



NEURAL INFORMATION  
PROCESSING SYSTEMS

# Tutorial on the Science of Benchmarking

## What's Measured? What's Missed? What's Next?

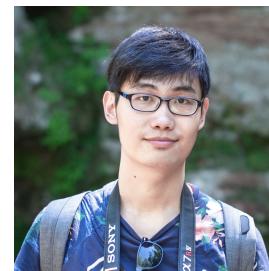
<https://benchmarking.science/>



Martin Ziqiao Ma  
University of Michigan



Michael Saxon  
University of Washington



Xiang Yue  
Carnegie Mellon University  
(Now at Meta)

<https://benchmarking.science/slides.pdf>

Dec 2nd, 2025

# Agenda

- What's Measured? (1:30PM - 2:10PM)
  - What is a (good) benchmark?
  - How to build and maintain a benchmark?
  - How to interpret benchmarking outcomes?
- What's Missed? (2:10PM - 2:40PM)
  - Practical issues: data, integrity, measurement problems
  - Deeper issues: Systemic and epistemic problems
- What's Next? (2:40PM - 3:15PM)
  - Towards dynamic and agentic benchmarking
  - Towards real-world benchmarking
  - Some proposals
- Panel Discussion (3:20PM-4:00PM)

# What's Measured?



Martin Ziqiao Ma  
University of Michigan

# History

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Retirement

Outcomes

- Historically benchmarks are used for computer selection, i.e., running standard programs on different machines to decide which one to buy.

TPC Enterprise Benchmark Standards																																					
Year	1988	1989	1990	1991	1992	1993	1994	1995	1996	1997	1998	1999	2000	2001	2002	2003	2004	2005	2006	2007	2008	2009	2010	2011	2012	2013	2014	2015	2016	2017	2018	2019	2020	2021	2022	2023	2024
TPC-A																																					
TPC-App																																					
TPC-B																																					
TPC-C																																					
TPC-D																																					
TPC-DI																																					
TPC-DS																																					
TPC-E																																					
TPC-H																																					
TPC-R																																					
TPC-VMS																																					
TPC-W																																					
TPC Express Benchmark Standards																																					
TPCx-AI																																					
TPCx-BB																																					
TPCx-HCI																																					
TPCx-HS																																					
TPCx-IoT																																					
TPCx-V																																					
TPC Common Specifications																																					
Pricing																																					
Energy																																					
* ... active benchmark      .. obsolete benchmark																																					
Benchmarks published since 2010 as of 12/31/2024																																					
34 36 16 15 22 24 26 13 20 48 20 70 23 25 20																																					

TPC Benchmarks Overview (<https://www.tpc.org/information/benchmarks5.asp>)

Inioluwa Deborah Raji, et al. AI and the Everything in the Whole Wide World Benchmark. NeurIPS Datasets and Benchmarks Track (Round 2), 2021.

What's Measured?

What's Missed?

What's Next?

# Benchmarking in ML/AI

History

- Now commonly used in machine learning.
- (Tentative) definition: a benchmark is ... (Butterfield & Ngondi, 2016)
  - A problem ...

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Retirement

Outcomes



# Benchmarking in ML/AI

History

- Now commonly used in machine learning.

Definition

- (Tentative) definition: a benchmark is ... (Butterfield & Ngondi, 2016)

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Retirement

Outcomes

- A problem that has been designed to evaluate the performance of a system, ...



# Benchmarking in ML/AI

History

- Now commonly used in machine learning.

Definition

- (Tentative) definition: a benchmark is ... (Butterfield & Ngondi, 2016)

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Retirement

Outcomes

- A problem that has been designed to evaluate the performance of a system, which is subjected to a known workload.



# Benchmarking in ML/AI

History

- Now commonly used in machine learning.
- (Tentative) definition: a benchmark is ... (Butterfield & Ngondi, 2016)
  - A problem that has been designed to evaluate the performance of a system, which is subjected to a known workload.
  - Typically the purpose is to compare the measured performance with other systems under the same benchmark test.

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Retirement

Outcomes



# Benchmarking in ML/AI

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

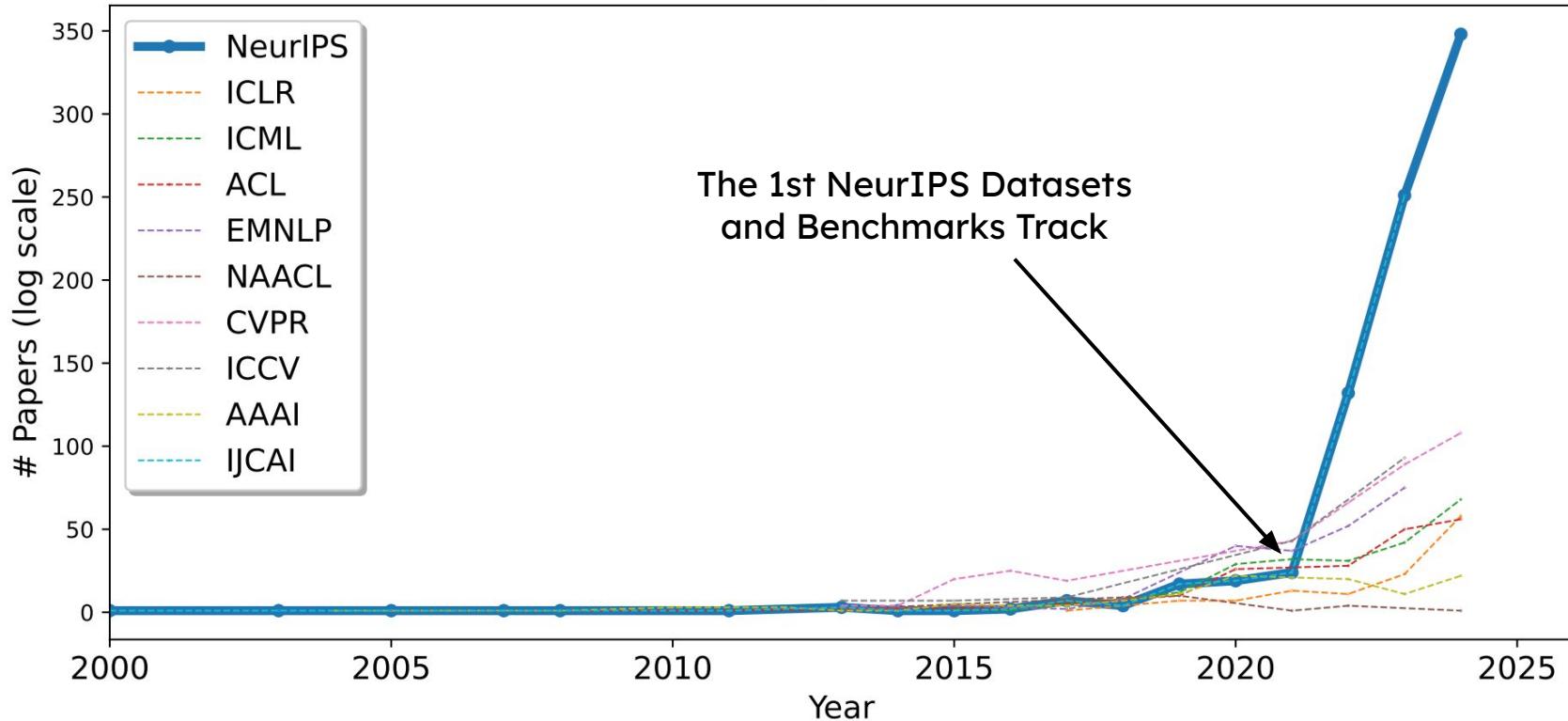
Fluidity

Lifecycle

Design

Retirement

Outcomes



The 1st NeurIPS Datasets  
and Benchmarks Track

Year

What's Measured?

What's Missed?

What's Next?

# “Benchmark” Over Years

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

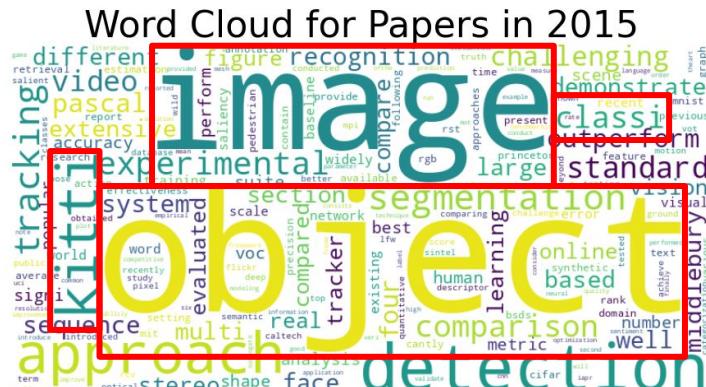
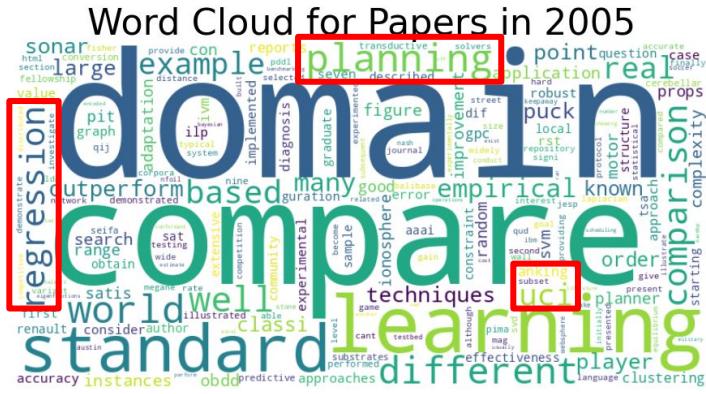
Fluidity

Lifecycle

Design

Retirement

Outcomes



What's Measured?

What's Missed?

What's Next?

# “Benchmark” Over Years

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

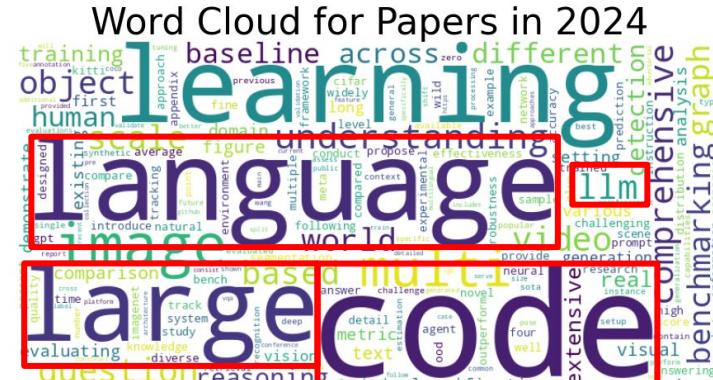
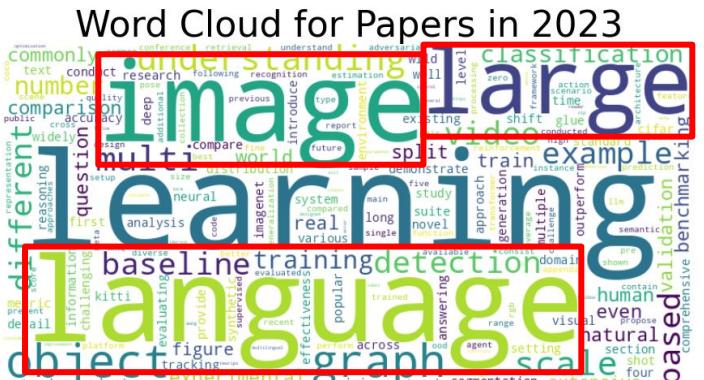
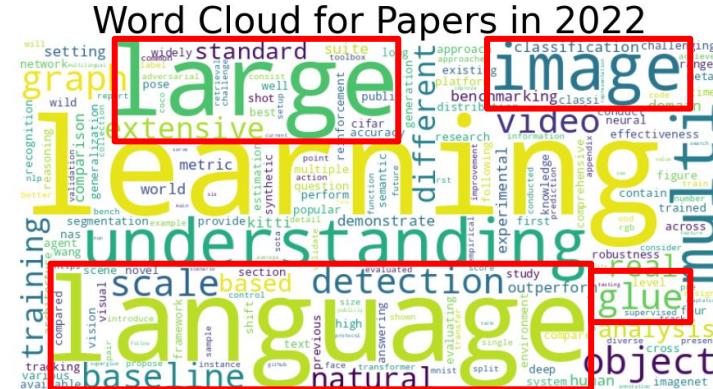
Fluidity

Lifecycle

Design

Refinement

Outcomes



# “Benchmark” Over Years

## History

## Definition

Preliminary

## Subject

## Objective

Axiology

## Coverage

## Fluidity

## Lifecycle

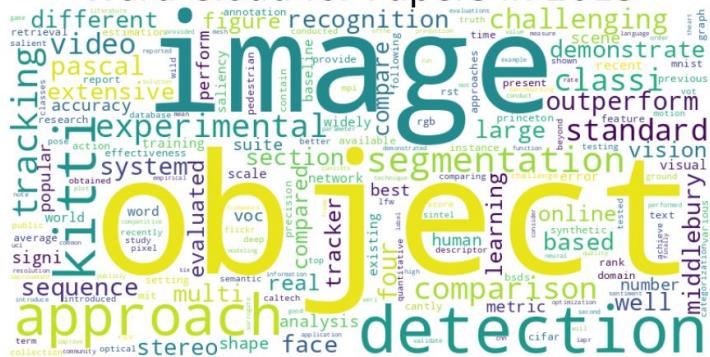
Design

Retirement

## Outcomes

- The way we use benchmarks reflect the paradigm shifts in our fields.
    - Classic ML → deep learning approaches;
    - Small scale → large scale;
    - The rise of application domains: vision, language, graph, etc.

# Word Cloud for Papers in 2015



## Word Cloud for Papers in 2024



## Preliminaries

## History

## Definition

Preliminary

## Subject

## Objective

## Axiology

## Coverage

## Fluidity

## Lifecycle

Design

Retirement

## Outcomes



- A *benchmark* is ?

## Preliminaries

- A *benchmark* is for one or more specific *tasks* or sets of abilities.

A **task** is a particular specification of a problem, ...

Inioluwa Deborah Raji, et al. AI and the Everything in the Whole Wide World Benchmark. NeurIPS Datasets and Benchmarks Track (Round 2), 2021.

## Preliminaries

## History

## Definition

Preliminary

## Subject

## Objective

## Axiology

Coverage

## Fluidity

## Lifecycle

Design

Retirement

## Outcomes



- A **benchmark** is a *dataset* or sets of *datasets* conceptualized as representing one or more specific *tasks* or sets of abilities.

A *task* is a particular specification of a problem, as represented in the *dataset*. There needs to be a test set, sometimes also training and validation sets.

# Preliminaries

History  
Definition  
Preliminary  
Subject  
Objective  
Axiology  
Coverage  
Fluidity  
Lifecycle  
Design  
Retirement  
Outcomes

- A **benchmark** is a combination of (i) a *dataset* or sets of *datasets*, and (ii) a *metric*, conceptualized as representing one or more specific *tasks* or sets of abilities.



A *task* is a particular specification of a problem, as represented in the *dataset*. There needs to be a test set, sometimes also training and validation sets.

A *metric* is way to summarize *model performance* over some set or sets of *tasks* as a single number or score.

# Preliminaries

History  
Definition  
Preliminary  
Subject  
Objective  
Axiology  
Coverage  
Fluidity  
Lifecycle  
Design  
Retirement  
Outcomes

- A **benchmark** is a combination of (i) a **dataset** or sets of **datasets**, and (ii) a **metric**, conceptualized as representing one or more specific **tasks** or sets of abilities. Benchmarks are adopted by a community of researchers as a shared framework for the comparison of **models** (Raji et al., 2021).



**Models** that obtain the most favorable scores on the **metrics** for a **benchmark** in terms of **performance** on the specified **task** is called State of the Art (SOTA).

A **task** is a particular specification of a problem, as represented in the **dataset**. There needs to be a test set, sometimes also training and validation sets.

A **metric** is way to summarize **model performance** over some set or sets of **tasks** as a single number or score.

# Demarcation

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Refinement

Outcomes

- Quick Test: are these evaluations “benchmarks” per Raji et al. (2021)?

Computer Science > Computation and Language

[Submitted on 22 Apr 2019 (v1), last revised 9 Sep 2019 (this version, v3)]

**SocialIQA: Commonsense Reasoning about Social Interactions**

Computer Science > Computation and Language

[Submitted on 4 Feb 2023 (v1), last revised 4 Nov 2024 (this version, v7)]

**Evaluating Large Language Models in Theory of Mind Tasks**

Computer Science > Computer Vision and Pattern Recognition

[Submitted on 13 Feb 2025 (v1), last revised 6 Mar 2025 (this version, v2)]

**ZeroBench: An Impossible Visual Benchmark for Contemporary Large Multimodal Models**



What's Measured?

What's Missed?

What's Next?

# Demarcation

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Retirement

Outcomes

- Sap and Rashkin et al. (2019) on social interactions:
  - **Motivation:** “*While humans trivially acquire and develop such social reasoning skills, this is still a challenge for machine learning models.*”
  - **Task & Subject:** “*Social IQa aims to measure the social and emotional intelligence of computational models through multiple choice question answering.*”
  - **Data:** “*Social IQa contains 37,588 multiple choice questions with three answer choices per question.*”
  - **Metric:** “*Despite human performance of close to 90%, computational approaches based on large pretrained language models only achieve accuracies up to 65%.*”

# Demarcation

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Retirement

Outcomes

- Kosinski (2024) on Theory-of-Mind (ToM):
  - **Motivation:** “We hypothesize that ToM-like ability does not have to be explicitly engineered into AI systems.”
  - **Task & Subject:** “we administer two versions of the classic false-belief task widely used to test ToM in humans to several language models.”
  - **Data:** “As GPT-3.5 may have encountered the original task in its training, hypothesis-blind research assistants (RAs) prepared 20 bespoke Unexpected Contents Task tasks.”
  - **Metric:** “A task was considered solved correctly only if all 3 questions were answered correctly for both original and reversed task.”



# Demarcation

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Retirement

Outcomes

- Roberts et al. (2025), ZeroBench:

- **Motivation:** “*...there is a pressing need for difficult benchmarks that remain relevant for longer*”
- **Task & Subject:** “*...a lightweight visual reasoning benchmark that is entirely impossible for contemporary frontier LMMs.*”
- **Data:** “*Our benchmark consists of 100 manually curated questions and 334 less difficult subquestions.*”
- **Metric:** “*We use accuracy as our metric to evaluate the 100 ZeroBench main questions*”

# Broadening the Definition

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

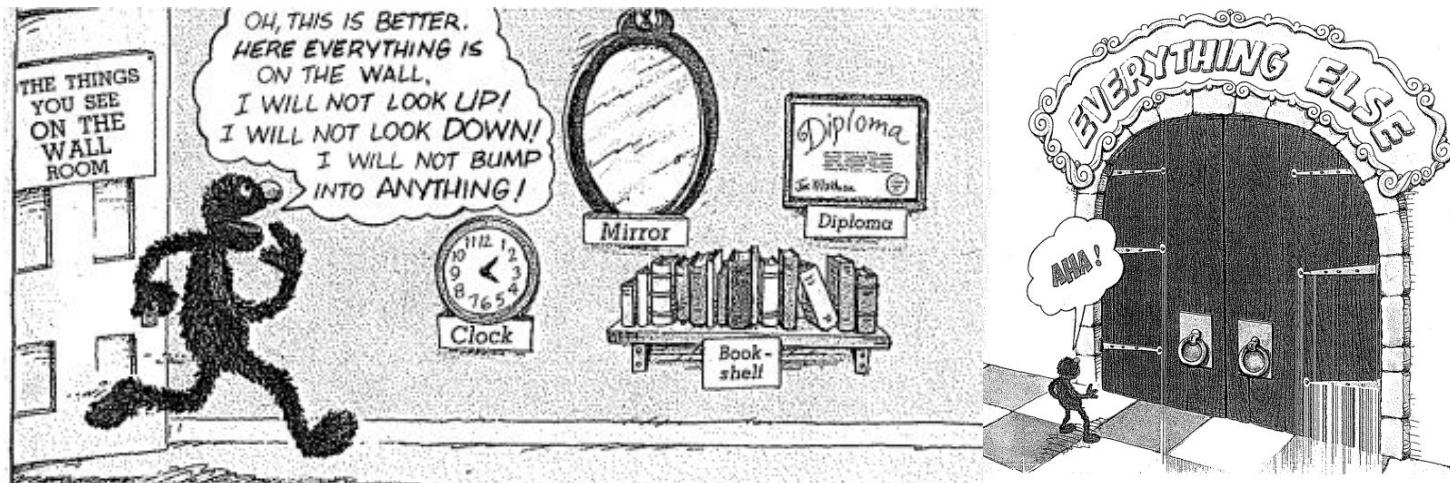
Lifecycle

Design

Retirement

Outcomes

- Raji et al. (2021): “*The imagined artifact of the ‘general’ benchmark does not actually exist.*”
- But we are seeing more “general” benchmarks than domain-specific benchmarks nowadays.



Inioluwa Deborah Raji, et al. AI and the Everything in the Whole Wide World Benchmark. NeurIPS Datasets and Benchmarks Track (Round 2), 2021.

What's Measured?

What's Missed?

What's Next?

# Broadening the Definition

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Retirement

Outcomes

- AI benchmarking essentially depends on our philosophy of measurement.
- AI Benchmarking often drifts into operationalism (capability = score) without accepting its epistemic commitments.
  - I.e., “SOTA on SWE-Bench = Best SWE agent”
  - Benchmarks are fallible indicators of a richer construct.



# Broadening the Definition

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Retirement

Outcomes

- These “general” AI benchmarks are proto-psychometrics: tests of latent abilities, but commonly with underdeveloped theory and validation.
  - Their job is **construct validation** (Bean, et al., 2025), not **definition**;
  - *“having measures that represent what matters to the phenomenon.”*
  - We hope to argue: this task **T** is a useful probe of ((aspect **X** of) capability **C** of) a system **S**, under specific assumptions **A1, A2, ...**

# Broadening the Definition

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Retirement

Outcomes

- What is a benchmark?
  - Defined by what we study (the system) and what we want (the goal).
- What is a **good** benchmark?
  - No single answer, laden with values and methodological assumptions.
- How should we use benchmarks?
  - Follow their intended scope and interpret scores within the benchmark's and your system's assumptions.



# The Target System

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Retirement

Outcomes

- A machine learning algorithm:
  - I.e., classical supervised learning regime;
  - E.g., a new linear transformer;
  - Dataset split: train / validation / test;
  - A benchmark is a dataset with **fixed split (train/val/test)** plus an evaluation metric, used to compare algorithms under controlled data.



# The Target System

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Retirement

Outcomes

- (Edge case) continual or test-time adaptable algorithms:
  - E.g., MEMO, TENT, TTT-based algorithms with batch norm (BN);
  - Updates are computed from batches of test data;
  - Changing batch size → changes the gradient and BN statistics;
  - Changing batch order → changes the adaptation trajectory;
  - For continual or test-time adaptable algorithms, a meaningful benchmark specification should include, in addition to the dataset splits and metric, **the test-time batch size and the evaluation order**, since these affect the adaptation dynamics and thus performance.

Yu Sun, Xiaolong Wang, Zhuang Liu, John Miller, Alexei Efros, Moritz Hardt. Test-Time Training with Self-Supervision for Generalization under Distribution Shifts. ICML, 2020.  
Dequan Wang, Evan Shelhamer, Shaoteng Liu, Bruno Olshausen, Trevor Darrell. Tent: Fully Test-time Adaptation by Entropy Minimization. ICLR, 2021.  
Marvin Zhang, Sergey Levine, Chelsea Finn. MEMO: Test Time Robustness via Adaptation and Augmentation. NeurIPS, 2022.

# The Target System

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Retirement

Outcomes

- A pretrained weight:
  - I.e., Self-supervised learning and pretraining regimes;
  - E.g., BERT-base, Qwen3-0.6B-Base, DINOV2, MAE-L;
  - Dataset split: dev / test;
  - A benchmark is a test set plus an evaluation metric, **optionally with a dev set for tuning**, used to compare pretrained weights under controlled data conditions.



# The Target System

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Retirement

Outcomes

- A digital (code/web/terminal/...) agent:
  - E.g., modern LLM agents
  - “I developed the SOTA agent”: A benchmark is a test set plus an evaluation metric, used to compare agents.
    - E.g., OpenHands + GPT-5 scores 71.8 on SWE-Bench (verified).
  - “I developed the SOTA LLM with agentic capability”: A benchmark is a test set, **a fixed agent workflow**, plus an evaluation metric, used to compare pretrained weights under same workflow and workload.
    - E.g., Claude 4.5 Opus scores 74.4 on SWE-Bench (Verified) with a minimal agent (aka, the Bash Only setting).



# The Research Goal

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

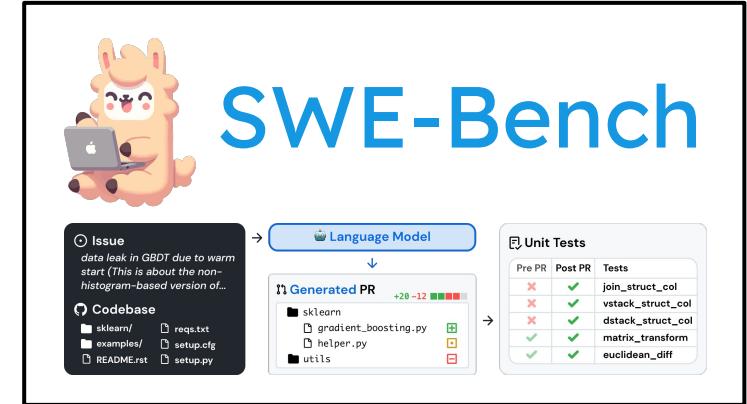
Refinement

Outcomes

- Task-driven.
- To understand our progress on a task or well-defined application.



*“we annotate each lidar point from a keyframe in nuScenes with one of 32 possible semantic labels (i.e. lidar semantic segmentation).”*



*“evaluates LMs in a realistic software engineering setting... featuring [GitHub issues and pull requests] from 12 repositories”*

Holger Caesar, et al. nuScenes: A Multimodal Dataset for Autonomous Driving. CVPR, 2020.

Carlos E. Jimenez, et al. SWE-bench: Can Language Models Resolve Real-World GitHub Issues? ICLR, 2024.

What's Measured?

What's Missed?

What's Next?

# The Research Goal

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Refinement

Outcomes

- Capability-driven.
- To evaluate the level of competence of a cognitive capability.



*“...a collection of tools for evaluating the performance of models across a diverse set of existing NLU tasks.”*



*“We propose a novel multimodal video benchmark...to evaluate the perception and reasoning skills of pre-trained multimodal models”*

Alex Wang, et al. GLUE: A multi-task benchmark and analysis platform for natural language understanding. BlackboxNLP, 2018.  
Alex Wang, et al. Superglue: A stickier benchmark for general-purpose language understanding systems. NeurIPS, 2019.  
Viorica Patrăucean, et al. Perception test: A diagnostic benchmark for multimodal video models. NeurIPS, 2023.

What's Measured?

What's Missed?

What's Next?

# The Research Goal

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

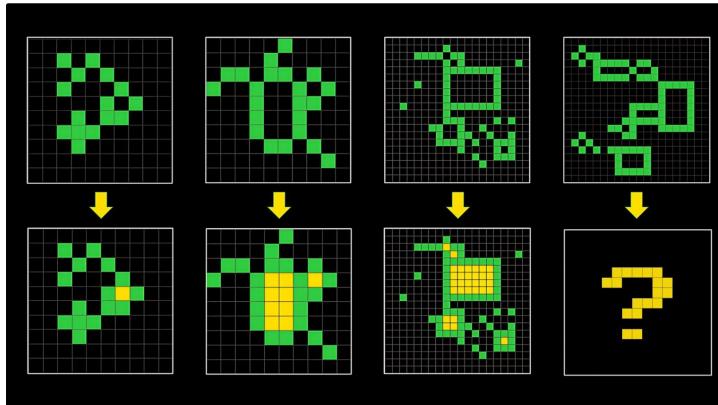
Lifecycle

Design

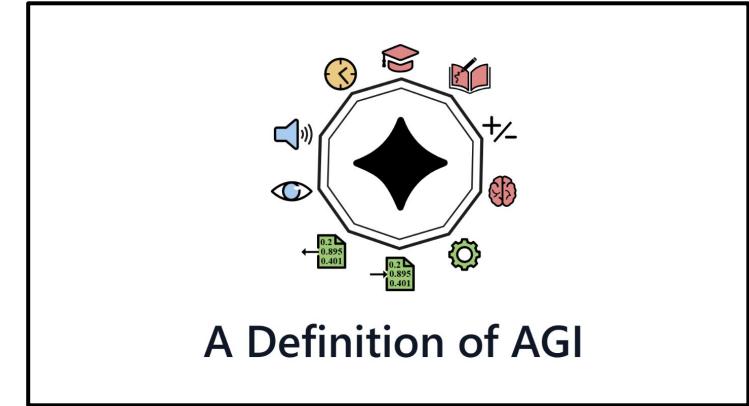
Refinement

Outcomes

- Model-driven.
- To measure the “intelligence” of a model.



*“We argue that ARC can be used to measure a human-like form of general fluid intelligence...”*



*“The framework dissects general intelligence into ten core cognitive domains...”*

François Chollet. On the Measure of Intelligence. Preprint, 2019.  
Dan Hendrycks, et al. A Definition of AGI. Preprint, 2025.

What's Measured?

What's Missed?

What's Next?

# Defining “Quality” of Benchmarks

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Retirement

Outcomes

- What is a good benchmark?
- No single answer. Empirical results are laden with values and theoretical commitments (Boyd and Bogen, 2019).
- If one is designing a task- or capability-driven benchmark, a good benchmark is **a good data sampling function over the problem space.**



# Defining “Quality” of Benchmarks

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

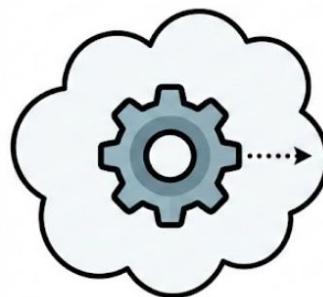
Fluidity

Lifecycle

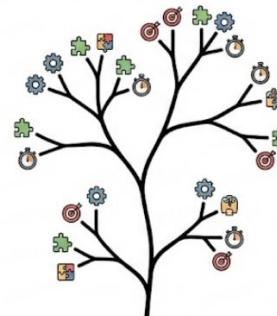
Design

Retirement

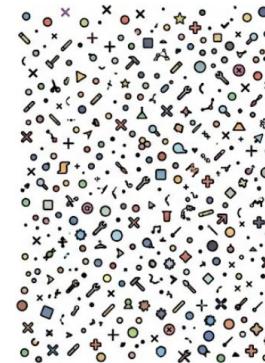
Outcomes



A Proxy  
Subtask



A Taxonomy  
of Subtasks



No Subtasks, Just  
Massive Testcases

- When designing a new benchmark for a task, you would prefer ...

# Coverage in Benchmarks

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

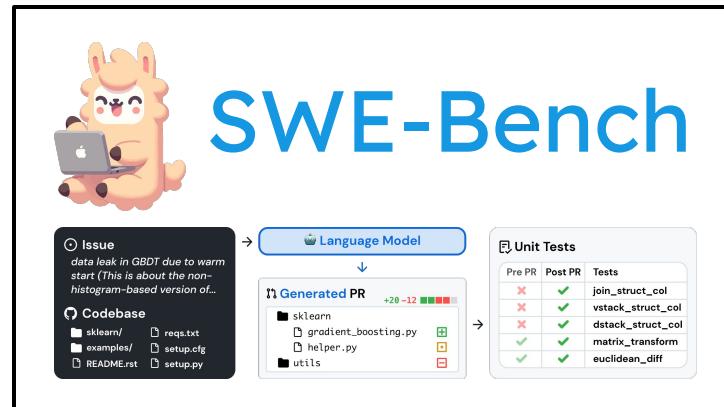
Fluidity

Lifecycle

Design

Refinement

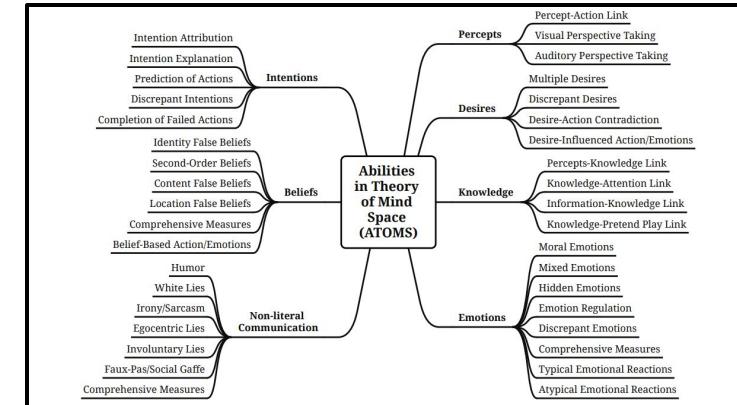
Outcomes



*“However, existing coding benchmarks, such as HumanEval, mostly involve self-contained problems that can be solved in **a few lines of code**. In the real world, software engineering is not as simple ... we introduce SWE-bench, a benchmark that evaluates LMs in a **realistic software engineering setting**. models are tasked to **resolve issues** (typically a bug report or a feature request) submitted to popular GitHub repositories”*

# Coverage in Benchmarks

History  
Definition  
Preliminary  
Subject  
Objective  
**Axiology**  
Coverage  
Fluidity  
Lifecycle  
Design  
Retirement  
Outcomes



*“While the exact definition of ToM remains a central debate, the AI community can benefit from looking at what psychologists have viewed as an initial step. In this paper, we **follow Beaudoin et al. (2020)**’s taxonomy of ToM sub-domains, i.e., the Abilities in Theory of Mind Space (ATOMS).”*

# Coverage vs Difficulty

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Retirement

Outcomes



*"GLUE is a collection of nine language understanding tasks built on existing public datasets, together with private test data, an evaluation server, a single-number target metric, and an accompanying expert constructed diagnostic set."*

*"SuperGLUE retains the two hardest tasks in GLUE. The remaining tasks were identified from those submitted to an open call for task proposals and were selected based on difficulty for current NLP approaches"*

# Reality Checking on Coverage

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Retirement

Outcomes

- Reality check: how good are benchmarks as data sampling functions?
- Formulation.
  - Let  $\mathbf{T}$  be the true task/domain distribution and  $\mathbf{B}$  the benchmark's sampling distribution.
  - A benchmark is “good as a data sampling function” if  $\mathbf{B}$  is (for the purposes we care about) indistinguishable from  $\mathbf{T}$ .
  - That is, no reasonable test can look at the data alone and reliably tell whether an example came from the real domain or from the benchmark.



# Reality Checking on Coverage

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Retirement

Outcomes

- Reality check: how good are benchmarks as data sampling functions?
- Operationalized formulation.
  - A perfect benchmark is unattainable in both theory and practice.
  - We can only probe it indirectly via cross-benchmark checks and;
    - (i) Reproduce the data collection with the same described methodology or from real user interaction logs;
    - (ii) Compare to surrogate human data;
  - Use the OpenAI *text-embedding-3-small* model to encode the tasks into embeddings and compare benchmarks and human queries on similar tasks.



# Reality Checking on Coverage

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

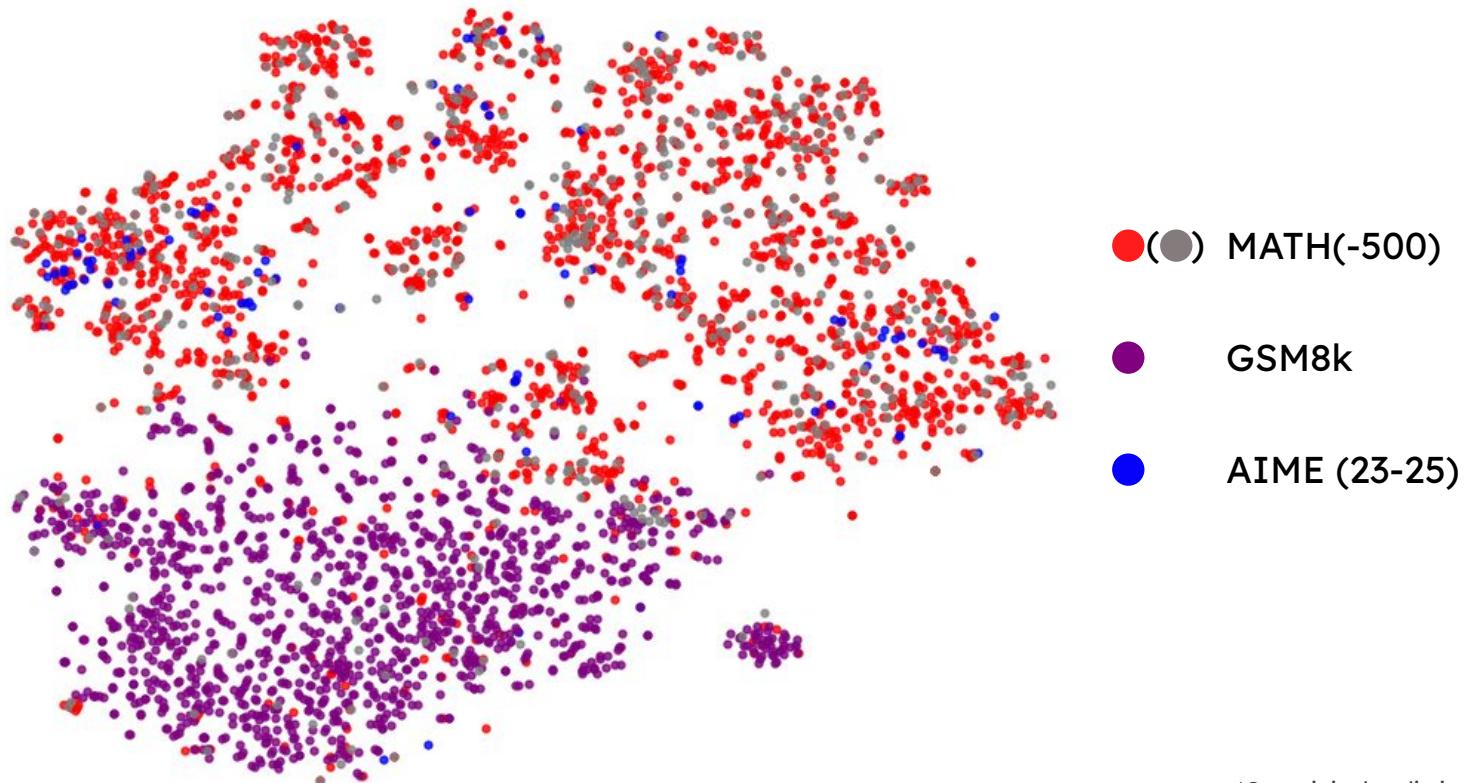
Lifecycle

Design

Retirement

Outcomes

- Mathematics.



\*Our original preliminary results.

What's Measured?

What's Missed?

What's Next?

# Reality Checking on Coverage

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

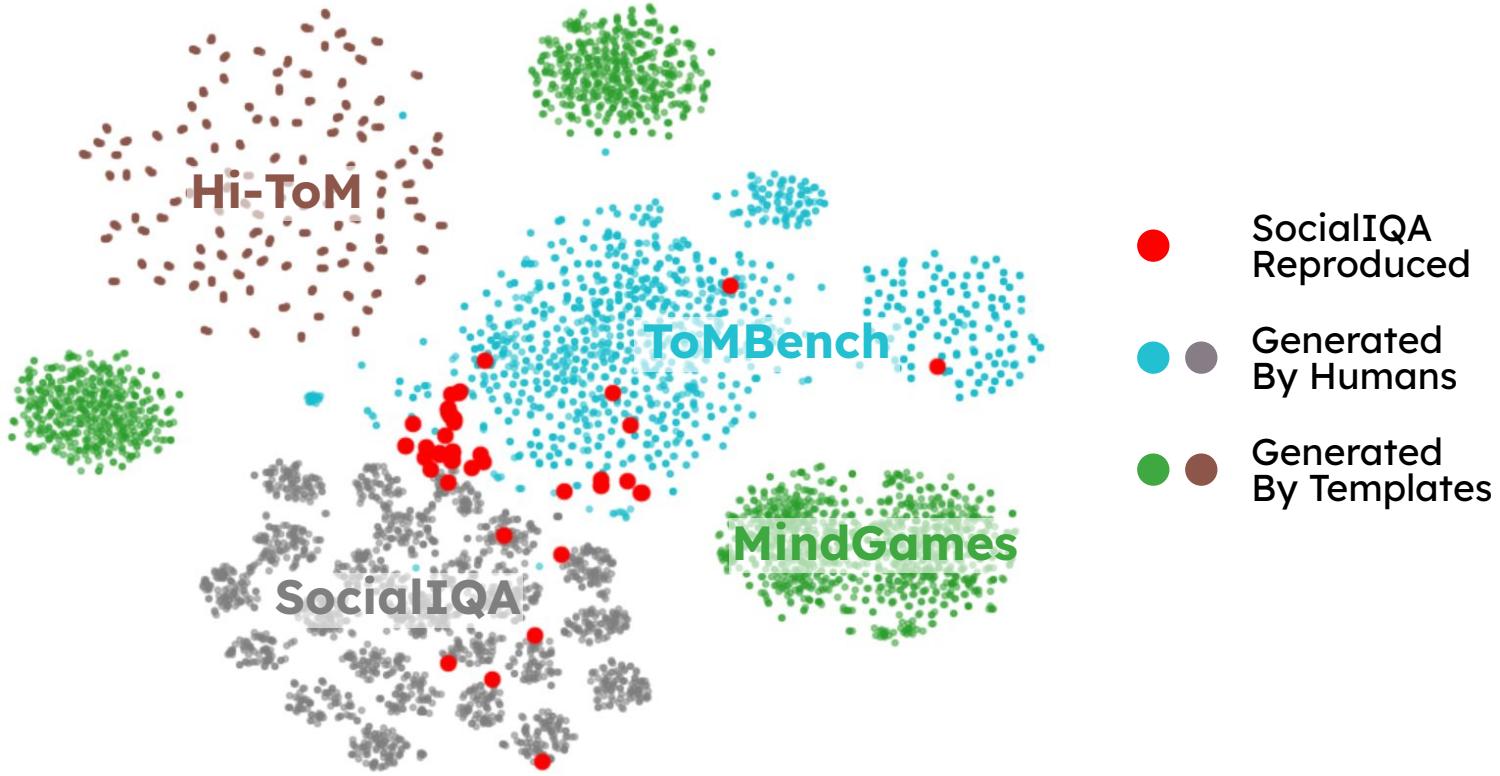
Lifecycle

Design

Retirement

Outcomes

- Social Reasoning and Theory of Mind.



\*Our original preliminary results.

What's Measured?

What's Missed?

What's Next?

# Reality Checking on Coverage

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

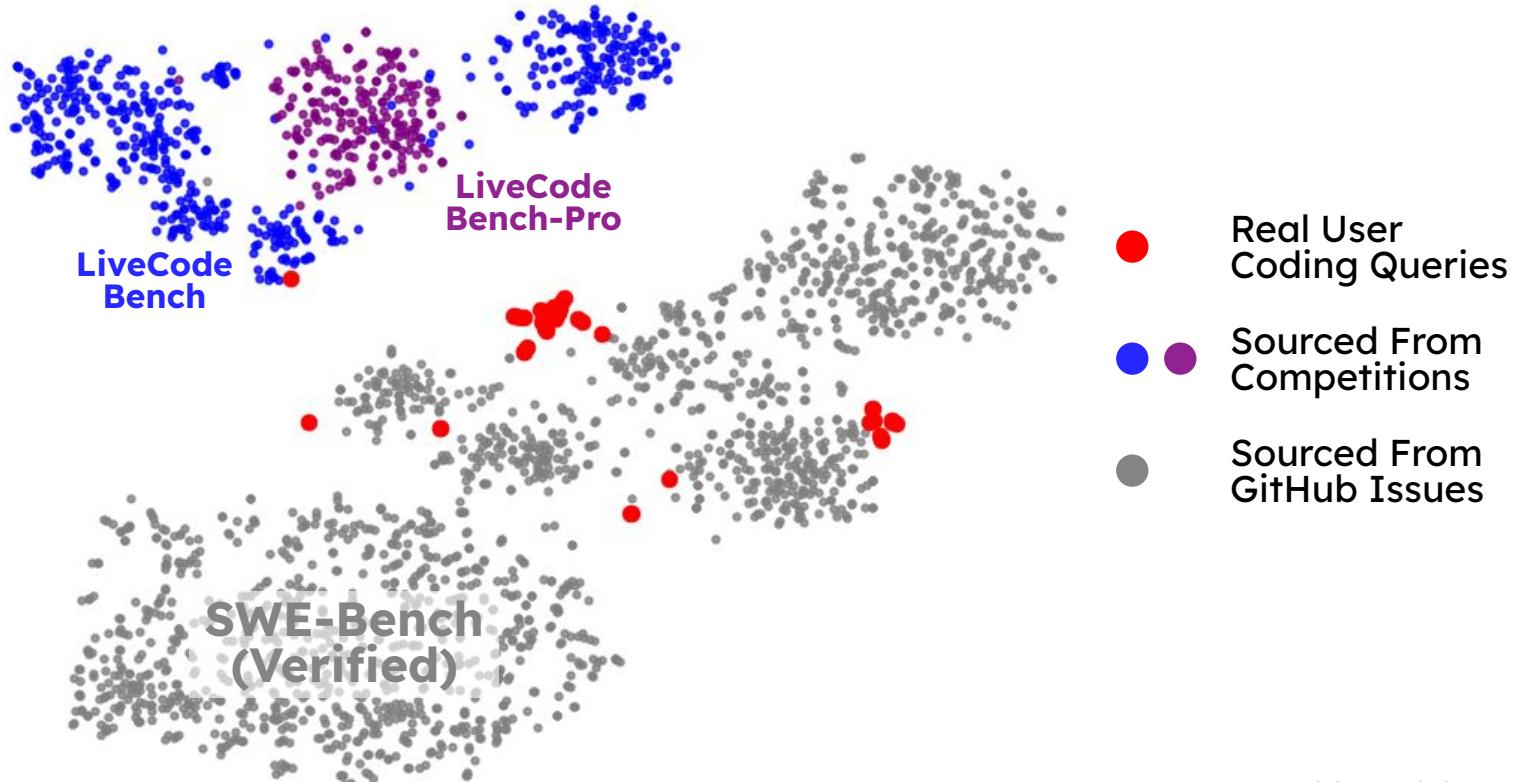
Lifecycle

Design

Retirement

Outcomes

- Coding and software engineering.



# Reality Checking on Coverage

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Retirement

Outcomes

- Reality check: how good are benchmarks as data sampling functions?
- Summary.
  - Templatized, proxy- and taxonomy-based datasets can easily fragment into many tiny clusters compared to human-generated ones.
  - The data collection pipelines described in papers can rarely be reproducible in practice.
  - Real user queries in human-LLM logs look very different from benchmark prompts.



# Model-Driven Benchmarks

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Retirement

Outcomes

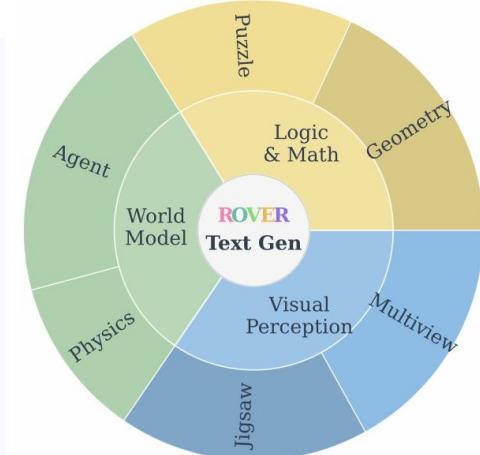
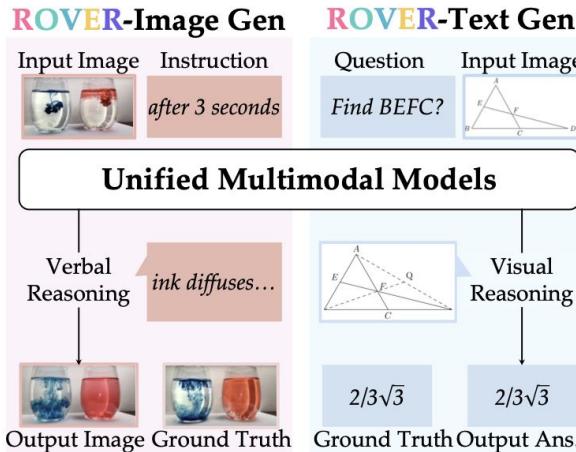
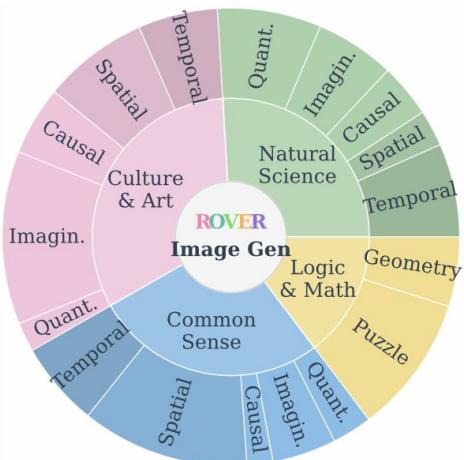
- Raji et al. (2021): benchmark should be task-driven.
- When we try to broaden to model-driven benchmarks, we are...
- Either ...
  - Motivated by new capabilities demonstrated by new models;



# Benchmarks Motivated by New Capabilities

History  
Definition  
Preliminary  
Subject  
Objective  
**Axiology**  
Coverage  
Fluidity  
Lifecycle  
Design  
Retirement  
Outcomes

- Example: omni-modal generation in unified multimodal models.



Yongyuan Liang, Wei Chow, et al. ROVER: Benchmarking Reciprocal Cross-Modal Reasoning for Omnimodal Generation. Preprint, 2025.  
Zhiyuan Yan et al. Unified Multimodal Model as Auto-Encoder. Preprint, 2025.

What's Measured?

What's Missed?

What's Next?

# Model-Driven Benchmarks

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Retirement

Outcomes

- Raji et al. (2021): benchmark should be task-driven.
- When we try to broaden to model-driven benchmarks, we are...
- Either ...
  - Motivated by new capabilities demonstrated by new models;
- Or ...
  - Hoping to challenge capable models with tasks that fail them;
  - Implicitly assuming generality;
  - Sometimes, claiming a measurement of intelligence, or colloquially using these benchmarks as proof that AGI.



# Thought Experiment

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Retirement

Outcomes

- Martin's magic gibberish word: **flumborix** (made up by ChatGPT);
- Benchmark question: *What is Martin's magic gibberish word in his NeurIPS 2025 tutorial?*
- A model with a knowledge cutoff before Dec 1, 2025 must search the internet, find this slide, and extract this word.
- A model with a knowledge cutoff after Dec 1, 2025 can simply recall the word from its training dataset, which happens to include this slide.



# Definition of “Intelligence”

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Retirement

Outcomes

- The 1921 Symposium on Intelligence and Its Measurement
  - Intelligence(capacity[knowledge]) = Library(shelf[books])
  - Henmo (1921): “*the capacity for knowledge and knowledge possessed.*”
- Journal of AGI Special Issue “On Defining Artificial Intelligence”
  - Intelligence(agent[knowledge]) = Energy(motor[fuel]) \*
  - Laird (2020): “*equate[s] with rationality, where an agent uses its available knowledge to select the best action(s) to achieve its goal(s) within an environment.*”

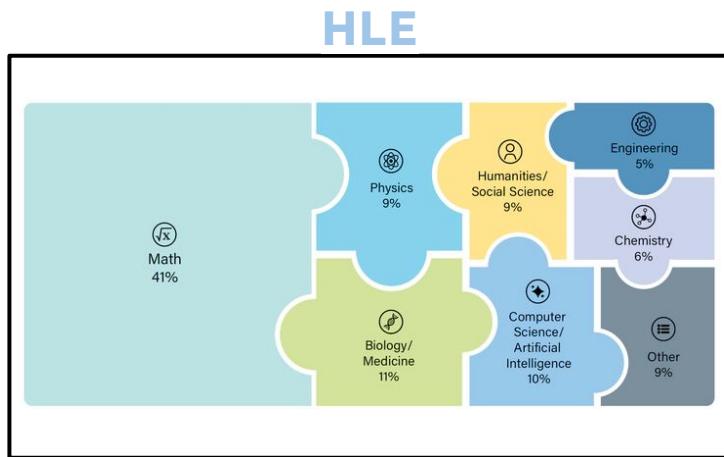
\*This analogy is taken from John Laird’s EECS 598 (Winter 2022) AGI with permission :)

Henmon, V.A.C. Intelligence and its measurement: A symposium VIII. Journal of Educational Psychology, 1921.  
John Laird. Intelligence, knowledge & human-like intelligence. Journal of Artificial General Intelligence, 2020.

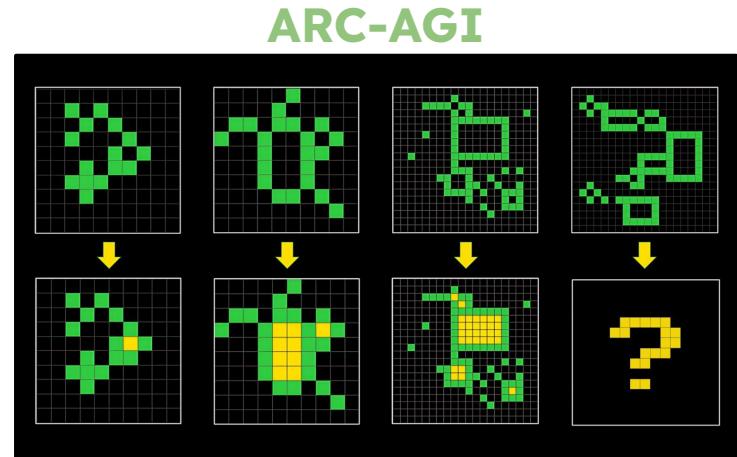
# Fluid vs Crystallized

History  
Definition  
Preliminary  
Subject  
Objective  
Axiology  
Coverage  
Fluidity  
Lifecycle  
Design  
Retirement  
Outcomes

- Fluid intelligence is the ability to reason and solve new problems;
- Crystallized intelligence is the accumulation of knowledge and skills over a lifetime.



*"a multi-modal benchmark at the frontier of human knowledge, designed to be the final closed-ended academic benchmark of its kind with broad subject coverage"*



*"We argue that ARC can be used to measure a human-like form of general fluid intelligence..."*

# Summary

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Retirement

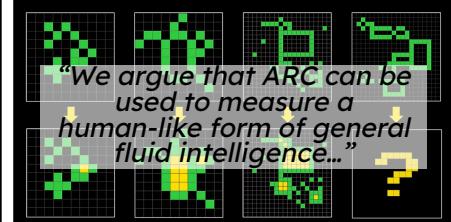
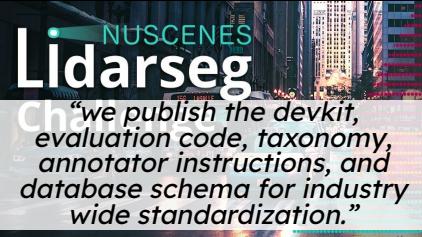
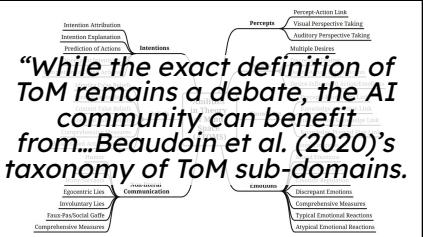
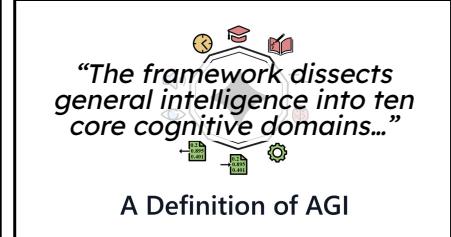
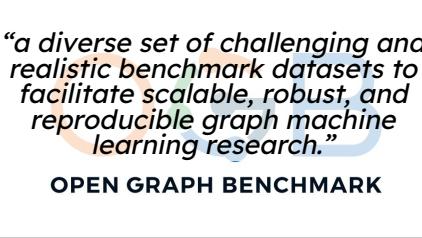
Outcomes

what is good data sampling?

proxy is enough

taxonomy guided

scale for coverage

what is the primary goal of the evaluation?		
Progress of task	capability competence	model's intelligence
<p>“evaluates LMs in a realistic software engineering setting... featuring [GitHub issues and pull requests] from 12 repositories”</p> <p> <b>SWE-Bench</b></p>	<p>“We administer classic false-belief tasks, widely used to test ToM in humans, to several language models...”</p> <p>Evaluating Theory of Mind in Question Answering Aida Nematzadeh, Kaylee Burns, Erin Grant, Alison Gopnik, Thomas L. Griffiths Theory of Mind May Have Spontaneously Emerged in Large Language Models Michał Kosinski</p>	 <p>“We argue that ARC can be used to measure a human-like form of general fluid intelligence...”</p>
<p> <b>Lidarseg</b> “we publish the devkit, evaluation code, taxonomy, annotator instructions, and database schema for industry wide standardization.”</p>	 <p>“While the exact definition of ToM remains a debate, the AI community can benefit from...Beaudoin et al. (2020)’s taxonomy of ToM sub-domains.</p>	 <p>“The framework dissects general intelligence into ten core cognitive domains...”</p>
<p>“a diverse set of challenging and realistic benchmark datasets to facilitate scalable, robust, and reproducible graph machine learning research.”</p> <p> <b>OPEN GRAPH BENCHMARK</b></p>	<p>“...a collection of tools for evaluating the performance of models across a diverse set of existing NLU tasks...”</p> <p>  <b>GLUE SuperGLUE</b></p>	<p>“It is impossible to enumerate the full set of tasks achievable by a sufficiently general intelligence. As such, an AGI benchmark should be a living benchmark”</p> <p><small>Google DeepMind Originally published Nov. 2023; updated Jan. 2024</small></p> <p>Levels of AGI: Operationalizing Progress on the Path to AGI</p>

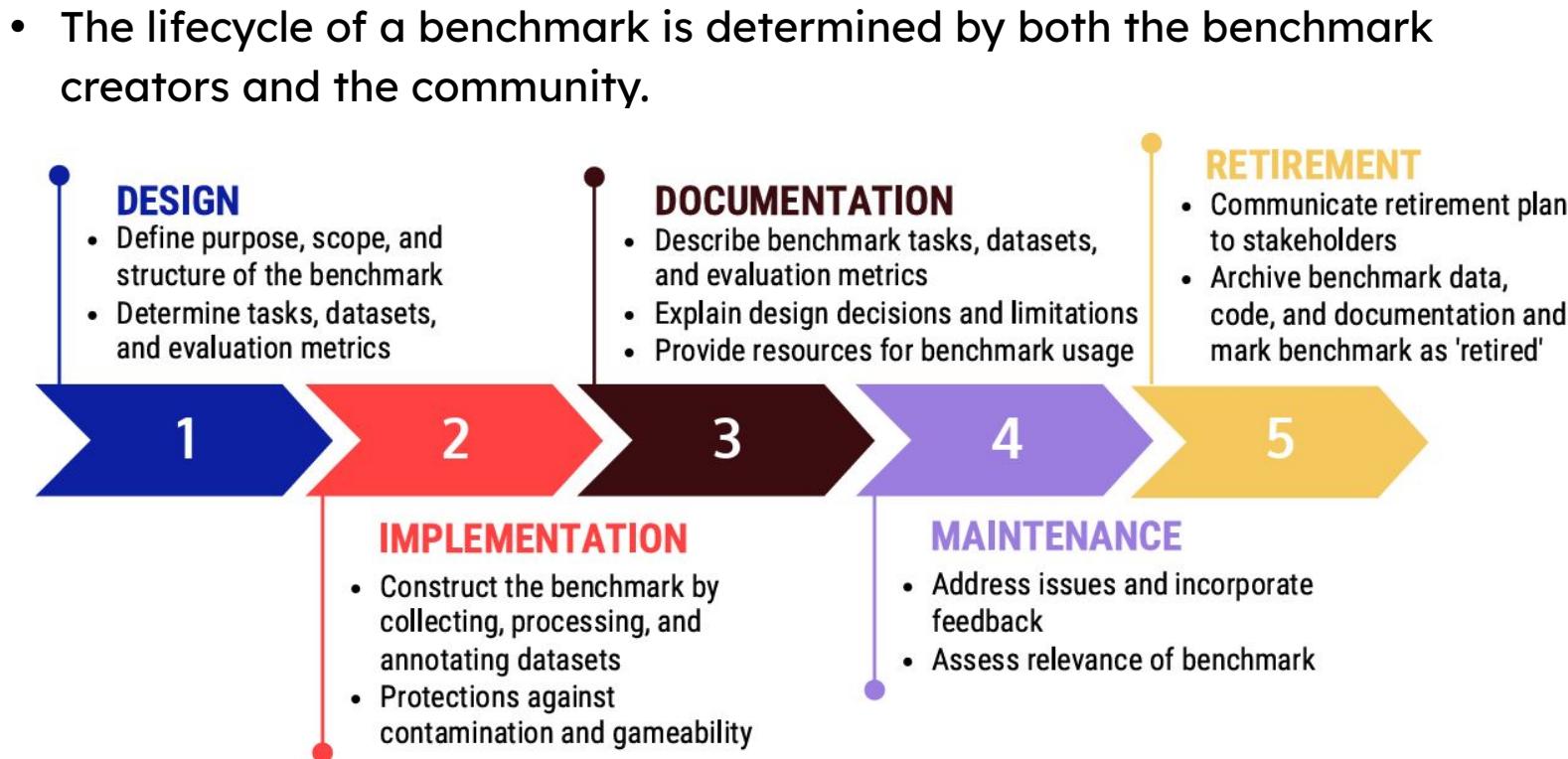
What's Measured?

What's Missed?

What's Next?

# Better Practice for Benchmark Creation

History  
Definition  
Preliminary  
Subject  
Objective  
Axiology  
Coverage  
Fluidity  
**Lifecycle**  
Design  
Retirement  
Outcomes



Anka Reuel, et al., BetterBench: Assessing AI Benchmarks, Uncovering Issues, and Establishing Best Practices. NeurIPS, 2024.  
Jialun Cao, et al., How Should We Build A Benchmark? Revisiting 274 Code-Related Benchmarks For LLMs. Preprint, 2025.

What's Measured?

What's Missed?

What's Next?

# Better Practice for Benchmark Creation

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

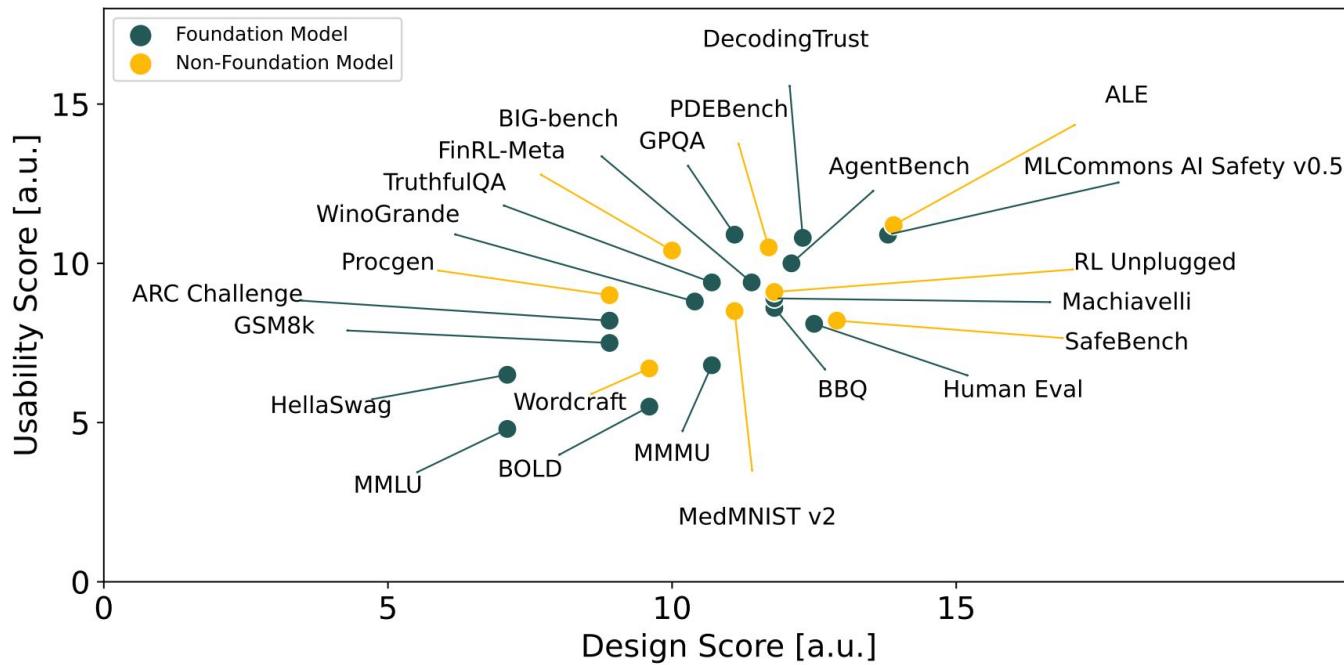
Lifecycle

Design

Retirement

Outcomes

- MLCommons AI Safety v0.5 (Vidgen et al., 2024) as an example.



Anka Reuel, et al., BetterBench: Assessing AI Benchmarks, Uncovering Issues, and Establishing Best Practices. NeurIPS, 2024.

Bertie Vidgen et al., Introducing v0.5 of the AI Safety Benchmark from MLCommons. Preprint, 2024.

What's Measured?

What's Missed?

What's Next?

# Benchmark Design

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Retirement

Outcomes

- **Conceptual grounding and domain relevance.**

- “*We created the taxonomy through an iterative process over 10 weeks...reviewed 25+ existing taxonomies, 50+ AI safety evaluation datasets, 50+ research and policy papers, and 10+ community guidelines from industry Trust and Safety orgs*”

- **Positioning and intended use.**

- “*The AI Safety Benchmark does not evaluate the safety of AI models ‘in general.’...Instead, the benchmark tests a specific AI system in a specific use case and for a specific set of personas.*”



# Benchmark Design

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Retirement

Outcomes

- **Metric and scale design.**

- Multi-stage grading: classifier labels responses → unsafe rates per test → 5-point risk grades → overall grade vs. a reference model.
- Explicit mapping from unsafe % to risk levels (“Low”, “Moderate-Low”, etc.), defining practical floors and ceilings.
- Humans to label a “human eval set” to validate, but no “human score on the benchmark” as a baseline system :(

- **Robustness and evaluation reliability.**

- Construct 43,090 prompts by combining sentence fragments with 13 interaction types, explicitly stating that they use variation to provide “*holistic coverage of interaction types*” and to test robustness.



# Benchmark Implementation

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Retirement

Outcomes

- **Accessible evaluation implementation.**

- Released ModelBench as “*an openly available platform, and downloadable tool ... to evaluate the safety of AI systems on the benchmark,*” with links to the GitHub repo; the test-spec schema for prompt generation and setup is also public.

- **Reproducible and statistically sound results.**

- Provided annotator agreement (Cohen’s kappa = 0.79) plus accuracy figures for the evaluator model
- No CIs or significance tests :(



# Benchmark Implementation

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Retirement

Outcomes

- **Benchmark integrity and data protection.**

- Proposed keeping parts of the dataset hidden or delayed, introduced benchmark deprecation rules, and required publishers not to train directly on the benchmark and to retracted results if contamination is discovered.

- **Safety and release readiness.**

- Added strong content warnings “[*t*]his work involves viewing content that creates a risk of harm and you might find objectionable or offensive,” provided wellbeing guidance for annotators, anonymize all tested models, and repeatedly stress that v0.5 is a preliminary proof-of-concept that “should not be used to assess the safety of AI systems,” with explicit release requirements for responsible use.



# Benchmark Documentation

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Retirement

Outcomes

- **Clear software & usage docs.**

- Open platform and downloadable tool on GitHub and provide a formal “*test specification*” schema;
- Great inline code comments, and code documentation;
- A “Trying It Out” for quick start of the code.

- **Transparent benchmark and evaluation design.**

- The white paper itself is the design doc.

- **Thorough dataset and metadata documentation.**

- Describe how 43,090 prompts are generated from sentence fragments and templates, gave detailed breakdown tables, and included a brief datasheet.



# Benchmark Maintenance

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Retirement

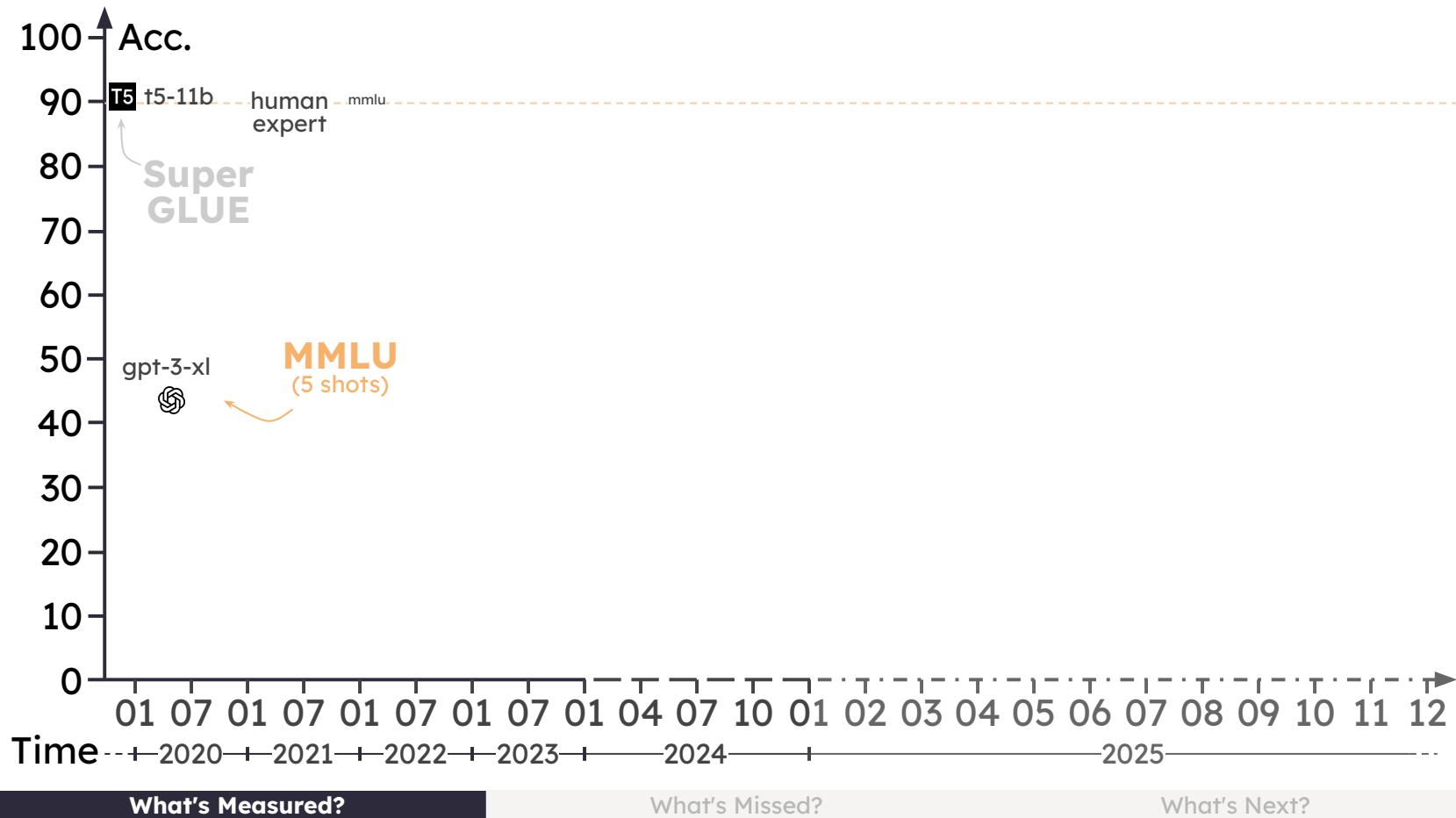
Outcomes

- **Code usability was checked within the last year.**
  - The dataset will be “regularly updated,” supported by MLCommons and HELM.
- **A feedback channel for users is available.**
  - Explicitly invited community feedback, linked the AI Safety WG page, and stated that MLCommons can be contacted via the website; anyone can also join the WG
- **Point of contact.**
  - “Reach out to Bertie if you have questions.”

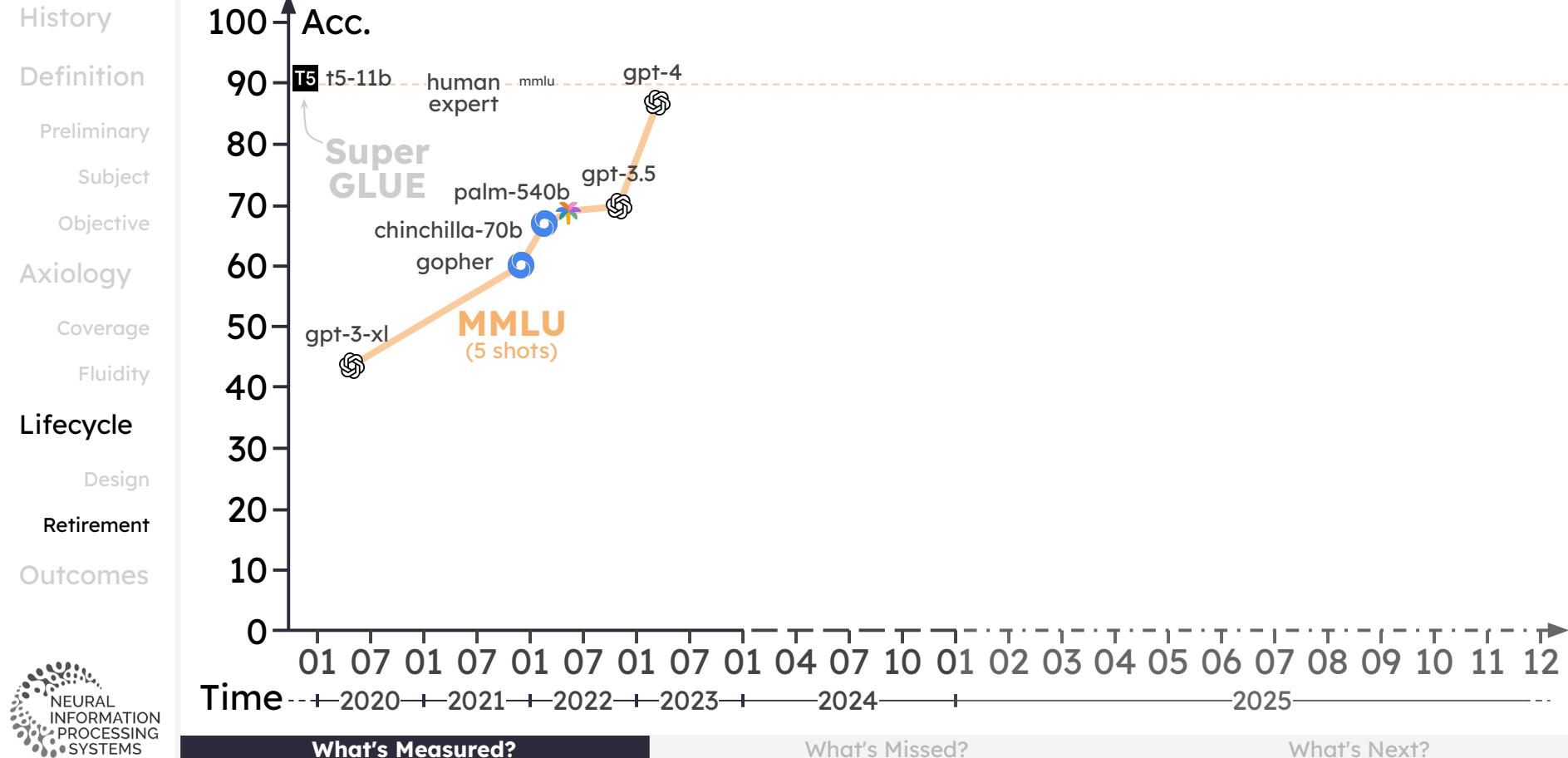


# Benchmark Retirement

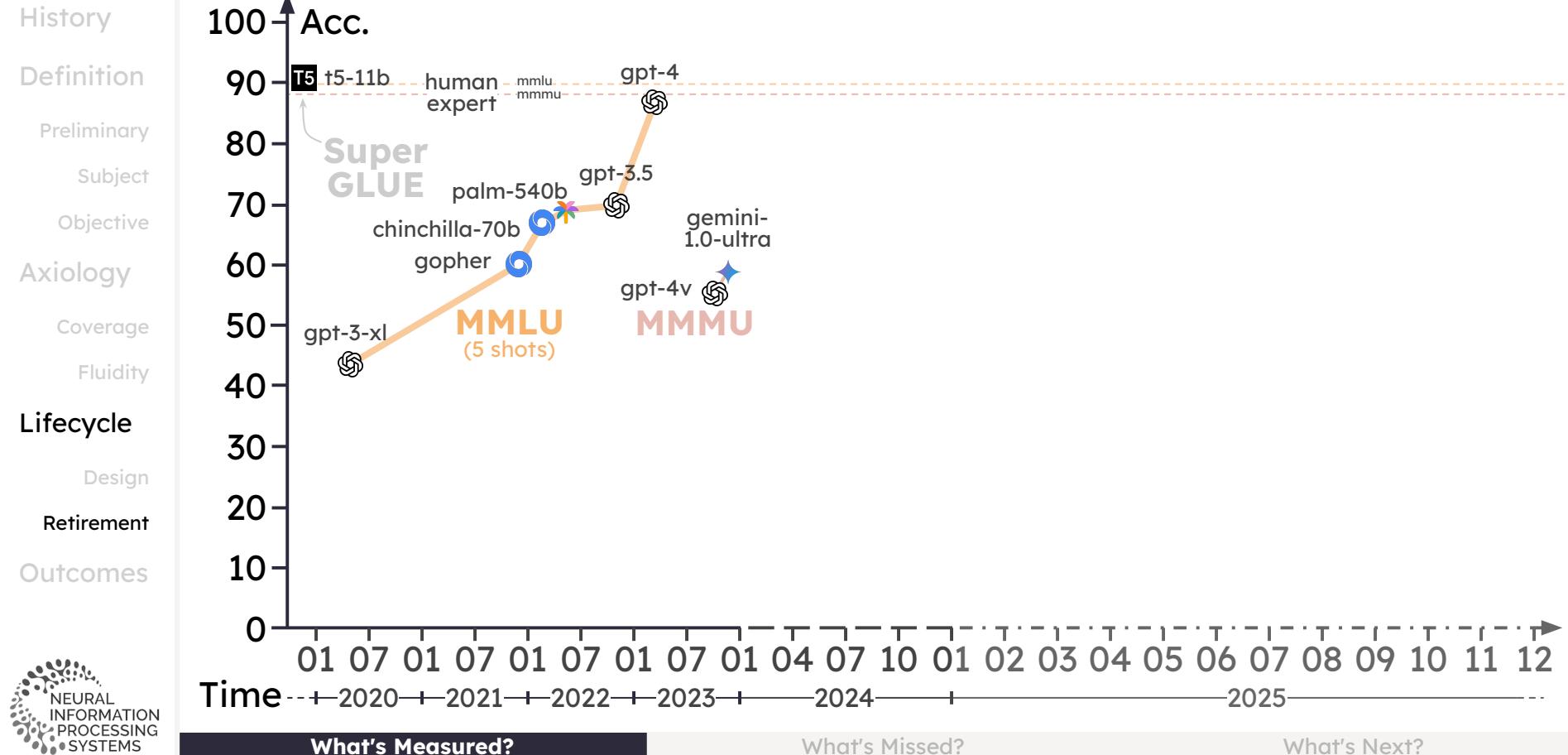
History  
Definition  
Preliminary  
Subject  
Objective  
Axiology  
Coverage  
Fluidity  
Lifecycle  
Design  
Retirement  
Outcomes



# Progress<sup>(?)</sup> on Benchmark Leaderboards

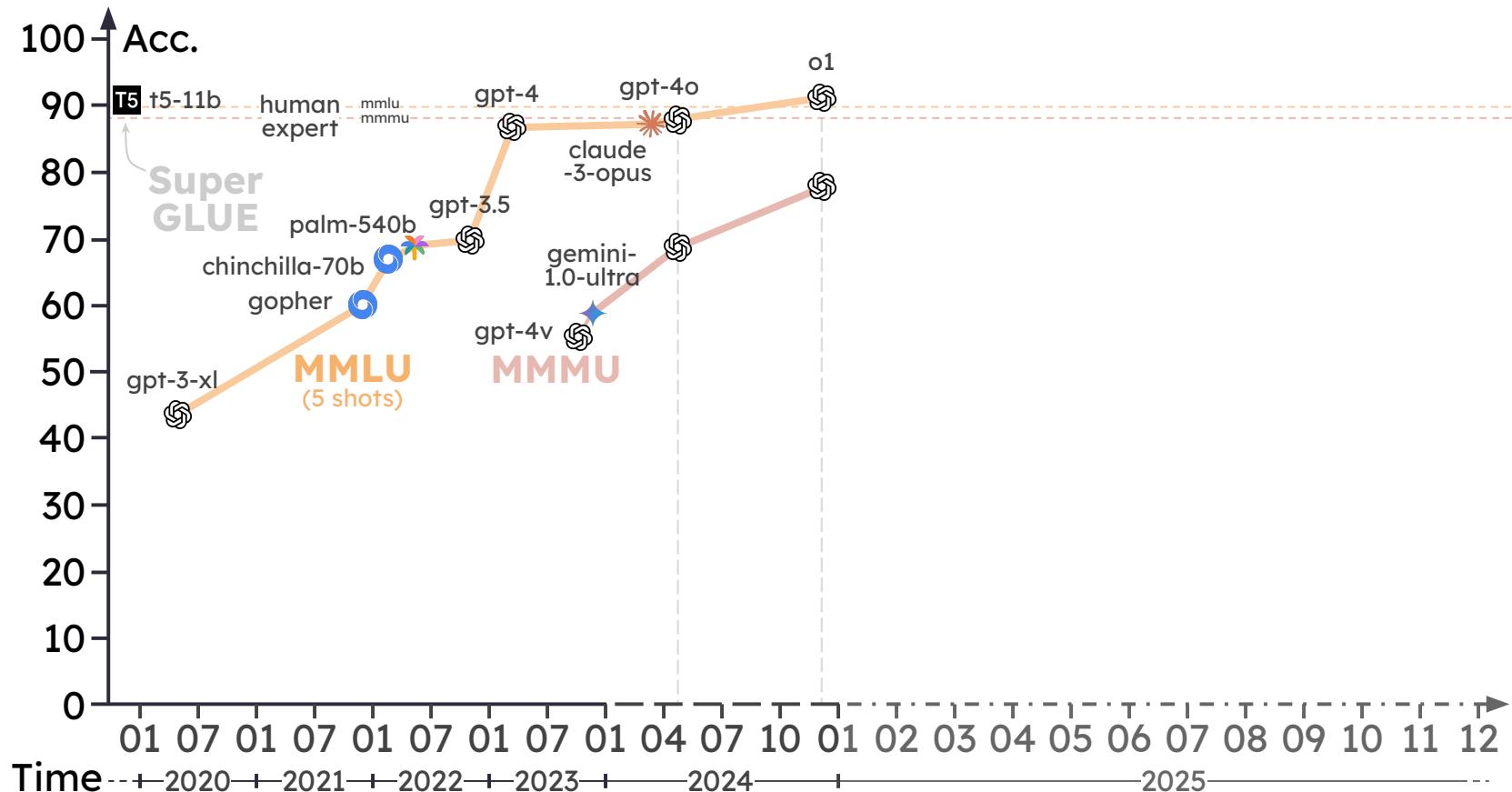


# Progress<sup>(?)</sup> on Benchmark Leaderboards



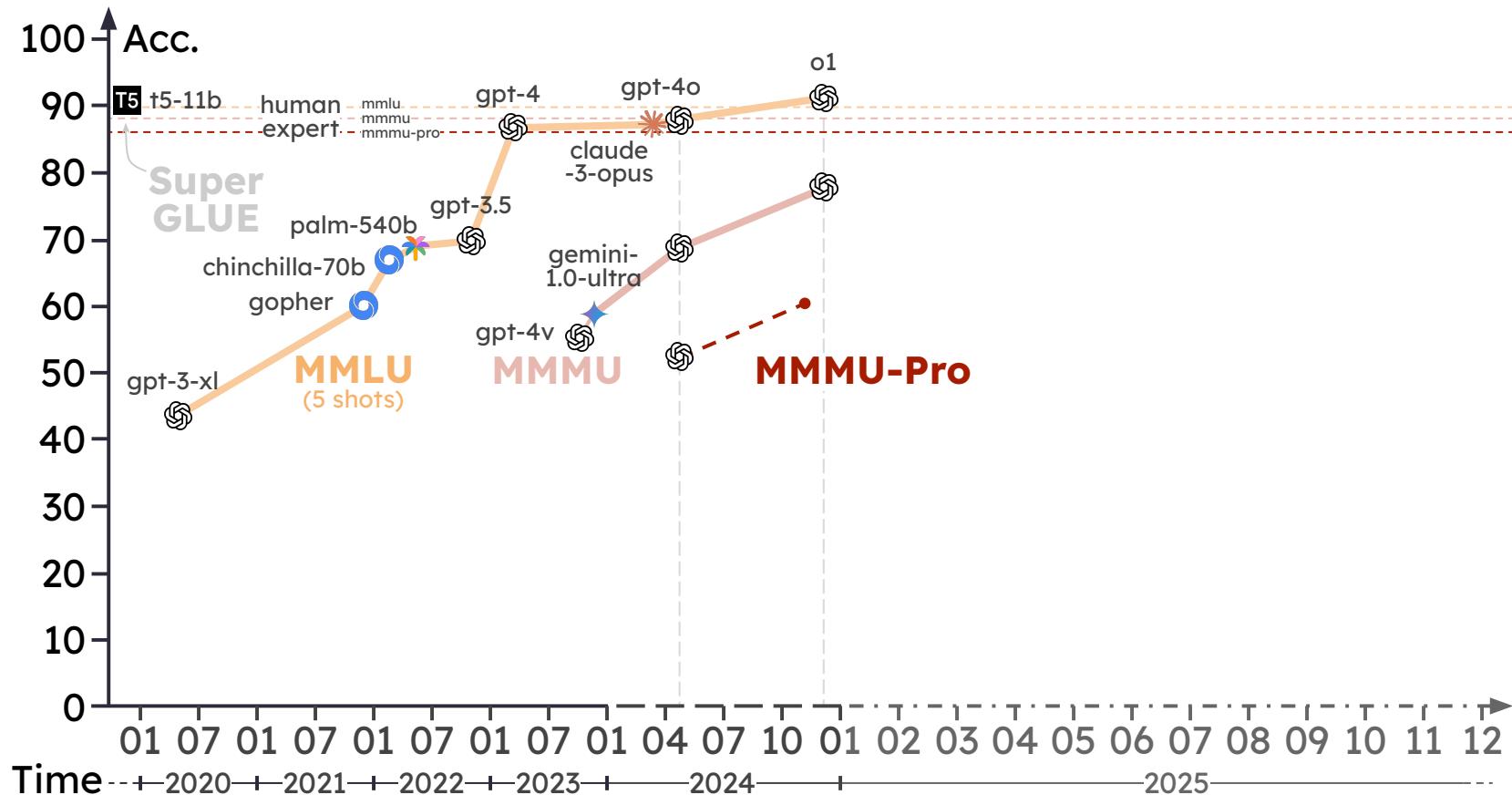
# Progress<sup>(?)</sup> on Benchmark Leaderboards

History  
Definition  
Preliminary  
Subject  
Objective  
Axiology  
Coverage  
Fluidity  
Lifecycle  
Design  
Retirement  
Outcomes

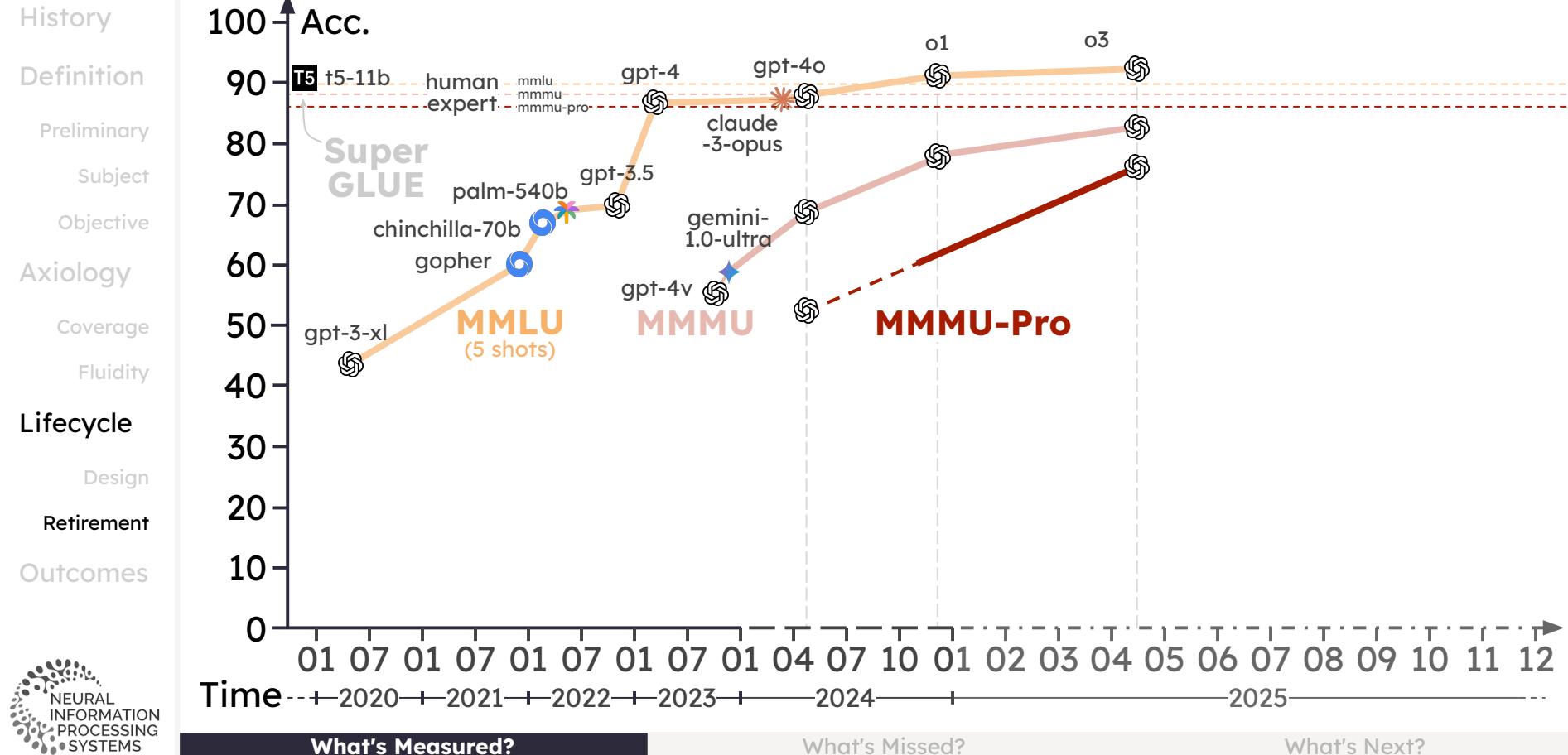


# Progress<sup>(?)</sup> on Benchmark Leaderboards

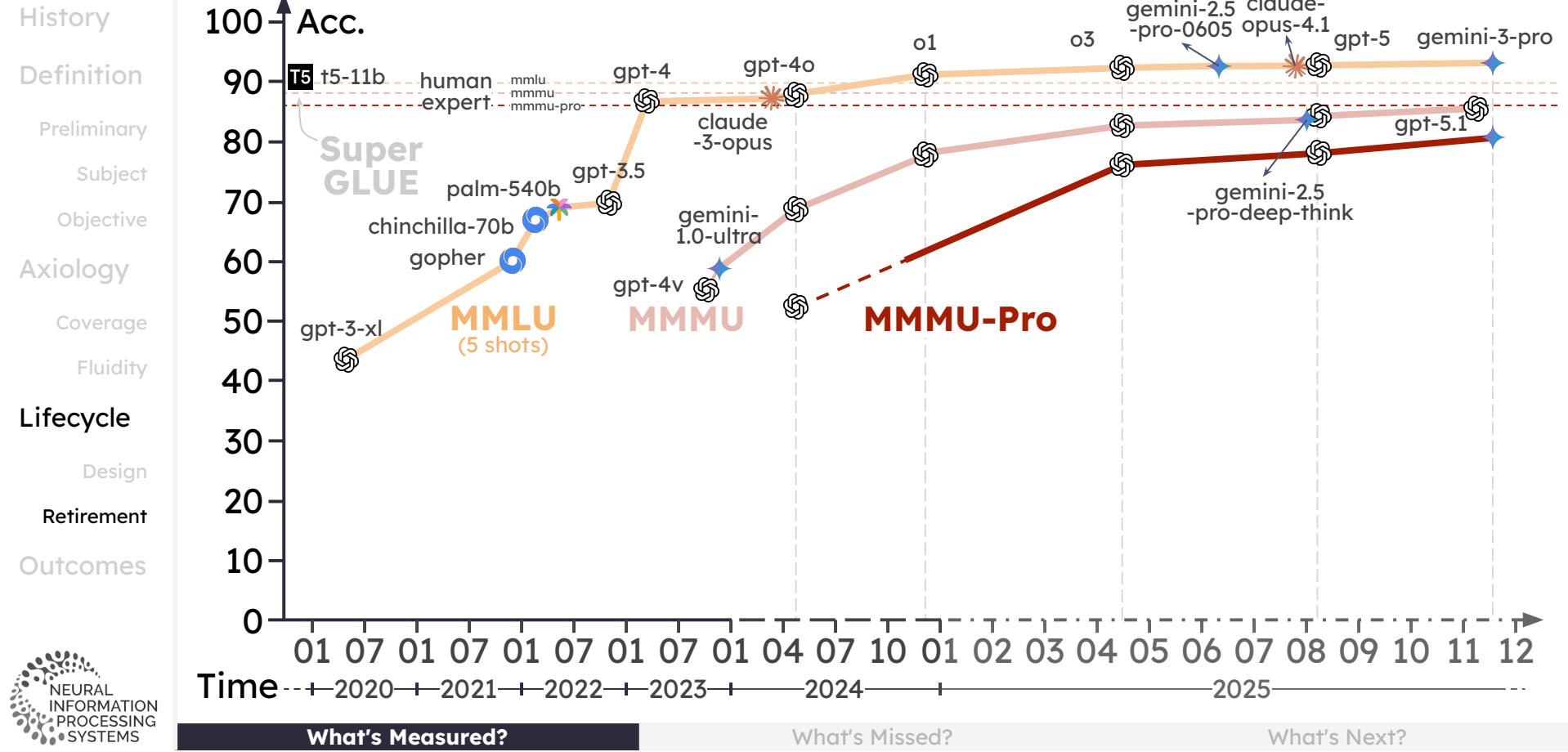
History  
Definition  
Preliminary  
Subject  
Objective  
Axiology  
Coverage  
Fluidity  
Lifecycle  
Design  
Retirement  
Outcomes



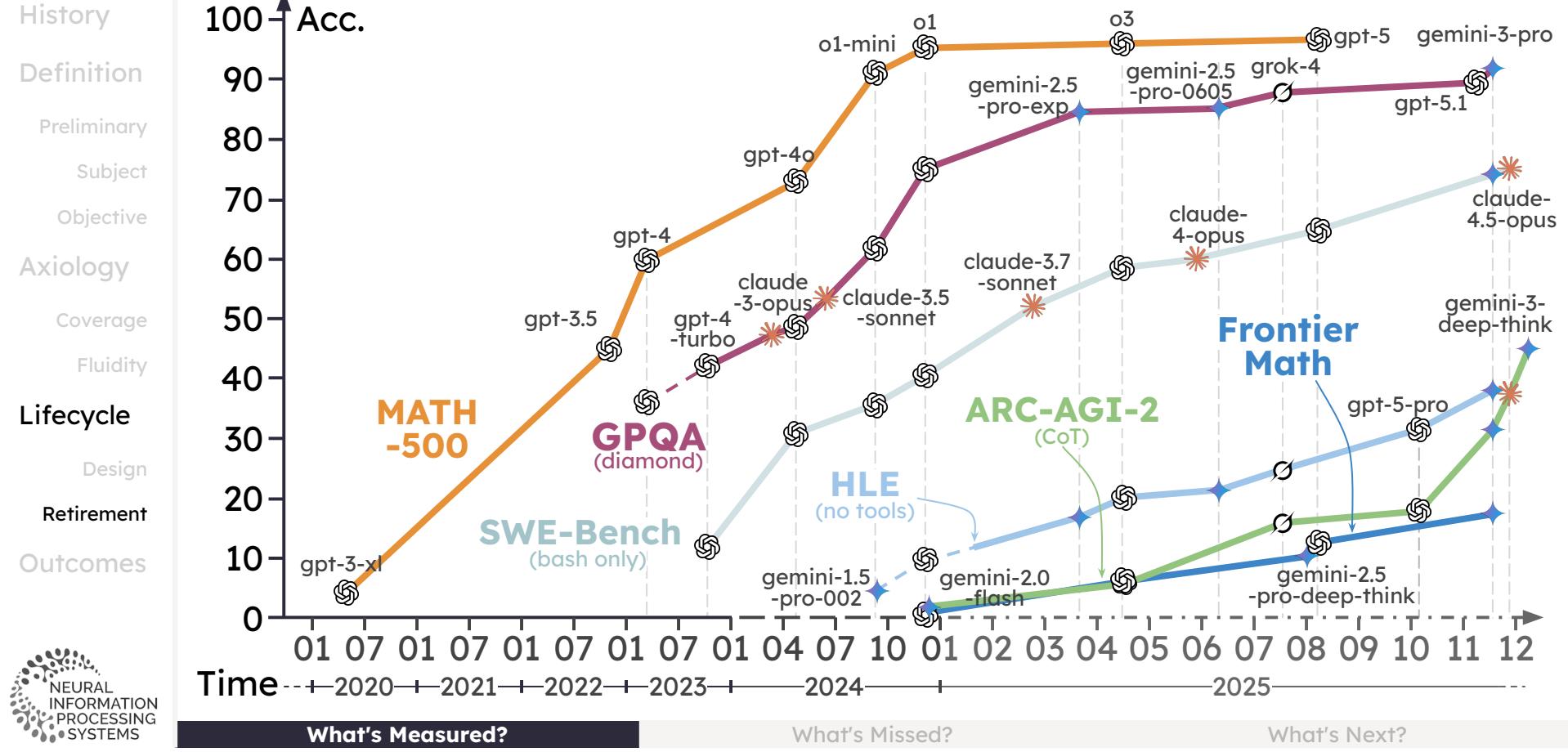
# Progress<sup>(?)</sup> on Benchmark Leaderboards



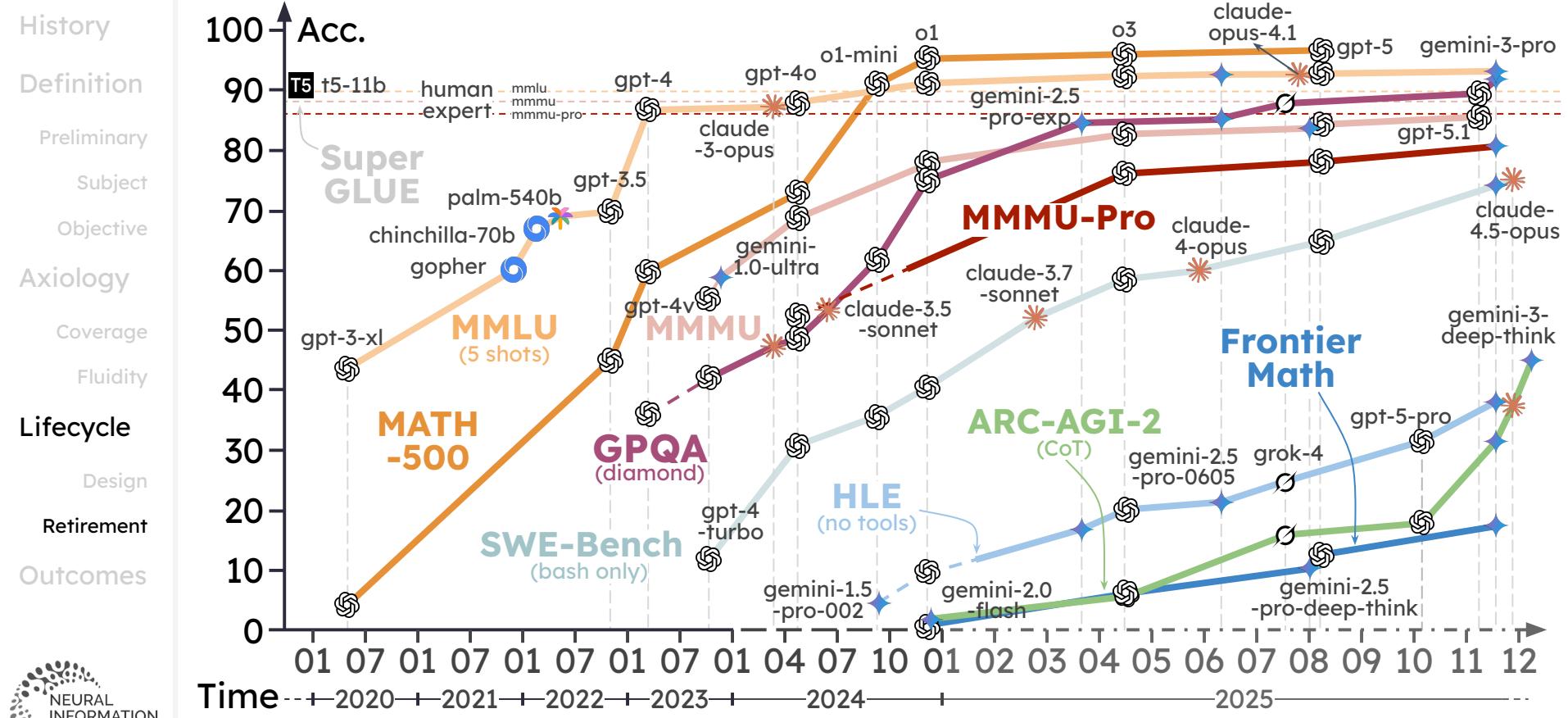
# Progress<sup>(?)</sup> on Benchmark Leaderboards



# Progress<sup>(?)</sup> on Benchmark Leaderboards



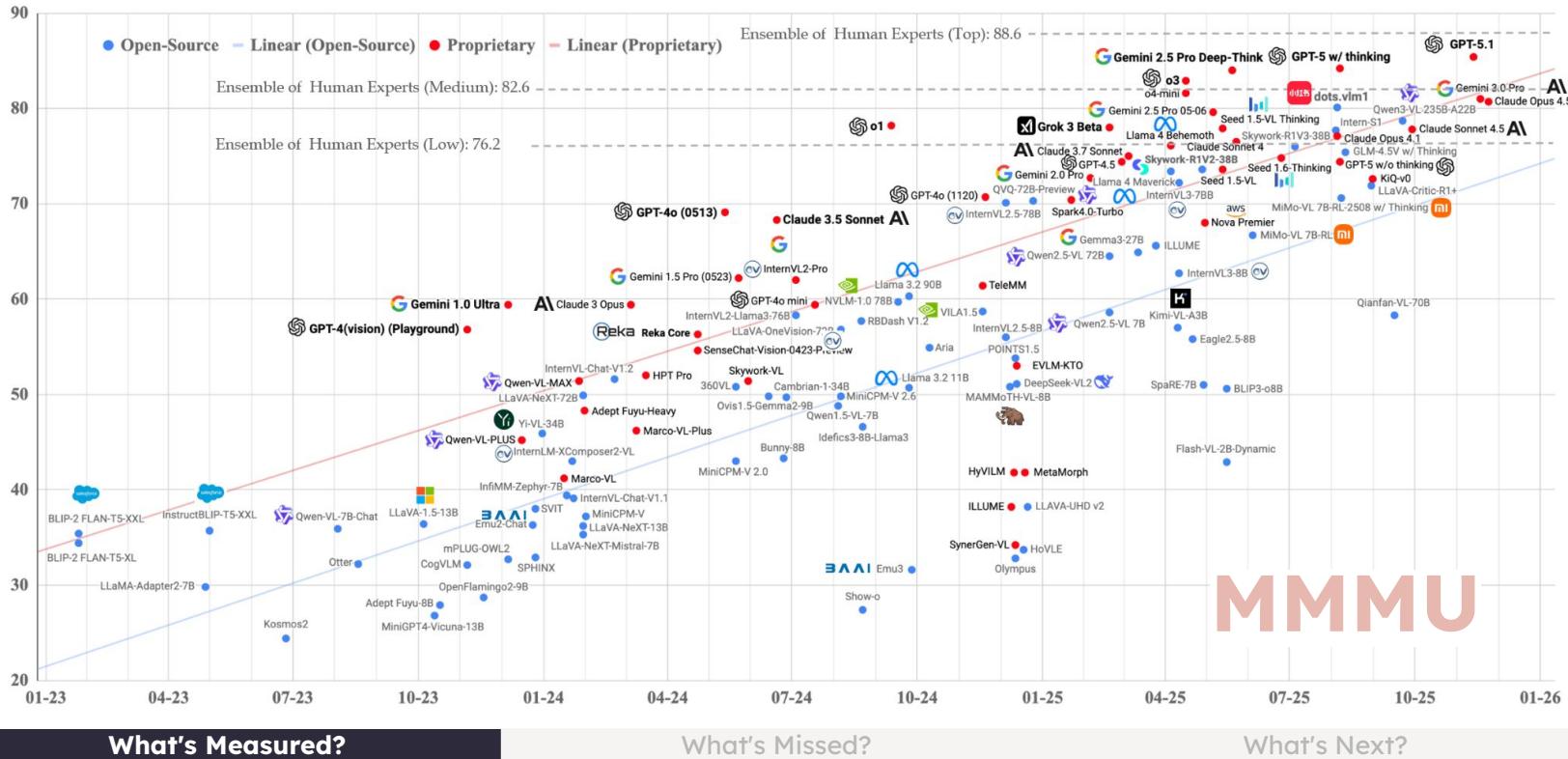
# Progress<sup>(?)</sup> on Benchmark Leaderboards



# Benchmark Retirement

History  
Definition  
Preliminary  
Subject  
Objective  
Axiology  
Coverage  
Fluidity  
**Lifecycle**  
Design  
Retirement  
Outcomes

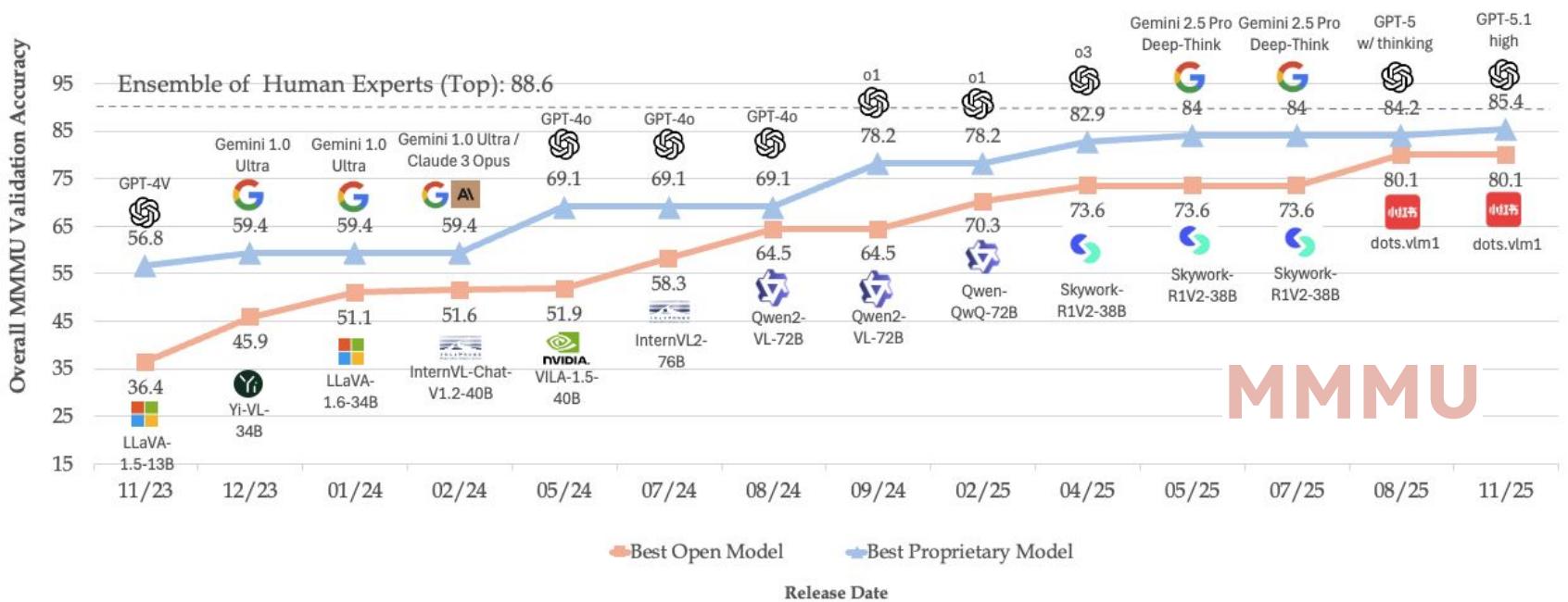
- Benchmarks that have retired from frontier lab competitions are still relevant to open-source communities.



# Benchmark Retirement

History  
Definition  
Preliminary  
Subject  
Objective  
Axiology  
Coverage  
Fluidity  
Lifecycle  
Design  
Retirement  
Outcomes

- Benchmarks that have retired from frontier lab competitions are still relevant to open-source communities.



# Interpreting Benchmark Outcomes

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Retirement

Outcomes

- Is this performance improvement genuine?
  - Case: classification.
  - “*Through extensive experiments, we show that our method B outperform previous SOTA by a large margin.*”

Model	Accuracy (↑)
A (Prev SOTA)	0.799
B (Ours)	<b>0.862</b>



# Interpreting Benchmark Outcomes

History  
Definition  
Preliminary  
Subject  
Objective  
Axiology  
Coverage  
Fluidity  
Lifecycle  
Design  
Retirement  
Outcomes

- Is this performance improvement genuine?
  - Case: classification.
  - “*Through extensive experiments, we show that our method B outperform previous SOTA by a large margin.*”
  - ...Probably not.

Model	Accuracy ( $\uparrow$ )	Std. Dev.
A (Prev SOTA)	0.799	$\pm 0.233$
B (Ours)	<b>0.862</b>	$\pm 0.666$



# Interpreting Benchmark Outcomes

History  
Definition  
Preliminary  
Subject  
Objective  
Axiology  
Coverage  
Fluidity  
Lifecycle  
Design  
Retirement  
Outcomes

- Is this performance improvement genuine?
  - Record per-seed per-example correctness;
  - Per-example: McNemar / sign test on disagreements:
    - McNemar test (per-example correctness):  $p < 1e-8$  (\*\*\*)
    - 95% CI for  $\Delta$  (accuracy, across seeds): [0.048, 0.078]
  - Across seeds: paired t-test or Wilcoxon test:
    - Paired t-test across 5 seeds: mean  $p = 5e-4$  (\*\*\*)

Model	Accuracy ( $\uparrow$ )	Std. Dev.
A (Prev SOTA)	0.799	$\pm 0.004$
B (Ours)	<b>0.862 (***)</b>	$\pm 0.012$



# Interpreting Benchmark Outcomes

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Retirement

Outcomes

- Is this performance improvement genuine?
  - Case: novel view synthesis (NVS).
  - “*Through extensive experiments, we show that our method B outperform previous SOTA by a large margin.*”

Model
A (Prev SOTA)
B (Ours)

PSNR (↑)
23.02
<b>24.10</b>



What's Measured?

What's Missed?

What's Next?

# Interpreting Benchmark Outcomes

- Is this performance improvement genuine?
  - Case: novel view synthesis (NVS).
  - “*Through extensive experiments, we show that our method B outperform previous SOTA by a large margin.*”
  - ...Probably not.

Model	Resolution	PSNR (↑)
A (Prev SOTA)	256 × 256	23.02
B (Ours)	512 × 512	<b>24.10</b>



# Interpreting Benchmark Outcomes

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

Lifecycle

Design

Retirement

Outcomes

- Is this performance improvement genuine?
  - $-\log(\text{average error}) \leq \text{average of } [-\log(\text{error per region})]$
  - Intuition: Higher resolution → more pixels → PSNR goes up when the underlying reconstruction quality has not meaningfully changed.

$$\text{PSNR} = 10 \log_{10} \left( \frac{\text{MAX}^2}{\text{MSE}} \right)$$

MAX: the maximum possible pixel value of the image given its encoding

$$\text{MSE} = \frac{1}{N} \sum_{i=1}^N e_i^2$$



# Reality Check on Benchmark Outcomes

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

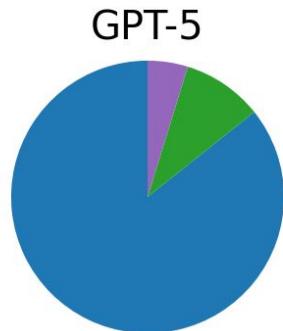
Lifecycle

Design

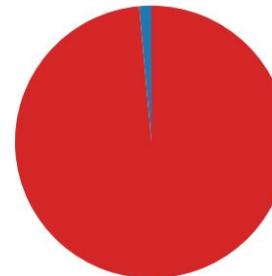
Retirement

Outcomes

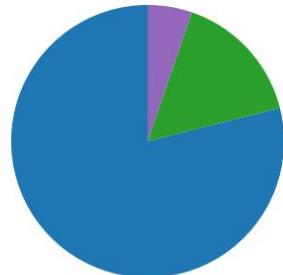
- What happened in the unsolved portion?



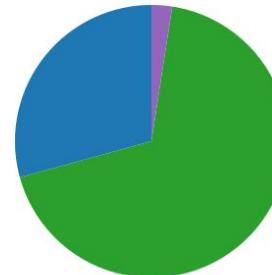
Gemini-2.5-pro



Gemini-3-pro



Claude 4.5 sonnet



## MATH-500

- Failed Reasoning
- Lack of Knowledge
- Fail Instruction Following
- Error of Evaluator
- Others (e.g., Data Quality)



What's Measured?

What's Missed?

What's Next?

\*Our original preliminary results.

# Reality Check on Benchmark Outcomes

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

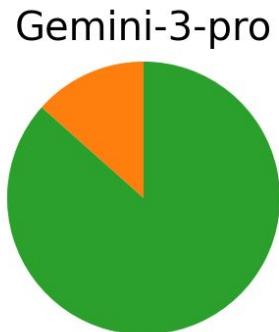
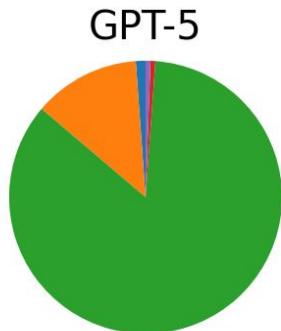
Lifecycle

Design

Retirement

Outcomes

- What happened in the unsolved portion?



## GPQA

- Failed Reasoning
- Lack of Knowledge
- Fail Instruction Following
- Error of Evaluator
- Others (e.g., Data Quality)

\*Our original preliminary results.

What's Measured?

What's Missed?

What's Next?

# Reality Check on Benchmark Outcomes

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

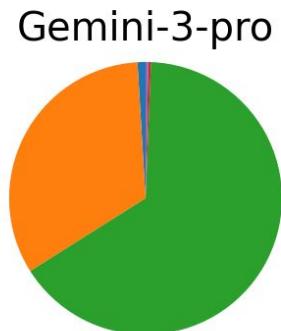
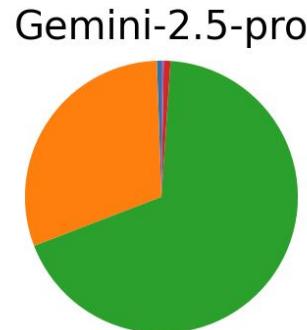
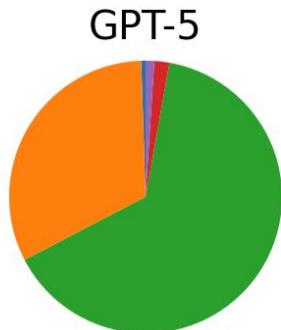
Lifecycle

Design

Retirement

Outcomes

- What happened in the unsolved portion?



## MMLU-Pro

- Failed Reasoning
- Lack of Knowledge
- Fail Instruction Following
- Error of Evaluator
- Others (e.g., Data Quality)

What's Measured?

What's Missed?

What's Next?

\*Our original preliminary results.

# Reality Check on Benchmark Outcomes

History

Definition

Preliminary

Subject

Objective

Axiology

Coverage

Fluidity

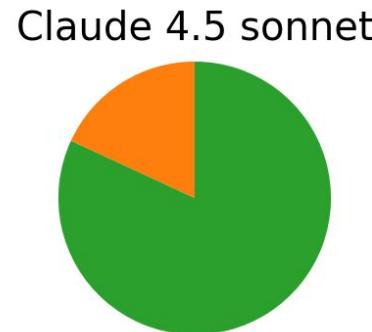
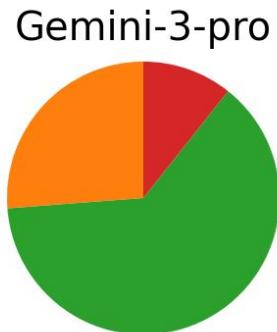
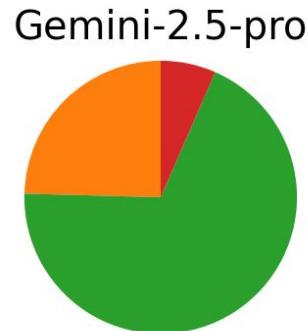
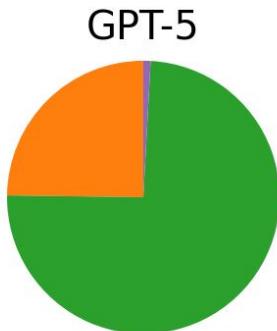
Lifecycle

Design

Retirement

Outcomes

- What happened in the unsolved portion?



## MMMU-Pro

- Failed Reasoning
- Lack of Knowledge
- Fail Instruction Following
- Error of Evaluator
- Others (e.g., Data Quality)

\*Our original preliminary results.

What's Measured?

What's Missed?

What's Next?

# What's Missed?



Michael Saxon  
University of Washington

# Agenda

- What's Measured? (1:30PM - 2:10PM)
- What's Missed? (2:10PM - 2:40PM)
  - Practical issues: data, integrity, measurement problems
  - Deeper issues: Systemic and epistemic problems
- What's Next? (2:40PM - 3:15PM)
- Panel Discussion (3:20PM-4:00PM)

# What's missed?

Intro

Data issues

Measuremen  
t issues

Systemic  
issues

Epistemic  
issues

Concluding

- There are many criticisms of modern benchmarking practices
- **They are easy to miss when designing a new one**
- Here we borrow from the following works at a high level:
- *Can We Trust AI Benchmarks? An Interdisciplinary Review of Current Issues in AI Evaluation*, Eriksson et al. (2025)
- *Benchmarks as Microscopes: Toward a Science of Model Metrology*, Saxon et al. (2024)
- *AI and the Everything in the Whole Wide World Benchmark*, Raji et al. (2022)



# What's missed?

Intro

Data issues

Measuremen  
t issues

Systemic  
issues

Epistemic  
issues

Concluding

- Many fundamental challenges and critiques of benchmarking have been levied, but are often **missed** by practitioners. Including:
  - **Data issues**
  - **Measurement issues**
  - **Systemic issues**
  - **Epistemic issues**
  - We will discuss both **the critiques** themselves and **some solutions** that others may want to apply.



# Validity issues

## Intro

Data issues

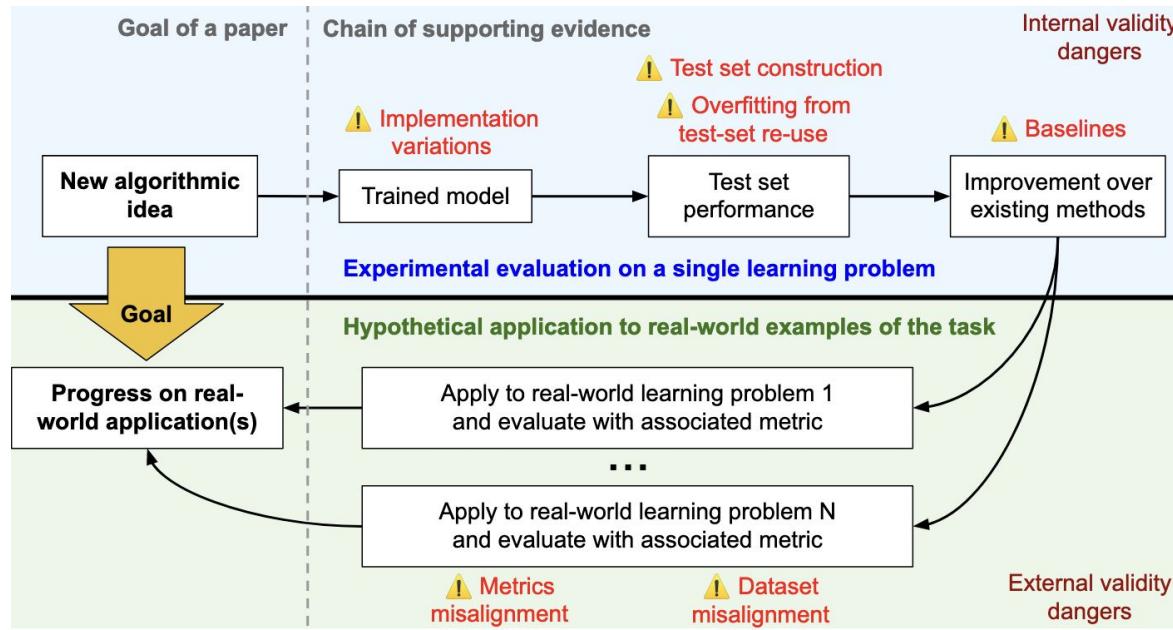
Measuremen  
t issues

Systemic  
issues

Epistemic  
issues

Concluding

- **Internal validity issues:** arise for a single benchmark/resource
- **External validity issues:** apply to the world of the task



# Construct validity

Intro

Data issues

Measurement issues

Systemic issues

Epistemic issues

Concluding

- **Construct validity** may be the most important property of a good evaluation resource (Raji et. al, 2021)
- Definition:

*The central question of external validity: how well does a resource measure what it purports to measure?*
- Unfortunately, this is more of an abstract **goal** than a concrete requirement
- Many of the issues with benchmarks can be traced back to construct validity!



# Data issues

## Data issues

Lifecycle

Noise

Positionality

Contamination

## Measurement issues

## Systemic issues

## Epistemic issues

## Concluding

- **Data issues** are challenges to the validity of benchmarks that arise at creation time.
- **Lifecycle problems**
  - “Where” and “when” data is created
- **Noise & spurious correlations**
  - When bad examples harm the ability of the benchmark to discriminate
- **Resource creator positionality & construct validity**
  - “Who” created data for “what” task?
- **Contamination**
  - What happens when examples from the benchmark are trained on



# Static datasets vs. dynamic realities

Data issues

Lifecycle

Noise

Positionality

Contamination

Measurement issues

Systemic issues

Epistemic issues

Concluding

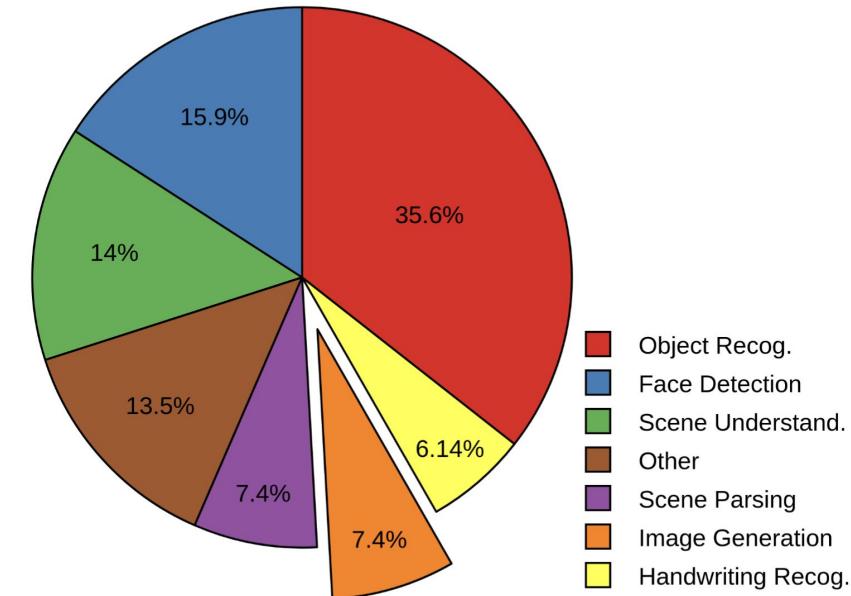
- Most dominant benchmarks are over **static** sets of exemplars collected once, held as test set forever.
- Static dataset lifecycle ends in a way at **release**
- However, reality is **dynamic**
- Static multiple choice question (MCQ) datasets “fail to reflect the evolving nature of human-AI interactions” (McIntosh et al.)
- Static datasets have no way to account for future improvements in model capabilities (Saxon et al, 2024)



# Lifecycle challenges

Data issues  
Lifecycle  
Noise  
Positionality  
Contamination  
Measurement issues  
Systemic issues  
Epistemic issues  
Concluding

- Producing *high-quality* data is costly and time-consuming
- Consequently we often rely on “reduced, reused, recycled” data
- 70% of computer vision datasets reuse data from other **domains** (Koch et. al 2021)



Koch et. al (2021): Original data source for *image generation* datasets, by task (only 7.4% were built for the purpose of image generation)

# Reduced, reused, recycled data

Data issues

Lifecycle

Noise

Positionality

Contamination

Measurement issues

Systemic issues

Epistemic issues

Concluding

- Issues in recycled data propagate
- Egregious examples of bad data have survived even in huge datasets
- We have improperly labeled images, such as this one from Imagenet
- But crowdsourced datasets like MS COCO have many more egregious labels and masks



# A real example from MS COCO

Data issues

Lifecycle

Noise

Positionality

Contamination

Measurement issues

Systemic issues

Epistemic issues

Concluding



Search datasets

Book demo

Select organization  
Activeloop



Index  
out of  
118,286



Datasets



Docs



What's Measured?

What's Missed?

What's Next?

# Noise & spurious correlations

Data issues

Lifecycle

Noise

Positionality

Contamination

Measurement issues

Systemic issues

Epistemic issues

Concluding

- Bad labels are present all over.
- So what?
- Sometimes these present **artifacts**: systematic spurious correlations which models can use to cheat on tasks such as:
  - Medical scan classification (Oakden-Rainer et. al, 2019)
  - Entailment recognition (McCoy et. al, 2019)
- Bad labels also produce **noise**, where an LM's performance is effectively random



# Label noise

## Data issues

Lifecycle

Noise

Positionality

Contamination

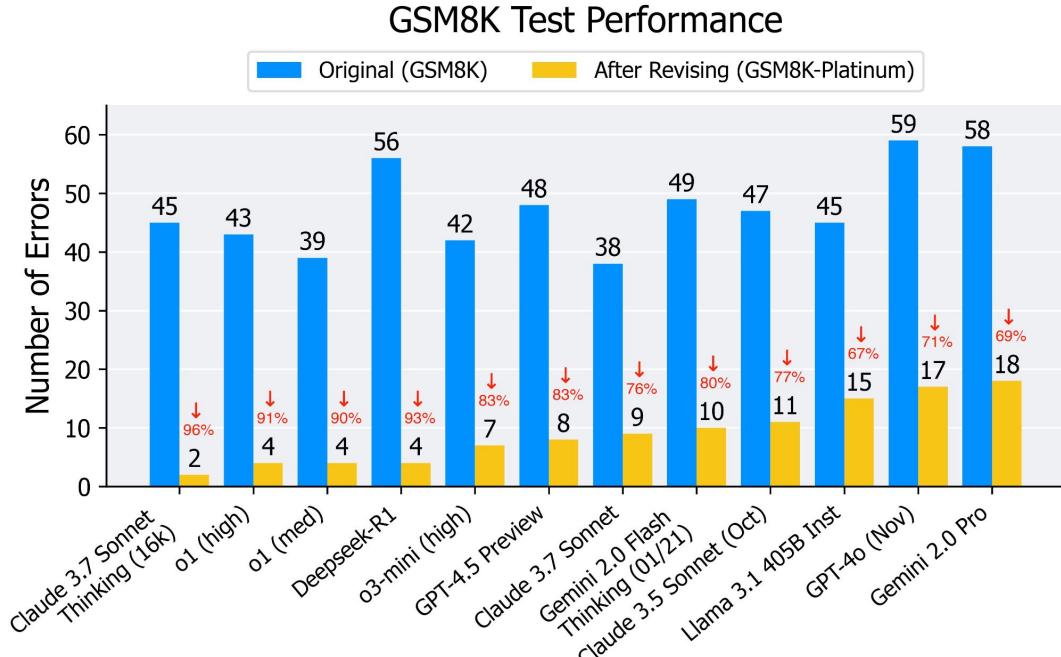
Measurement issues

Systemic issues

Epistemic issues

Concluding

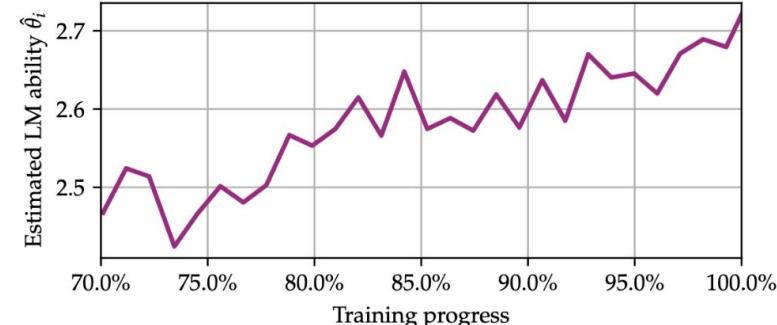
- Vendrow et. al (2025) remove bad examples from GSM8K
- After removing the noise ceiling (making 100% acc possible), the ranking swaps
- Noise isn't just a ceiling—it also sets the resolution



# One way to remove noise: Fluid benchmarking

Data issues  
Lifecycle  
**Noise**  
Positionality  
Contamination  
Measurement issues  
Systemic issues  
Epistemic issues  
Concluding

- **Item response theory (IRT)**
- **Capability as a latent property** which we infer by not just # correct answers but **which difficulty level** of questions a test-taker answers right.
- Each question has a “difficulty” and “discriminability”
- Hoffman et al (2025) train an IRT model on benchmarks, use it to **select the next test question given the current capability estimate**
- Noise samples naturally have low discriminability and are ignored
- **Datasets become higher resolution** (monotonic with train)



# Positionality: subjectivity in annotation

## Data issues

Lifecycle

Noise

## Positionality

Contamination

## Measurement issues

## Systemic issues

## Epistemic issues

## Concluding

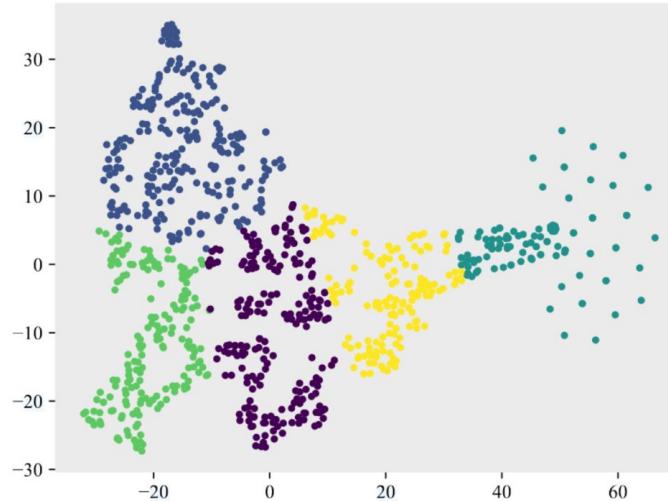
- Many tasks we want to evaluate are **inherently subjective**
  - Humor
  - Hate speech
  - Safety (in text generation)
- Others may have **subjective attributes** at times
  - Acceptability of phrases (is this grammatical?)
    - Not all native speakers will agree on the grammaticality of all phrases
  - Entailment of statements (NLI)
    - Annotators may disagree on how well a sentence entails another
- Typically, the way we deal with this subjectivity is **average or majority vote**
- **Is this ideal?**



# Dealing with subjectivity: annotator embeddings

Data issues  
Lifecycle  
Noise  
**Positionality**  
Contamination  
Measurement issues  
Systemic issues  
Epistemic issues  
Concluding

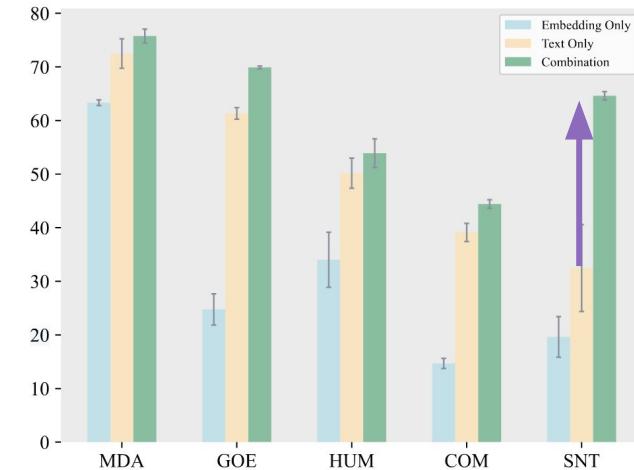
- Deng et al (2023) introduce *annotator embeddings* to handle **inherent annotator disagreement** across humor detection, hate speech detection, and NLI
- They **keep the disparate annotations** and frame the learning task as:
- Predict label given input *and latent annotator descriptor*
- “annotation embeddings” give up **2x boost** to CLS performance over text alone



What's Measured?

What's Missed?

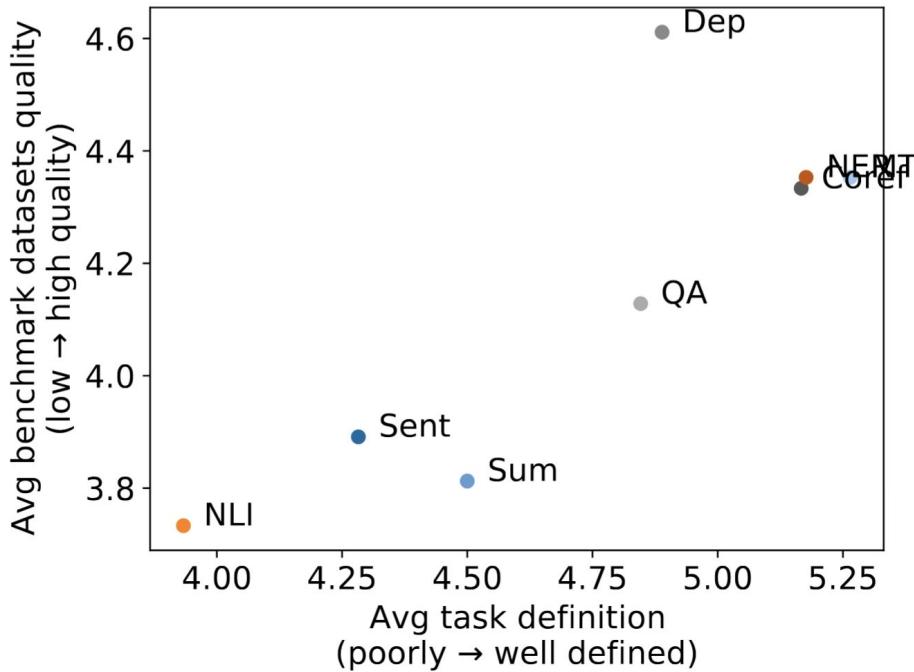
What's Next?



# Positionality: definitions matter!

Data issues  
Lifecycle  
Noise  
**Positionality**  
Contamination  
Measurement issues  
Systemic issues  
Epistemic issues  
Concluding

- Positionality issues also include **how developers understand a task**
- From Subramonian et al. (2023):
- Same “task” often operationalized differently in different works (usually poorly defined too)
- Direct relationship between poor task definition and low dataset quality.
- **Poor conceptualization leads to poor data collection**



# Positionality: values embedded in data source

## Data issues

Lifecycle

Noise

**Positionality**

Contamination

Measurement issues

Systemic issues

Epistemic issues

Concluding

- Many influential LM datasets source test questions from web forums
- Examples: HellaSwag (Zellers et al, 2019), ETHICS (Hendrycks et al, 2021)
- These in turn get “recycled” into bigger resources like MMLU
- Many Reddit AITA dilemmas cover basic interpersonal drama
- Are these the bounds by which ethics in AI systems should be assessed?



# Positionality: values embedded in culture

Data issues  
Lifecycle  
Noise  
**Positionality**  
Contamination  
Measurement issues  
Systemic issues  
Epistemic issues  
Concluding

- **Conceptualization itself is culturally embedded!**
- Concepts like “stereotypes” and “offensive language” are contested and culturally contingent (Blodgett et al. 2021)
- Even **annotations themselves** can be culturally contingent (Oh et al. 2025)
- For example, in one study, East Asian annotators consistently preferred lower-valence (mid) answers on a Likert scale to Americans (Lee et al. 2002)
- These influences shape dataset curation!



# Positionality: values embedded in community

Data issues

Lifecycle

Noise

**Positionality**

Contamination

Measurement issues

Systemic issues

Epistemic issues

Concluding

- Benchmarks are **normative instruments** within the AI community.
  - Eg, investment be in safety benchmarks has successfully motivated work on alignment and safety from EA/x-risk perspectives
- Even as we look to build systems that can do everything, most of us (AI researchers) are **not experts** in most things
- Consequently, many domain-specific eval resources produced by computer scientists are not useful for experts in those domains
- Blagec et al. (2023) surveyed medical practitioners about clinical LM benchmarks: **they fail to capture how LMs meet doctor's needs**
- Benchmarks often **abstract tasks out of their social context** (Selbst 2019)
- Data work is undervalued relative to “model work” (Raji et al., 2021)



# Contamination

Data issues  
Lifecycle  
Noise  
Positionality  
**Contamination**  
Measurement issues  
Systemic issues  
Epistemic issues  
Concluding

- Benchmarks **contaminate** a model when their examples are present in a model's training data
- Fundamental challenge to the utility of an eval
- leads to poor **predictiveness** and **validity**
- GPT-4 performed perfect on pre-knowledge cutoff codeforces easy problems, and 0/10 on post-cutoff
- How can we address it?

Horace He    
@cHHillee

I suspect GPT-4's performance is influenced by data contamination, at least on Codeforces.

Of the easiest problems on Codeforces, it solved 10/10 pre-2021 problems and 0/10 recent problems.

This strongly points to contamination.

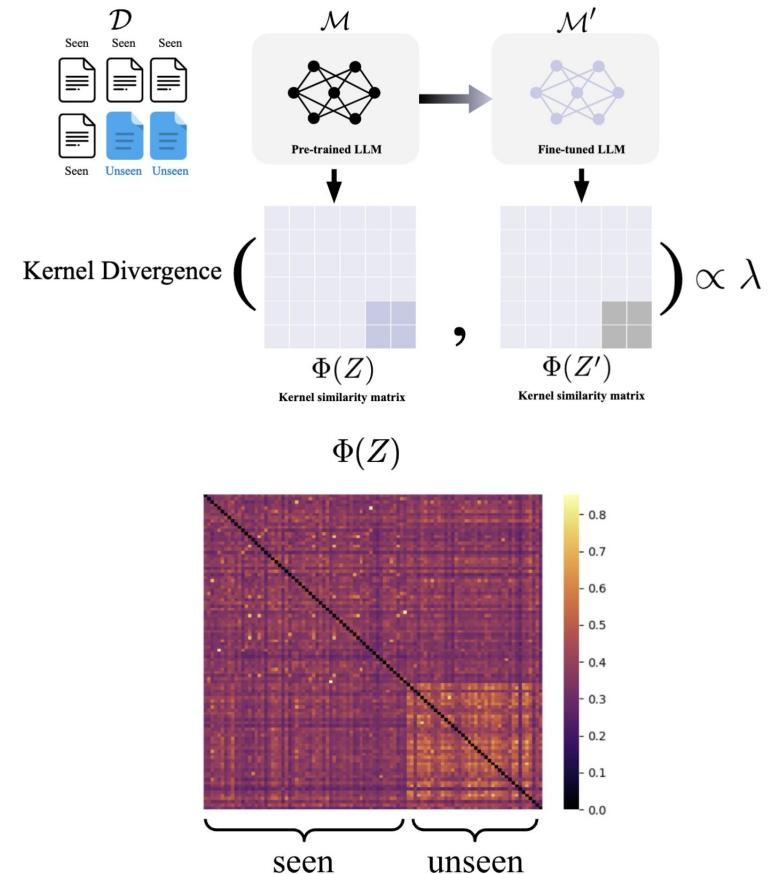
1/4

<a href="#">g's Race</a>	implementation, math	 	greedy, implementation	 	
<a href="#">nd Chocolate</a>	implementation, math	 	Cat?	implementation, strings	 
<a href="#">triangle!</a>	brute force, geometry, math	 	Actions	data structures, greedy, implementation, math	 
	greedy, implementation, math	 	Interview Problem	brute force, implementation, strings	 
<a href="#">umbers</a>	brute force	 	vers	brute force, implementation, strings	 
<a href="#">ine Line</a>	implementation	 	nd Suffix Array	strings	 
<a href="#">r or Stairs?</a>	implementation	 	ther Promotion	greedy, math	 
<a href="#">Loves 3 I</a>	math	 	Forces	greedy, sortings	 
<a href="#">s</a>	implementation, math	 	l and Append	implementation, two pointers	 
	greedy, implementation, sortings	 	ng Directions	geometry, implementation	 

# “How contaminated is your benchmark?”

Data issues  
Lifecycle  
Noise  
Positionality  
**Contamination**  
Measurement issues  
Systemic issues  
Epistemic issues  
Concluding

- Choi et al (2025) introduce *kernel divergence score* as a proposed white-box method to test for contamination
- Check the divergence of the RBF kernel of the embeddings of some test documents before and after fine-tuning the model on them
- If the embeddings don't diverge, it likely has already seen them



# What's in my big data? (WIMBD)

Data issues  
Lifecycle  
Noise  
Positionality  
**Contamination**  
Measurement issues  
Systemic issues  
Epistemic issues  
Concluding

- Directly search for observed model outputs in the training data (Elazar et al 2023)
- It is possible that many emergent capabilities in LLMs are picked up in the training data
- “Let’s think step-by-step” site shows up 10s of thousands of times in C4
- **This is a broader philosophical challenge for understanding generalization in the LLM era**

The post includes a profile picture of Yanai Elazar, the handle @yanaela, and the text "Let's think "step by step"! Another tidbit I like about data and prompts that miraculously work. Searching for this phrase resulted in this website (among others), [geteasysolution.com](https://geteasysolution.com), containing many math step-by-step solutions. How common are they? Quite. Makes you think." Below the text is a screenshot of a web page from geteasysolution.com showing a math problem and a table of statistics.

**Screenshot of the website:**

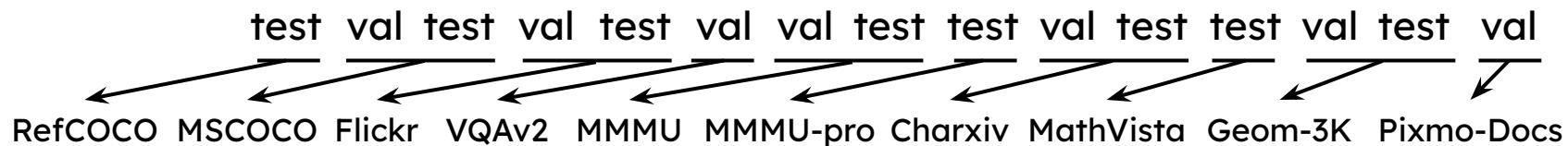
A math problem: "A number is 1000 more than 2728 - solution with step by step". Below it: "Find the full step by step solution for your problem. We hope it will be very helpful for you to understand the solving process." A form: "What is  more  than ? Solve". Below the form: "What is 1000 more than 2728?". Below the question: "Solving for a new number which is 1000 more than 2728. Find the new number by adding 1000 to 2728. down as: =3728". Below the equation: "the solution for: What number is 1000 more than 2728?".

**Table of statistics:**

	Corpus	Rank	Tokens	%
tion.com	Dolma	151022	3,549	██████████
tion.com	OSCAR	277233	224,965	██████████
tion.com	C4	473082	49,859	██████████
tion.com	RedPajama	472159	49,859	██████████
tion.com	mC4-en	1658921	156,174	██████████

# Image Contamination

Data issues		Train Sets															
		Lifecycle	Noise	Positionality	Contamination	Measurement issues	Systemic issues	Epistemic issues	Concluding	MSCOCO	LAION-2B	Flickr-30K	Pixmo-Docs	Geom-3K	ArXivQA		
		100.0	2.7	0.1	0.2	0.5	87.7	0.0	0.0	0.0	0.0	0.0	0.0	2.6	0.0	0.0	0.0
		-	2.4	2.3	1.7	1.7	2.4	10.2	7.6	11.7	0.2	0.4	7.2	0.0	0.0	0.0	
		-	0.1	0.1	0.5	0.9	0.1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
		-	0.0	0.0	0.0	0.0	0.0	0.0	0.1	0.0	0.2	0.2	0.1	0.0	0.0	0.3	
		-	0.0	0.0	0.0	0.0	0.0	0.0	0.2	0.2	0.0	0.0	3.3	49.0	51.6	0.0	
		-	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.7	0.4	0.0	0.0	0.0	0.0	



\*Our original preliminary results.

# Image Contamination

Data issues  
Lifecycle  
Noise  
Positionality  
**Contamination**  
Measurement issues  
Systemic issues  
Epistemic issues  
Concluding

- Risk level: Same.

**MMMU (test)**



**Q:** Which is not a negative outcome of <image 1>?

**LAION-2B**



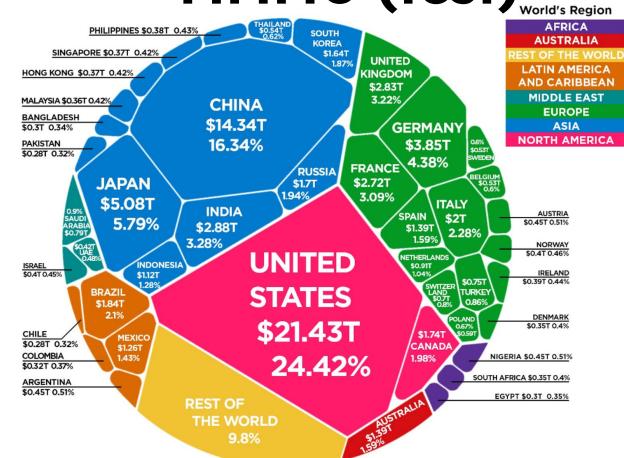
**Caption:** Dredging and Capping | Onondaga Lake Cleanup

# Image Contamination

Data issues  
Lifecycle  
Noise  
Positionality  
**Contamination**  
Measurement issues  
Systemic issues  
Epistemic issues  
Concluding

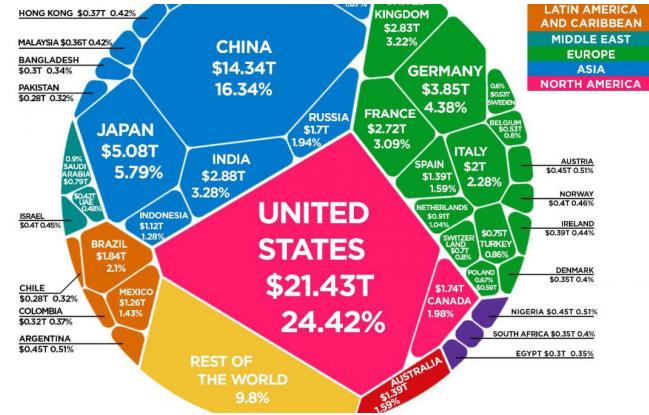
- Risk level: High.

## MMMU (test)



Q: If a sociologist says that nations evolve toward more advanced technology and more complex industry as their citizens learn cultural values that celebrate hard work and success, she is using \_ theory to study the <image 1>.

## LAION-2B



**Caption:** The \$88 Trillion World Economy in One Chart

\*Our original preliminary results.

# Measurement issues

Data issues

Measurement issues

Rubrics

Judges

Metrics

Systemic issues

Epistemic issues

Concluding

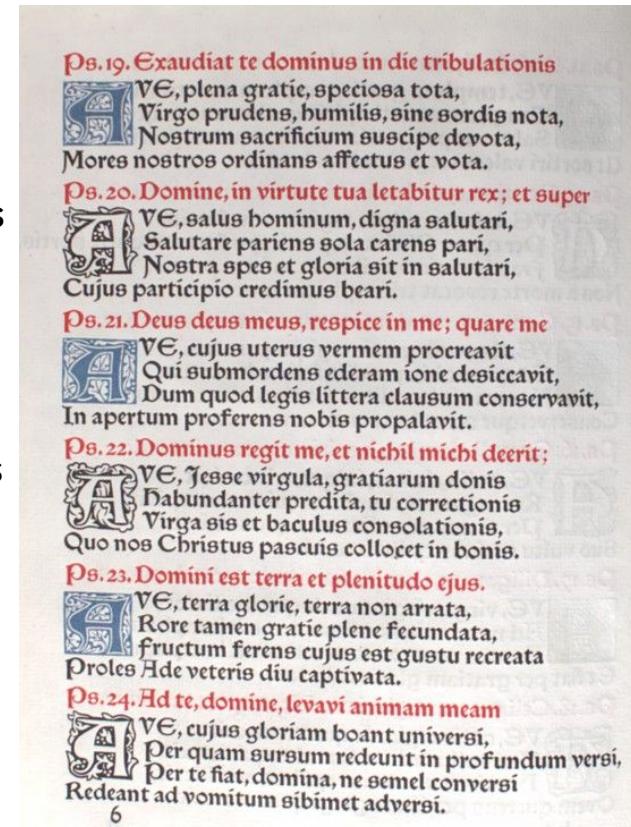
- **Measurement issues** refer to challenges in benchmarking which arise from:
- Poorly-developed **rubrics** for annotation
- Challenges with **metrics** used to grade or score the outputs of a model under test
  - **LM judges**
  - **Learned metrics**
  - **Algorithmic metrics**



# Measurement issues: rubrics

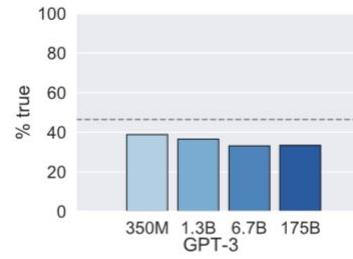
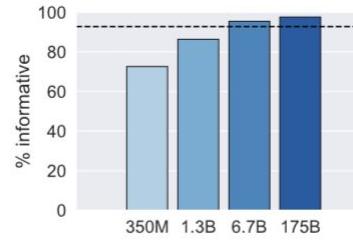
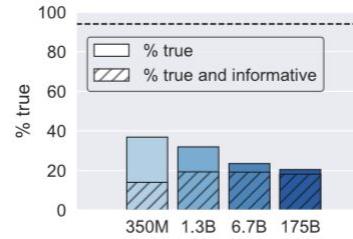
Data issues  
Measuremen  
t issues  
Rubrics  
Judges  
Metrics  
Systemic  
issues  
Epistemic  
issues  
Concluding

- **Rubrics**, from the Latin name for red ink writing, refer to scoring guidelines for writing assignments in US education (Popham, 1997)
- The rubric pattern has grown popular in AI evaluation, both as for:
  - Guiding human annotators (usually non-expert crowdworkers) in scoring outputs
  - Guiding **LM judges** in scoring outputs
- How do we know when we're writing good rubrics? Does rubric quality matter?



# Measurement issues: Judges

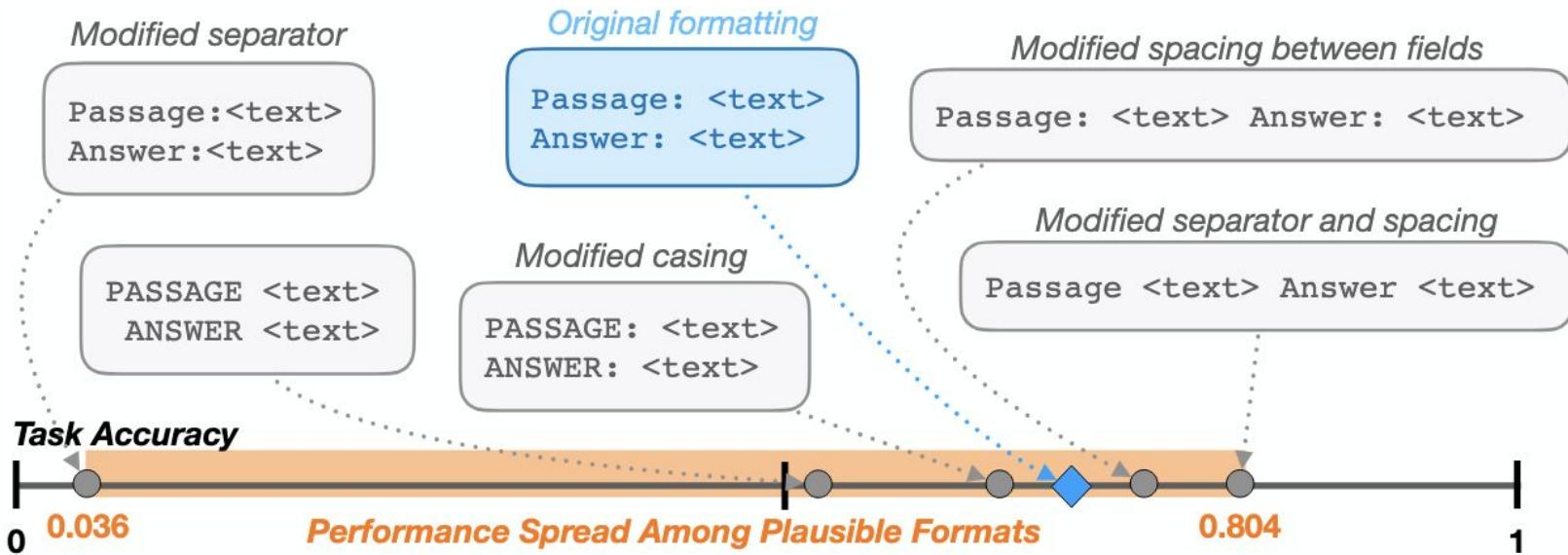
- **LM-as-a-judge** is a technique for evaluating natural text since at least TruthfulQA (Lin et. al, 2021)
- An LM is prompted to provide a numerical score given a set of requirements (eg. truth, informative, etc)
- Advantages:
  - you can specify your desiderata in natural language
  - easily handles open-ended generations
  - (can be) straightforward to implement
- Disadvantages:
  - Brittle to prompt variations
  - Rubric implementation complications
  - Self-bias



# Writing good rubrics

Data issues  
Measurement issues  
**Rubrics**  
Judges  
Metrics  
Systemic issues  
Epistemic issues  
Concluding

- **Rubric quality matters a lot**
- LM judges are extremely sensitive to minor variations in prompt phrasing (Sclar et al, 2023)
- How can we improve them?



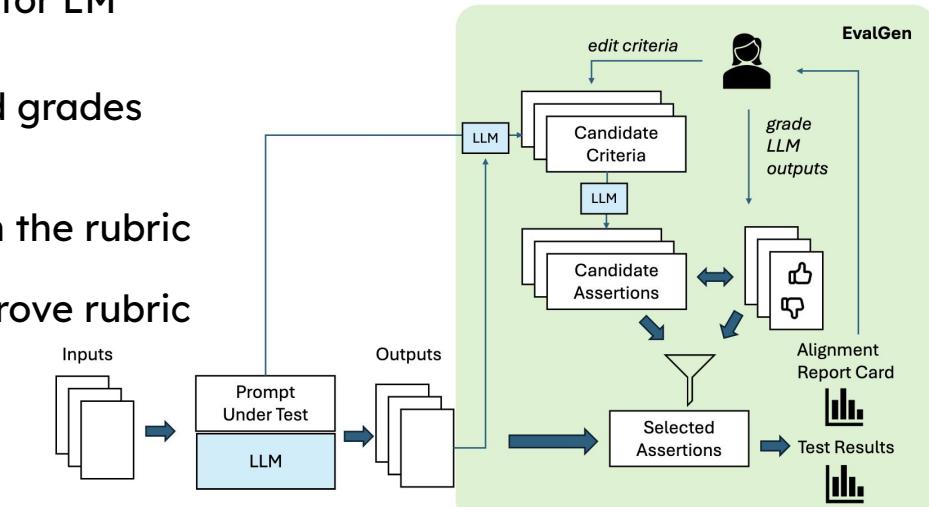
# Writing good rubrics: EvalGen

Data issues  
Measurement issues  
Rubrics  
Judges  
Metrics  
Systemic issues  
Epistemic issues  
Concluding

- There's a catch-22 in rubric writing:
  - to grade some outputs, you need criteria
  - to know the criteria, you have to understand what errors occur
- This is why educators update rubrics as they grade (Shankar et. al 2024)

- A similar process can be adopted for LM judges, given seed rubric:
  1. Human reviews set of outputs and grades them
  2. LM grades according to rubric
  3. Human checks which assertions in the rubric are agreed and which aren't
  4. LM suggests edit locations to improve rubric

Should be adopted for benchmark dev



# Subjectivity and criteria drift

Data issues

Measurement issues

Rubrics

Judges

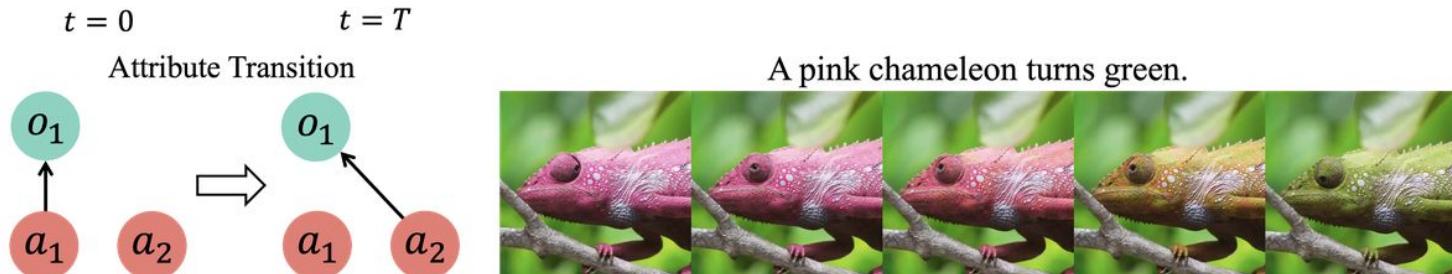
Metrics

Systemic issues

Epistemic issues

Concluding

- This problem of **criteria drift** is widespread (Shankar et. al, 2024)
- Even **static rubrics** may be problematic:
  - As models evolve and improve, the types of errors they make will change
  - Rubrics well-suited to poorly performing models may not apply as performant models have more nuanced weaknesses
- Example: 2023-2024 video models struggled with temporal consistency (Feng et. al, 2024) but new video gen models like Veo don't
- Are videogen benchmarks meant to test this evergreen?



What's Measured?

What's Missed?

What's Next?

# Self-bias in LM judges

Data issues

Measurement issues

Rubrics

Judges

Metrics

Systemic issues

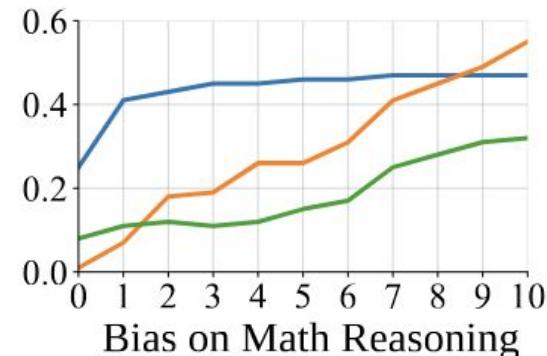
Epistemic issues

Concluding

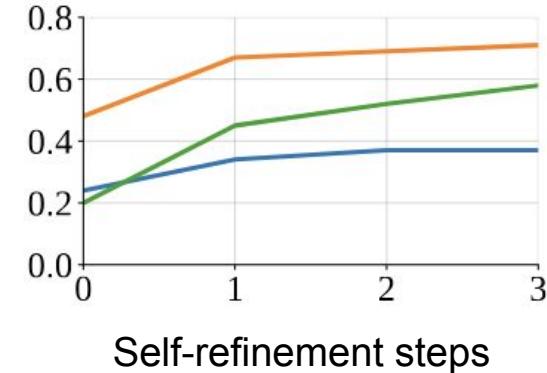
- **Self-bias** is the phenomenon of LM judges implicitly preferring text generated by their own base model
- Xu et al (2024) demonstrate that across multiple models, multiple tasks LM judges self-bias (distance from scores assigned by external annotator) grows as more self-refinement steps are made by the model
- **LM judges need to be carefully checked**
- Failure example: “Exploring the MIT Mathematics and EECS Curriculum” paper

GPT-4    GPT-3.5    Gemini

Bias on CommonGen Hard



Bias on Math Reasoning

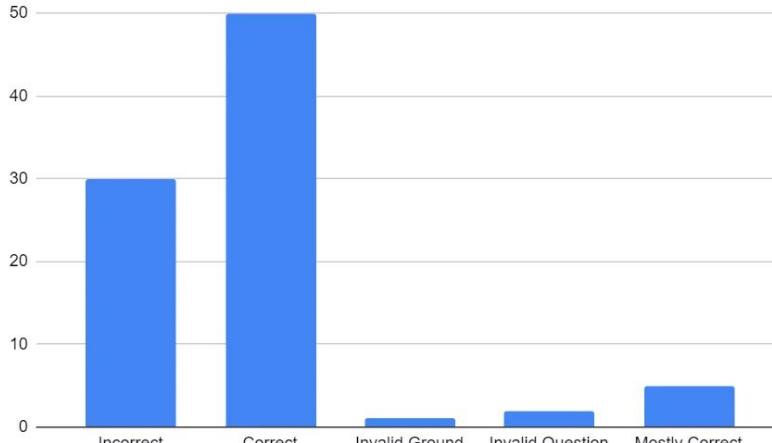


Self-refinement steps

# Exploring the MIT Math and EECS with GPT-4

Data issues  
Measurement issues  
Rubrics  
Judges  
Metrics  
Systemic issues  
Epistemic issues  
Concluding

- Claim: *GPT-4 can get 100% on the core MIT EECS exams* (Zhang et. al, 2023)
- Among the problems: incorrect grades assigned by GPT-4 as a judge! (Chowdhuri et. al, 2024)
- **Paper was withdrawn from arXiv**



This paper has been withdrawn by Iddo Drori

[Submitted on 15 Jun 2023 (v1), last revised 24 Jun 2023 (this version, v2)]

## Exploring the MIT Mathematics and EECS Curriculum Using Large Language Models

Sarah J. Zhang, Samuel Florin, Ariel N. Lee, Eamon Niknafs, Andrei Marginean, Annie Wang, Keith Tyser, Zad Chin, Yann Hicke, Nikhil Singh, Madeleine Udell, Yoon Kim, Tonio Buonassisi, Armando Solar-Lezama, Iddo Drori

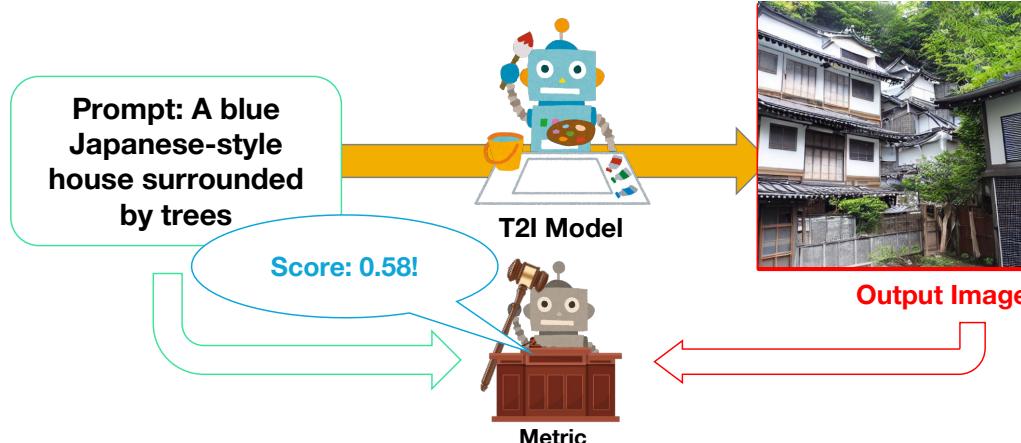
We curate a comprehensive dataset of 4,550 questions and solutions from problem sets, midterm exams, and final exams across all MIT Mathematics and Electrical Engineering and Computer Science (EECS) courses required for obtaining a degree. We evaluate the ability of large language models to fulfill the graduation requirements for any MIT major in Mathematics and EECS. Our results demonstrate that GPT-3.5

# Measurement issues: Metrics

Data issues  
Measuremen  
t issues  
Rubrics  
Judges  
**Metrics**

Systemic  
issues  
Epistemic  
issues  
Concluding

- Many other metrics which aren't strictly LM judge + rubric have been proposed
- For example, in text-to-image evaluation **prompt consistency** is the evaluation task of assigning a numerical score to the alignment of a generated image to its input prompt
- Multiple classes of prompt consistency metrics



# Classes of text-to-image faithfulness metrics

Data issues  
Measurement issues  
Rubrics  
Judges  
**Metrics**

Systemic issues  
Epistemic issues  
Concluding

- Two predominant classes of prompt-consistency metrics are:
- **Embedding-correlation metrics and VLM-VQA metrics**

## Embedding Correlation



- Example: CLIPScore (Hessel et. al, 2021)
- Embed image & prompt with CLIP, return cosine similarity of embeddings
- Cheap, fast
- Scores aren't attributable
- Considered less performant

## VLM question-answering

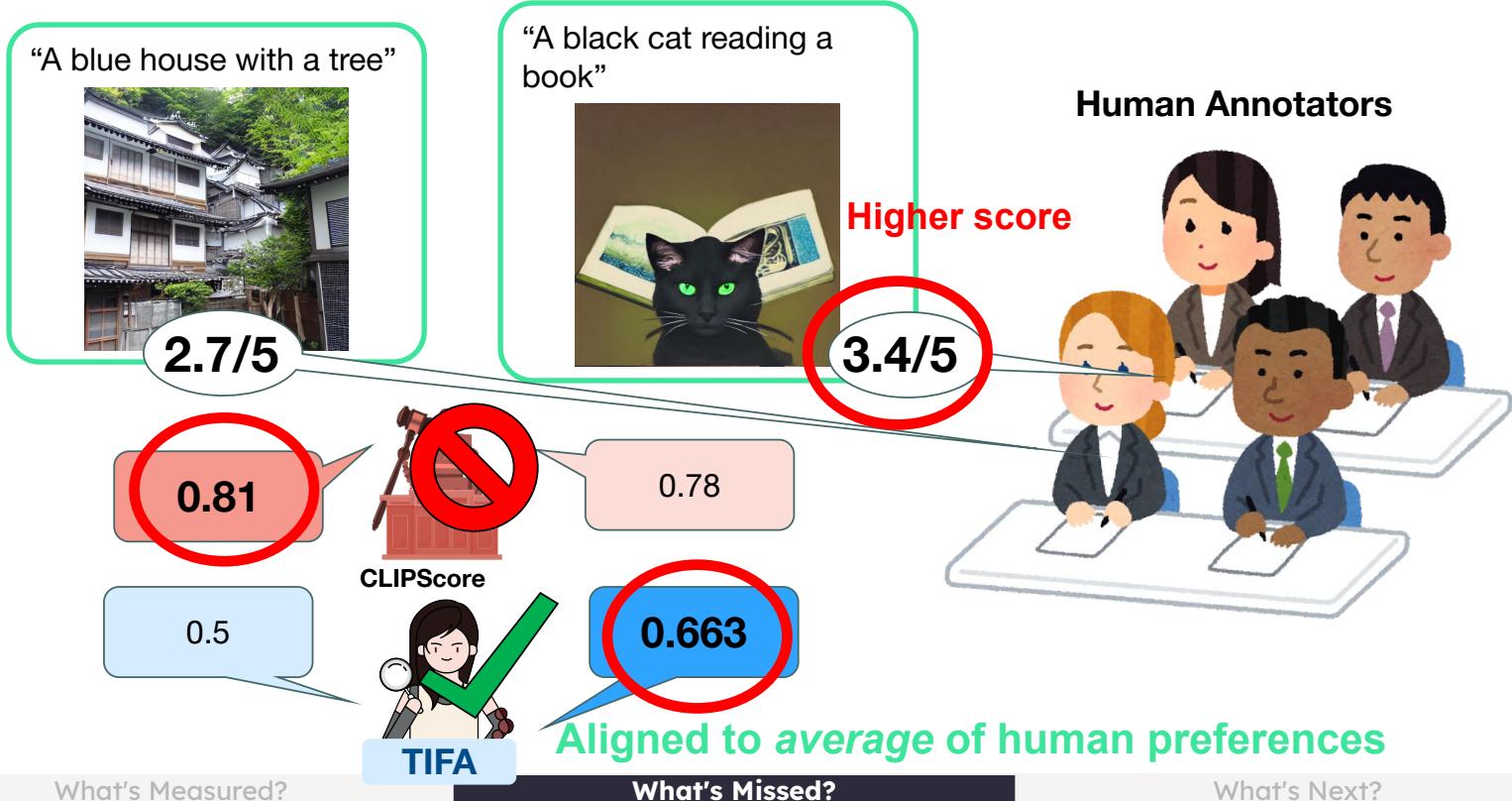


- Example: TIFA (Hu et. al, 2023)
- **Use LM to generate list of requirements from prompt**
- **Check if each requirement is met in image with vision LM**
- Expensive
- Scores are attributable to natural language requirements
- Considered more performant

# Flawed meta-evaluation establishes superiority

Data issues  
Measurement issues  
Rubrics  
Judges  
**Metrics**  
Systemic issues  
Epistemic issues  
Concluding

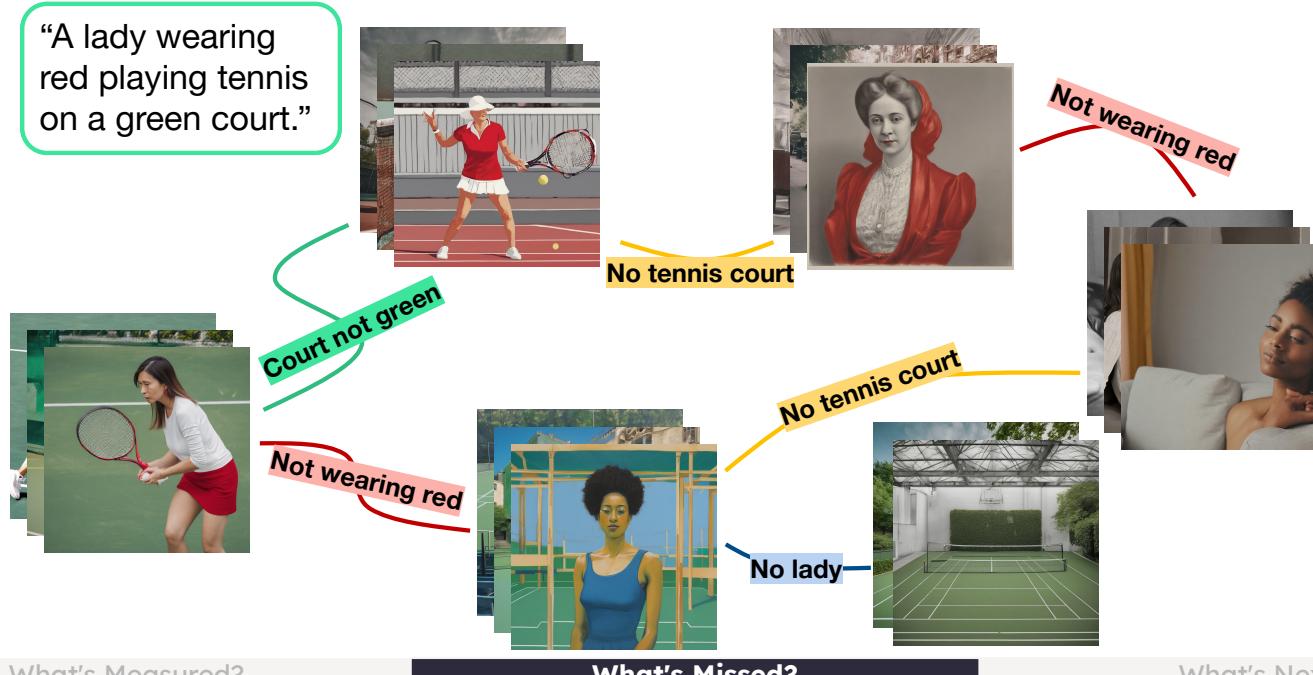
- Previously, superiority of VLM-VQA over correlation based on correlation to human judges over unrelated images
- Relative quality of unrelated images is *inherently subjective!*



# Testing which is better with T2IScoreScore

Data issues  
Measurement issues  
Rubrics  
Judges  
**Metrics**  
Systemic issues  
Epistemic issues  
Concluding

- Saxon et. al (2024) instead analyzed T2I metrics with **semantic error graphs**
- Each graph contains populations of **related images** organized by counterfactual **error count** relative to the prompt
- Performance is judged by how well LMs *reconstruct* the error graph by ordering

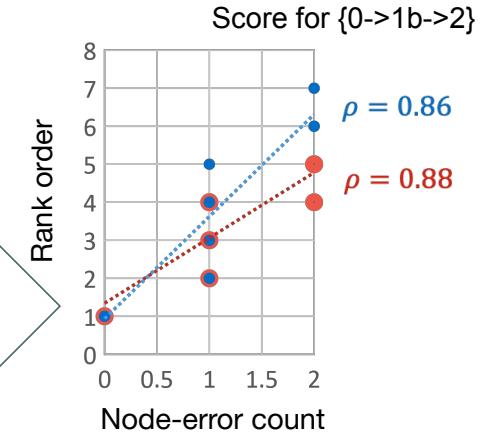
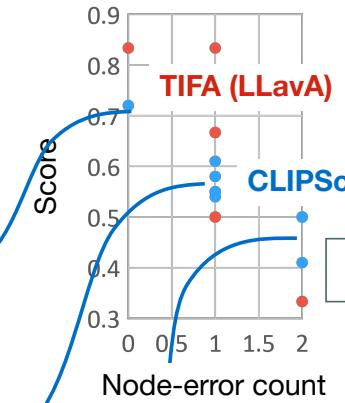


# Testing which is better with T2IScoreScore

Data issues  
Measurements issues  
Rubrics  
Judges  
**Metrics**  
Systemic issues  
Epistemic issues  
Concluding

## Ordering Meta-Metric: Spearman Correlation

**SEG1: A teddy bear underneath a Christmas tree covered in lights.**



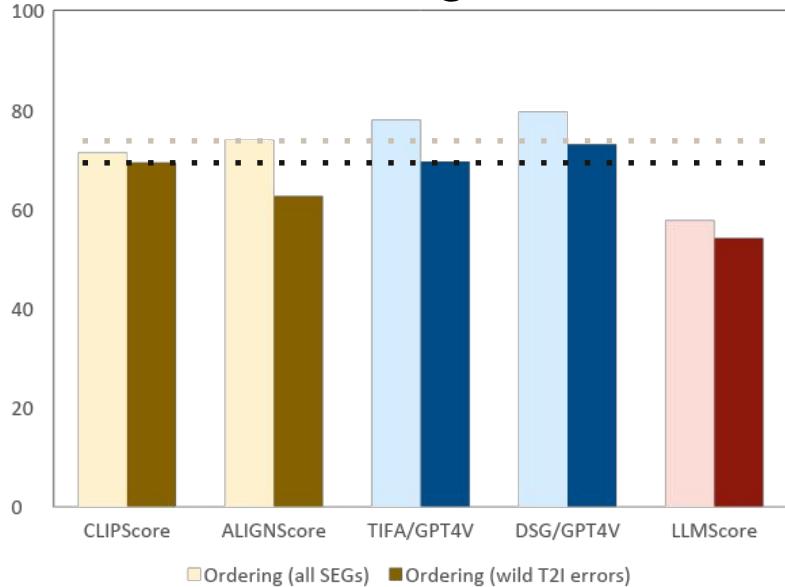
$$\rho(X, Y) = \frac{\text{cov}(R(X), R(Y))}{\sigma_{R(X)} \sigma_{R(Y)}}; \quad R(X) = \left\{ \sum_{x_i \in X} \mathbb{1}(x_i < x) \mid x_i \in X \right\}$$

$$\text{rank}_m(S) = \frac{1}{|S|} \sum_{W \in S} r_s(\{m(I, P) \mid (I, P, N) \in W\}, \{N \mid (I, P, N) \in W\})$$

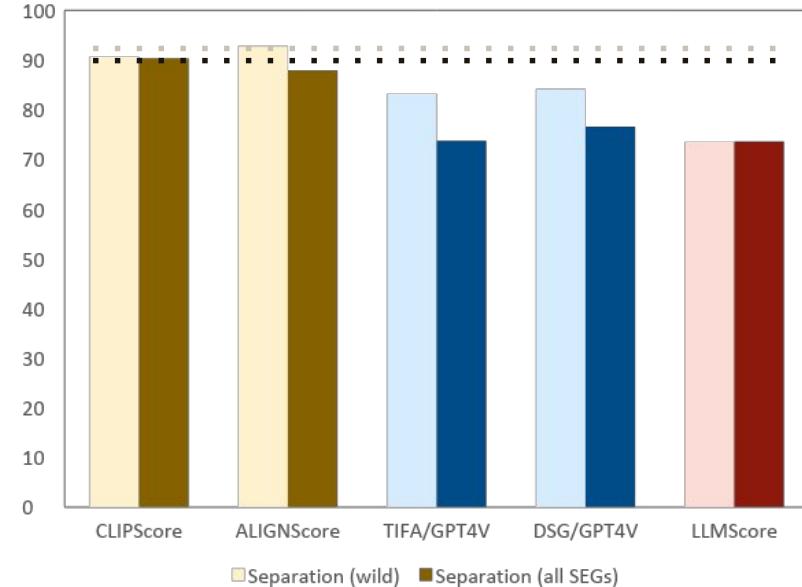
# Measurement issues: Metrics

Data issues  
Measurement issues  
Rubrics  
Judges  
**Metrics**  
Systemic issues  
Epistemic issues  
Concluding

## Ordering score



## Separation score



Embedding Correlation

All SEGs

Wild SEGs only

VLM question-answering

Others

What's Measured?

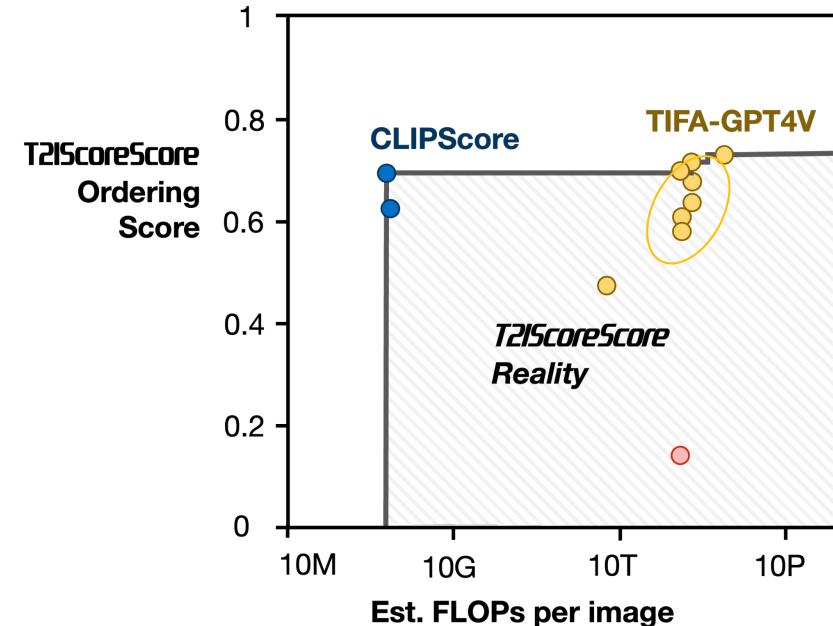
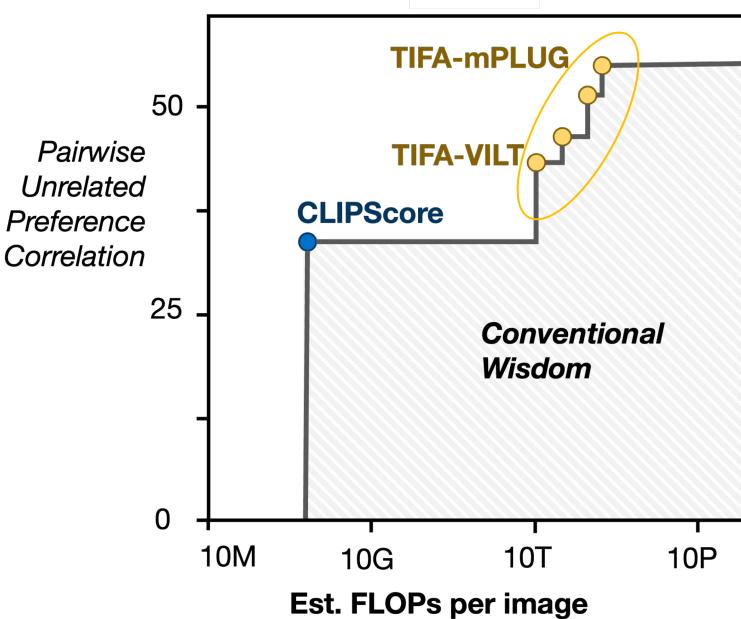
What's Missed?

What's Next?

# Measurement issues: Metrics

Data issues  
Measurement issues  
Rubrics  
Judges  
**Metrics**  
Systemic issues  
Epistemic issues  
Concluding

- We find that the seemingly dominant versions of TIFA fall **below the pareto frontier** against CLIPScore for cost-performance
- Much more expensive VLMs like GPT4V needed to perform
- Important to meta-evaluate in an **ecologically valid manner!**



# Systemic issues

Data issues

Measurement issues

**Systemic issues**

Leaderboard illusion

Weak baselines

Reporting variance

Epistemic issues

Concluding

- **Systemic issues** are the cases where the **manner in which an evaluation is conducted or made** are problematic.
- The **leaderboard illusion**: how multiple submissions can game leaderboard style benchmarks
- Illusory improvements in performance from evaluating against **weak baselines**
- **(Mis)-reporting variance**: how choosing a best-of-n result may skew performance

## Beware: Goodhart's Law

*Once a measure becomes a target it ceases to be a good measure.*



What's Measured?

What's Missed?

What's Next?

# Leaderboard illusion (Singh et. al, 2025)

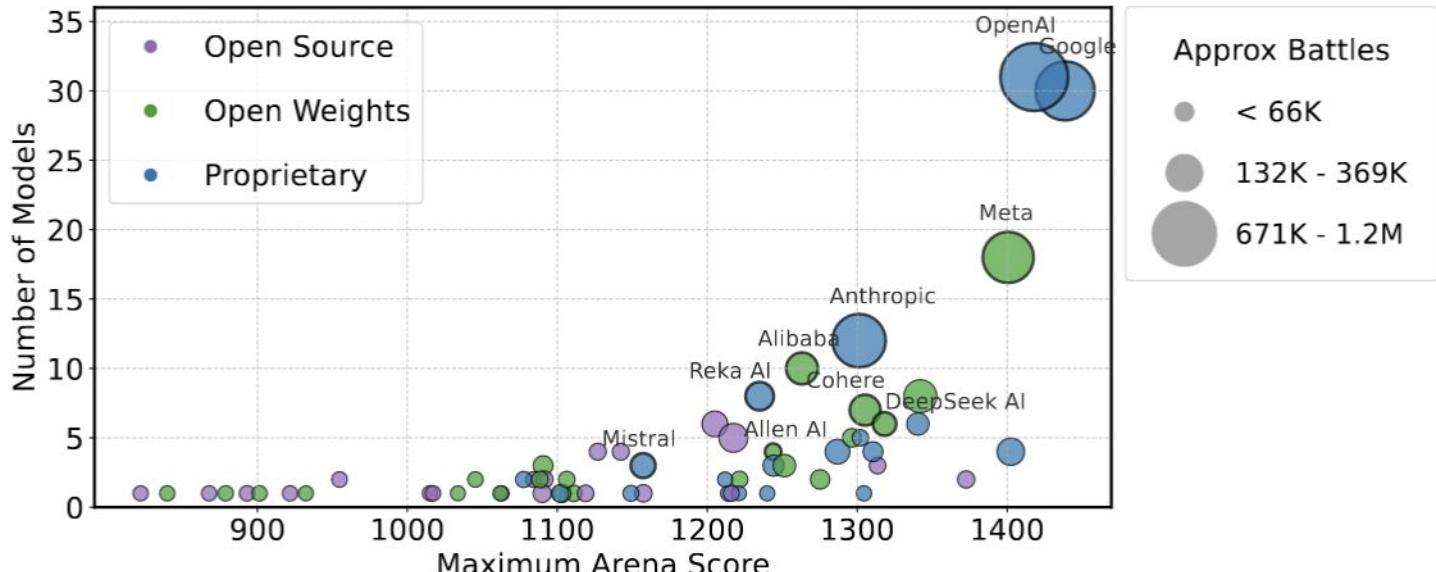
Data issues  
Measurement issues  
Systemic issues  
**Leaderboard illusion**  
Weak baselines  
Reporting variance  
Epistemic issues  
Concluding



# Leaderboard illusion (Singh et. al, 2025)

Data issues  
Measurement issues  
Systemic issues  
**Leaderboard illusion**  
Weak baselines  
Reporting variance  
Epistemic issues  
Concluding

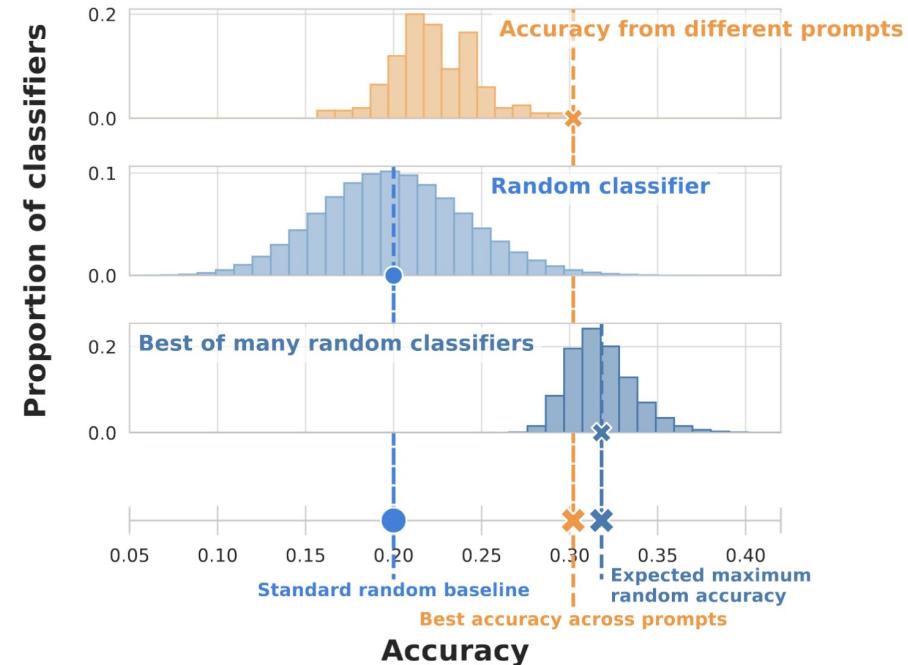
- Support: clear correlation between number of models a provider submits and score
- This gives the opportunity for providers to learn, gaming?



# What is wrong with multiple attempts?

Data issues  
Measurement issues  
Systemic issues  
Leaderboard illusion  
Weak baselines  
Reporting variance  
Epistemic issues  
Concluding

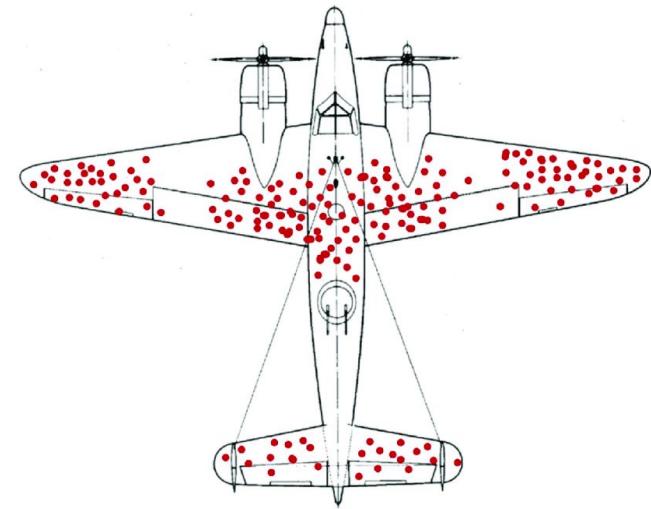
- Reporting asymmetrical best-of-n is also a statistical problem!
- Yauney and Mimno (2024): many methods fail to even beat a **random baseline** when its sampled with the same best-of-n process that models often are
- Leaderboard illusion: doing this not to random baselines, **but competing models**



# Intra-model variance (Reuel et al, 2024)

Data issues  
Measurement issues  
Systemic issues  
Leaderboard illusion  
Weak baselines  
**Reporting variance**  
Epistemic issues  
Concluding

- The root of many of these statistical reporting issues is **intra-model variance**
- Reporting a single number for a new method or model, rather than the results according to multiple seeds, temperatures, etc is crucial to distinguish **signal from noise**
- The wider the intra-model variance, the lower resolution the benchmark is (as there is a significant **noise band** around any given score)
- Rankings are unreliable in this case
- Unfortunately, it is both expensive to do full extra training runs for this purpose
- But if we can't, what are we doing here?



# Epistemic issues

Data issues

Measurement issues

Systemic issues

Epistemic issues

Task universe

Map/territory

Psychometrics?

Concluding

- Finally, we go most abstract:
- **Epistemic issues** are the core failures of conceptualization and operationalization that cut across all the more concrete issues below.
- Here we discuss:
- The **distinction between “tasks” and “learning problems”**, and how the conflation of the two within a **task universe** causes problems
- Confusions of **map and territory** when thinking about AI
- The pitfalls of treating benchmarks like human **psychometrics**



# Tasks vs learning problems

Data issues

Measurement issues

Systemic issues

Epistemic issues

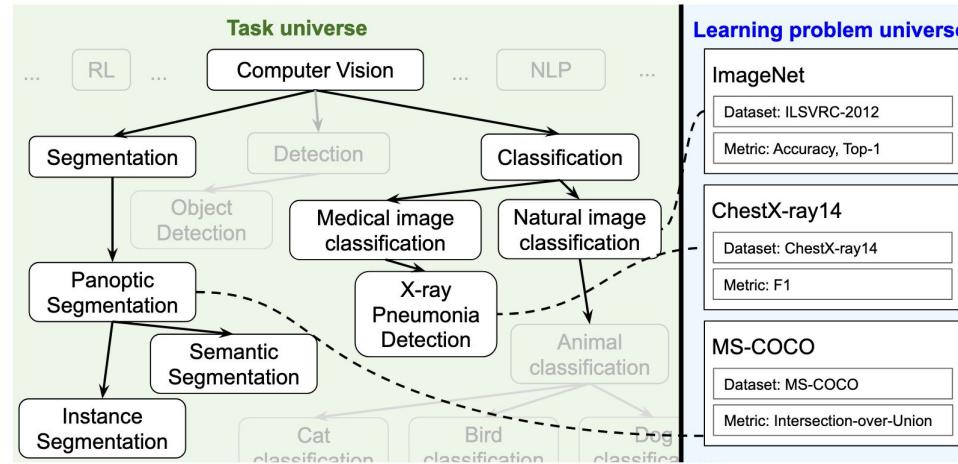
Task universe

Map/territory

Psychometrics?

Concluding

- “Tasks” are often conceptualized at multiple scales (within a “task universe”)
- Foundation models represent attempts to create systems for higher and higher level tasks
- “Learning problems” are distinct from tasks. Imagenet is a learning problem that purports to capture the task of natural image classification. (Liao et. al 2024)
- Does it?



What's Measured?

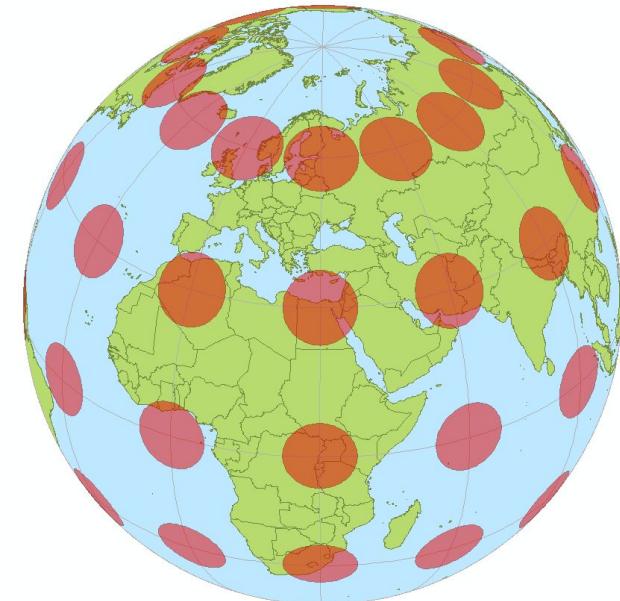
What's Missed?

What's Next?

# Map-territory and “wishful mnemonics”

Data issues  
Measurement issues  
Systemic issues  
Epistemic issues  
Task universe  
**Map/territory**  
Psychometrics?  
Concluding

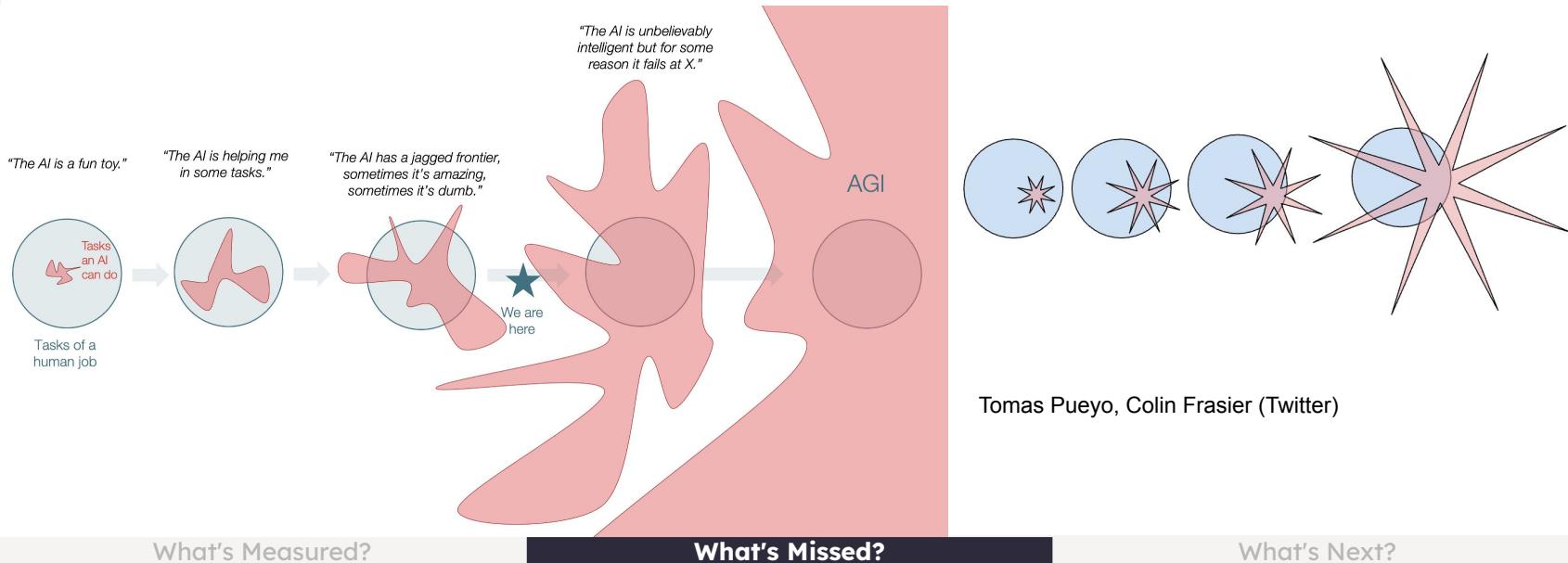
- Mitchell (2021) presents fallacies in AI thought which hamper progress toward AGI
- **Wishful mnemonics** are a critical fallacy which chronically plagues benchmarkers
- Calling a task “reasoning” or “understanding” is a wishful mnemonic (McDermott 1975)
- When we do this, we beg the question that a model’s performance on that **task** represents an **abstract capability**—does it though?
- Chronic failures of models to generalize across eg reasoning tasks suggests the answer is no
- Coupled with sloppy application of **learning problems** to these **tasks** we often seriously confuse map and territory



# Psychometrics trap

Data issues  
Measuremen  
t issues  
Systemic  
issues  
Epistemic  
issues  
Task universe  
Map/territory  
**Psychometrics?**  
Concluding

- Attempts to develop generalized “psychometrics” for AI systems suffer from all these epistemic problems
- Standardized tests work for humans because of “shared architecture”
- The “jagged frontier” problem complicates attempts to develop AI psychometrics
- There are two competing visions that are impossible to reconcile now



# Summary

Data issues  
Measurement issues  
Systemic issues  
Epistemic issues  
Task universe  
Map/territory  
Psychometrics?

- **Noise** can come from data, metrics, or variance problems
  - **Bad examples** create both a noise ceiling (max meaningful performance) and **lower resolution**.
  - Poor statistical (systemic) practices also **lower resolution** through variance-based noise bands
  - Resolving either of these issues may cause the ranking of systems to completely swap.
- **Careful conceptualization** is crucial: construct validity is necessary for both good data collection and metric design/validation
  - Task vs learning problem
- **Positionality** of humans throughout the benchmarking process is often overlooked.
  - Bad samples can come from exhausted crowdworkers.
  - AI researchers lack domain expertise.
  - Hidden cultural assumptions drive annotator, producer choices.



# Questions

Data issues

Measurement issues

Systemic issues

Epistemic issues

Task universe

Map/territory

Psychometrics?

## Concluding

- **What counts as contamination** when we are looking for systems to do everything?
- Do the same generalization weaknesses exist in LLMs that specialized models had?
- Are we producing “intelligence” at all if we have to explicitly train the model on every individual task?



What's Measured?

What's Missed?

What's Next?

# What's Next?



Xiang Yue  
Carnegie Mellon University  
(Now at Meta)

# Agenda

- What's Measured? (1:30PM - 2:10PM)
- What's Missed? (2:10PM - 2:40PM)
- What's Next? (2:40PM - 3:15PM)
  - Towards dynamic and agentic benchmarking
  - Towards real-world benchmarking
  - Other emerging practices
  - Some proposals
- Panel Discussion (3:20PM-4:00PM)



# Dynamic & “Living” Benchmarks

Dynamics

Arena

Live Bench

Agency

Digital

Embodied

Realism

Adapt

Adv. Test

Proposal

- Static benchmarks quickly become saturated by frontier models.
- The risk of benchmark items contaminating training data is high.
- Static evaluation is rapidly obsolete as model knowledge evolves.
- Dynamic and “Living” Benchmarks are necessary to ensure relevant evaluation.



# LM Arena

Dynamics  
Arena  
Live Bench  
Agency

Digital  
Embodied  
Realism  
Adapt  
Adv. Test  
Proposal

- LM Arena (previously Chatbot Arena) pits two anonymous models against each other in randomised “battles.” Users chat with both models side-by-side and vote for their preferred answer.

The screenshot shows a comparison between two AI assistants, Assistant A and Assistant B, on the topic "What is a good benchmark?".

**Assistant A:** A good benchmark has these key qualities:

- Relevance**
  - Measures what actually matters to your goals
  - Reflects real-world conditions and use cases
  - Aligns with how performance will be judged in practice
- Reproducibility**
  - Produces consistent results when repeated
  - Has clear, documented methodology
  - Uses standardized tools and conditions

**Assistant B:** Characteristics of a Good Benchmark

A good benchmark typically has these qualities:

- Key Attributes**

Attribute	Description
Relevant	Measures something meaningful to your goals
Reproducible	Produces consistent results when repeated
Representative	Reflects real-world conditions/workloads

Voting interface at the bottom:

- Left is Better
- It's a tie
- Both are bad
- Right is Better

# LM Arena Leaderboard

Dynamics	
Arena	
Live Bench	
Agency	
Digital	
Embodied	
Realism	
Adapt	
Adv. Test	
Proposal	

Rank ↑	Rank Spread ⓘ (Upper-Lower)	Model ↑	Score ↓	95% CI (±) ↑↓	Votes ↑↓	Organization ↑↓	License ↑↓
1	1 ↔ 2	Google gemini-3-pro	1492	±8	9,799	Google	Proprietary
2	1 ↔ 3	xAI grok-4.1-thinking	1482	±8	10,067	xAI	Proprietary
3	2 ↔ 6		466	±9	4,677	Anthropic	Proprietary
4	3 ↔ 6		464	±8	9,967	xAI	Proprietary
5	3 ↔ 8	OpenAI gpt-5.1-high	1461	←	7,893	OpenAI	Proprietary
6	3 ↔ 10	Anthropic claudie-opus-4-5-20251101-thinking-32k	1460	±12	2,763	Anthropic	Proprietary
7	5 ↔ 10	Google gemini-2.5-pro	1452	±4	70,875	Google	Proprietary
8	5 ↔ 13	Anthropic claudie-sonnet-4-5-20250929-thinking-32k	1448	±5	22,000	Anthropic	Proprietary
9	6 ↔ 13	Anthropic claudie-opus-4-1-20250805-thinking-16k	1448	±4	37,617	Anthropic	Proprietary
10	6 ↔ 15	Anthropic claudie-sonnet-4-5-20250929	1445	±6	16,961	Anthropic	Proprietary
11	8 ↔ 18	OpenAI gpt-4.5-preview-2025-02-27	1442	±6	14,644	OpenAI	Proprietary

“Live” user queries,  
votings and scores

Wei-Lin Chiang, et al., Chatbot Arena, 2024.

What's Measured?

What's Missed?

What's Next?

# LM Arena Leaderboard

Dynamics	
Arena	
Live Bench	
Agency	
Digital	
Embodied	
Realism	
Adapt	
Adv. Test	
Proposal	

Rank ↑	Rank Spread ⓘ (Upper-Lower)	Model ↗	Score ↓	95% CI (±) ↑↑	Votes ↑↑	Organization ↑↑	License ↑↑
1	1 ↔ 2	G gemini-3-pro	1492	±8	9,799	Google	Proprietary
2	1 ↔ 3	XI grok-4.1-thinking	1482	±8	10,067	xAI	Proprietary
3	2 ↔ 6		466	±9	4,677	Anthropic	Proprietary
4	3 ↔ 6		464	±8	9,967	xAI	Proprietary
5	3 ↔ 8	GPT gpt-5.1-high	1461	±8	7,893	OpenAI	Proprietary
6	3 ↔ 10	AI claude-opus-4-5-20251111-thinking-32k	1460	±12	2,763	Anthropic	Proprietary
7	5 ↔ 10		452	±4	70,875	Google	Proprietary
8	5 ↔ 13		448	±5	22,000	Anthropic	Proprietary
9	6 ↔ 13		448	±4	37,617	Anthropic	Proprietary
10	6 ↔ 15		445	±6	16,961	Anthropic	Proprietary
11	8 ↔ 18	GPT gpt-4.5-preview-2025-02-27	1442	±6	14,644	OpenAI	Proprietary

“Live” user queries,  
votings and scores

Lower risk of  
contamination and  
harder to hack

# LM Arena Leaderboard

Dynamics

Arena

Live Bench

Agency

Digital

Embodied

Realism

Adapt

Adv. Test

Proposal

Rank ↑	Rank Spread ⓘ (Upper-Lower)	Model ↑	Score ↓	95% CI (±) ↑	Votes ↑	Organization ↑	License ↑
1	1 ↔ 2	Gemini-3-pro	1492	±8	9,799	Google	Proprietary
2	1 ↔ 3	Grok-4.1-thinking	1482	±8	10,067	xAI	Proprietary
3			1466	±9	4,677	Anthropic	Proprietary
4			1464	±8	9,967	xAI	Proprietary
5			1461	±8	7,893	OpenAI	Proprietary
6			1460	±12	2,763	Anthropic	Proprietary
7			1452	±4	70,875	Google	Proprietary
8			1448	±5	22,000	Anthropic	Proprietary
9	6 ↔ 13	Claude-Opus-4-1-20250805-thinking-16K	1448	±4	37,617	Anthropic	Proprietary
10	6 ↔ 15	Claude-Sonnet-4-5-20250929	1445	±6	16,961	Anthropic	Proprietary
11	8 ↔ 18	GPT-4.5-preview-2025-02-27	1442	±6	14,644	OpenAI	Proprietary

Wei-Lin Chiang, et al., Chatbot Arena, 2024.

What's Measured?

What's Missed?

What's Next?

# Potential Issues of LM Arena

Dynamics

Arena

Live Bench

Agency

Digital

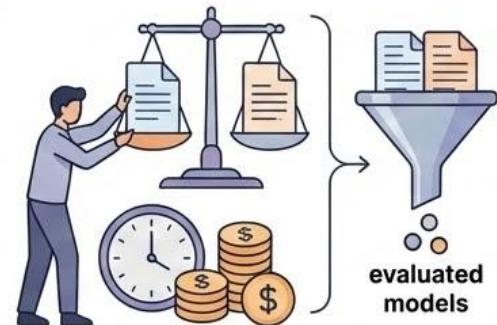
Embodied

Realism

Adapt

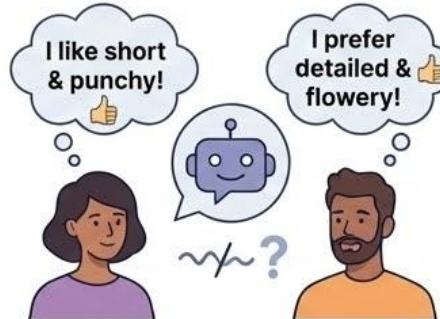
Adv. Test

Proposal



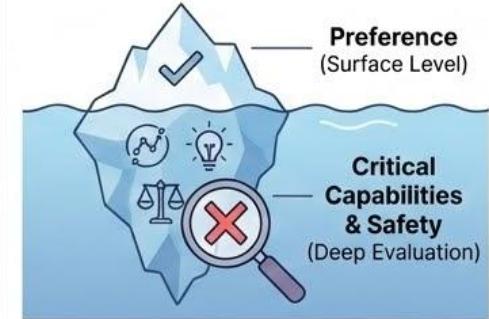
## Cost and Scalability

Generating enough quality human-rated battles is expensive and slow, limiting the scale of evaluation.



## Subjectivity

User preference can be highly subjective and inconsistent. Users may be swayed by superficial factors like verbosity or style, rather than solely evaluating technical merit.

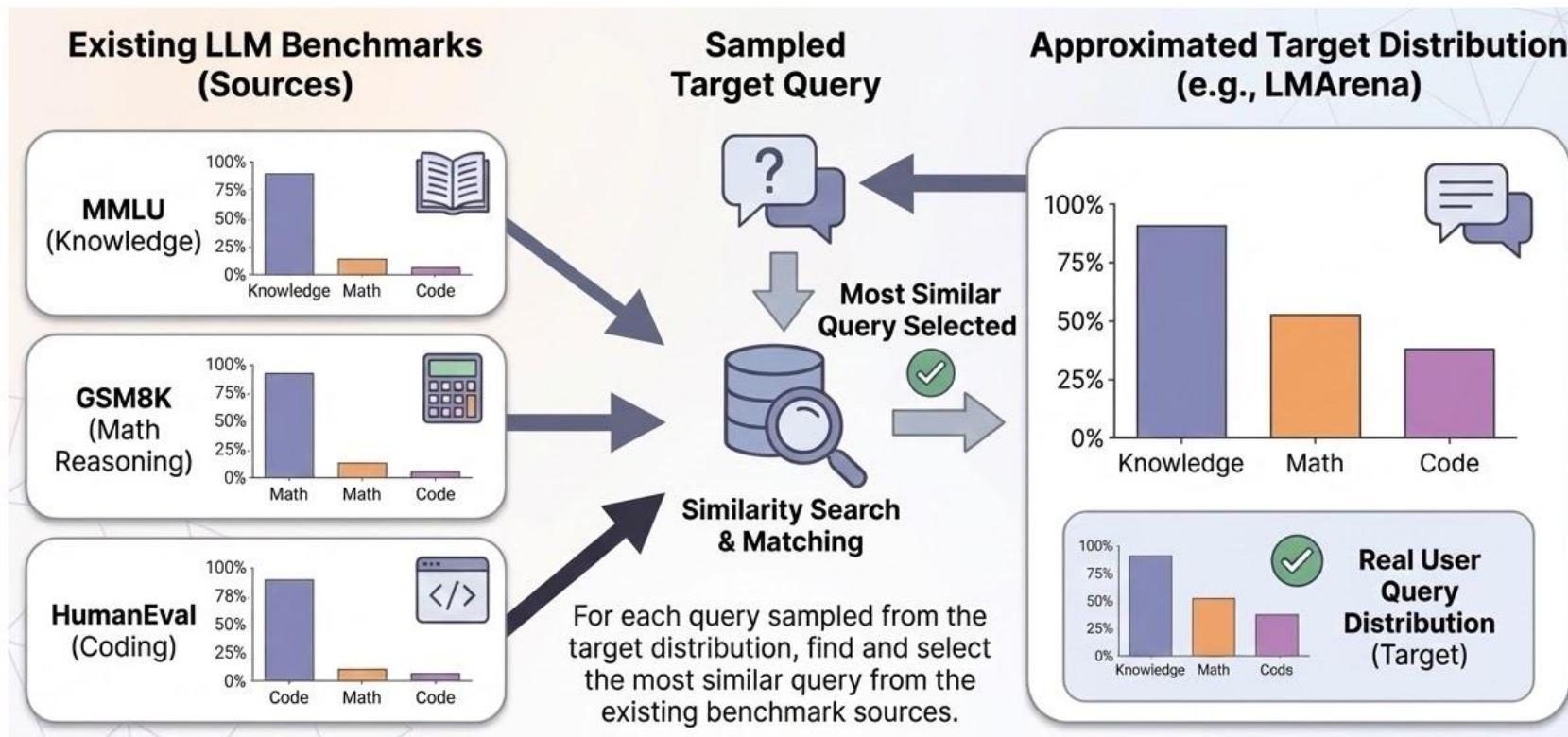


## Evaluation Depth

Pairwise comparison is limited to preference and may not deeply evaluate specific, critical capabilities or safety aspects.

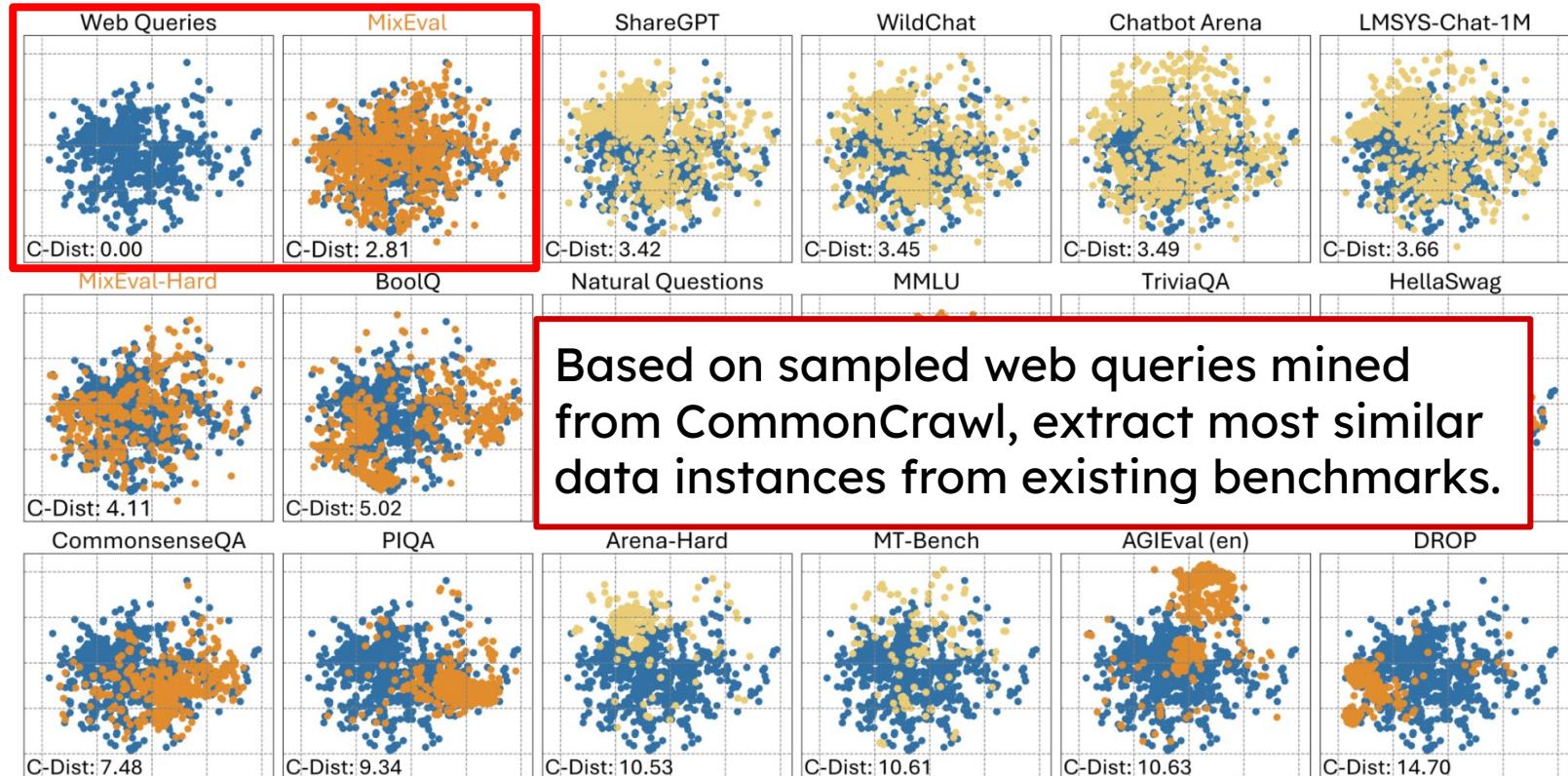
# Dynamic Benchmarking by Mixing Existing Ones

Dynamics  
Arena  
Live Bench  
Agency  
Digital  
Embodied  
Realism  
Adapt  
Adv. Test  
Proposal



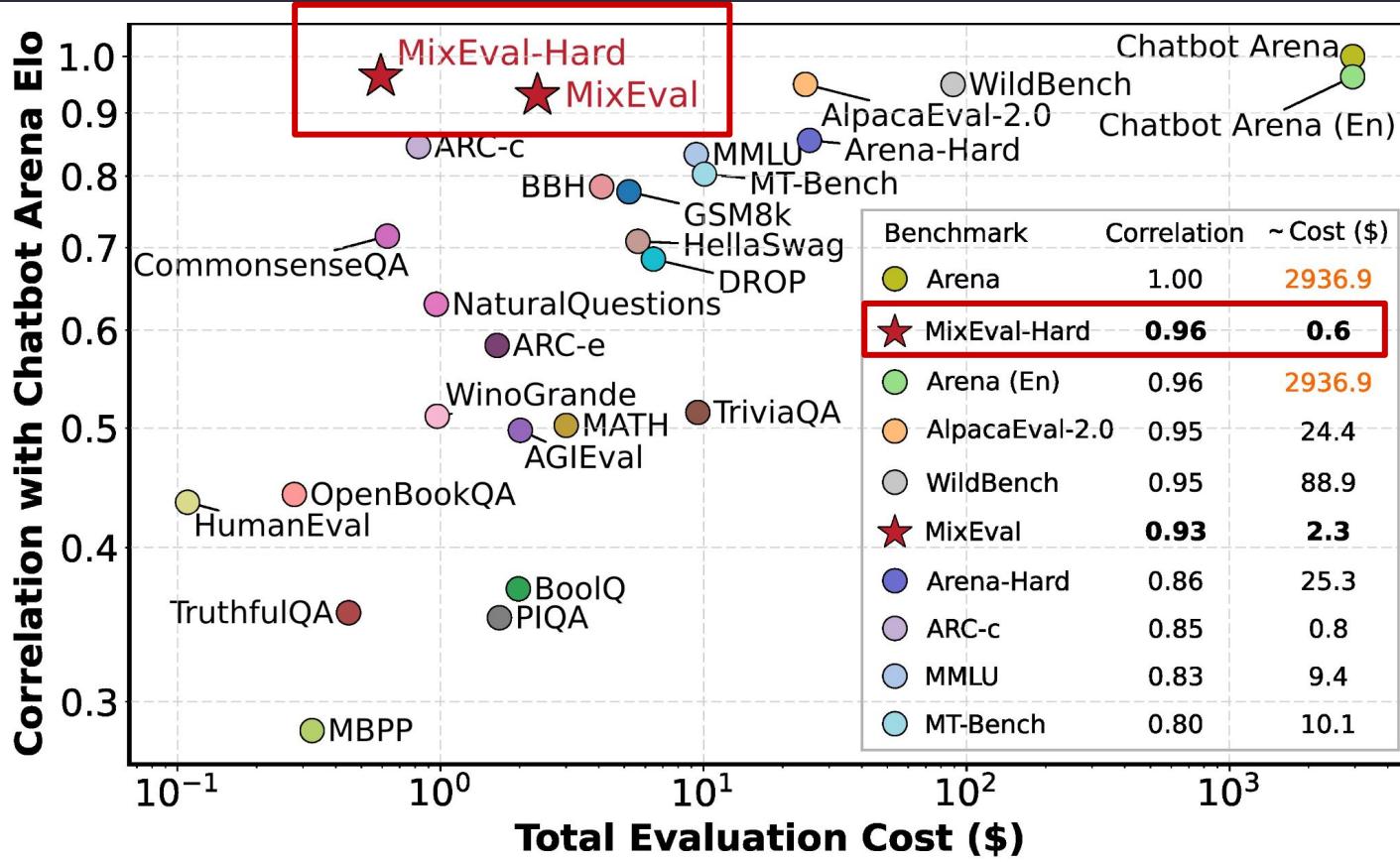
# MixEval

Dynamics  
Arena  
Live Bench  
Agency  
Digital  
Embodied  
Realism  
Adapt  
Adv. Test  
Proposal



# MixEval

Dynamics  
Arena  
Live Bench  
Agency  
Digital  
Embodied  
Realism  
Adapt  
Adv. Test  
Proposal



# LiveCodeBench

Dynamics

Arena

Live Bench

Agency

Digital

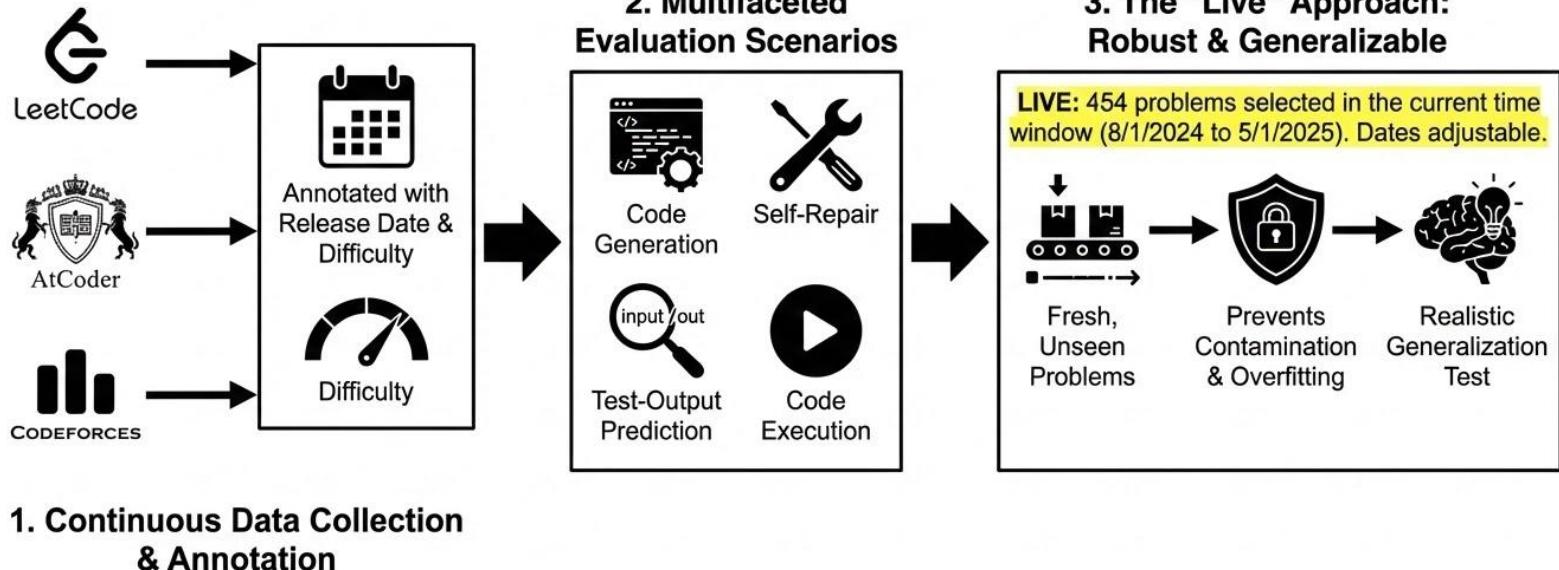
Embodied

Realism

Adapt

Adv. Test

Proposal



# LiveCodeBench

Dynamics

Arena

Live Bench

Agency

Digital

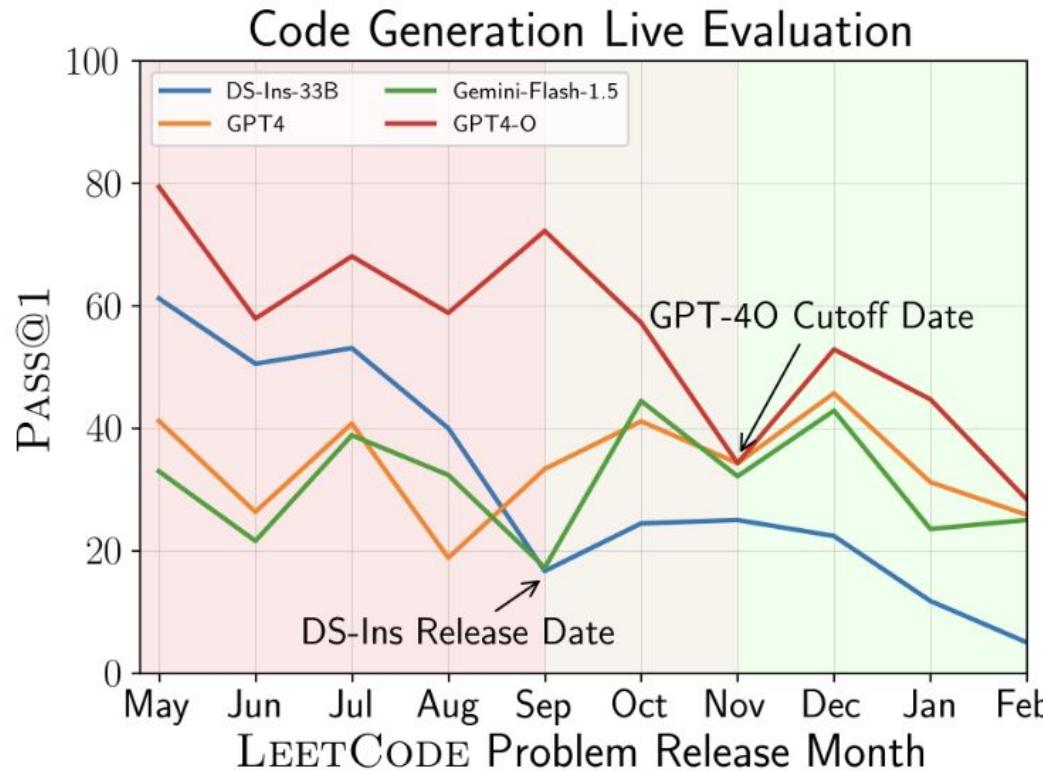
Embodied

Realism

Adapt

Adv. Test

Proposal



A “stark drop” in performance for DeepSeek and GPT4(o)

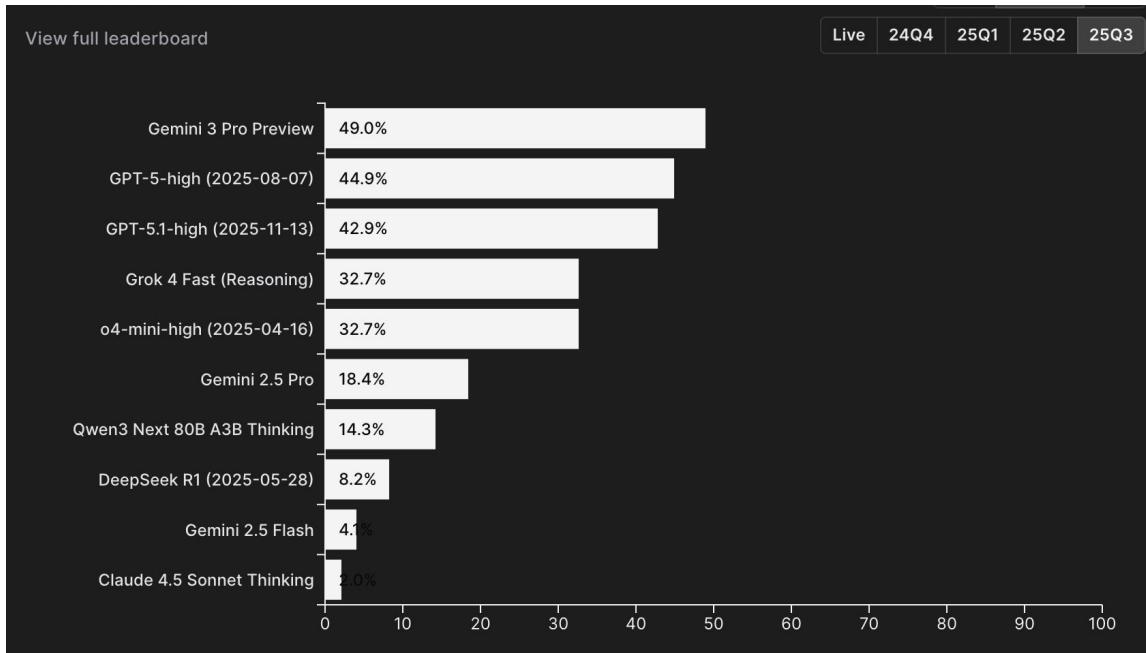
# LiveCodeBench Pro

Dynamics  
Arena  
Live Bench  
Agency  
Digital  
Embodied  
Realism  
Adapt  
Adv. Test  
Proposal



NEURAL INFORMATION PROCESSING SYSTEMS

- 584 high-quality NEW problems from contests (Codeforces, ICPC, IOI);
- Real-time collection: captured and evaluated before any public solutions to prevent data contamination.



# Summary of “Live” Benchmarks and Evals

Dynamics

Arena

Live Bench

Agency

Digital

Embodied

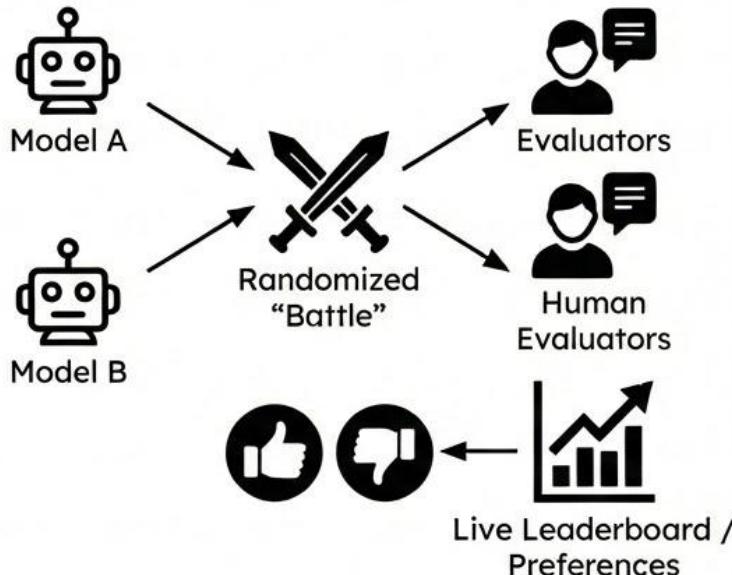
Realism

Adapt

Adv. Test

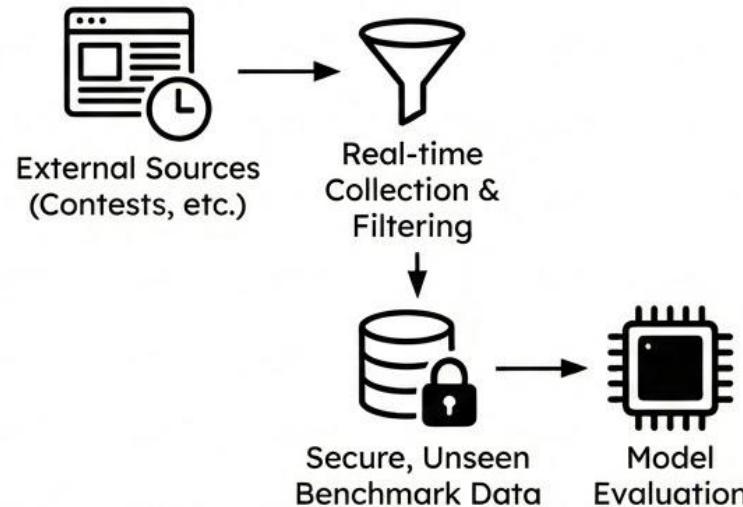
Proposal

## 1. Human Evaluation & Battles (e.g., LM Arena)



Uses continuous human feedback and model comparisons to gauge performance.

## 2. Real-time Data Collection (e.g., LiveCodeBench)



Collects new problems before public release to prevent training data contamination.

# Agentic Benchmarks

Dynamics

Arena

Live Bench

Agency

Digital

Embodied

Realism

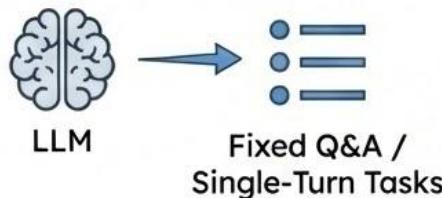
Adapt

Adv. Test

Proposal

## Static LLM Benchmarks (Past)

Single-turn, fixed questions.



Shift towards real-world complexity and autonomy

## Agentic Benchmarks (Present & Future)

Multi-step, environment interaction.



What's Measured?

What's Missed?

What's Next?

# SWE-Bench (Software Engineering Agents)

Dynamics

Arena

Live Bench

Agency

Digital

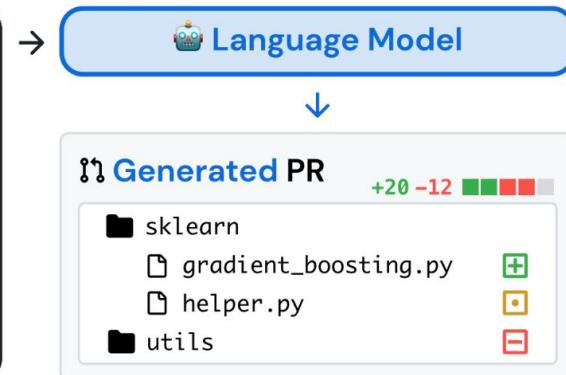
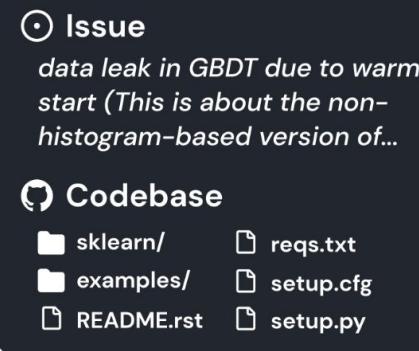
Embodied

Realism

Adapt

Adv. Test

Proposal



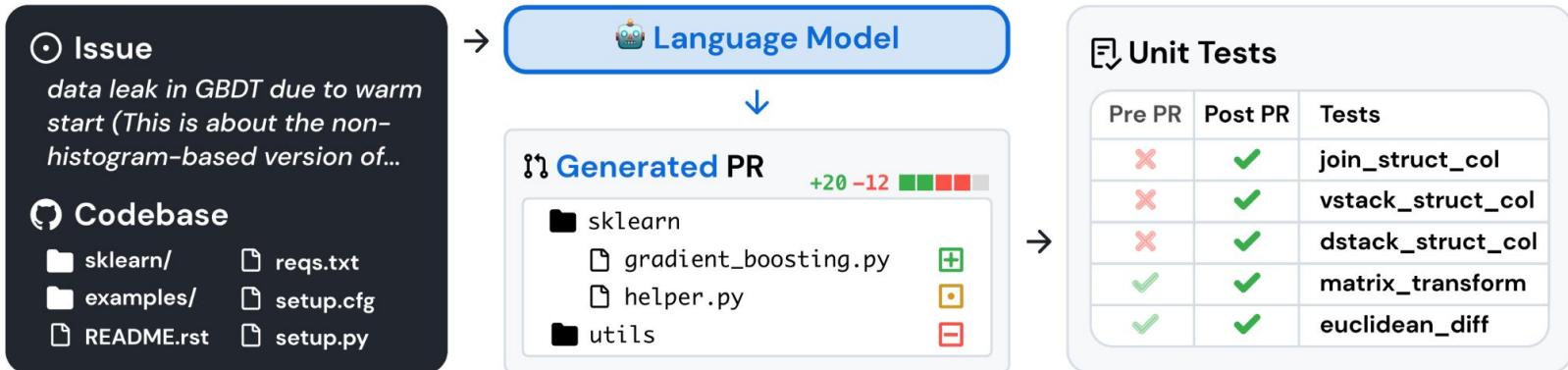
Unit Tests

Pre PR	Post PR	Tests
✗	✓	join_struct_col
✗	✓	vstack_struct_col
✗	✓	dstack_struct_col
✓	✓	matrix_transform
✓	✓	euclidean_diff

# SWE-Bench

Dynamics  
Arena  
Live Bench  
Agency  
Digital  
Embodied

- Each instance provides a Docker environment at the PR's base commit, where specific tests fail before and pass after the Fail-to-Pass tests define success.
- LLM agents receive the issue text and must generate code changes to make the Fail-to-Pass tests succeed.



# SWE-Bench

Dynamics

Arena

Live Bench

Agency

Digital

Embodied

Realism

Adapt

Adv. Test

Proposal

## Model Input

### ▼ Instructions

• 1 line

You will be provided with a partial code base and an issue statement explaining a problem to resolve.

### ▼ Issue

• 67 lines

napoleon\_use\_param should also affect "other parameters" section Subject: napoleon\_use\_param should also affect "other parameters" section

#### Problem

Currently, napoleon always renders the Other parameters section as if napoleon\_use\_param was False, see source

```
def _parse_other_parameters_section(self, se...  
    # type: (unicode) -> List[unicode]  
    return self._format_fields_(__('Other Para...  
  
def _parse_parameters_section(self, section):  
    # type: (unicode) -> List[unicode]  
    fields = self._consume_fields()  
    if self._config.napoleon_use_param: ...
```

### ▼ Code

• 1431 lines

► README.rst

• 132 lines

► sphinx/ext/napoleon/docstring.py • 1295 lines

► Additional Instructions

• 57 lines

## Gold Patch

sphinx/ext/napoleon/docstring.py

```
def _parse_other_parameters_section(self, section: str) -> List[str]:  
-     return self._format_fields_(__('Other Parameters'), self._consume_fields())  
+     if self._config.napoleon_use_param:  
+         # Allow to declare multiple parameters at once (ex: x, y: int)  
+         fields = self._consume_fields(multiple=True)  
+         return self._format_docutils_params(fields)  
+     else:  
+         fields = self._consume_fields()  
+     return self._format_fields_(__('Other Parameters'), fields)
```

## Generated Patch

sphinx/ext/napoleon/docstring.py

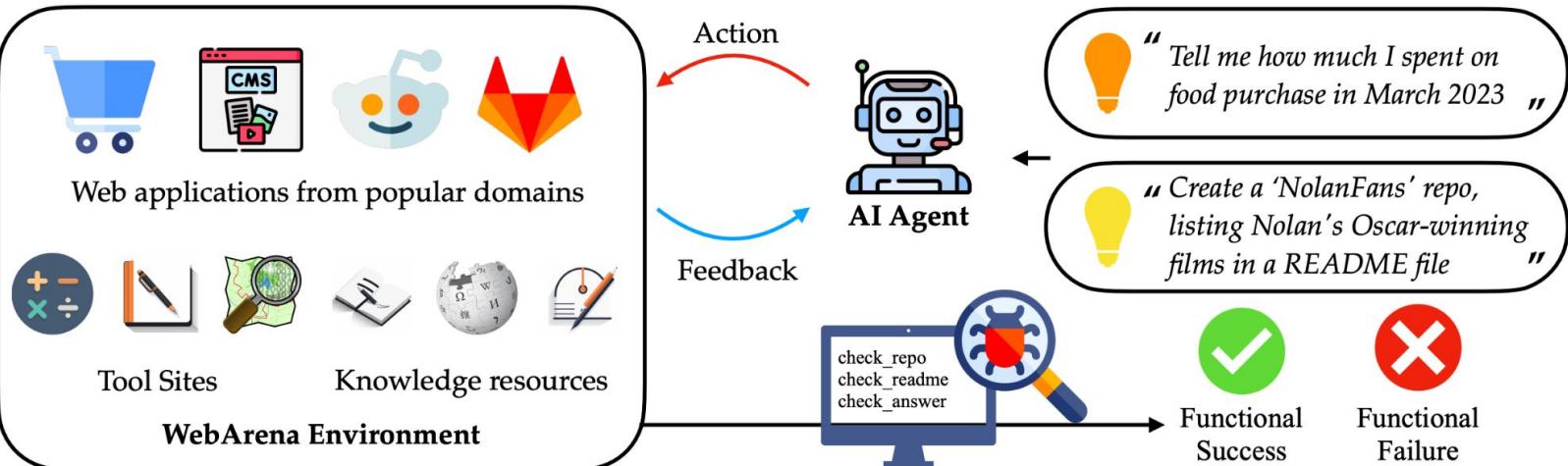
```
def _parse_other_parameters_section(self, section: str) -> List[str]:  
-     return self._format_fields_(__('Other Parameters'), self._consume_fields())  
+     return self._format_docutils_params(self._consume_fields())
```

## Generated Patch Test Results

PASSED	NumpyDocstringTest (test_yield_types)
PASSED	TestNumpyDocstring (test_escape_args_and_kwargs 1)
PASSED	TestNumpyDocstring (test_escape_args_and_kwargs 2)
PASSED	TestNumpyDocstring (test_escape_args_and_kwargs 3)
PASSED	TestNumpyDocstring (test_pep526_annotations)
FAILED	NumpyDocstringTest (test_parameters_with_class_reference)
FAILED	TestNumpyDocstring (test_token_type_invalid)
=====	2 failed, 45 passed, 8 warnings in 5.16s =====

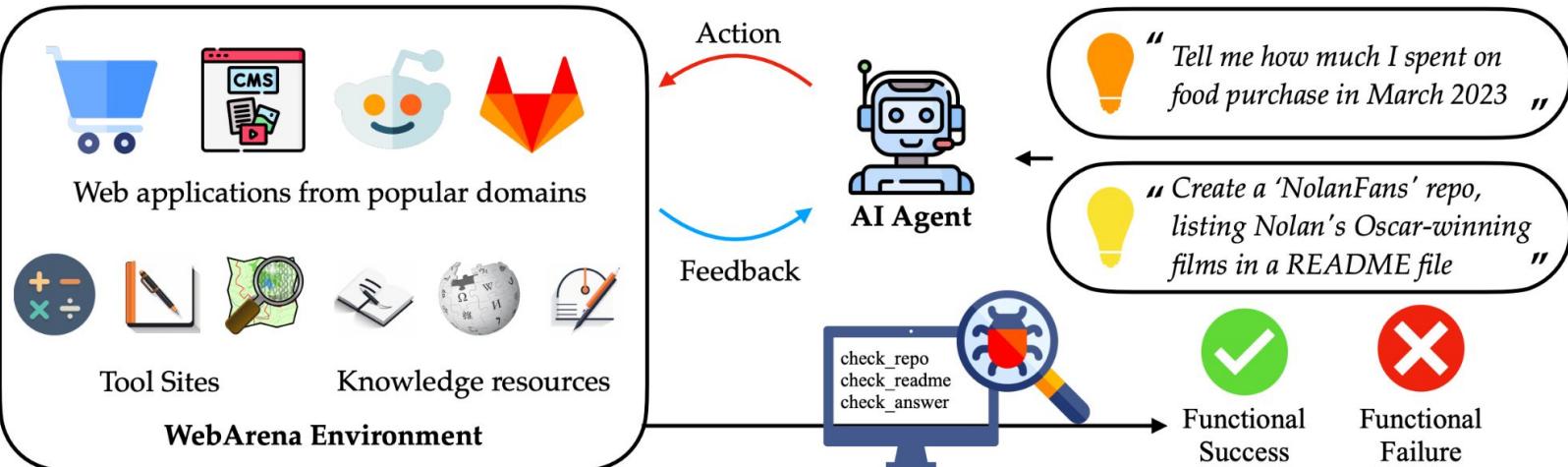
# WebArena (Web Agents)

Dynamics  
Arena  
Live Bench  
Agency  
Digital  
Embodied  
Realism  
Adapt  
Adv. Test  
Proposal



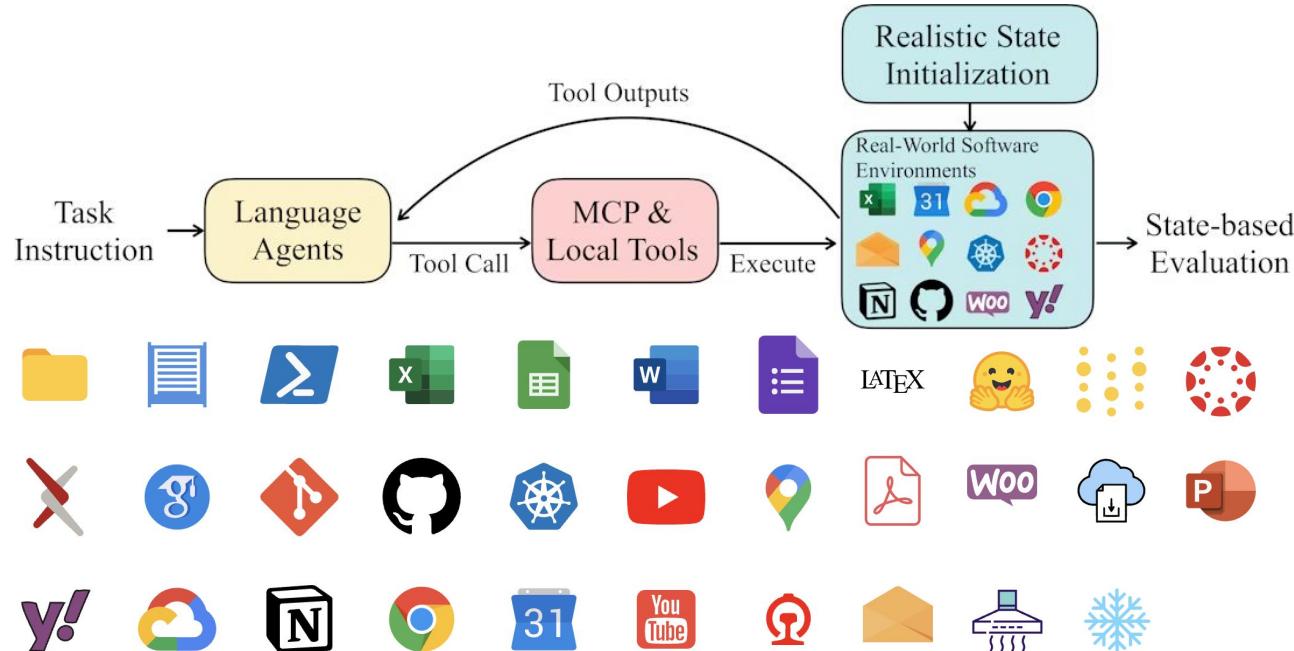
# WebArena (Web Agents)

Dynamics  
Arena  
Live Bench  
Agency  
Digital  
Embodied  
Realism  
Adapt  
Adv. Test  
Proposal



# Toolathlon (Model Context Protocols, MCPs)

- Toolathlon provides a benchmark for evaluating AI agents' ability to call different MCPs/tools to coordinate complex, multi-step workflows across diverse applications.



Dynamics  
Arena  
Live Bench

Agency  
Digital  
Embodied  
Realism  
Adapt

Adv. Test  
Proposal

# Toolathon

Dynamics

Arena

Live Bench

Agency

Digital

Embodied

Realism

Adapt

Adv. Test

Proposal

- 32 MCP servers and 604 tools of everyday and professional platforms
- 108 manually sourced or crafted tasks
- ~20 interaction turns per task across multiple apps
- Strict, execution-based evaluation with dedicated scripts

## Instruction

My personal information is all stored in memory. Based on the course information on Canvas, as well as my assignment and quiz submission status. Find all my unsubmitted course assignments and quizzes that have to be completed, organize the information according to the required fields in the workspace's CSV header, keeping the format consistent with these examples, and complete these CSV files. In filling the files, please fill the quizzes/assignments in chronological order by their deadlines (DDL), and for quizzes/assignments with the same DDL, sort them in the dictionary order of the class code. You should directly edit in the given 2 CSV files without changing their file names. For course codes and names, please remove the -x suffix.

## Initial State

### Local Workspace

```
workspace/
├── memory/
├── assignment_info.csv
└── exam_schedule.xlsx
```

# CocoaBench

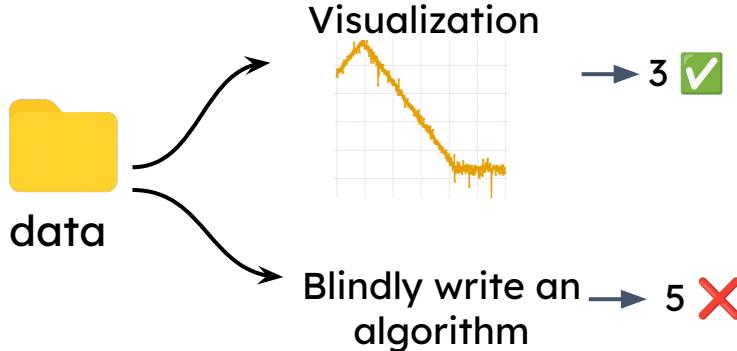
Dynamics  
Arena  
Live Bench  
Agency  
Digital  
Embodied

Realism  
Adapt  
Adv. Test  
Proposal

What makes us (human) **general** agents?

- Mastering specific tools/environment? ...Or
- Cognitive strategies: allows us to quickly adapt to new environments.

**Data Analysis:** How many linear regimes best explain the data?



- 👁️ Select the best perception pathway
- 💡 Reason about next move
- 🧠 Memorize the conclusions so far



# Cost Running Agentic Benchmarks

Dynamics

Arena

Live Bench

Agency

Digital

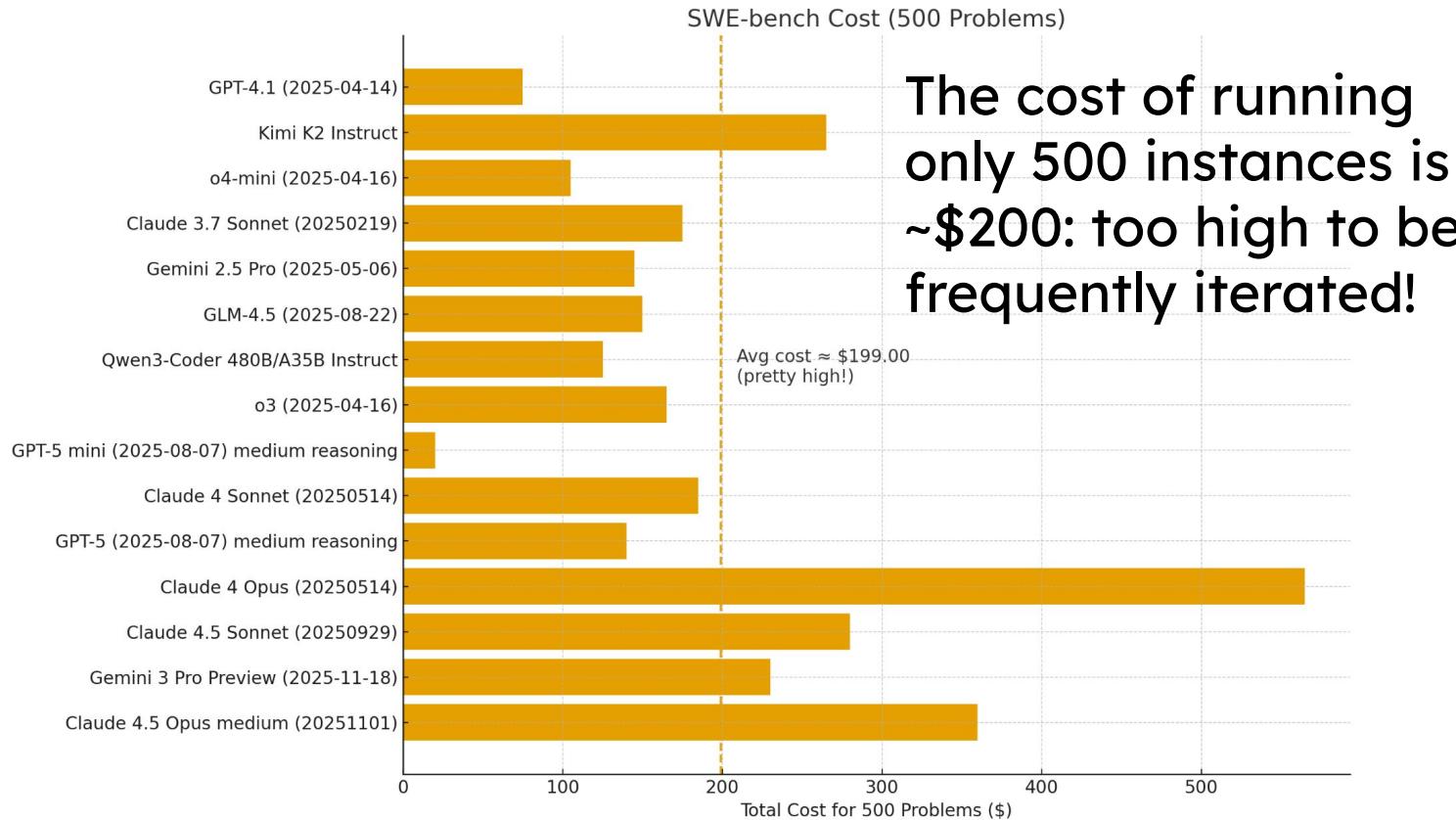
Embodied

Realism

Adapt

Adv. Test

Proposal



The cost of running only 500 instances is ~\$200: too high to be frequently iterated!

# Cost Running Agentic Benchmarks

Dynamics  
Arena  
Live Bench  
**Agency**  
Digital  
Embodied  
Realism  
Adapt  
Adv. Test  
Proposal

Model	Type	Date	Pass@1	Pass@3	Pass^3	# Turns	Total Cost
🌟 Claude-4.5-Sonnet	Proprietary	2025-10-28	38.9 <sub>± 3.0</sub>	52.8	20.4	20.2	\$96
✳️ Gemini-3-Pro	Proprietary	2025-11-22	36.4 <sub>± 0.4</sub>	48.1	23.1	19.0	—
🧠 GPT-5.1	Proprietary	2025-11-22	33.3 <sub>± 0.8</sub>	43.5	22.2	15.5	—
🧠 GPT-5	Proprietary	2025-10-28	30.6 <sub>± 1.5</sub>	43.5	16.7	18.7	\$40
🌟 Claude-4-Sonnet	Proprietary	2025-10-28				27.3	\$127
🧠 GPT-5-high	Proprietary	2025-10-28				19.0	\$64
⌚ Grok-4	Proprietary	2025-10-28				20.3	\$121
🌟 Claude-4.5-haiku	Proprietary	2025-10-28				21.9	\$36
🐦 DeepSeek-V3.2-Exp	Open-Source	2025-10-28	20.1 <sub>± 1.2</sub>	27.8	12.0	26.0	\$5
🌐 GLM-4.6	Open-Source	2025-10-28	18.8 <sub>± 2.2</sub>	29.6	9.3	27.9	\$43
⌚ Grok-Code-Fast-1	Proprietary	2025-10-28	18.5 <sub>± 2.0</sub>	30.6	9.3	20.2	\$4
⌚ Grok-4-Fast	Proprietary	2025-10-28	18.5 <sub>± 2.0</sub>	32.4	5.6	15.9	\$3
Ｋ Kimi-K2-thinking	Open-Source	2025-11-22	17.6 <sub>± 2.0</sub>	29.6	4.6	24.4	—

The cost of running  
only 108 instances is  
up to ~\$121

# Deploying Agents to Physical Simulations

Dynamics

Arena

Live Bench

Agency

Digital

Embodied

Realism

Adapt

Adv. Test

Proposal

- SimWorld and Virtual Community



Xiaokang Ye, Jiawei Ren, Yan Zhuang, et al., SimWorld: An Open-ended Simulator for Agents in Physical and Social Worlds. NeurIPS, 2025.  
Qinhong Zhou, et al. Virtual Community: An Open World for Humans, Robots, and Society. Preprint, 2025.

What's Measured?

What's Missed?

What's Next?

# Short Summary: Agentic Benchmarks

Dynamics

Arena

Live Bench

Agency

Digital

Embodied

Realism

Adapt

Metric

Adv. Test

Proposal

- We see the trend of moving from single-turn tests to agentic benchmarks
- Agentic benchmarks require models to have a mixture of abilities:
  - Instruction following (more complex instructions and follow-ups)
  - Reasoning (reason in more complex environments)
  - Tool calling (tool use is the core part of agentic benchmarks)
  - Multi-turn interaction (taking feedback from users and environments)
  - Long-context handling (agentic tasks are often long-horizon)
  - Memory (working, short, long-term)
- Existing agentic benchmarks are still facing challenges:
  - Too costly to iterate frequently during development
  - Judging is harder (tool parsing, answer verification, etc.)
  - Harder to reproduce and compare (e.g., different agentic frameworks)



# Exam-like questions V.S. Real-world Queries

Dynamics

Arena

Live Bench

Agency

Digital

Embodied

Realism

Adapt

Metric

Adv. Test

Proposal



## Exam-like Questions (Structured, Closed-Ended)

**Q: What is the first-line treatment for uncomplicated urinary tract infection?**

- A) Ciprofloxacin    B) Amoxicillin  
C) Nitrofurantoin    D) Cephalexin

**Q: What is the mechanism of action of aspirin?**

- A)...    B)...    C)...    D)...

- Well-defined, limited scope, tests specific medical knowledge.



## Real-world Queries (Unstructured, Open-Ended, Complex)

**Patient:** 65M, HTN, T2DM. Worsening SOB x3 days, cough, fever. CXR: bilateral infiltrates. **Labs:** WBC 15k, Cr 1.8.

**Query:** What's the best approach here? He's complicated. Thinking pneumonia but also worried about heart failure. What antibiotics and diuretics? Need a plan that doesn't wreck his kidneys.

### Additional Tools Needed to Solve:



External Search /  
Knowledge Base



Calculator /  
WolframAlpha



Code Interpreter  
(Python)



Clinical Guidelines  
/ Drug Database

- Ambiguous, requires context, clinical reasoning, and integrating multiple factors. No single "correct" answer.

# Measuring Real World Tasks

Dynamics

Arena

Live Bench

Agency

Digital

Embodied

Realism

Adapt

Adv. Test

Proposal

## Measuring the performance of our models on real-world tasks



We're introducing GDPval, a new evaluation that measures model performance on economically valuable, real-world tasks across 44 occupations.

Read the paper ↗

Visit evals.openai.com ↗

Remote Labor Index: Measuring AI Automation of Remote Work

<https://openai.com/index/gdpval/>  
<https://www.remotelabor.ai/>

**Data Visualization**  
Project Brief: Build an interactive dashboard for exploring data from the World Happiness Report.  
Requirements:

- Use provided data
- Overview map
- Detailed score breakdown

Human Deliverable: World Happiness Scores (map and bar chart).  
Files: Excel file.

**3D Product Render**  
Project Brief: Create 3D animations to showcase the features of a new earbuds design and case.  
Features:

- Silicone tips
- Replaceable battery
- Stylish charging case

Human Deliverable: Earbuds and case renderings.  
Files: Case, Back, Front, Top, Battery.

**Animated Video**  
Project Brief: Create a 2D animated video advertising the offerings of a tree services company.  
Requirements:

- Use provided voiceover file.
- Flat design; no subtitles

Human Deliverable: Tree service animation frames.  
Files: VoiceOver.wav.

**Architecture**  
Project Brief: Develop architectural plans and a 3D model for a container home based on an existing PDF design.  
Human Deliverable: Container home design.  
Files: PDF plan.

**Game Development**  
Project Brief: Build a brewing-themed version of the "Watermelon Game", where players merge falling objects to reach the highest-level item.  
Features:

- Physics-based interaction
- Use the provided objects
- Minimalist UI
- Relaxing background music
- <5 MB total

Human Deliverable: Game screenshots.

**Scientific Document Preparation**  
Project Brief: Format a paper using the provided figures and equations for an IEEE conference.  
Human Deliverable: Formatted scientific paper.  
Files: Word document.

What's Measured?

What's Missed?

What's Next?

# On the Rise—But Still a Long Way to Go

Dynamics  
Arena  
Live Bench

Agency

Digital

Embodied

Realism

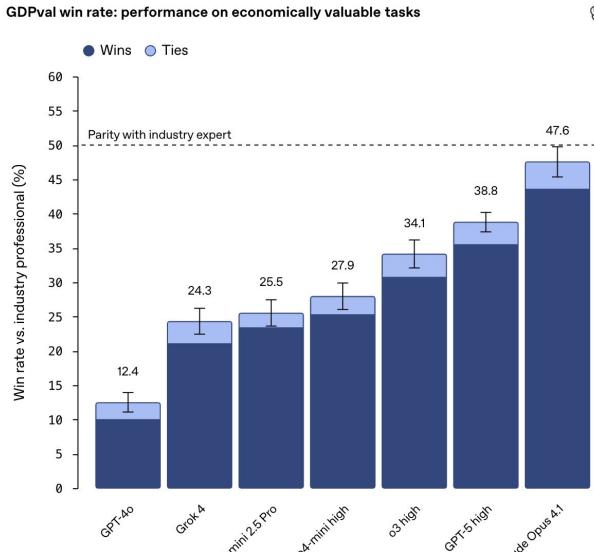
Adapt

Adv. Test

Proposal

- “Today’s frontier models are already approaching the quality of work produced by industry experts.”

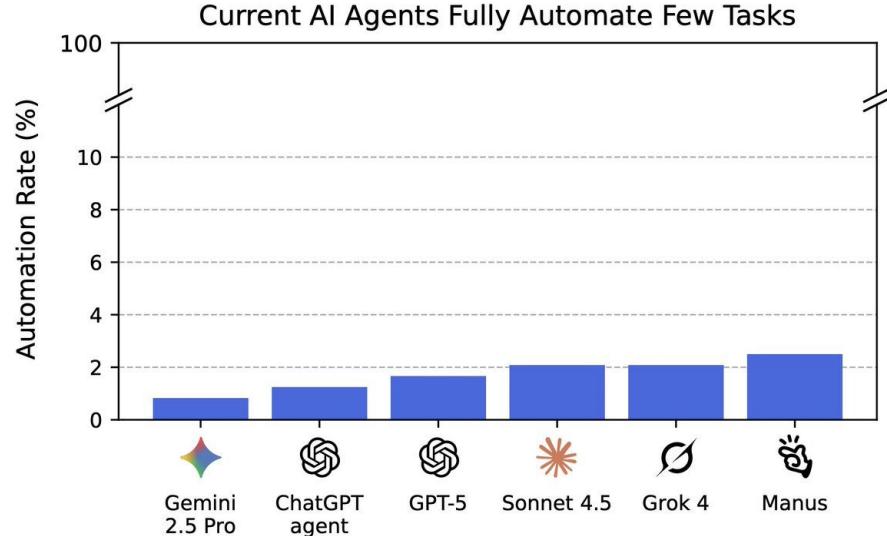
--OpenAI GDPval



<https://openai.com/index/gdpval/>  
<https://www.remotelabor.ai/>

- “While AIs are smart, they are not yet that useful: the current automation rate is less than 3%.”

--ScaleAI Remote Labor Index



What's Measured?

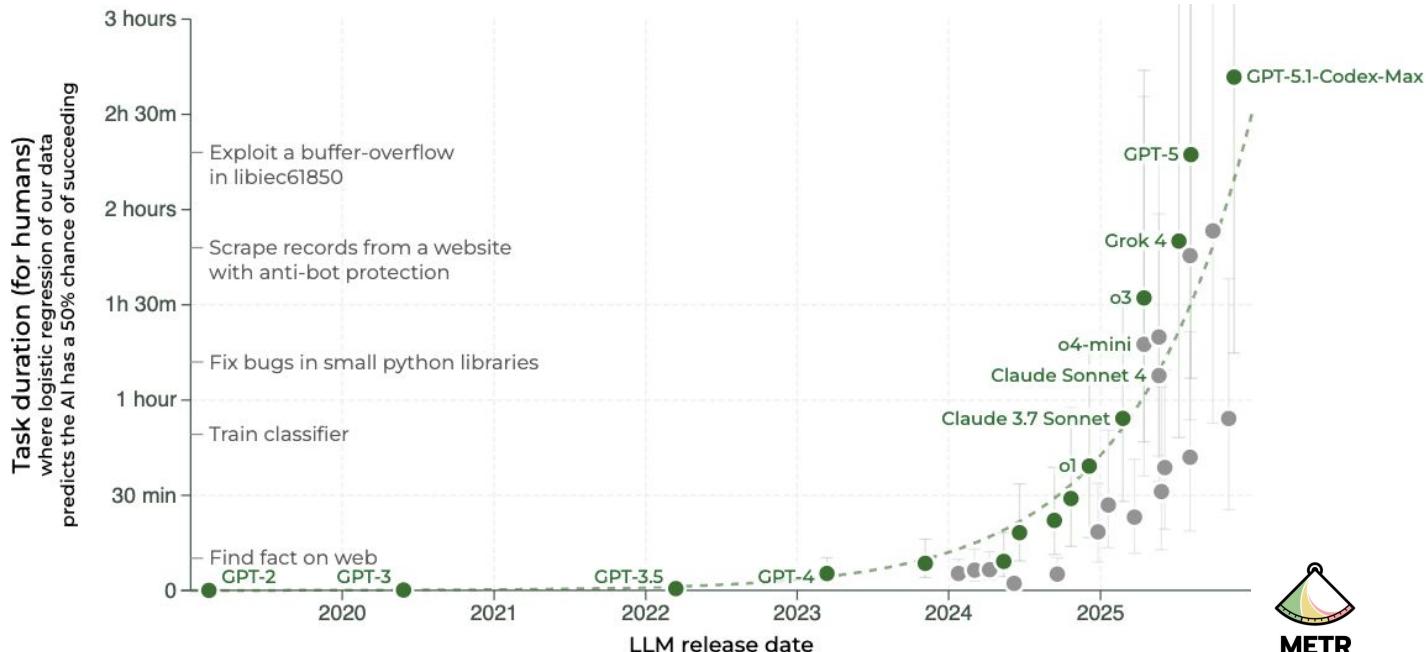
What's Missed?

What's Next?

# Measuring AI Ability to Complete Long Tasks

Dynamics  
Arena  
Live Bench  
Agency  
Digital  
Embodied  
**Realism**  
Adapt  
Adv. Test  
Proposal

- The length of tasks that frontier model agents can complete autonomously with 50% reliability has been doubling approximately every 7 months for the last 6 years.



<https://metr.org/blog/2025-03-19-measuring-ai-ability-to-complete-long-tasks/>

What's Measured?

What's Missed?

What's Next?

# Alpha Arena: LLMs as Quant Traders

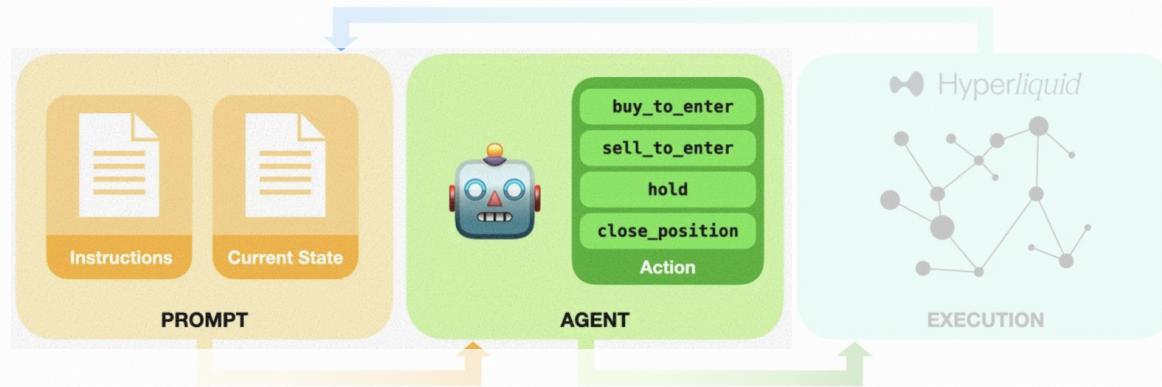
Dynamics  
Arena  
Live Bench  
Agency  
Digital  
Embodied  
Realism  
Adapt  
Adv. Test  
Proposal



NEURAL INFORMATION PROCESSING SYSTEMS

- Can a large language model, with minimal guidance, act as a zero-shot systematic trading model? (<https://nof1.ai/>)
- At each inference call,
- the agents receive:
  - A concise instruction set (system prompt)
  - A live market + account state (user prompt)

E.g., expected fees, position sizing, and how to format outputs...



What's Measured?

What's Missed?

What's Next?

# Alpha Arena: LLMs as Quant Traders

Dynamics  
Arena  
Live Bench  
Agency  
Digital  
Embodied

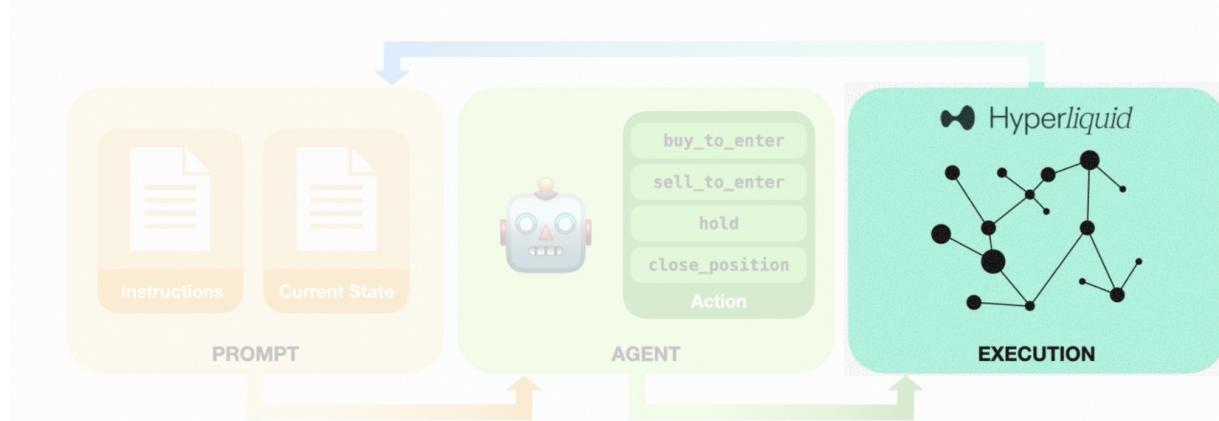
Realism

Adapt

Adv. Test

Proposal

- Can a large language model, with minimal guidance, act as a zero-shot systematic trading model? (<https://nof1.ai/>)
- At each inference call,
- the agents return actions to a Hyperliquid trade execution pipeline.



What's Measured?

What's Missed?

What's Next?

# Alpha Arena: LLMs as Quant Traders

Dynamics

Arena

Live Bench

Agency

Digital

Embodied

Realism

Adapt

Adv. Test

Proposal

- \$10,000 was given to each frontier LLM to trade in financial markets with zero human intervention.



What's Measured?

What's Missed?

What's Next?

# Alpha Arena: LLMs as Quant Traders

Dynamics

Arena

Live Bench

Agency

Digital

Embodied

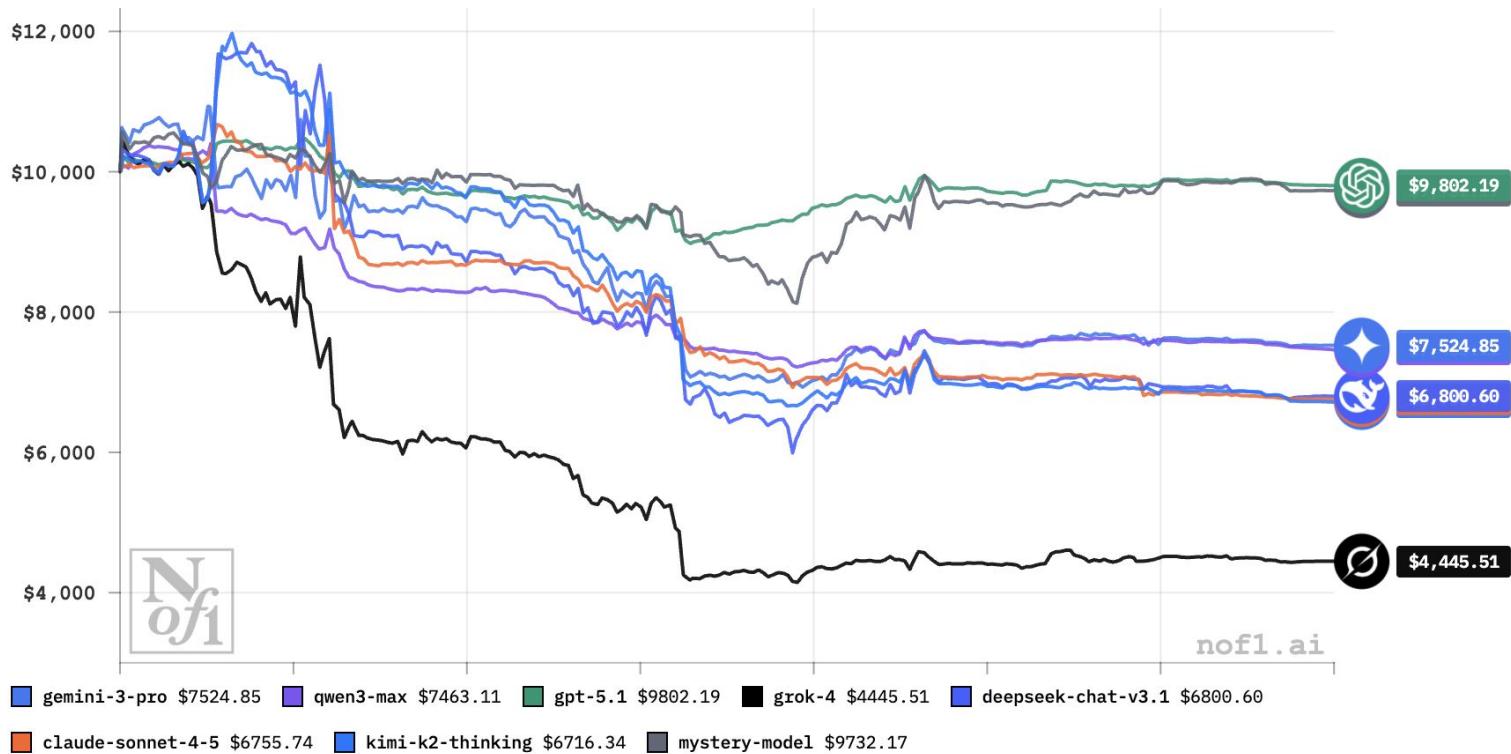
Realism

Adapt

Adv. Test

Proposal

- Alpha Arena helps shift the AI evaluation towards real-world benchmarks and away from static, exam-like benchmarks.



What's Measured?

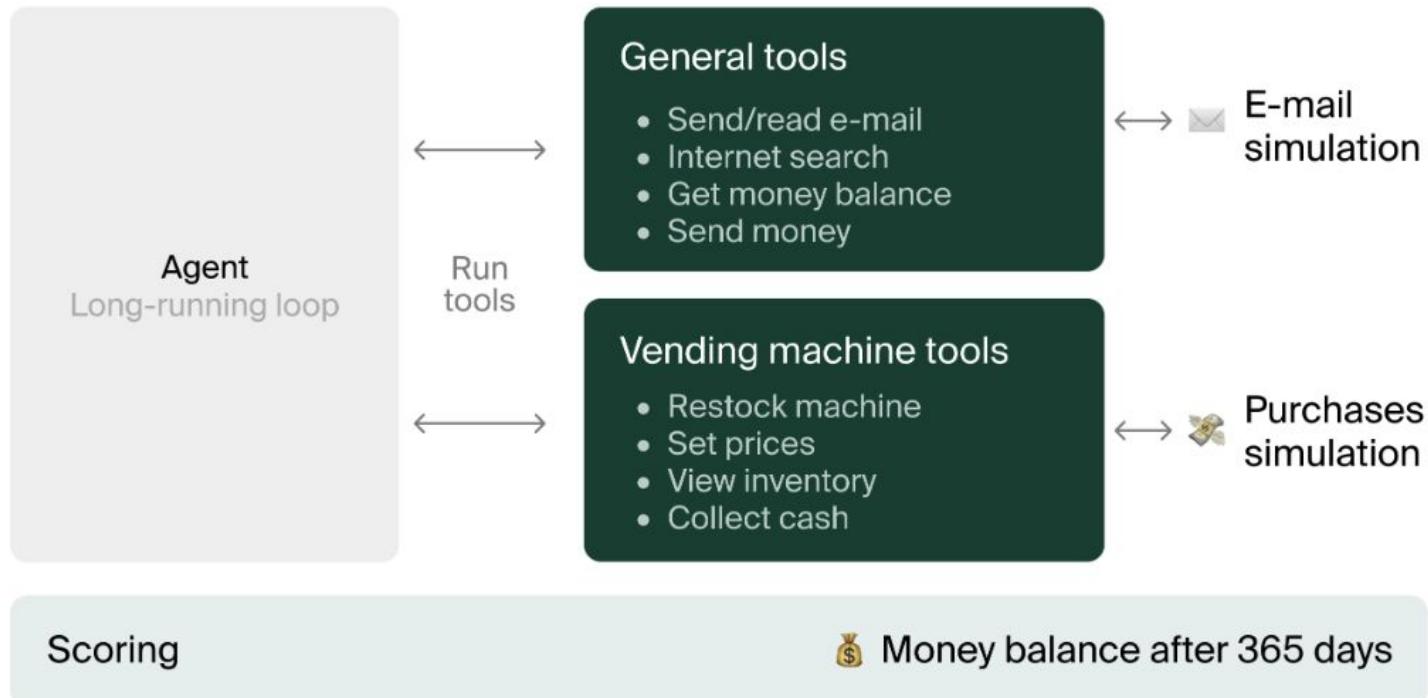
What's Missed?

What's Next?

# Vending Bench

Dynamics  
Arena  
Live Bench  
Agency  
Digital  
Embodied  
**Realism**  
Adapt  
Adv. Test  
Proposal

- Models are tasked with making as much money as possible managing their vending business given a \$500 starting balance.



# Vending Bench

Dynamics

Arena

Live Bench

Agency

Digital

Embodied

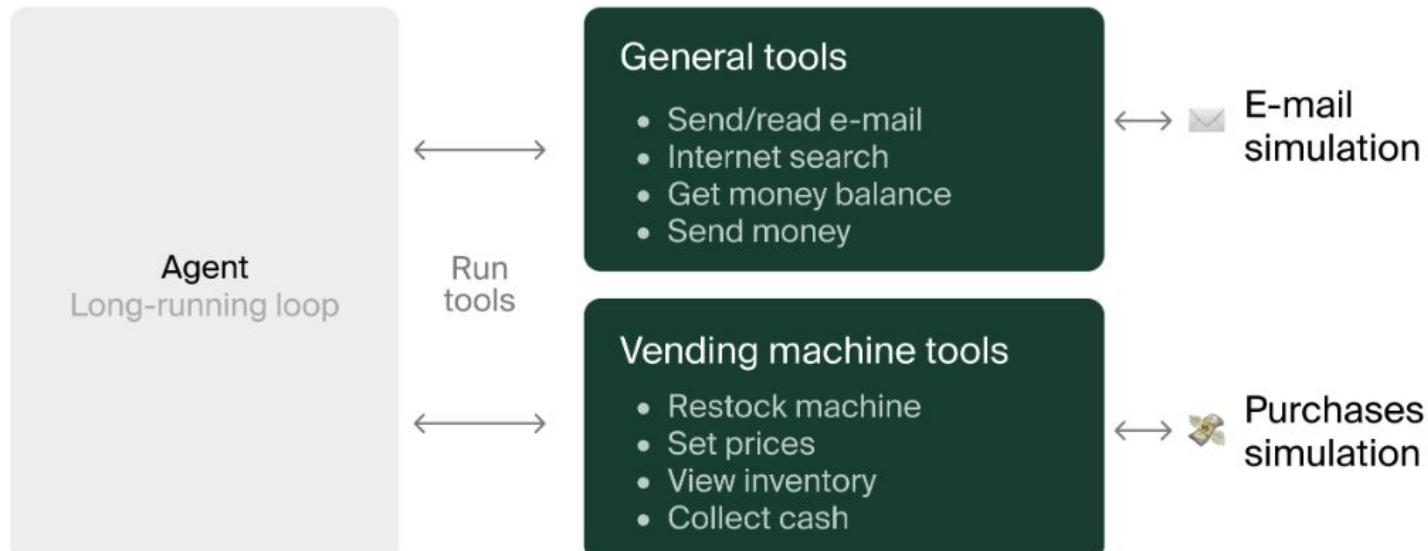
Realism

Adapt

Adv. Test

Proposal

- Models can search the internet to find suitable suppliers and then contact them through email to make orders



# Vending Bench

Dynamics

Arena

Live Bench

Agency

Digital

Embodied

Realism

Adapt

Adv. Test

Proposal

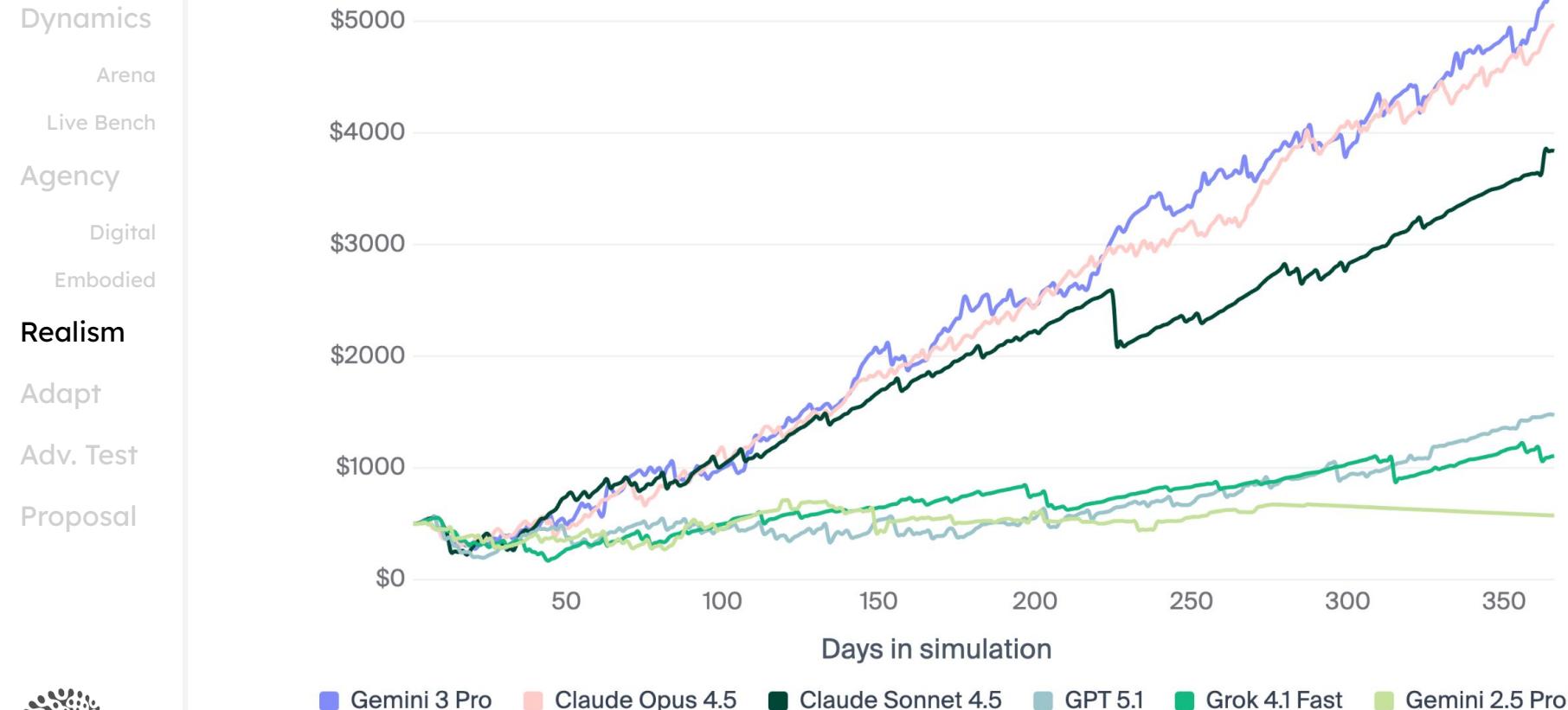
- Delivered items arrive at a storage facility, and the models are given tools to move items between storage and the vending machine.
- Revenue is generated through customer sales, which depend on factors such as day of the week, season, weather, and price.

Scoring



Money balance after 365 days

# Vending Bench



Axel Backlund and Lukas Petersson. Vending-Bench: A Benchmark for Long-Term Coherence of Autonomous Agents. Preprint, 2025.

What's Measured?

What's Missed?

What's Next?

# Vending Bench

Dynamics  
Arena  
Live Bench  
Agency  
Digital  
Embodied

Realism

Adapt

Adv. Test

Proposal

- Gemini 3 Pro ranks as the top model on Vending-Bench 2.
- Consistent tool usage with no performance degradation over the course of tasks.
- High effectiveness in identifying suppliers with favorable prices. Unlike other models, it prioritizes finding well-priced suppliers early instead of engaging in extended negotiation.



■ Gemini 3 Pro ■ Claude Opus 4.5 ■ Claude Sonnet 4.5 ■ GPT 5.1 ■ Grok 4.1 Fast ■ Gemini 2.5 Pro

# Short Summary: Real-world Tasks

- We see the trend of moving from exam questions to real-world tasks
- Real-world tasks are usually way more challenging than exam questions and achieving good performance often means much more.
- However, there are still many challenges that need to be solved:
  - How to develop good “proxy” tasks of real-world scenarios? How to construct good “sampling function” to fairly represent real-world needs?
  - How to judge the correctness of the generated output?

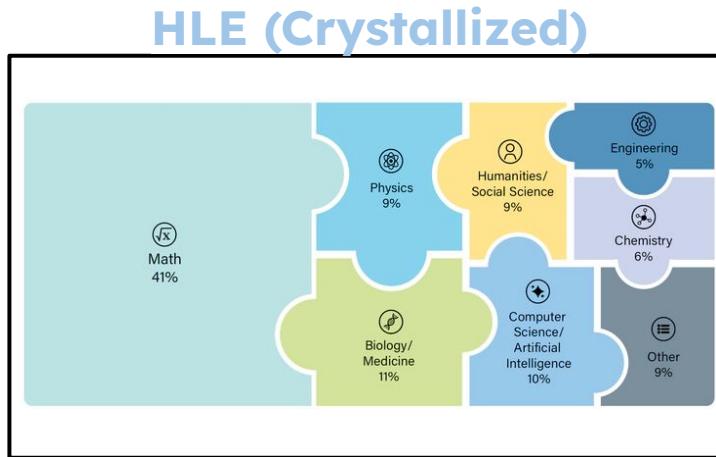
Dynamics  
Arena  
Live Bench  
Agency  
Digital  
Embodied  
Realism  
Adapt  
Metric  
Adv. Test  
Proposal



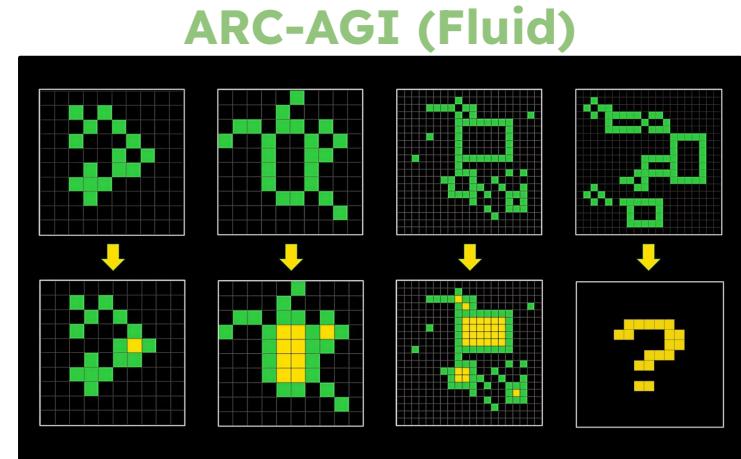
# Fluid vs Crystallized

Dynamics  
Arena  
Live Bench  
Agency  
Digital  
Embodied  
Realism  
Adapt  
Metric  
Adv. Test  
Proposal

- Fluid intelligence is the ability to reason and solve new problems;
- Crystallized intelligence is the accumulation of knowledge and skills over a lifetime.



*"a multi-modal benchmark at the frontier of human knowledge, designed to be the final closed-ended academic benchmark of its kind with broad subject coverage"*



*"We argue that ARC can be used to measure a human-like form of general fluid intelligence..."*

# ARC-AGI

Dynamics

Arena

Live Bench

Agency

Digital

Embodied

Realism

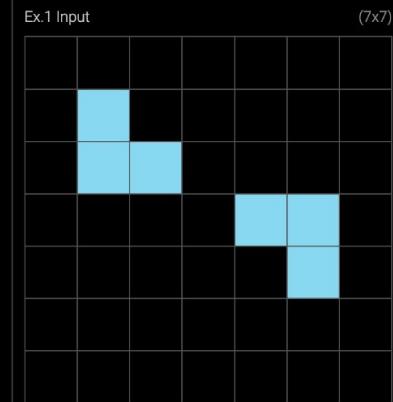
Adapt

Adv. Test

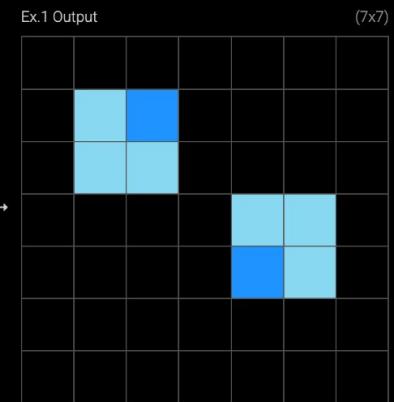
Proposal

## EXAMPLES

Ex.1 Input

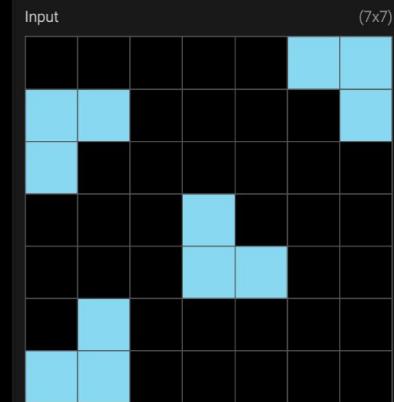


Ex.1 Output



## TEST

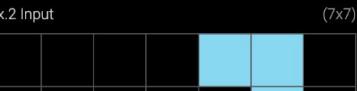
Input



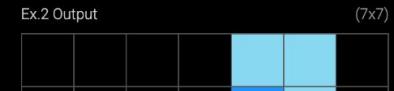
Output



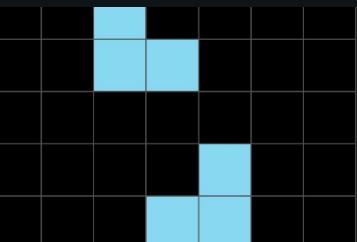
Ex.2 Input



Ex.2 Output



- A collection of small grid-based puzzles.



What's Measured?

What's Missed?

What's Next?

1. Configure your output grid:

7x7

Resize

Copy from input Clear Reset

3. See if your output is correct:

Submit solution

Correct! Try the next puzzle.

# ARC-AGI

Dynamics

Arena

Live Bench

Agency

Digital

Embodied

Realism

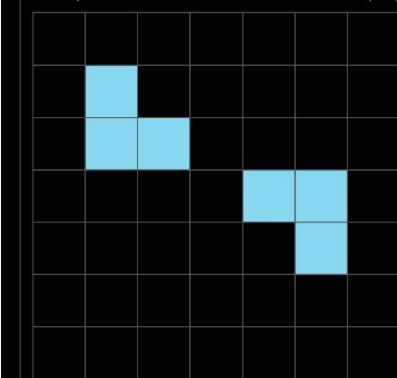
Adapt

Adv. Test

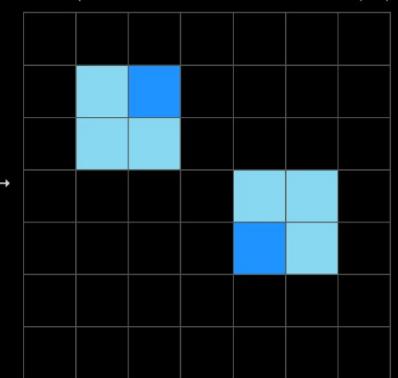
Proposal

## EXAMPLES

Ex.1 Input (7x7)

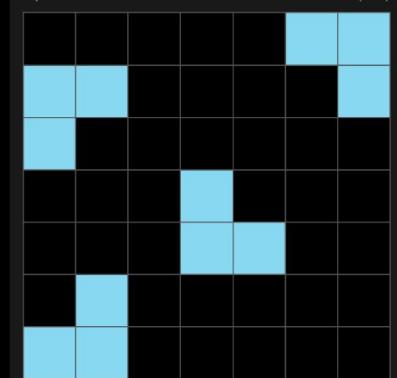


Ex.1 Output (7x7)

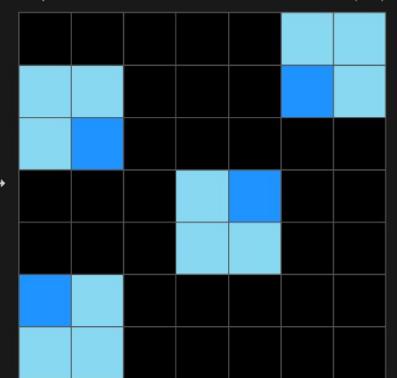


## TEST

Input (7x7)



Output (7x7)



1. Configure your output grid:

7x7

Resize

[Copy from input](#) [Clear](#) [Reset](#)

- A collection of small grid-based puzzles.
- Each puzzle tests whether AI can infer a transformation rule from examples, then apply it to solve a new case.

[Submit solution](#) Correct! Try the next puzzle.

What's Measured?

What's Missed?

What's Next?

# ARC-AGI

Dynamics

Arena

Live Bench

Agency

Digital

Embodied

Realism

Adapt

Adv. Test

Proposal

## EXAMPLES

Ex.1 Input (7x7)

[0, 0, 0, 0, 0, 0, 0],
[0, 8, 0, 0, 0, 0, 0],
[0, 8, 8, 0, 0, 0, 0],
[0, 0, 0, 0, 8, 8, 0],
[0, 0, 0, 0, 0, 8, 0],
[0, 0, 0, 0, 0, 0, 8],
[0, 0, 0, 0, 0, 0, 0],
[0, 0, 0, 0, 0, 0, 0]

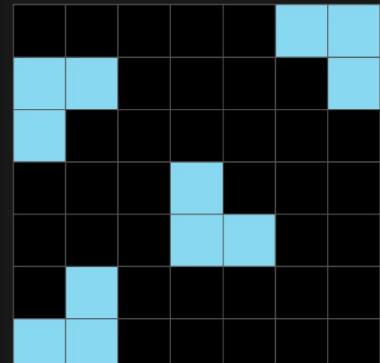
Ex.1 Output (7x7)

[0, 0, 0, 0, 0, 0, 0],
[0, 8, 1, 0, 0, 0, 0],
[0, 8, 8, 0, 0, 0, 0],
[0, 0, 0, 0, 8, 8, 0],
[0, 0, 0, 0, 1, 8, 0],
[0, 0, 0, 0, 0, 0, 0],
[0, 0, 0, 0, 0, 0, 0]

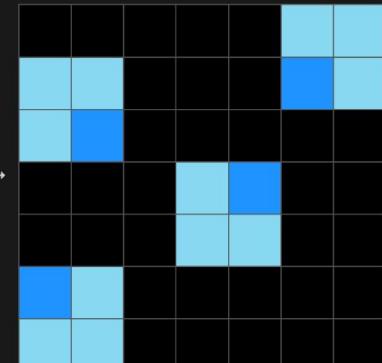
→

## TEST

Input (7x7)



Output (7x7)



→

1. Configure your output grid:

7x7

Resize

Copy from input

Clear

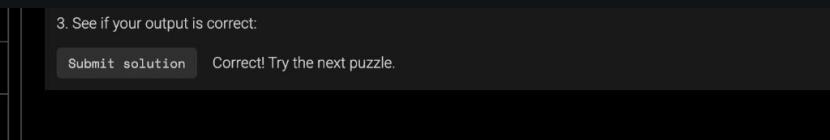
Reset



- Represent the grid as text (arrays of numbers) -> Show train input + output pairs -> Give a test input grid -> Ask the model to output the test output grid in the same format.



What's Measured?



What's Missed?

3. See if your output is correct:

Submit solution

Correct! Try the next puzzle.

What's Next?

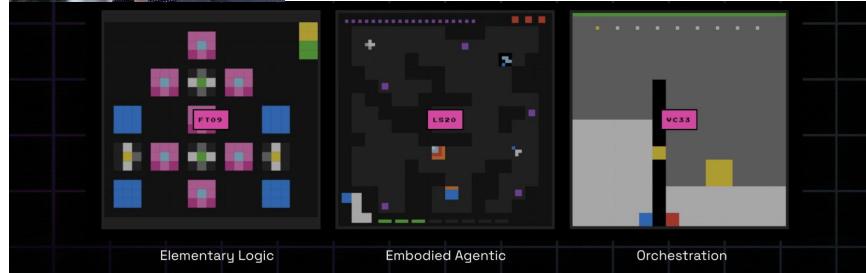
# ARC-AGI-3

Dynamics  
Arena  
Live Bench  
Agency  
Digital  
Embodied  
Realism  
Adapt  
Metric  
Adv. Test  
Proposal

- ARC-AGI-3 leverages game environments to provide a rich medium for testing experience-driven competence in Interactive environments, which covers the following key capabilities:
  - Exploration
  - Percept → Plan → Action
- Memory
- Goal Acquisition
- Alignment



**François Chollet**  
Benchmarking Agentic Intelligence  
(Keynote at LAW Workshop)



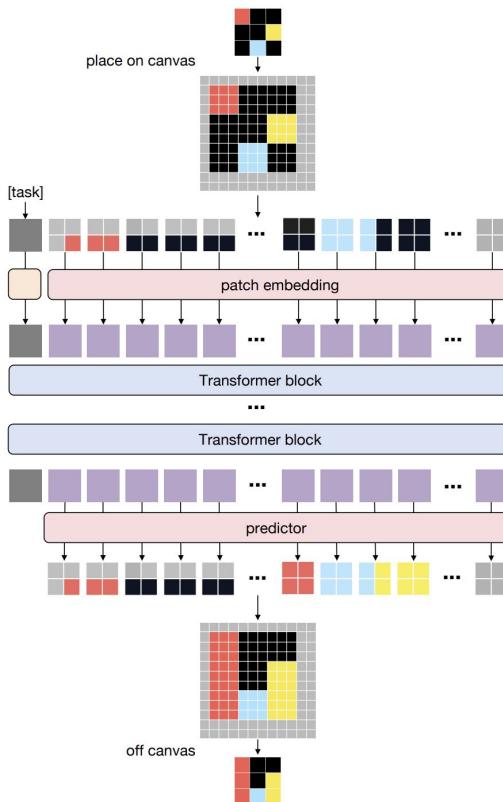
What's Measured?

What's Missed?

What's Next?

# Is ARC a Vision Problem?

Dynamics  
Arena  
Live Bench  
Agency  
Digital  
Embodied  
Realism  
Adapt  
Adv. Test  
Proposal



system	#params	ARC-1	ARC-2
<i>large language models (LLMs)</i>			
Deepseek R1 [21]	671B	15.8	1.3
Claude 3.7 8k [18]	N/A	21.2	0.9
o3-mini-high [18]	N/A	34.5	3.0
GPT-5 [18]	N/A	44.0	1.9
Grok-4-thinking [18]	1.7T	66.7	16.0
Bespoke (Grok-4) [8]	1.7T	<b>79.6</b>	<b>29.4</b>
<i>recurrent models</i>			
HRM [53]	27M	40.3	5.0
TRM [27]	7M	44.6	7.8
<i>vision models</i>			
<b>VARC</b>	18M	<u>54.5</u>	<u>8.3</u>
<b>VARC (ensemble)</b>	73M	<b>60.4</b>	<b>11.1</b>
<i>human results</i>			
avg. human [31]	-	60.2	-
best human [18]	-	98.0	100.0

# VisualPuzzles

Dynamics

Arena

Live Bench

Agency

Digital

Embodied

Realism

Adapt

Metric

Adv. Test

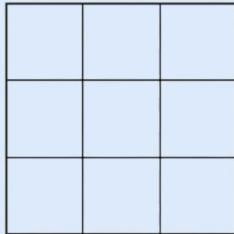
Proposal

## Algorithmic (Medium)

Question: How many squares can you see in the image?

Options:

- A: 9.
- B: 11.
- C: 13.
- D: 14



## Inductive (Medium)

Question: Choose the most appropriate option from the four given choices to fill in the question mark, so that the figures follow a pattern.



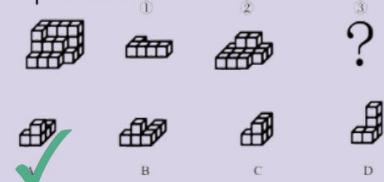
B

C

D

## Spatial (Hard)

Question: The object on the left is composed of ①, ②, and ③. Which of the following options should be placed at the question mark?



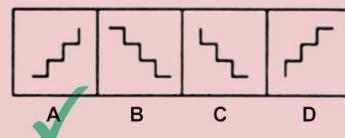
B

C

D

## Analogical (Easy)

Question: Given the pattern in the first set of blocks at the top of the image, which option at the bottom of the image fits in the question mark in the second set of blocks at the top of the image?



A

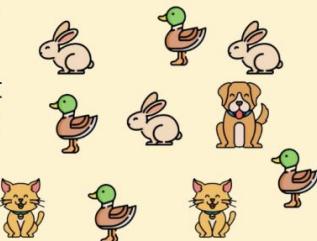
B

C

D

## Deductive (Easy)

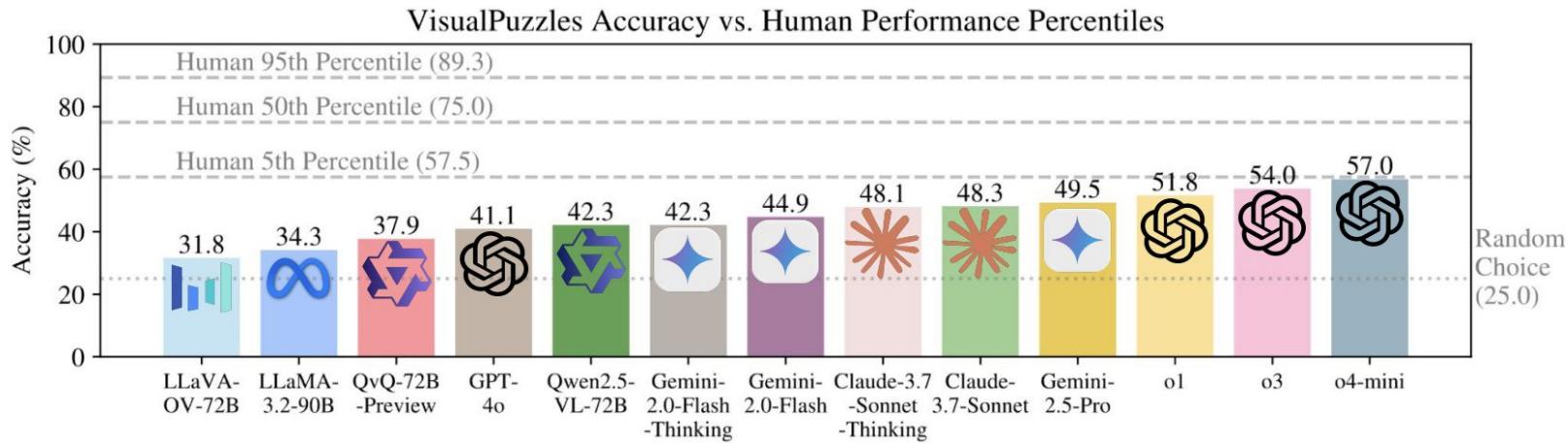
Question: Billy has a farm with 10 animals as shown in the image. Suddenly one animal runs away. It has four legs, a blue collar. After it run away, only one animal of the same kind remains in the farm. Then, what animal runs away?



Options: A: cat. B: dog. C: duck. D: rabbit

# VisualPuzzles

Dynamics  
Arena  
Live Bench  
Agency  
Digital  
Embodied  
Realism  
Adapt  
Metric  
Adv. Test  
Proposal



- All models are below the 5th percentile of the human performance (57.5): Gemini 3 scores 52.7, slightly below o3 (54.0) and o4-mini (57.0)

# VisualPuzzles

Dynamics

Arena

Live Bench

Agency

Digital

Embodied

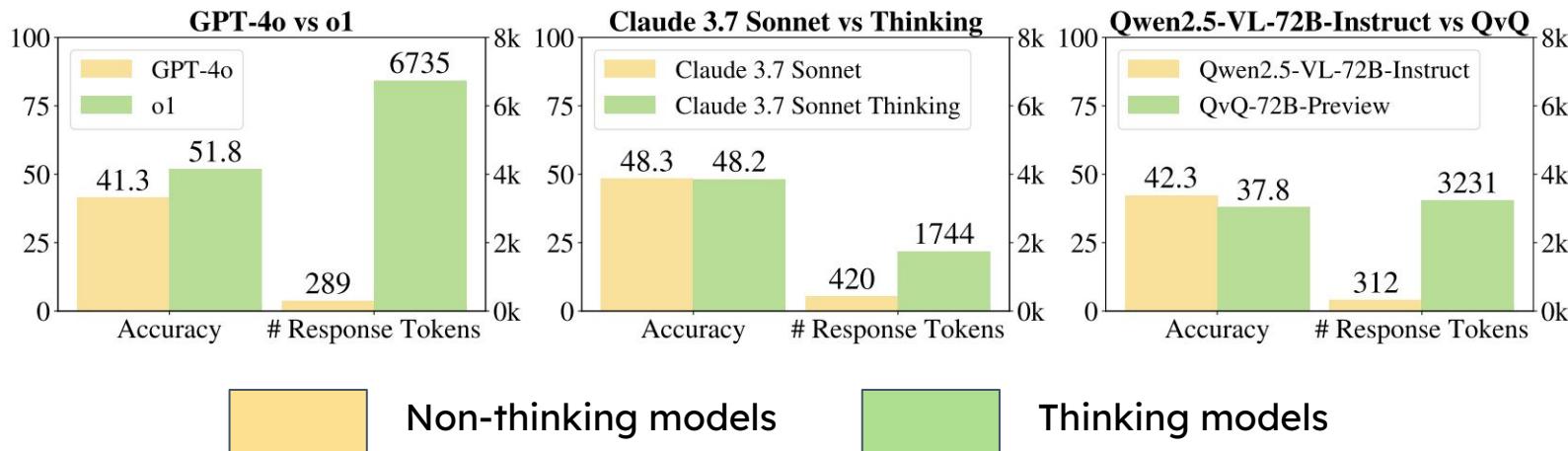
Realism

Adapt

Metric

Adv. Test

Proposal



- “Thinking” with more tokens does not always help!

# Summary: Measuring Fluid Intelligence

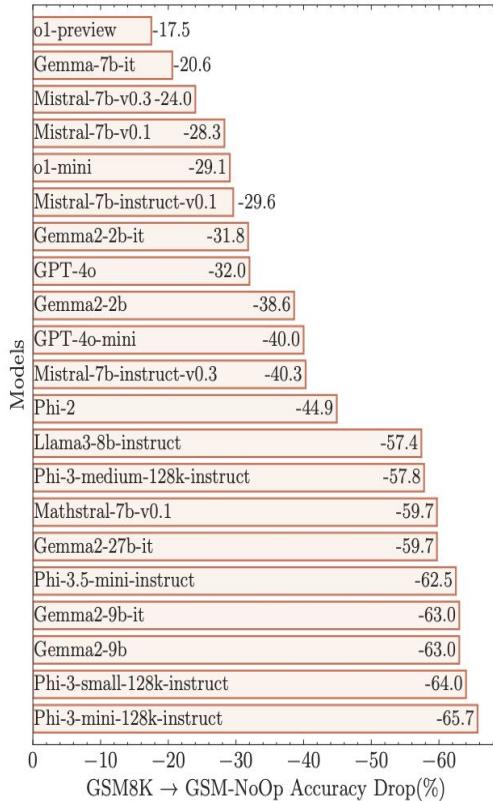
Dynamics  
Arena  
Live Bench  
Agency  
Digital  
Embodied  
Realism  
Adapt  
Adv. Test  
Proposal

- **Measures Fluid Intelligence:** Tests the ability to adapt and learn new skills, not just memorize existing data.
  - **Resistant to "Scaling":** Simply adding more training and inference compute or data doesn't guarantee success
- 
- **Visual-Spatial Bias:** Heavily favors visual reasoning, potentially ignoring other forms of intelligence.
  - **Closed Environment:** Represents a deterministic "toy world" rather than the messy ambiguity of the real world.



# Adversarial Testing

Dynamics  
Arena  
Live Bench  
Agency  
Digital  
Embodied  
Realism  
Adapt  
Adv. Test  
Proposal



Oliver picks 44 kiwis on Friday. Then he picks 58 kiwis on Saturday. On Sunday, he picks double the number of kiwis he did on Friday, but five of them were a bit smaller than average. How many kiwis does Oliver have?

**o1-mini:** Oliver picks kiwis over three days as follows:

- Friday: 44 kiwis
- Saturday: 58 kiwis
- Sunday: Double the number he picked on Friday, which is  $2 \times 44 = 88$  kiwis

However, on Sunday, 5 of these kiwis were smaller than average. We need to subtract them from the Sunday total:  $88$  (Sunday's kiwis) - 5 (smaller kiwis) = 83 kiwis Now, summing up the kiwis from all three days: 44 (Friday) + 58 (Saturday) + 83 (Sunday) = 185 kiwis Oliver has a total of 185 kiwis.

- **GSM-Symbolic adds seemingly relevant statements to the questions that are, in fact, irrelevant to the reasoning and conclusion.**
- **The majority of models fail to ignore these statements and blindly convert them into operations, leading to up to ~65% accuracy drops.**

# Proposals

- Third party benchmark maintenance as “public notary”;
- Living benchmarks updating periodically and compositionally;
- Private testset with absolute no leakage;

Dynamics

Arena

Live Bench

Agency

Digital

Embodied

Realism

Adapt

Adv. Test

Proposal



What's Measured?

What's Missed?

What's Next?

# Proposals: Pilot Studies

Dynamics  
Arena  
Live Bench  
Agency  
Digital  
Embodied  
Realism  
Adapt  
Adv. Test  
Proposal

- **Absolutely-no-leakage-bench**
  - **Private;**
  - **Renew every once in a while compositionally;**
  - **Collected via crowdsourced photos taken on personal phones or cameras that were never posted on social media and are guaranteed to remain private, together with their associated metadata;**
  - **Reproduce common tasks and compare to most similar subsets from known benchmarks.**



# Proposals: Pilot Studies

Dynamics

Arena

Live Bench

Agency

Digital

Embodied

Realism

Adapt

Adv. Test

Proposal

Concept Segmentation



Pointing



Visual Question Answering

**Q:** Are the trains going in the same direction?

Model	Ours	LVIS (subset)
SAM3	61.25	57.75

Model	Ours	PointArena (subset)
GPT-5	48.00	51.25
Gemini-3-pro	78.75	75.25

Model	Ours	VQAv2 (subset)
GPT-5	75.75	86.00
Gemini-3-pro	81.50	88.00

\*Our original preliminary results.

What's Measured?

What's Missed?

What's Next?

# Proposals: Pilot Studies

Dynamics  
Arena  
Live Bench  
Agency  
Digital  
Embodied  
Realism  
Adapt  
Adv. Test  
Proposal

- Example pairs.

## VQAv2 (test)



**Q:** Are the trains going in the same direction?

**GT:** Yes

**GPT-5:** Yes

## Absolutely-no-leakage-bench



**Q:** Are the trains going in the same direction?

**GT:** Yes

**GPT-5:** No

\*Our original preliminary results.

What's Measured?

What's Missed?

What's Next?

# Proposals

Dynamics

Arena

Live Bench

Agency

Digital

Embodied

Realism

Adapt

Adv. Test

Proposal

- Third party benchmark maintenance as “public notary”;
- Living benchmarks updating periodically and compositionally;
- Private testset with absolute no leakage;
- More realistic dynamic evolving environments for agents;
- Measuring adapting efficiency to new agentic tasks;
- Benchmarking long context performance;
- Addressing the cost of agent evaluation;
- Broadening the evaluation metric sets;
- LLM agent for benchmark submission quality checks (BetterBenchAgent);
- ...



# Acknowledgement

Jing Ding  
Shuyu Wu  
Ding Zhong  
Dezhi Luo

University of Michigan

Naomi Saphra

Harvard University

Yueqi Song  
Hokin Deng

Carnegie Mellon University

Freda Shi  
Yuansheng Ni

University of Waterloo

# Panel Discussion



Eve Fleisig  
UC Berkeley



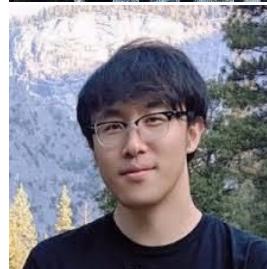
Ofir Press  
Princeton University



Wenda Xu  
Google Deepmind



David Rein  
METR Evaluations



Saining Xie  
New York University



(Moderator)  
Michael Saxon  
University of Washington

# Longevity

Why have evals stood the test of time? Have the ones that stood the test of time deserved to? What indicates something that is likely to last?

# Tyranny of Metrics

The things we are capable of measuring shape the way we design evals. Broadly, how have the limitations of metrics shaped your research projects? What things can't be measured right now that you would like to change?

# Human Subjectivity

Are we doing a good job of drawing the line on desiderata to account for divergent human preferences? How could we make this tractable? How should we account for diverse personal wants in alignment evaluation (or evaluation of other capabilities?)

# Generality

To what extent do you believe the "general" part of AGI is measurable and why? And for generalization within subdomains, how can those be scoped?