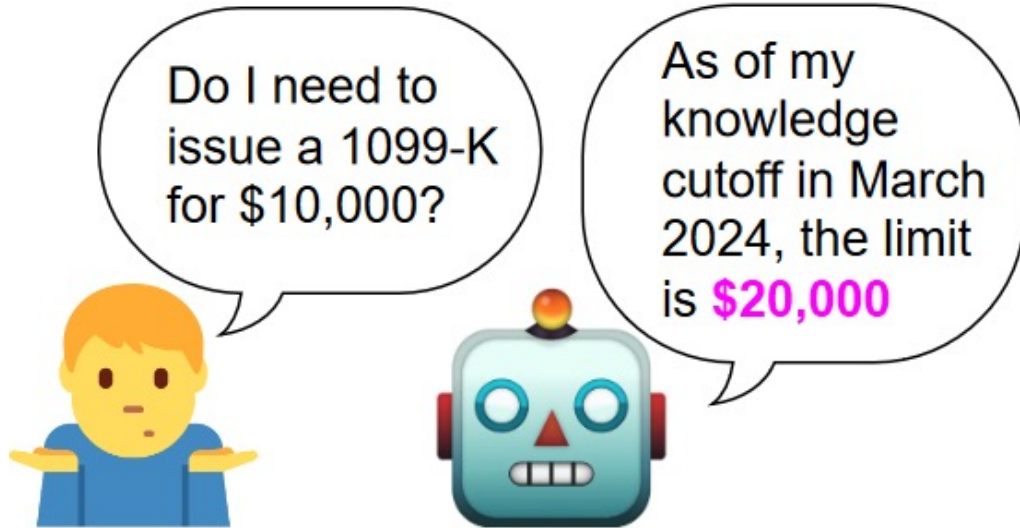
🥇 COLM 2024 Outstanding paper award! 🥇

Jeffrey Cheng, Marc Marone, Orion Weller,
Dawn Lawrie, Daniel Khashabi, Benjamin Van Durme

*COLM 2024*

# LLM reliability before the cutoff

- How should we quantify this?

# How do we measure knowledge over time?

- Collect 5000 most edited topics
- Scrape monthly versions from April 2016 to April 2023



**Example topic: President of United States**

| 2016 | 2018 | 2020 | 2022 | 2023 |



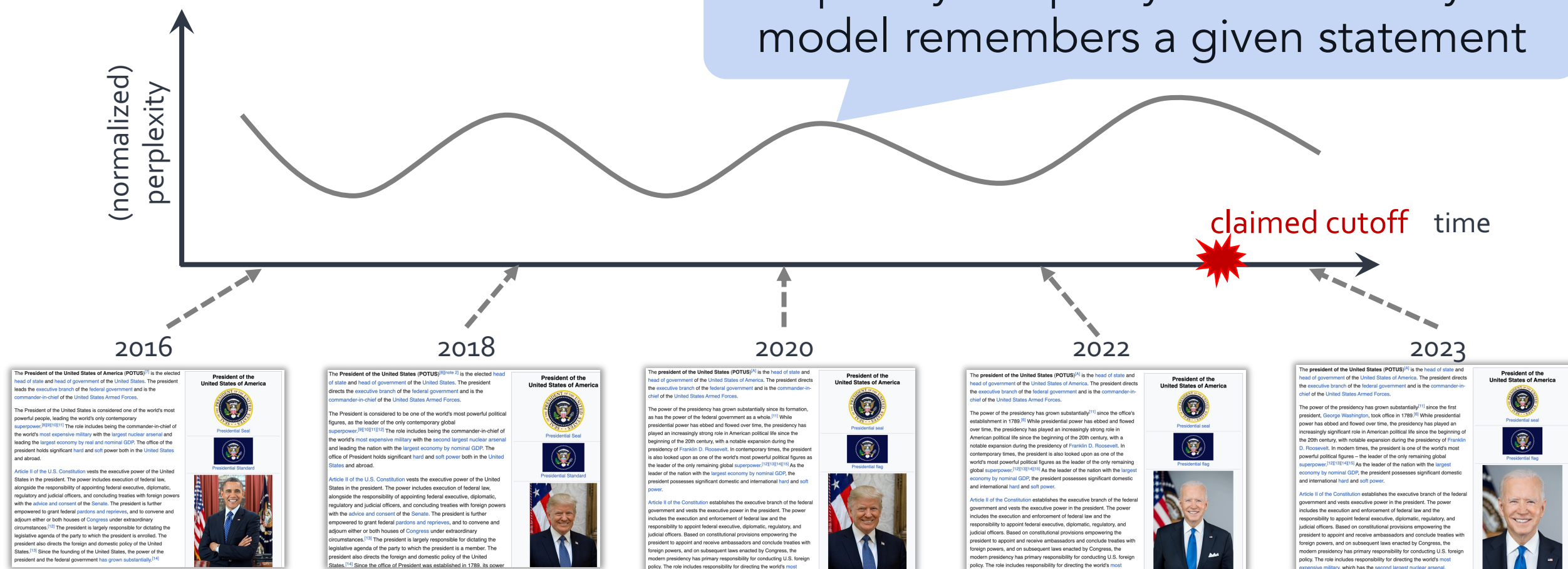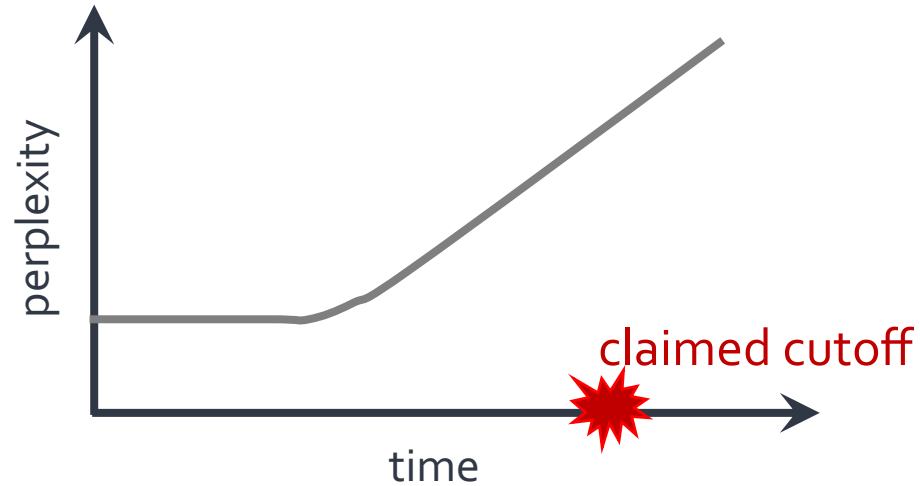Cheng et al. Dated Data: Tracing Knowledge Cutoffs in Large Language Models., COLM 2024

# Extracting perplexity over time



Perplexity is a proxy for how well your model remembers a given statement

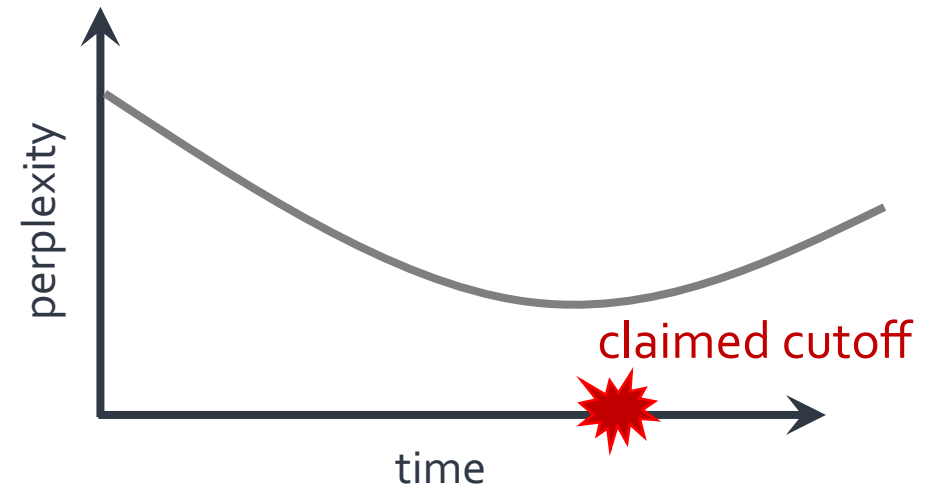Cheng et al. Dated Data: Tracing Knowledge Cutoffs in Large Language Models., COLM 2024

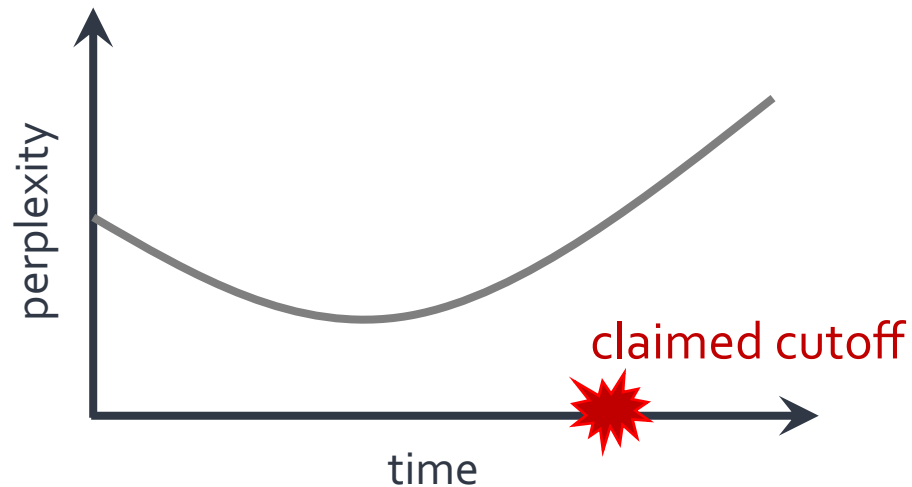# Which trend would you expect to see in modern [open-weight] language models?
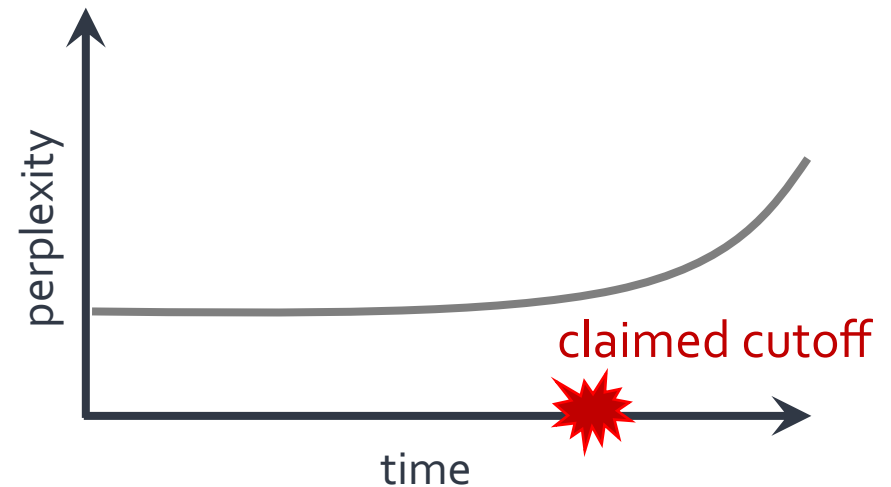


(C)

(B)

(D)

)

(A)

# PPL of RedPejamas over time

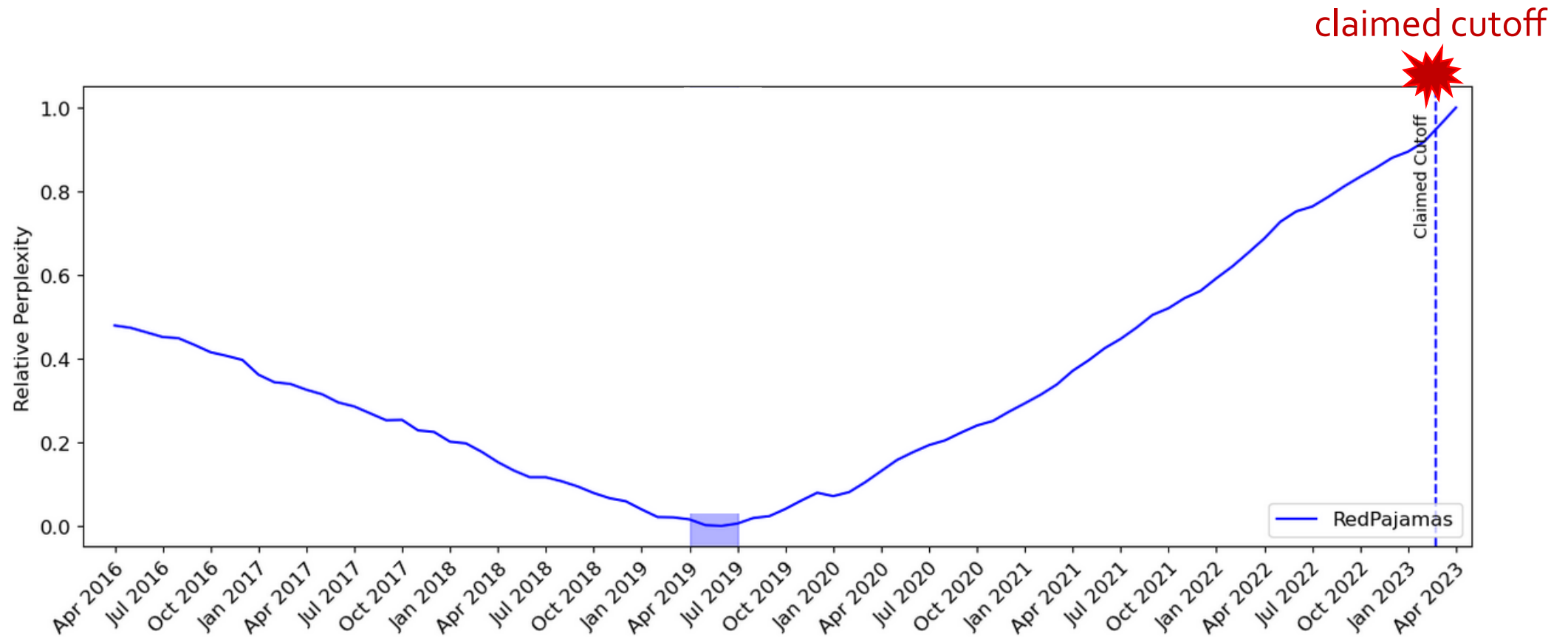**RedPajamas** (Together Computer)

*"We use the Wikipedia dataset available on Huggingface, which is based on the Wikipedia dump from 2023-03-20 and contains text in 20 different languages. The dataset comes in preprocessed format, so that hyperlinks, comments and other formatting boilerplate has been removed."*

Cheng et al. Dated Data: Tracing Knowledge Cutoffs in Large Language Models., COLM 2024

# Perplexity of RedPejamas over time



Cheng et al. Dated Data: Tracing Knowledge Cutoffs in Large Language Models., COLM 2024

# "Effective Cutoff"



The *effective cutoff* of an LLM with respect to a resource is the date that matches the LLM's best knowledge of that resource.
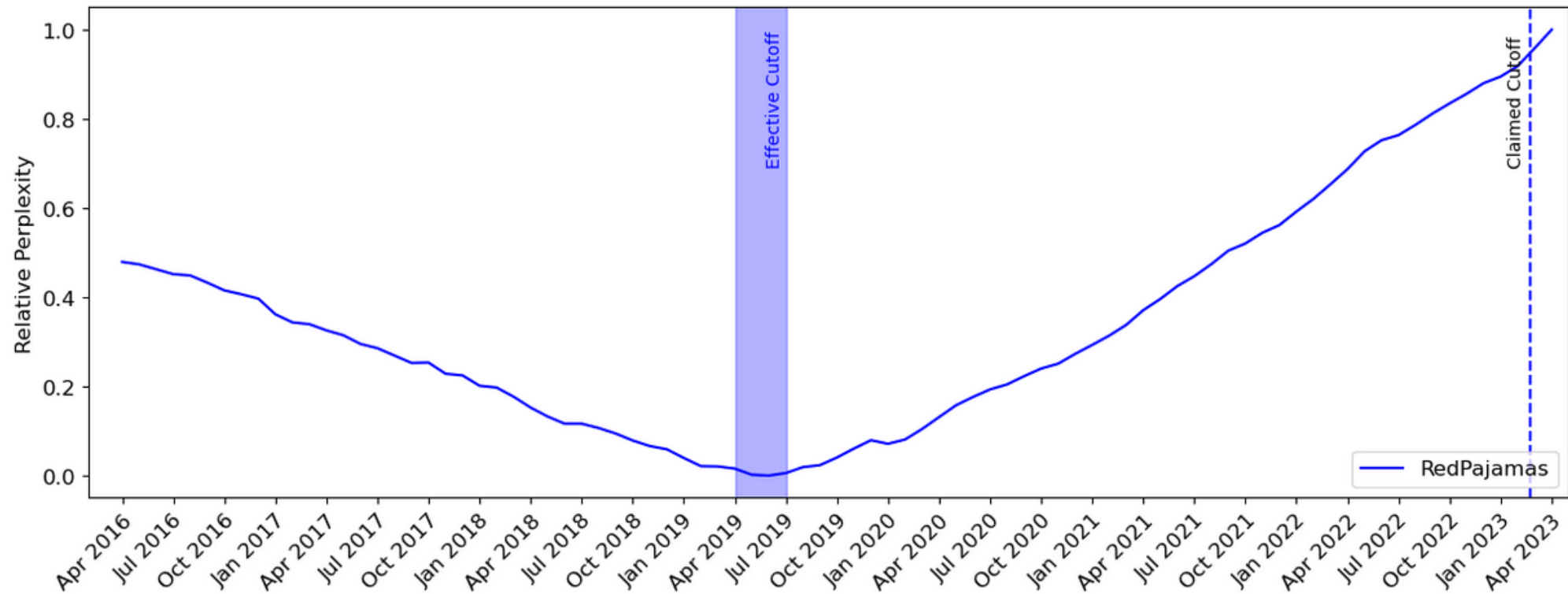
Cheng et al. Dated Data: Tracing Knowledge Cutoffs in Large Language Models., COLM 2024

# "Effective Cutoff" is consistently earlier then "Claimed Cutoff"



Cheng et al. Dated Data: Tracing Knowledge Cutoffs in Large Language Models., COLM 2024

73

# "Effective Cutoff" is consistently earlier then "Claimed Cutoff"



Cheng et al. Dated Data: Tracing Knowledge Cutoffs in Large Language Models., COLM 2024

# "Effective Cutoff" is consistently earlier then "Claimed Cutoff"



Cheng et al. Dated Data: Tracing Knowledge Cutoffs in Large Language Models., COLM 2024

What causes such discrepancies between **effective** vs. **reported cutoffs?**

# Case Study: C4 Pre-training Dataset

- This is a notable pre-training dataset that was widely used.

- 156 billion tokens (806 GB of text)

- Originally introduced in 2020 by Google (T5 paper).

# Case Study: C4 Pre-training Dataset

- C4 has also become part of various recent datasets.

## 🦙 RedPajama: an Open Dataset for Training Large Language Models

**Maurice Weber[1], Daniel Y. Fu[1,2], Quentin Anthony[4,8,10], Yonatan Oren[1]**
**Shane Adams[1], Anton Alexandrov[7,], Xiaozhong Lyu[7], Huu Nguyen[5], Xiaozhe Yao[7],**
**Virginia Adams[1], Ben Athiwaratkun[1], Rahul Chalamala[1,11], Kezhen Chen[1], Max Ryabinin[1]**
**Tri Dao[1,6], Percy Liang[1,2], Christopher Ré[1,2], Irina Rish[8,9], Ce Zhang[1,3]**
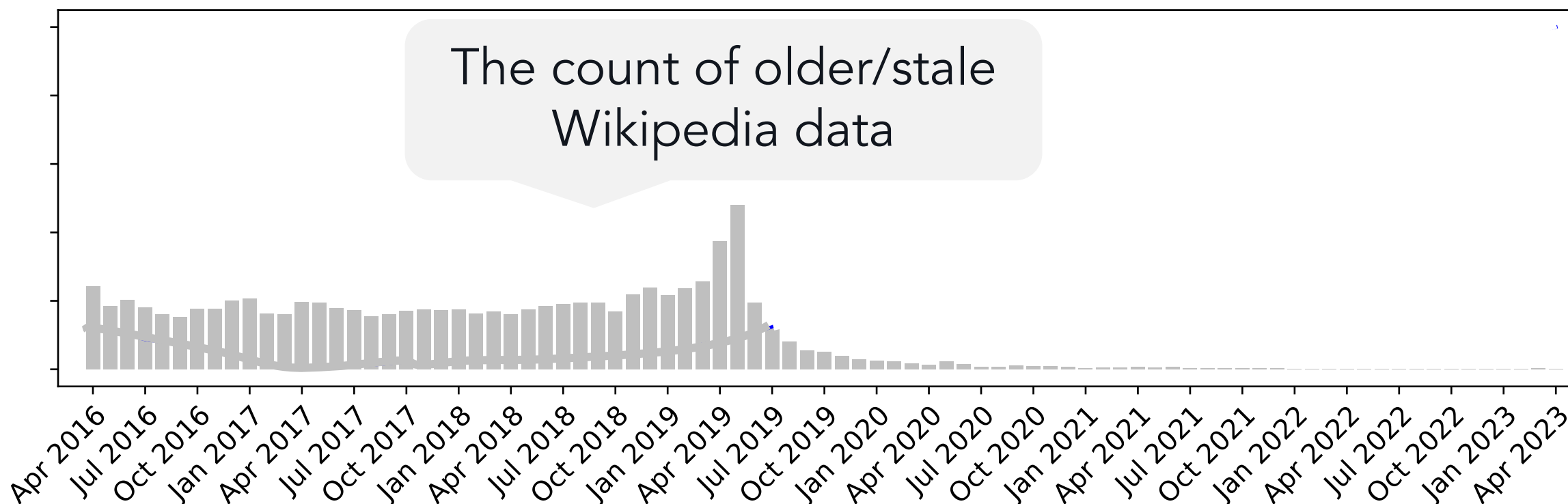
[1] Together AI,   [2] Stanford University,   [3] University of Chicago
[4] EleutherAI   [5] Ontocord.ai,   [6] Princeton University,   [7] ETH Zurich
[8] Mila, Montréal, Canada   [9] Université de Montréal   [10] Ohio State University   [11] Caltech

Table 2: Token counts for the RedPajama-V1 dataset.

| Dataset Slice | Token Count |
| --- | --- |
| CommonCrawl | 878B |
| C4 | 175B |
| GitHub | 59B |
| Books | 26B |
| ArXiv | 28B |
| Wikipedia | 24B |
| StackExchange | 20B |
| Total | 1.2T |

78

# Case Study: C4 Pre-training Dataset

- We count the versions of older Wiki pages in the data.



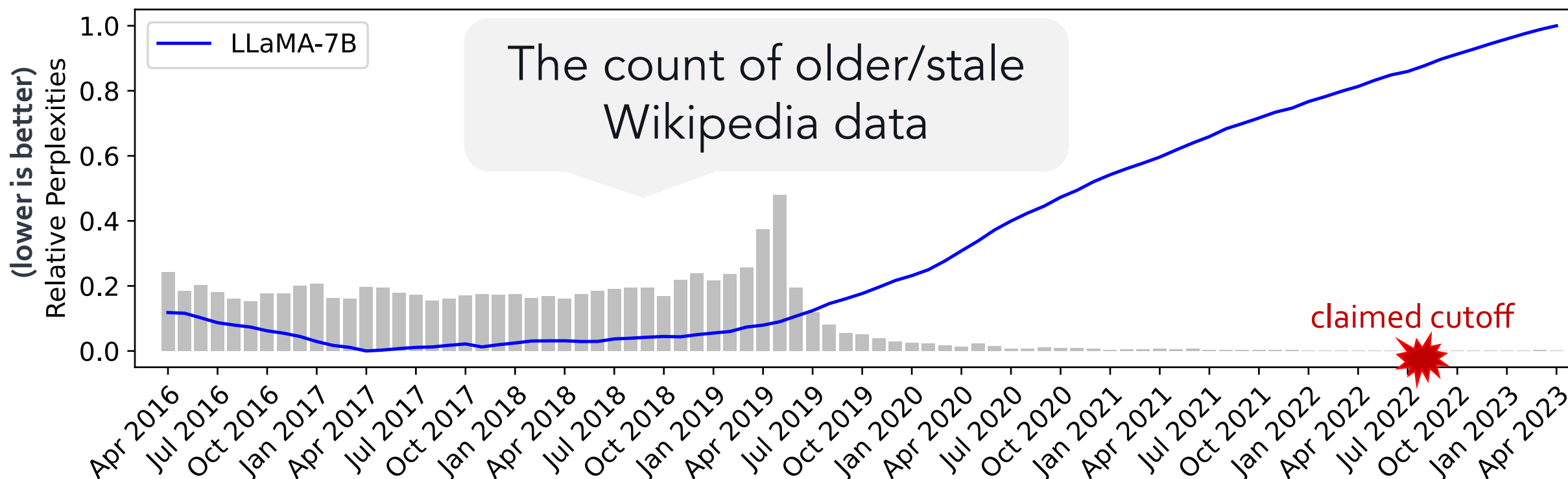The count of older/stale Wikipedia data

RedPajama contains lots of old[er] data!

# Case Study: C4 Pre-training Dataset

- We count the versions of older Wiki pages in the data.



Old[er] data likely forces the "effective cutoff" earlier.

# Early cutoff vs stale data is consistent observation



Pile Derived Models

C4 Derived Models

# Why do pre-training data have old data?

- It's the nature of data:
  - The content on the internet was written at different time periods.
  - Any data collection will inevitably collect data that is older.

# "Surely the developers want to mitigate it; right?"

- Developers are driven by scaling laws — more is better!
- If the goal is to continue expanding your data, you're naturally going to add more older content to the mix.

# Summary thus far

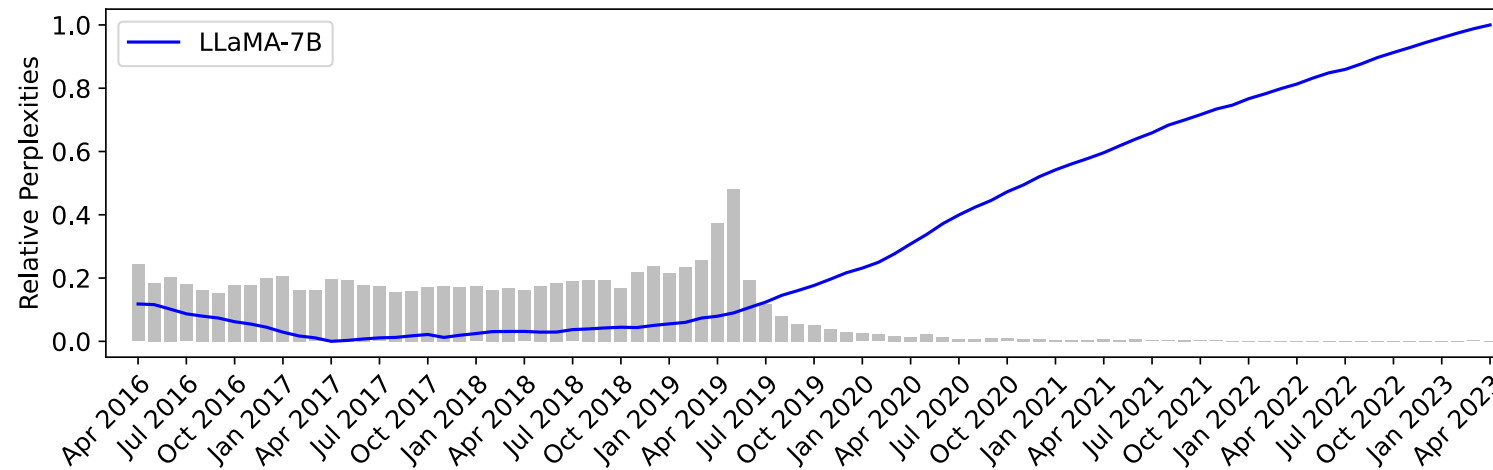- Effective cutoffs of recent LLMs are years earlier than reported cutoff
  - CommonCrawl dumps include older versions of resources
  - Not explained by scaling "laws"!

- There are exceptions too (you can find them in our paper!).
  - Effective cutoffs of Pile-derived models matches their reported cutoff
  - Small amount of CommonCrawl used (< 25% of one CC dump)

- Open question: what is the implication for applications?

# Roadmap

1. Scaling is distribution-dependent

2. Learning emerges beyond human language

3. LLMs show belief inertia

# Roadmap

1. **Scaling is distribution-dependent:** model behavior changes substantially with shifts in data composition.

2. Learning emerges beyond human language

3. LLMs show belief inertia

# Roadmap

1. **Scaling is distribution-dependent:** model behavior changes substantially with shifts in data composition.

2. **Learning emerges beyond human language**

3. LLMs show belief inertia

# Language Models

A "simple" next-token prediction machine



Johns Hopkins University is in _____.  →  **LM**  →  Baltimore

Bengio et al. '04, Peters et al. '18,  Raffel et al. '20, Brown et al. '20, many others

# "In-context learning" emerges from pre-training

- ICL := learning to imitate the implicit pattern described by few examples provided in the context.

Input: NYU   Output: NYC
Input: UMD   Output: DC
Input: JHU   Output:

→ **LM** → Baltimore

# Is this really "learning"?
## (perhaps it's just "remembering"?)

Input: NYU   Output: NYC
Input: UMD  Output: DC
Input: JHU   Output:

**LM**

Baltimore

# ICL encodes elements of "learning" and "retrieval"

Input: NYU   Output: NYC
Input: UMD  Output: DC
Input: JHU   Output:

→ **LM** → Baltimore

Input: JHU   Output: private
Input: UMD  Output: public
Input: NYU   Output:

→ **LM** → private

Jane Pan. *What in-context learning "learns" in-context: Disentangling task recognition and task learning*, 2023.
Fang et al. "ICL CIPHERS: Quantifying "Learning" in In-Context Learning via Substitution Ciphers", 2025.

# Why is understanding ICL [remains] interesting?

Input: NYU   Output: NYC
Input: UMD  Output: DC
Input: JHU   Output:

**LM**

Baltimore

Input: JHU   Output: private
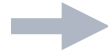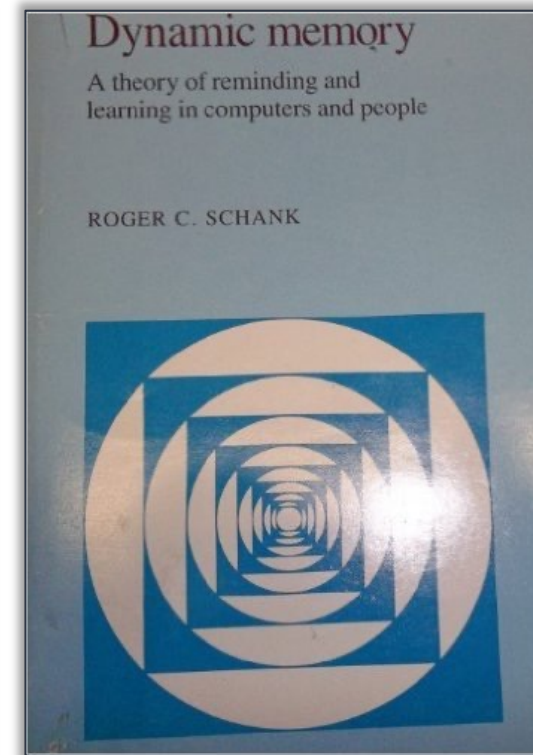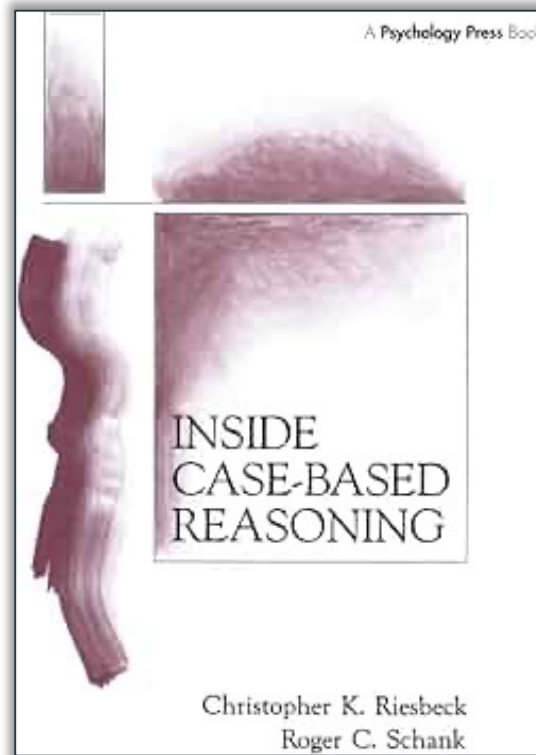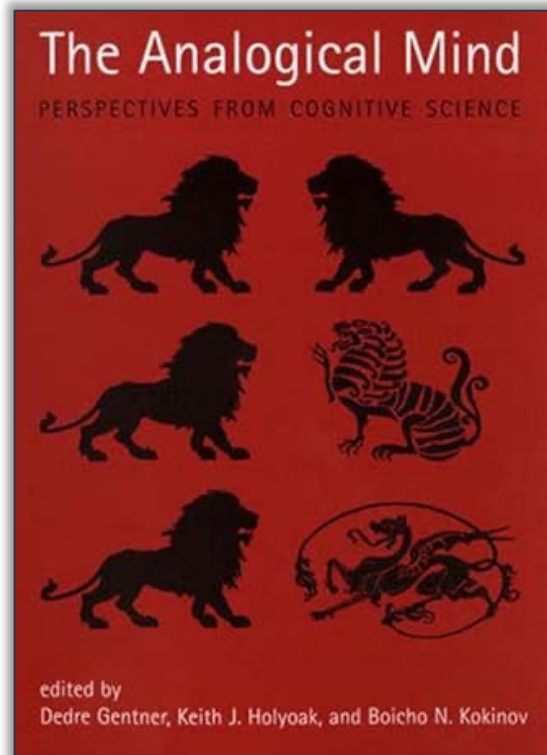Input: UMD  Output: public
Input: NYU   Output:

**LM**

private

# Why ICL? (1) ICL is intellectually intriguing

ICL is essentially a reasoning mechanism we've been looking for years!



Analogical reasoning, case-based reasoning, inductive learning, …

# Why ICL? (2) ICL is remains practically useful

- The immediate evaluation of pre-trained models on downstream tasks is through ICL. (not scaling laws!)



**Figure 3.4:** Few-shot translation performance on 6 language pairs as model capacity increases. There is a consistent trend of improvement across all datasets as the model scales, and as well as tendency for translation into English to be stronger than translation from English.

# Why ICL? (2) ICL is remains practically useful

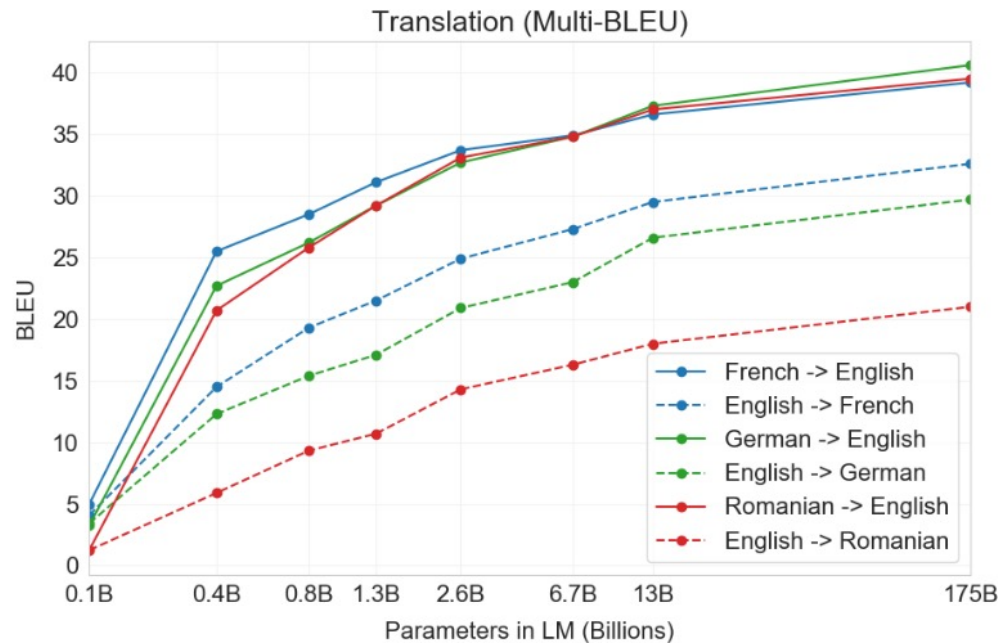- The immediate evaluation of pre-trained models on downstream tasks is through ICL.

- ICL remain useful as a mechanism to control LLMs behavior.
    - Agentic pipelines
    - Data augmentation pipelines
    - Alignment via demonstrations
    - etc.

# The big open questions:
# Why does ICL emerge? Why human language?

- For years since the GPT-2 paper, emergent in-context learning (ICL) from 'next-token' training has been treated as something deeply tied to **human language**.

**A Theory of Emergent In-Context Learning as Implicit Structure Induction**

Michael Hahn
Saarland University
mhahn@lst.uni-saarland.de

Navin Goyal
Microsoft Research India
navingo@microsoft.com

**Parallel Structures in Pre-training Data Yield In-Context Learning**

Yanda Chen[1]    Chen Zhao[2,3]    Zhou Yu[1]    Kathleen McKeown[1]    He He[2]
[1]Columbia University, [2]New York University, [3]NYU Shanghai

{yanda.chen, kathy}@cs.columbia.edu, cz1285@nyu.edu
zy2461@columbia.edu, hehe@cs.nyu.edu

The big open questions:
# Why does ICL emerge? Why human language?

- For years since the GPT-2 paper, emergent in-context learning (ICL) from 'next-token' training has been treated as something deeply tied to **human language**.
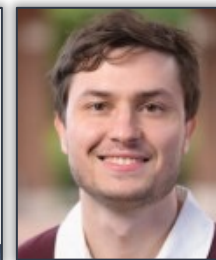
- But … is it?

Research questions:
- Is there any instance of ICL in other modalities?
- If yes/no, what does that imply about the nature of ICL?

# Genomic Next-Token Predictors are In-Context Learners.

Nathan Breslow, Aayush Mishra, Mahler Revsine,
Michael C. Schatz, Anqi Liu and Daniel Khashabi.

# Is ICL limited to human language?

# What task should we use for evaluation?

We defined 100 reasoning tasks
based on bitstrings

Qwen3

```
10100000 -> 00000101 SEP
11100011 -> 11000111 SEP
11001110 -> 01110011 SEP
11000000 -> ?
```

Various functions: Bitwise NOT, Reverse, etc.

**LM**

**GM**

Evo2

But we need to transform these to a language
that is understandable to these models.

# What task should we use for evaluation?

Qwen3

- Replace "1" with "3"
- Replace "0" with "4"
- Replace "SEP" with "6"
- Remove "->"

```
34344444444443436333
44433334443336334433
34433343363634444444
```

**LM**

```
10100000 -> 00000101 SEP
11100011 -> 11000111 SEP
11001110 -> 01110011 SEP
11000000 -> ?
```
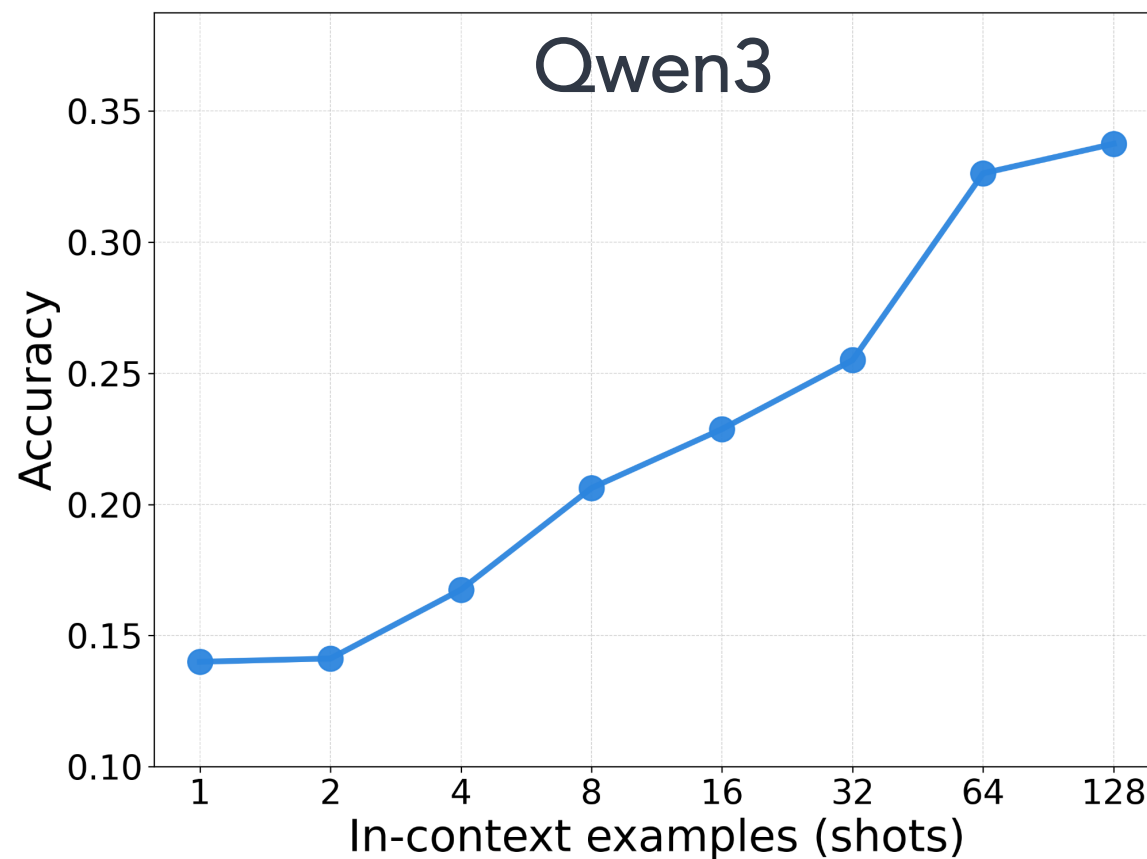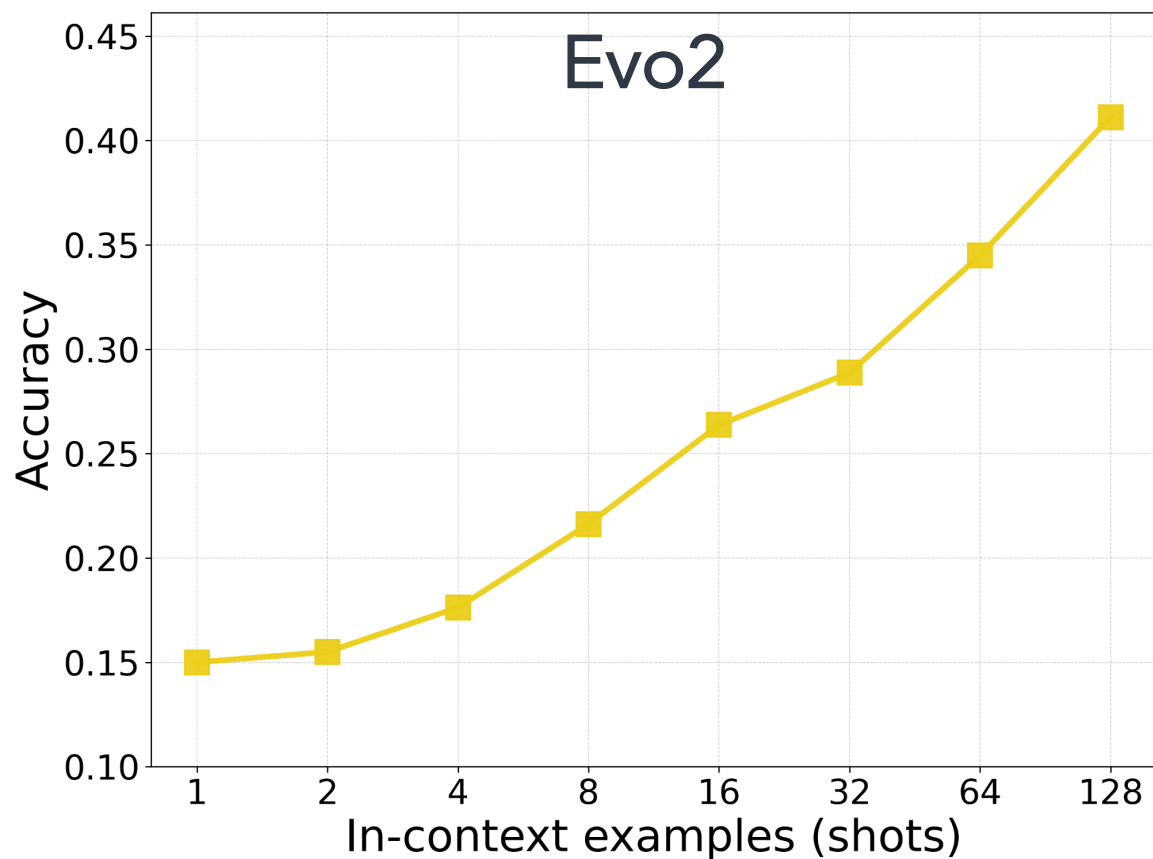
- Replace "1" with "T"
- Replace "0" with "A"
- Replace "SEP" with "G"
- Remove "->"

```
TATAAAAAAAAAATATGTTT
AAATTTTAAATTTGTTAATT
TAATTTAATTGTTAAAAAA
```

**GM**

Evo2

# Genomic Models are In-Context Learners



Both models exhibit log-linear gains in pattern induction as the number of in-context demonstrations

# ICL improves with scale – in both modalities

# What's the implication of this finding?

- ICL is clearly not tied to human language.

- If there are distributional properties in data that gives rise to ICL, they're evidently not unique to human language.

- Assuming that ICL is a manifestation of "reasoning", then reasoning is modality-agnostic.

**Hypothesis:** Pre-training on sequence data of **other modalities** will facilitate scaling **language** models.*

*there is some evidence on this, but they're narrow (e.g., task-specific fine-tuning)

# Roadmap

1. **Scaling is distribution-dependent:** model behavior changes substantially with shifts in data composition.

2. **Learning emerges beyond human language**

3. LLMs show belief inertia

# Roadmap

1. **Scaling is distribution-dependent:** model behavior changes substantially with shifts in data composition.

2. **Learning emerges beyond human language:** structure and abstraction arise even in non-linguistic distributions.

3. LLMs show belief inertia

# Roadmap

1. **Scaling is distribution-dependent:** model behavior changes substantially with shifts in data composition.

2. **Learning emerges beyond human language:** structure and abstraction arise even in non-linguistic distributions.

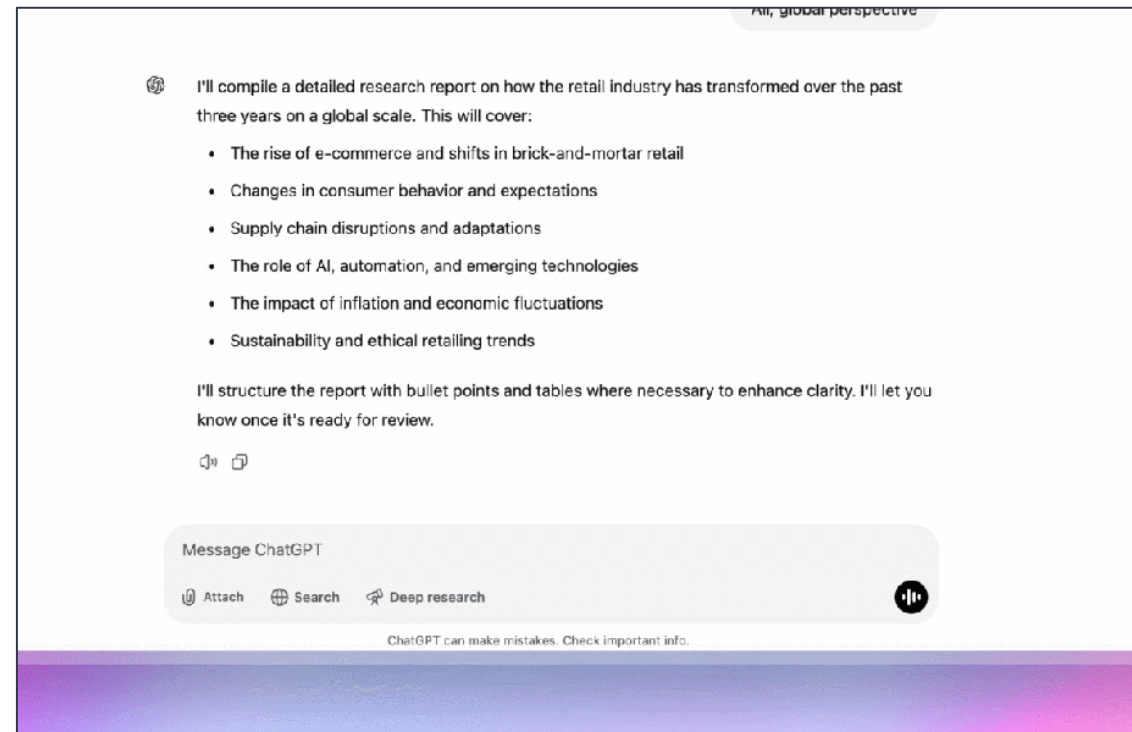3. **LLMs show belief inertia**

# From Passive Solvers to Active Agents

- We are increasing delegating more **freedom (autonomy)** to AI.
  - Freedom to think and act over a long horizon;
  - Freedom to change course and try a different solution, etc.

# Agents Live in Environments

- Agents don't live in a vacuum—they act inside dynamic environments.

- Their behavior is shaped by the feedback they receive
  - user instructions, tool outputs, search results, compiler errors, etc.

When do models decide
**when to stand firm** vs **when to change their mind?**

**Extreme thought experiment**: If we tell the agent that it made a mistake, a "perfect" agent would incorporate *all* corrective signals.

# Feedback Friction: LLMs Struggle to Fully Incorporate External Feedback

Dongwei Jiang, Alvin Zhang, Andrew Wang, Nicholas Andrews, Daniel Khashabi

*NeurIPS 2025*

# Setup: Interaction w/ a Feedback Model

- Goal: How well do LLMs incorporate external feedback?

Feedback Friction: LLMs Struggle to Fully Incorporate External Feedback, 2025

# Setup: Interaction w/ a Feedback Model

- Goal: How well do LLMs incorporate external feedback?



Feedback Friction: LLMs Struggle to Fully Incorporate External Feedback, 2025

# Setup: Interaction w/ a Feedback Model

- Goal: How well do LLMs incorporate external feedback?



*Michael had 58 golf balls. On Tuesday, he lost 23 golf balls. On Wednesday, he lost 2 more. How many golf balls did he have at the end of Wednesday?*

Solver → Response → Correct? — Yes → ✔

Correct? — No → Feedback Giver → Feedback → Solver

Feedback Friction: LLMs Struggle to Fully Incorporate External Feedback, 2025

# Setup: Interaction w/ a Feedback Model

- Goal: How well do LLMs incorporate external feedback?



*Michael had 58 golf balls. On Tuesday, he lost 23 golf balls. On Wednesday, he lost 2 more. How many golf balls did he have at the end of Wednesday?*

**Solver**

*Michael started with 58 golf balls. After losing 23 on Tuesday, he had 58 - 23 = 35 golf balls. After losing 2 more on Wednesday, he had 58 - 2 = 56 golf balls. The final answer is 56.*

Correct?

Yes ✔

No

**Feedback Giver**

**Feedback**

Feedback Friction: LLMs Struggle to Fully Incorporate External Feedback, 2025

118

# Setup: Interaction w/ a Feedback Model
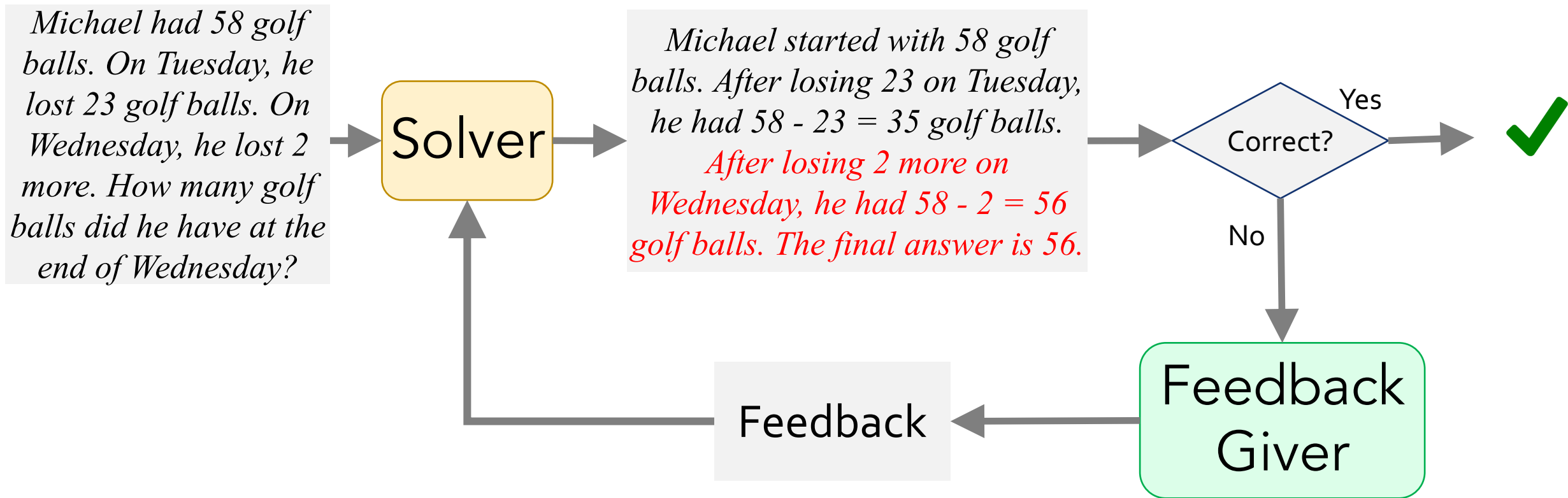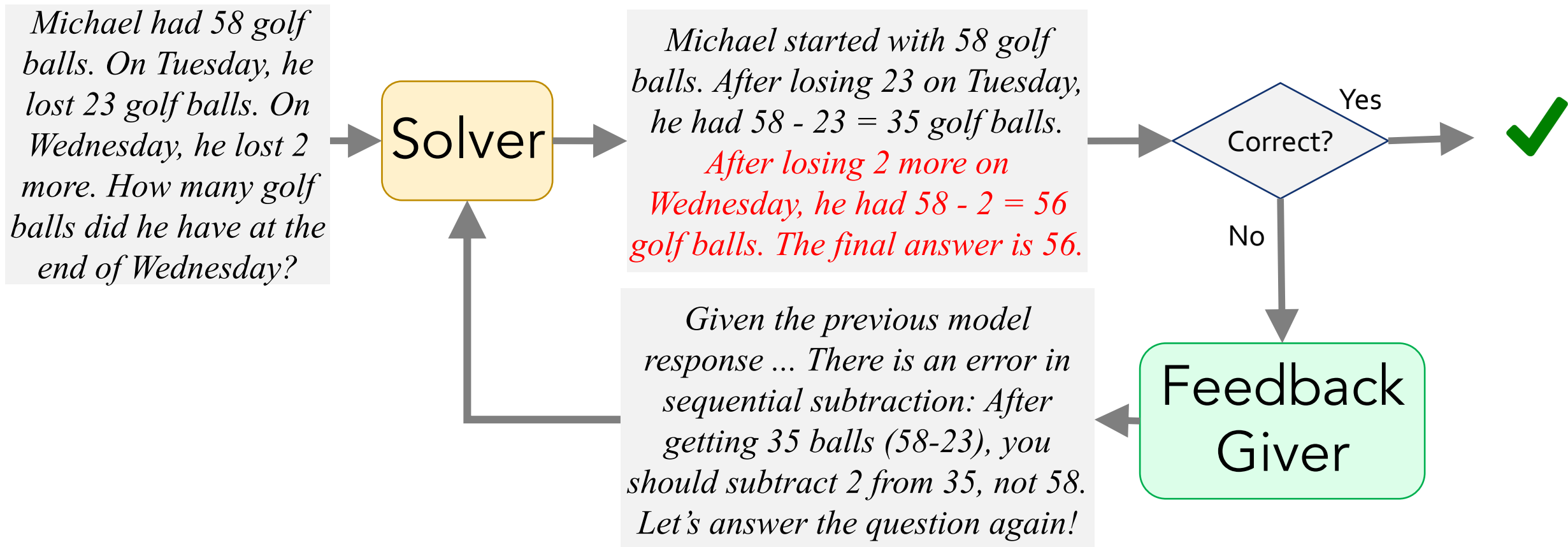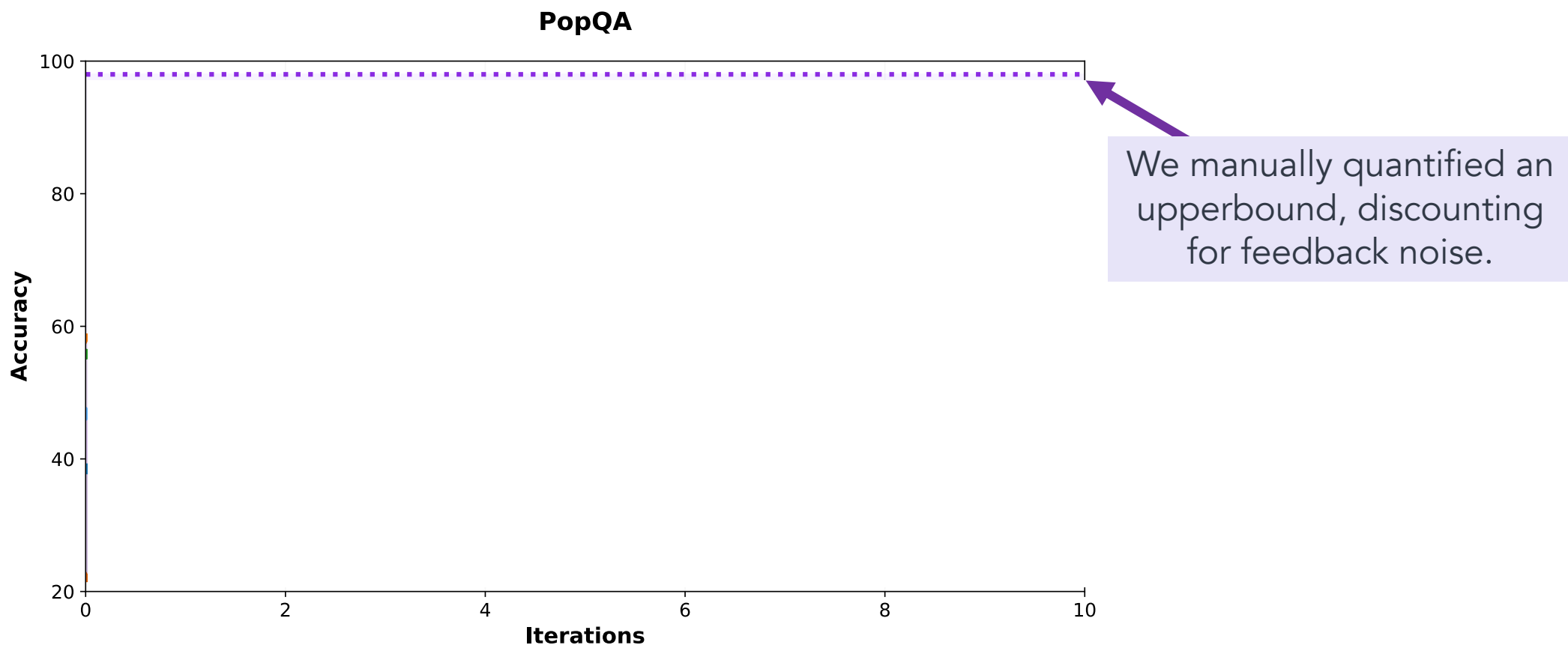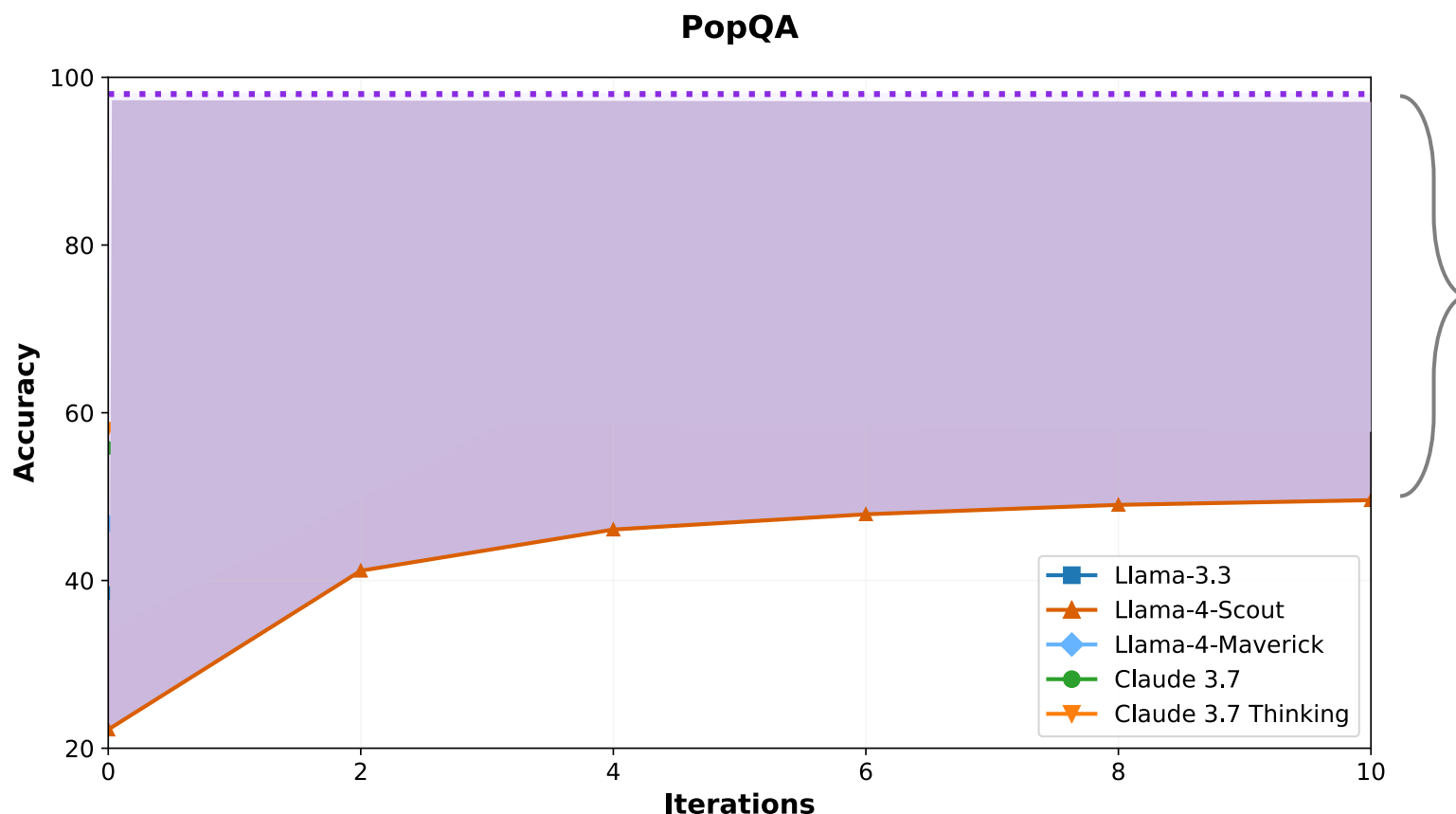
- Goal: How well do LLMs incorporate external feedback?



*Michael had 58 golf balls. On Tuesday, he lost 23 golf balls. On Wednesday, he lost 2 more. How many golf balls did he have at the end of Wednesday?*

Solver

*Michael started with 58 golf balls. After losing 23 on Tuesday, he had 58 - 23 = 35 golf balls. After losing 2 more on Wednesday, he had 58 - 2 = 56 golf balls. The final answer is 56.*

Correct?

Yes ✔

No

Feedback Giver

*Given the previous model response ... There is an error in sequential subtraction: After getting 35 balls (58-23), you should subtract 2 from 35, not 58. Let's answer the question again!*

Feedback Friction: LLMs Struggle to Fully Incorporate External Feedback, 2025

# Interaction w/ a Corrective Feedback: Results

**PopQA**


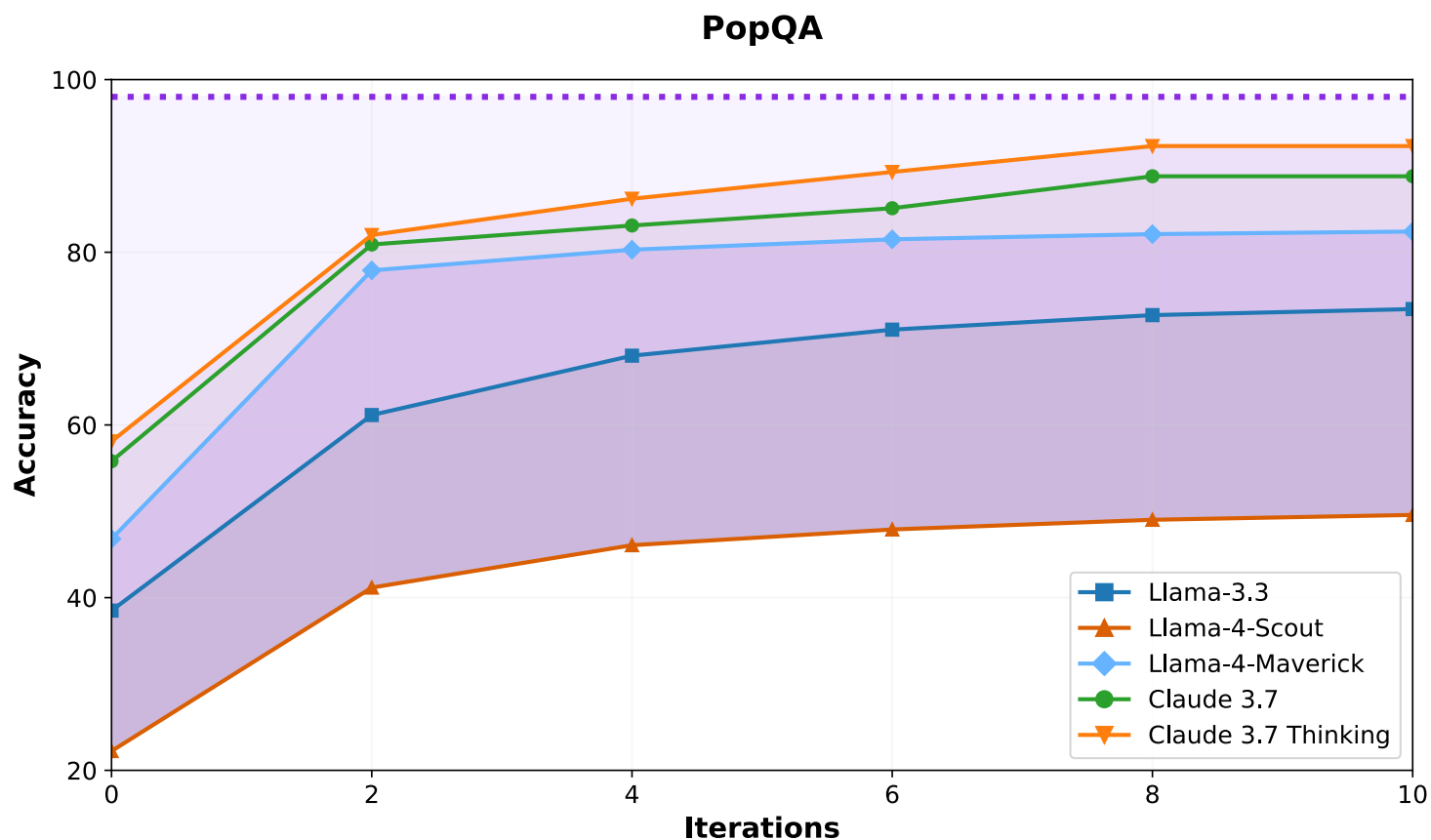
We manually quantified an upperbound, discounting for feedback noise.

- An ideal model should be able to fully incorporate all the constructive feedback.
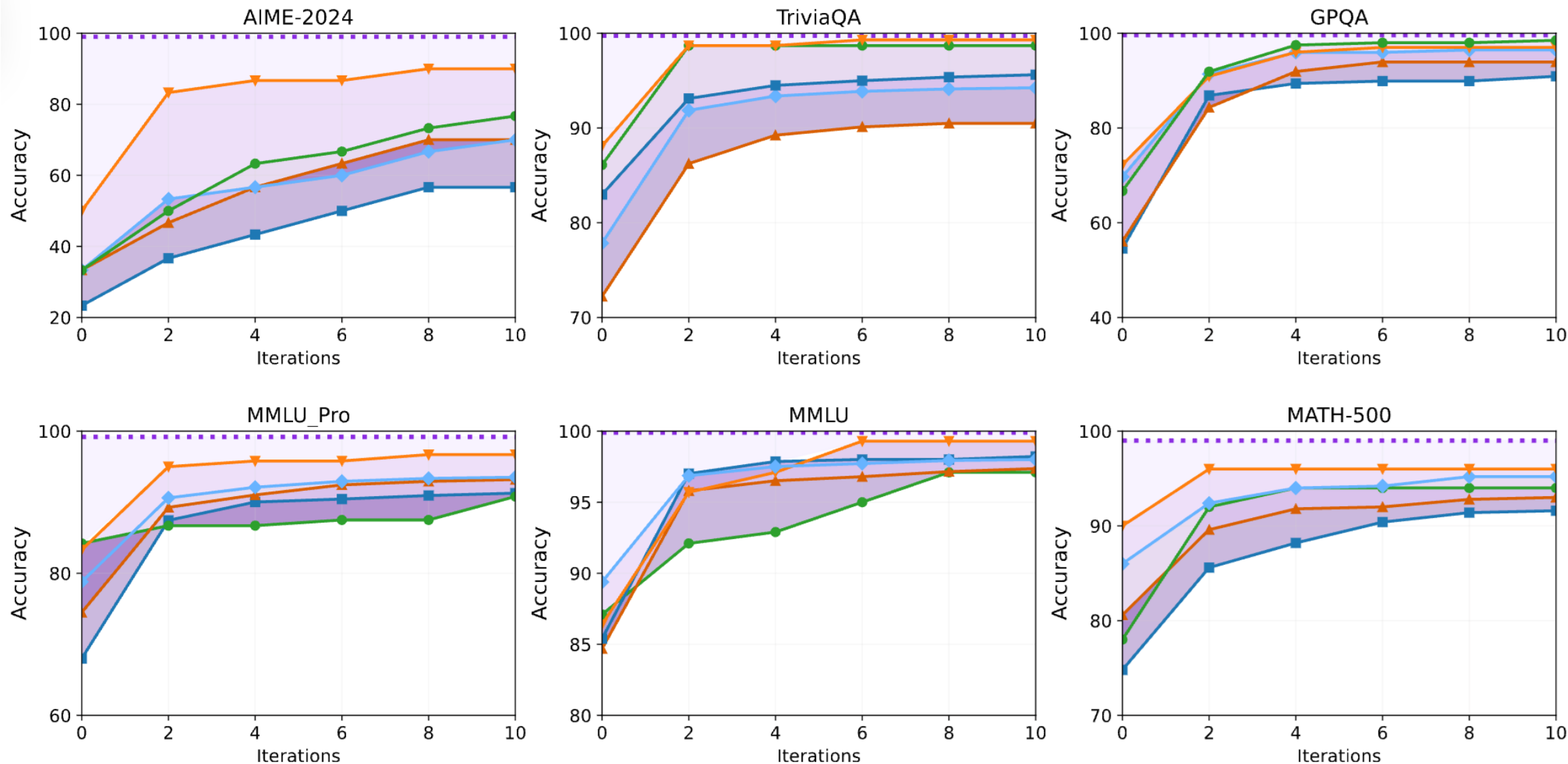
# Interaction w/ a Corrective Feedback: Results



**PopQA**

Legend:
- Llama-3.3
- Llama-4-Scout
- Llama-4-Maverick
- Claude 3.7
- Claude 3.7 Thinking

Models fail to fully integrate the constructive feedback.

Feedback Friction: LLMs Struggle to Fully Incorporate External Feedback, 2025

121

# Interaction w/ a Corrective Feedback: Results



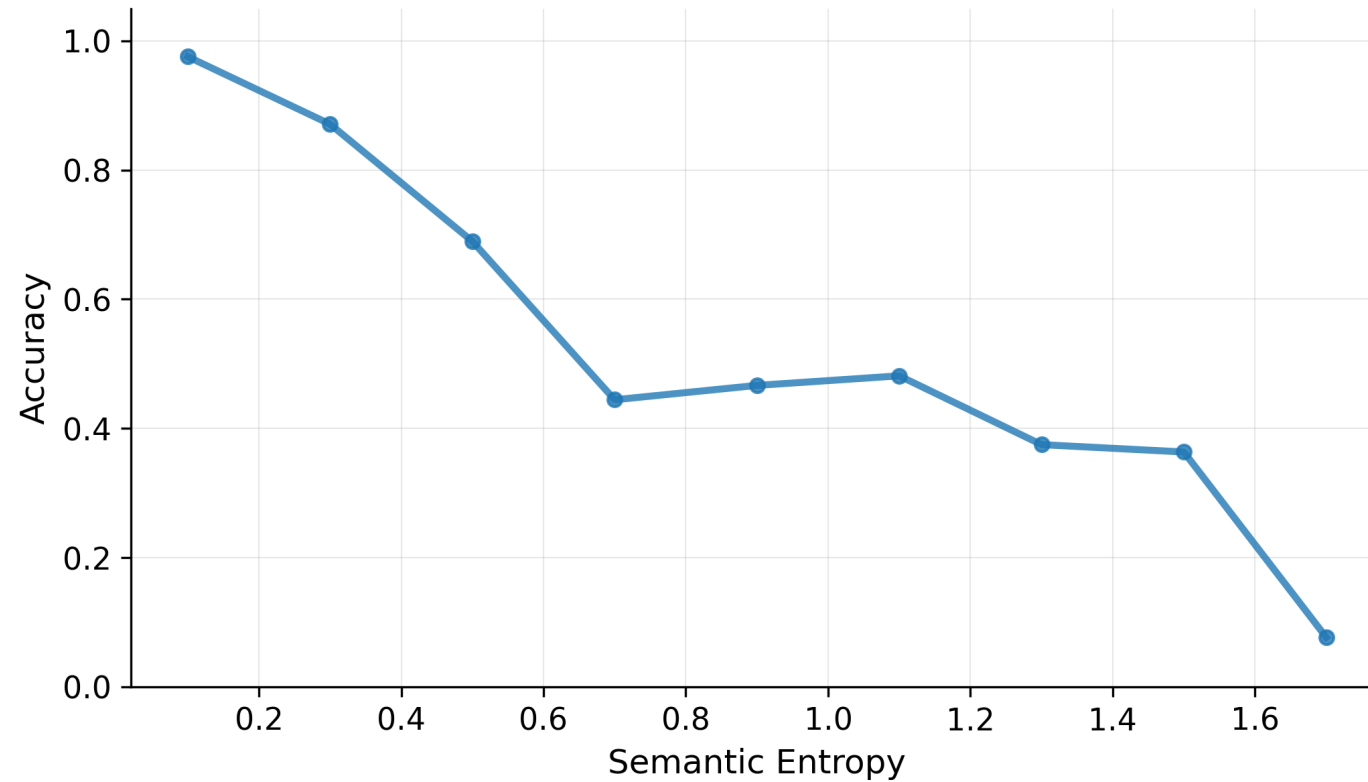Models fail to fully integrate the constructive feedback.

Feedback Friction: LLMs Struggle to Fully Incorporate External Feedback, 2025

Models fail to fully integrate the constructive feedback.

Feedback Friction: LLMs Struggle to Fully Incorporate External Feedback, 2025

123

# Model uncertainty may explain "feedback friction"

# Model uncertainty may explain "feedback friction"

# Model uncertainty may explain "feedback friction"



**Llama-3.3-70B | MATH**

Legend: Initial Accuracy · Final Accuracy · Absolute Improvement Rate

Higher semantic entropy (**more uncertainty**) correlates with greater receptiveness to feedback.

# Model uncertainty may explain "feedback friction"



Higher semantic entropy correlates with greater receptiveness to feedback.

128

# Summary: Feedback Friction

- Models don't always listen to feedback, even if it's constructive.
(Feedback Friction)

- One can trace this back to model uncertainty: when model is certain, it tends to ignore external feedback.

- "Certainty" may correlate with frequency of related data.

- It may also correlate with model accuracy, if the model is calibrated. But most models are not calibrated.

# Stability-Plasticity Tension

- That's where the behavioral tension here:
    - **Too much stability**—Resistant to even high-quality feedback.
    - **Too much plasticity**—Easily swayed by feedback.



Stability-Plasticity Tug-of-War

# Stability-Plasticity Tension: Too Much Plasticity

- LLMs can behave as *interlocutor-pleasers* in dialogue, even if they're initially correct. (conversational sycophancy)
- This makes them vulnerable to flawed feedback.

Published as a conference paper at ICLR 2024

## TOWARDS UNDERSTANDING SYCOPHANCY IN LANGUAGE MODELS

**Mrinank Sharma**, **Meg Tong**, **Tomasz Korbak, David Duvenaud**

**Amanda Askell, Samuel R. Bowman, Newton Cheng, Esin Durmus, Zac Hatfield-Dodds, Scott R. Johnston, Shauna Kravec, Timothy Maxwell, Sam McCandlish, Kamal Ndousse, Oliver Rausch, Nicholas Schiefer, Da Yan, Miranda Zhang,**

**Ethan Perez**

### Challenging the Evaluator: LLM Sycophancy Under User Rebuttal

**Sungwon Kim**
Johns Hopkins University
Baltimore, Maryland
skim434@jhu.edu

**Daniel Khashabi**
Johns Hopkins University
Baltimore, Maryland
danielk@jhu.edu

#### Abstract

Large Language Models (LLMs) often exhibit *sycophancy*, distorting responses to align with

Fanous et al., 2025; Laban et al., 2024). Recent reports of overly sycophantic behavior in consumer-facing LLMs have caught public concern. For example, therapists have cautioned against relying on

Other related effort:
* Are You Sure? Challenging LLMs Leads to Performance Drops in The FlipFlop Experiment, 2023
* Towards Understanding Sycophancy in Language Models, 2024
* Quantifying Multi-Turn Sycophancy in Language Models, 2025
* SycEval: Evaluating LLM Sycophancy, 2025

# Stability-Plasticity Tension

- Overall, no side always dominates.
- But that doesn't imply balance; it's constantly off-balance.



Stability-Plasticity Tug-of-War

- And again, "scaling laws" has nothing to say about these.

# Stability-Plasticity Tension: Open Questions

- What forces govern plasticity-stability?
  - Need to disentangle the factors (data mixtures?, context repetition? etc.)
    
    (Laban et al. 2025)

- How do we engineer (instill) or guarantee a desired balance?

- It's possible that there is no ideal here (?).
  - Perhaps a fundamental trade-off, similar to bias-variance trade-off (?)

- If so, what does that mean for future of agentic AI?
  - (safety, autonomy, reliability, etc.)

# Roadmap

1. **Scaling is distribution-dependent:** model behavior changes substantially with shifts in data composition.

2. **Learning emerges beyond human language:** structure and abstraction arise even in non-linguistic distributions.

3. **LLMs show belief inertia**

# Roadmap

1. **Scaling is distribution-dependent:** model behavior changes substantially with shifts in data composition.

2. **Learning emerges beyond human language:** structure and abstraction arise even in non-linguistic distributions.

3. **LLMs show belief inertia:** models often discount correct updates when they conflict with high-confidence prior beliefs.

# Data ↔ LLM behavior

- LMs are likely to remain brittle:
  - Diminishing returns from scaling.
  - There are numerous nuances that are not captured by scaling laws. (We saw a few of them — probably tip of the iceberg)

- Our understanding of data and its impact on behavior remains quite primitive.

- The heavy focus on "scaling laws" may be have been counterproductive, as it disincentivizing a deeper understanding of data.

# Ilya Sutskever on The State of AI

"We are no longer in the age of scaling, we're back to the age of research."

# Thanks to our wonderful collaborators on these projects!

**Students:**



**Other JHU collaborators:**



**External collaborators & mentors:**