

CS11-711 Advanced NLP

# Evaluation and Benchmarks

Sean Welleck

**Carnegie  
Mellon  
University**



<https://cmu-l3.github.io/anlp-spring2026/>

<https://github.com/cmu-l3/anlp-spring2026-code>

Some slides adapted from Spring 2025 by Seungone Kim

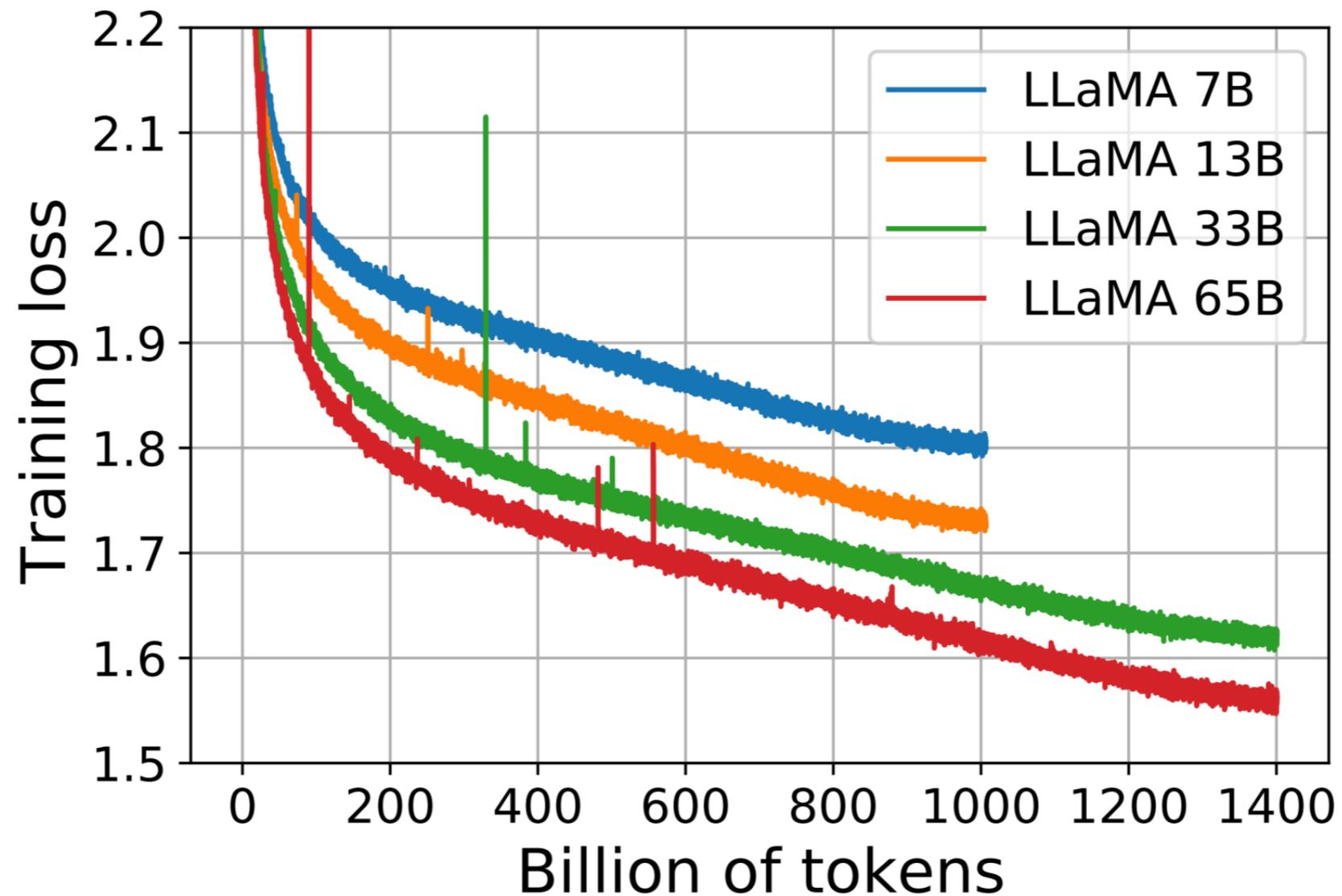
# Recap

- Learning: pre-training, fine-tuning
- Inference: decoding algorithms
- Today: How do we evaluate how good our system is?
  - How do we choose between two systems?

# Recap: Loss-Based Evaluation

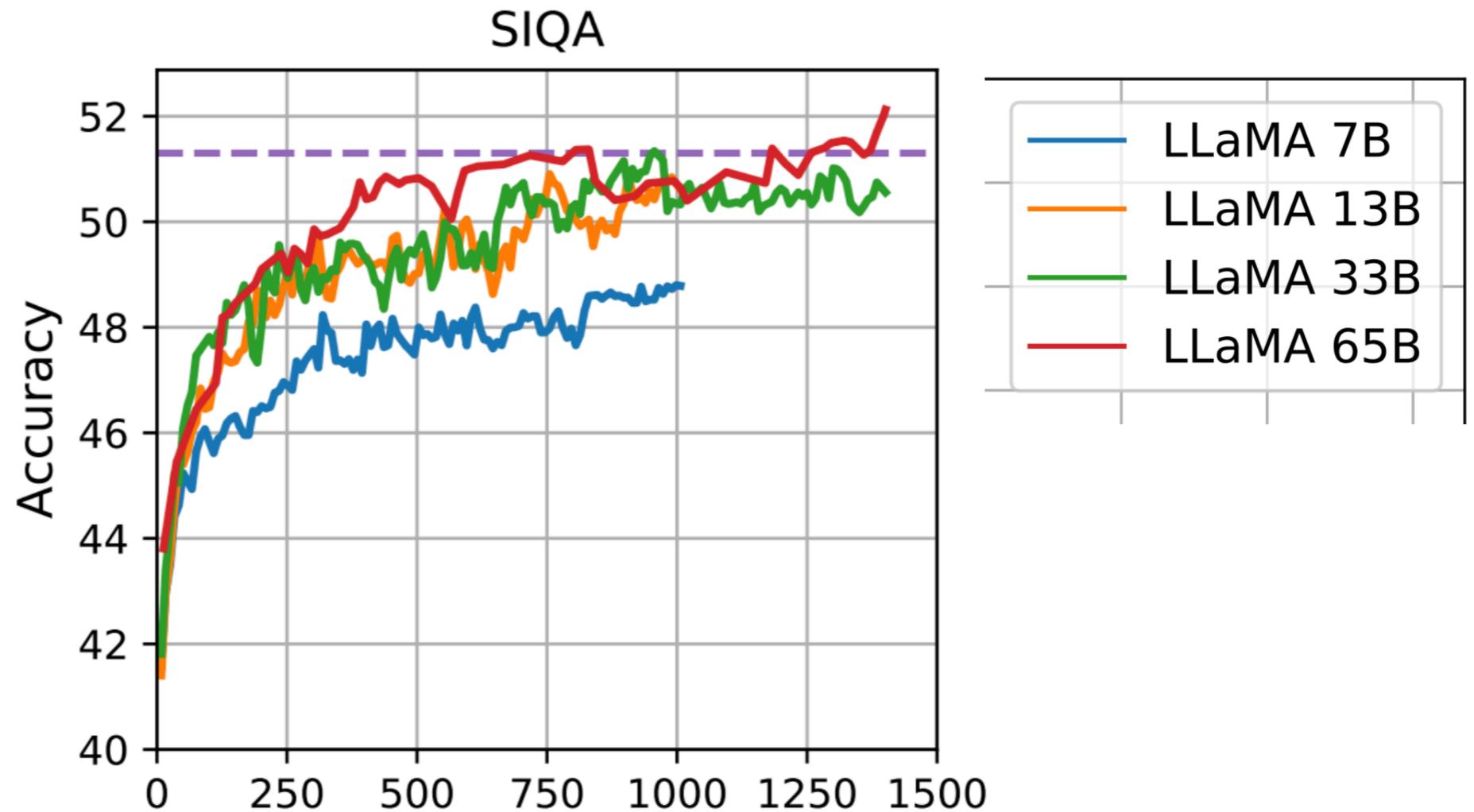
- Let  $p_\theta$  be a language model
- Let  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$  be a dataset
- Let  $\mathcal{L}(\theta; \mathcal{D})$  be the loss function
  - Typically  $-\log p_\theta(y | x)$
- Then we can use the training loss, validation loss, or test loss to evaluate the model

# Recap: Loss-Based Evaluation



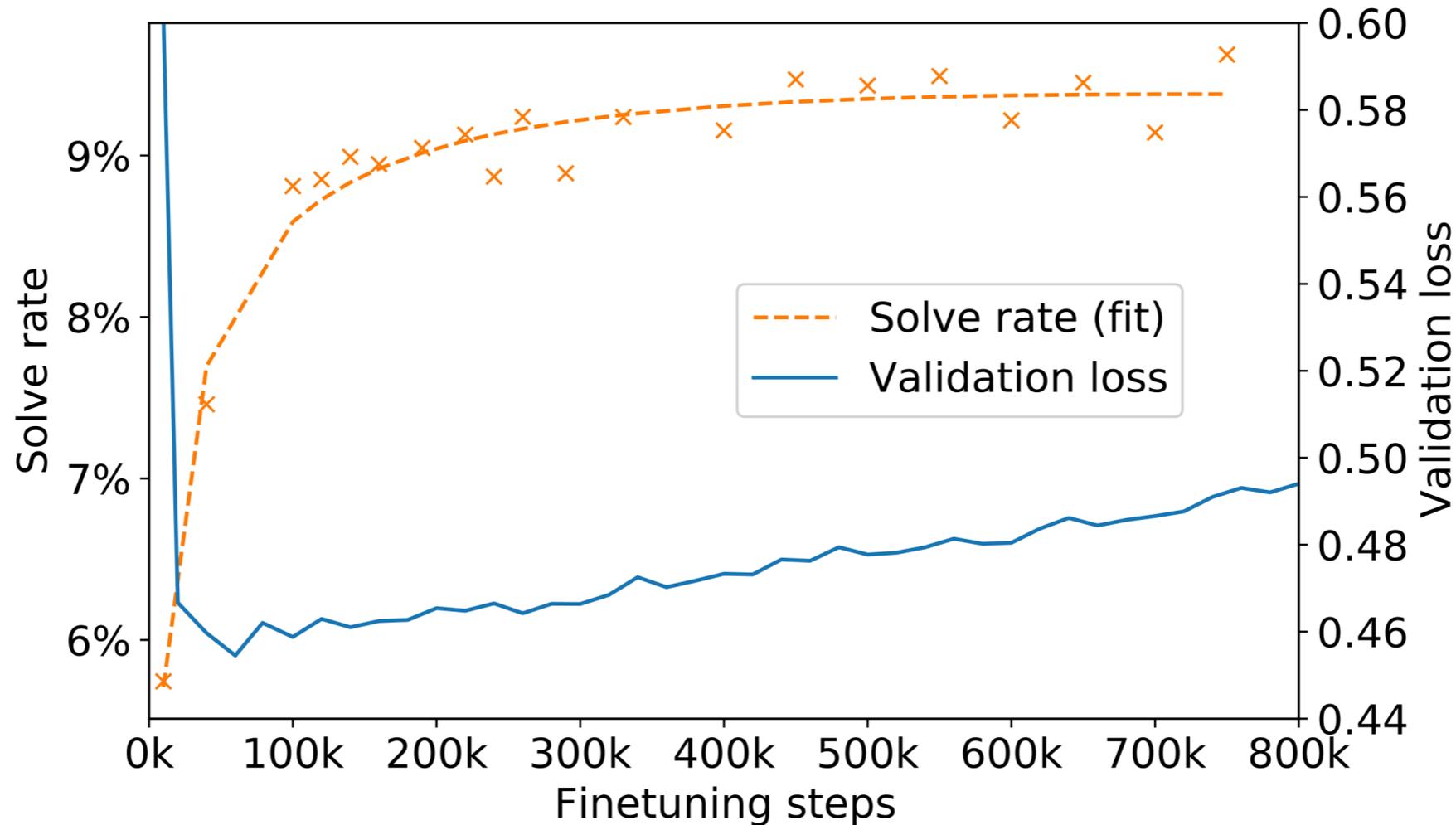
Claim: “Llama 33B is better than llama 13B”

# Loss isn't always enough



Claim: “Llama 33B is better than llama 13B” (?)

# Loss isn't always enough



Li et al. 2022 (DeepMind): *Competition-Level Code Generation with AlphaCode*

# Loss isn't always enough

- Loss measures how well the model fits the data distribution
- It does not directly measure how well a model performs tasks

# Task metrics and downstream evaluation

- Let  $g(p_\theta, x)$  denote running an inference algorithm
- Let  $m(y, \hat{y}) \rightarrow \mathbb{R}$  be a task metric (e.g., solve rate)
- Let  $p(D)$  be a distribution over task datasets
- We want the model to have high expected task performance

$$\mathbb{E}_{D \sim p(D)} \mathbb{E}_{x, y \sim D} \mathbb{E}_{\hat{y} \sim g(p_\theta, x)} [m(y, \hat{y})]$$

- Each dataset  $D$  represents a “downstream task”, i.e. a task that we use the model for downstream of its training

# How do we evaluate?

- In practice, we use a finite number of *evaluation benchmarks*. Each benchmark typically consists of:
  - A dataset  $D$  (sometimes only inputs  $x$ )
  - A task metric  $m$  (typically  $m(y, \hat{y})$ )
- We want the benchmarks to cover the breadth and depth of functionality that we want in our system.
  - For a machine translation system, we want various machine translation benchmarks
  - For a language model API, we want benchmarks related to the many capabilities that customers want from the API

# This lecture

- Intro/evaluation setup
- Properties of good benchmarks
  - Example benchmarks
  - Example task metrics
- Can I trust the benchmark score?

# Example: Llama 3 paper

Category	Benchmark	Llama 3 8B	Gemma 2 9B	Mistral 7B	Llama 3 70B	Mixtral 8x22B	GPT 3.5 Turbo	Llama 3 405B	Nemotron 4 340B	GPT-4 (0125)	GPT-4o	Claude 3.5 Sonnet
General	MMLU (5-shot)	69.4	<b>72.3</b>	61.1	<b>83.6</b>	76.9	70.7	87.3	82.6	85.1	89.1	<b>89.9</b>
	MMLU (0-shot, CoT)	<b>73.0</b>	72.3 <sup>△</sup>	60.5	<b>86.0</b>	79.9	69.8	88.6	78.7 <sup>△</sup>	85.4	<b>88.7</b>	88.3
	MMLU-Pro (5-shot, CoT)	<b>48.3</b>	–	36.9	<b>66.4</b>	56.3	49.2	73.3	62.7	64.8	74.0	<b>77.0</b>
	IFEval	<b>80.4</b>	73.6	57.6	<b>87.5</b>	72.7	69.9	<b>88.6</b>	85.1	84.3	85.6	88.0
Code	HumanEval (0-shot)	<b>72.6</b>	54.3	40.2	<b>80.5</b>	75.6	68.0	89.0	73.2	86.6	90.2	<b>92.0</b>
	MBPP EvalPlus (0-shot)	<b>72.8</b>	71.7	49.5	<b>86.0</b>	78.6	82.0	88.6	72.8	83.6	87.8	<b>90.5</b>
Math	GSM8K (8-shot, CoT)	<b>84.5</b>	76.7	53.2	<b>95.1</b>	88.2	81.6	<b>96.8</b>	92.3 <sup>◇</sup>	94.2	96.1	96.4 <sup>◇</sup>
	MATH (0-shot, CoT)	<b>51.9</b>	44.3	13.0	<b>68.0</b>	54.1	43.1	73.8	41.1	64.5	<b>76.6</b>	71.1
Reasoning	ARC Challenge (0-shot)	83.4	<b>87.6</b>	74.2	<b>94.8</b>	88.7	83.7	<b>96.9</b>	94.6	96.4	96.7	96.7
	GPQA (0-shot, CoT)	32.8	–	28.8	<b>46.7</b>	33.3	30.8	51.1	–	41.4	53.6	<b>59.4</b>
Tool use	BFCL	<b>76.1</b>	–	60.4	84.8	–	<b>85.9</b>	88.5	86.5	88.3	80.5	<b>90.2</b>
	Nexus	<b>38.5</b>	30.0	24.7	<b>56.7</b>	48.5	37.2	<b>58.7</b>	–	50.3	56.1	45.7
Long context	ZeroSCROLLS/QuALITY	81.0	–	–	90.5	–	–	<b>95.2</b>	–	<b>95.2</b>	90.5	90.5
	InfiniteBench/En.MC	65.1	–	–	78.2	–	–	<b>83.4</b>	–	72.1	82.5	–
	NIH/Multi-needle	98.8	–	–	97.5	–	–	98.1	–	<b>100.0</b>	<b>100.0</b>	90.8
Multilingual	MGSM (0-shot, CoT)	<b>68.9</b>	53.2	29.9	<b>86.9</b>	71.1	51.4	<b>91.6</b>	–	85.9	90.5	<b>91.6</b>

**Table 2 Performance of finetuned Llama 3 models on key benchmark evaluations.** The table compares the performance of

- “Key benchmarks” in the Llama 3 paper

# Example: Gemini 3.1 (Feb 2026)

Benchmark		Gemini 3.1 Pro Thinking (High)	Gemini 3 Pro Thinking (High)	Sonnet 4.6 Thinking (Max)	Opus 4.6 Thinking (Max)	GPT-5.2 Thinking (xhigh)	GPT-5.3-Codex Thinking (xhigh)
<b>Humanity's Last Exam</b> Academic reasoning (full set, text + MM)	No tools	<b>44.4%</b>	37.5%	33.2%	40.0%	34.5%	—
	Search (blocklist) + Code	<b>51.4%</b>	45.8%	49.0%	<b>53.1%</b>	45.5%	—
<b>ARC-AGI-2</b> Abstract reasoning puzzles	ARC Prize Verified	<b>77.1%</b>	31.1%	58.3%	68.8%	52.9%	—
<b>GPQA Diamond</b> Scientific knowledge	No tools	<b>94.3%</b>	91.9%	89.9%	91.3%	92.4%	—
<b>Terminal-Bench 2.0</b> Agentic terminal coding	Terminus-2 harness	<b>68.5%</b>	56.9%	59.1%	65.4%	54.0%	64.7%
	Other best self-reported harness	—	—	—	—	62.2% (Codex)	<b>77.3%</b> (Codex)
<b>SWE-Bench Verified</b> Agentic coding	Single attempt	<b>80.6%</b>	76.2%	79.6%	<b>80.8%</b>	80.0%	—
<b>SWE-Bench Pro (Public)</b> Diverse agentic coding tasks	Single attempt	<b>54.2%</b>	43.3%	—	—	55.6%	<b>56.8%</b>
<b>LiveCodeBench Pro</b> Competitive coding problems from Codeforces, ICPC, and IOI	Elo	<b>2887</b>	2439	—	—	2393	—
<b>SciCode</b> Scientific research coding		<b>59%</b>	56%	47%	52%	52%	—
<b>APEX-Agents</b> Long horizon professional tasks		<b>33.5%</b>	18.4%	—	29.8%	23.0%	—
<b>GDPval-AA Elo</b> Expert tasks		1317	1195	<b>1633</b>	1606	1462	—
<b>τ2-bench</b> Agentic tool use	Retail	90.8%	85.3%	91.7%	<b>91.9%</b>	82.0%	—
	Telecom	<b>99.3%</b>	98.0%	97.9%	<b>99.3%</b>	98.7%	—
<b>MCP Atlas</b> Multi-step workflows using MCP		<b>69.2%</b>	54.1%	61.3%	59.5%	60.6%	—
<b>BrowseComp</b> Agentic search	Search + Python + Browse	<b>85.9%</b>	59.2%	74.7%	84.0%	65.8%	—
<b>MMMU Pro</b> Multimodal understanding and reasoning	No tools	80.5%	<b>81.0%</b>	74.5%	73.9%	79.5%	—
<b>MMMLU</b> Multilingual Q&A		<b>92.6%</b>	91.8%	89.3%	91.1%	89.6%	—
<b>MRCR v2 (8-needle)</b> Long context performance	128k (average)	<b>84.9%</b>	77.0%	<b>84.9%</b>	84.0%	83.8%	—
	1M (pointwise)	26.3%	26.3%	Not supported	Not supported	Not supported	—

# Example: MMLU

- Multiple-choice questions in 57 subjects
  - Model generates A, B, C, or D
  - Metric: exact match; is the answer the same as the one in the dataset
    - $m(y, \hat{y}) = 1[y = \hat{y}]$

**Microeconomics** One of the reasons that the government discourages and regulates monopolies is that

- (A) producer surplus is lost and consumer surplus is gained. 
- (B) monopoly prices ensure productive efficiency but cost society allocative efficiency. 
- (C) monopoly firms do not engage in significant research and development. 
- (D) consumer surplus is lost with higher prices and lower levels of output. 

Figure 3: Examples from the Microeconomics task.

**Conceptual Physics** When you drop a ball from rest it accelerates downward at  $9.8 \text{ m/s}^2$ . If you instead throw it downward assuming no air resistance its acceleration immediately after leaving your hand is

- (A)  $9.8 \text{ m/s}^2$  
- (B) more than  $9.8 \text{ m/s}^2$  
- (C) less than  $9.8 \text{ m/s}^2$  
- (D) Cannot say unless the speed of throw is given. 

# Example: HumanEval

- HumanEval: LeetCode-style Python problems
  - Model generates code
  - Metric: execute the code and check whether it passes all test cases

$$\bullet m(y, \hat{y}) = \prod_{j=1}^{numtests} 1[\hat{y} \text{ passes test } j]$$

$$\bullet \text{pass @ } K(y, \{y_1, \dots, y_K\}) = \max(m(y, y_1), \dots, m(y, y_K))$$

```
def solution(lst):
    """Given a non-empty list of integers, return the sum of all of the odd elements
    that are in even positions.

    Examples
    solution([5, 8, 7, 1]) ==>12
    solution([3, 3, 3, 3, 3]) ==>9
    solution([30, 13, 24, 321]) ==>0
    """
    return sum(lst[i] for i in range(0, len(lst)) if i % 2 == 0 and lst[i] % 2 == 1)
```

# Example: GSM8k

- Grade school mathematics questions
  - Model generates a chain-of-thought and then an answer
  - Metric: check whether the answer is the same as the one in the dataset

**Problem:** Tina buys 3 12-packs of soda for a party. Including Tina, 6 people are at the party. Half of the people at the party have 3 sodas each, 2 of the people have 4, and 1 person has 5. How many sodas are left over when the party is over?

**Solution:** Tina buys 3 12-packs of soda, for  $3*12=36$  sodas

6 people attend the party, so half of them is  $6/2=3$  people

Each of those people drinks 3 sodas, so they drink  $3*3=9$  sodas

Two people drink 4 sodas, which means they drink  $2*4=8$  sodas

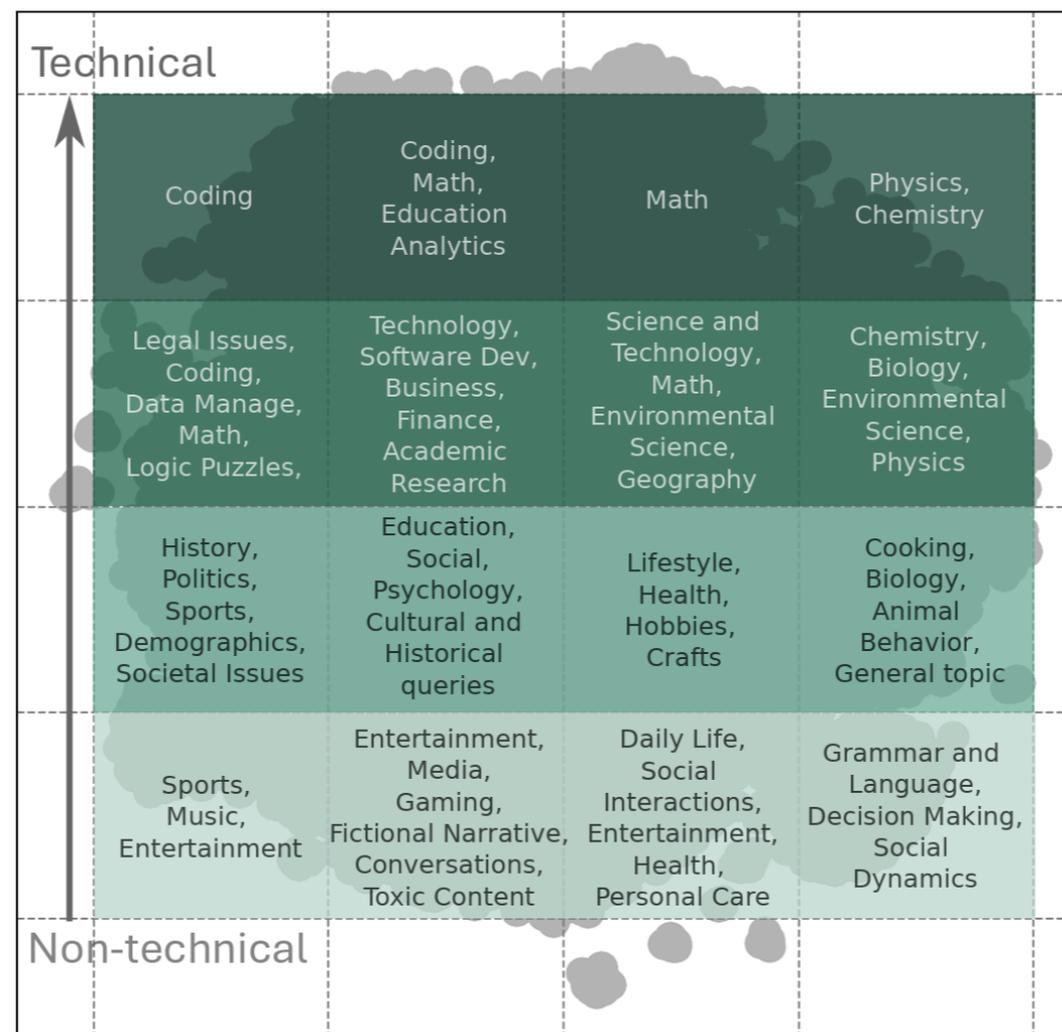
With one person drinking 5, that brings the total drank to  $5+9+8+3=25$  sodas

As Tina started off with 36 sodas, that means there are  $36-25=11$  sodas left

**Final Answer:** 11

# Discussion: properties of good benchmarks?

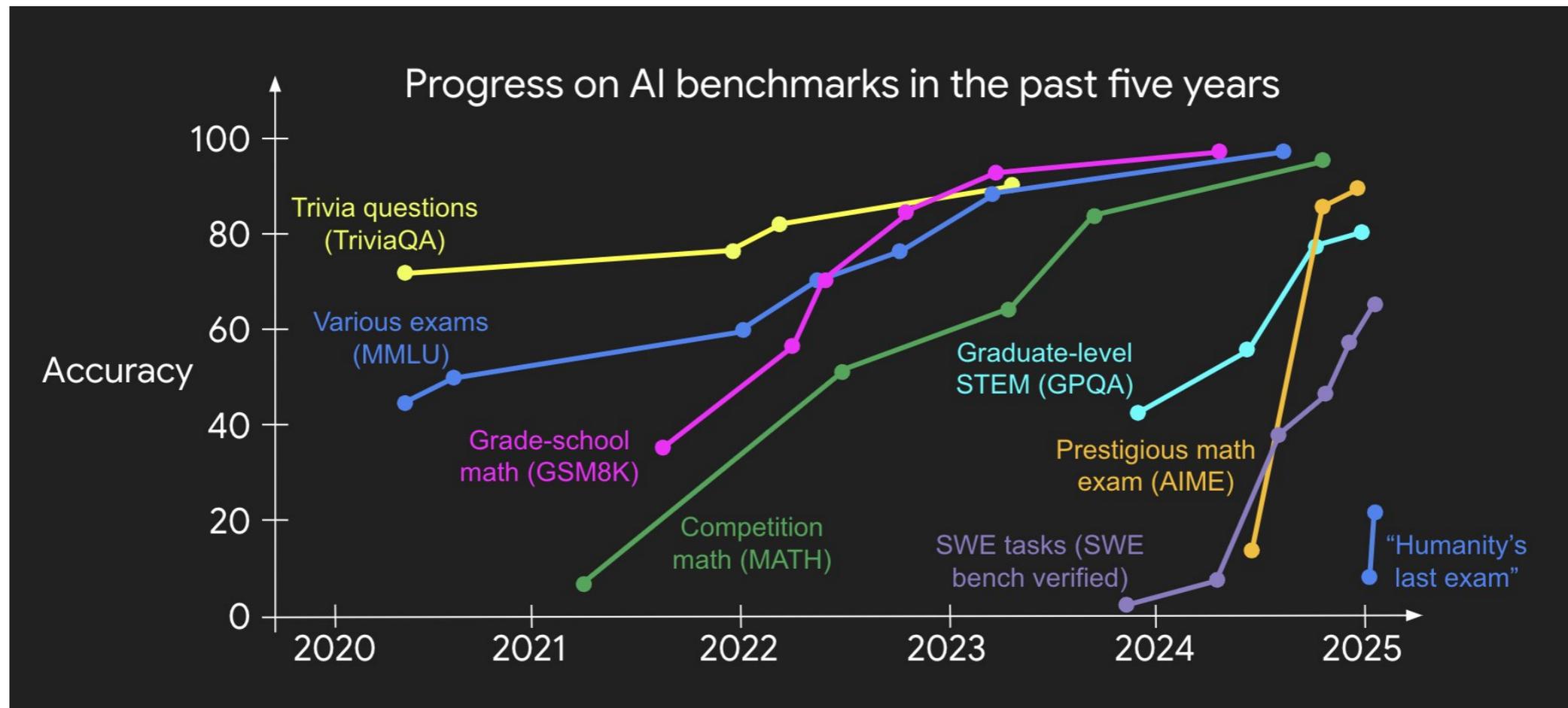
- Breadth/diversity
- Evaluates all desired input-output functionality



Prompt types from active users [Ni et al 2024]

# Discussion: properties of good benchmarks?

- Depth/difficulty
  - Examples are difficult enough to distinguish good and bad models.



# Discussion: properties of good benchmarks?

- Utility
  - Task is a proxy for desired functionality

## MATH Dataset (Ours)

**Problem:** Tom has a red marble, a green marble, a blue marble, and three identical yellow marbles. How many different groups of two marbles can Tom choose?

**Solution:** There are two cases here: either Tom chooses two yellow marbles (1 result), or he chooses two marbles of different colors ( $\binom{4}{2} = 6$  results). The total number of distinct pairs of marbles Tom can choose is  $1 + 6 = \boxed{7}$ .

**Problem:** The equation  $x^2 + 2x = i$  has two complex solutions. Determine the product of their real parts.

**Solution:** Complete the square by adding 1 to each side. Then  $(x + 1)^2 = 1 + i = e^{\frac{i\pi}{4}} \sqrt{2}$ , so  $x + 1 = \pm e^{\frac{i\pi}{8}} \sqrt[4]{2}$ . The desired product is then  $(-1 + \cos(\frac{\pi}{8}) \sqrt[4]{2})(-1 - \cos(\frac{\pi}{8}) \sqrt[4]{2}) = 1 - \cos^2(\frac{\pi}{8}) \sqrt{2} = 1 - \frac{(1 + \cos(\frac{\pi}{4}))}{2} \sqrt{2} = \boxed{\frac{1 - \sqrt{2}}{2}}$ .

Why do we need an LLM that is good at solving math word problems?

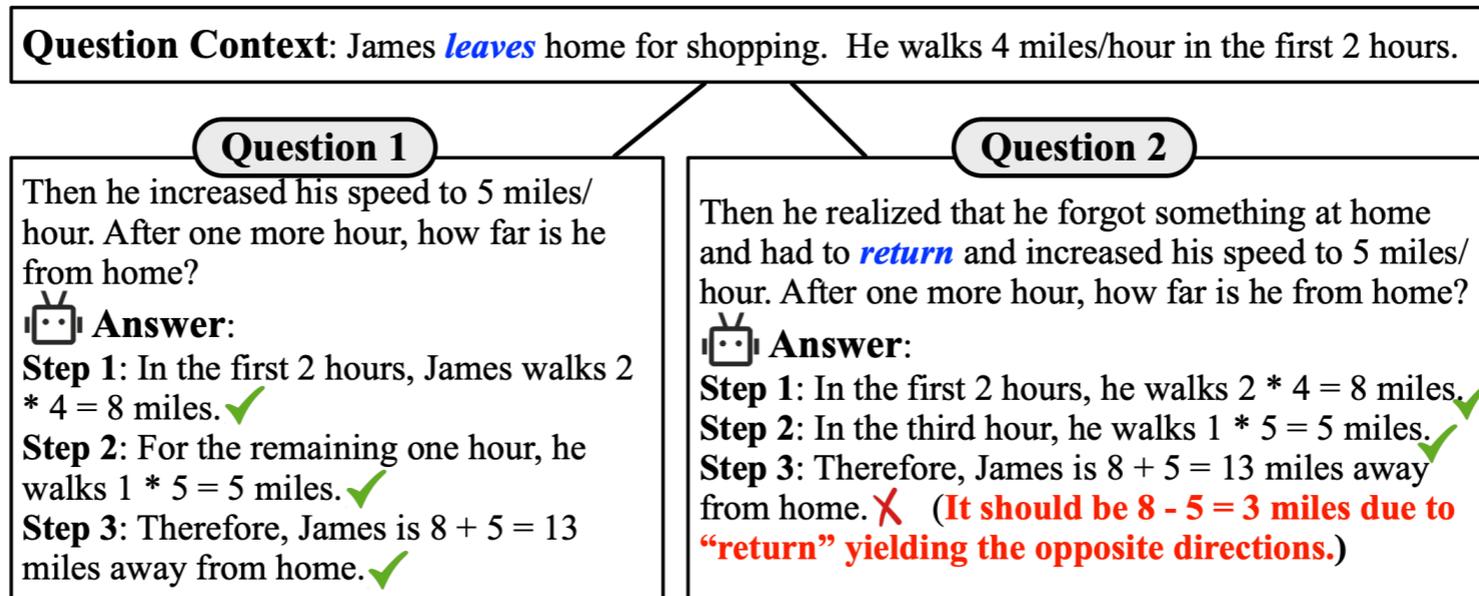
Foundation for **more complex tasks**  
(e.g., financial analysts)

Useful for **real users with similar questions**  
(e.g., tutoring high school math)

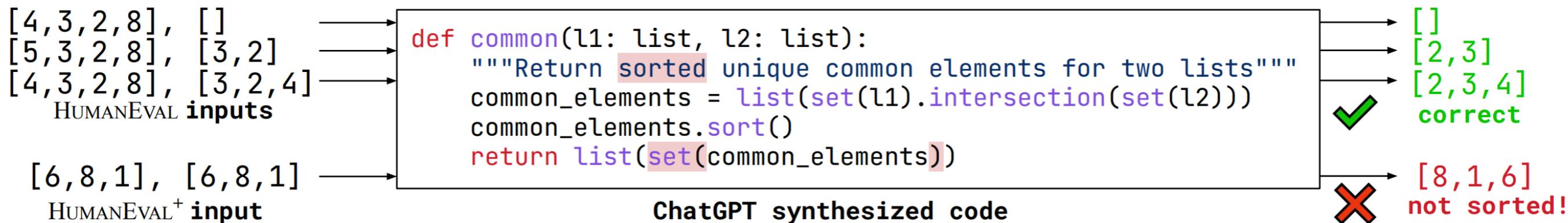
Medium for **addressing research questions**  
(e.g., Can an AI model reason?)

# Discussion: properties of good benchmarks?

- Robustness



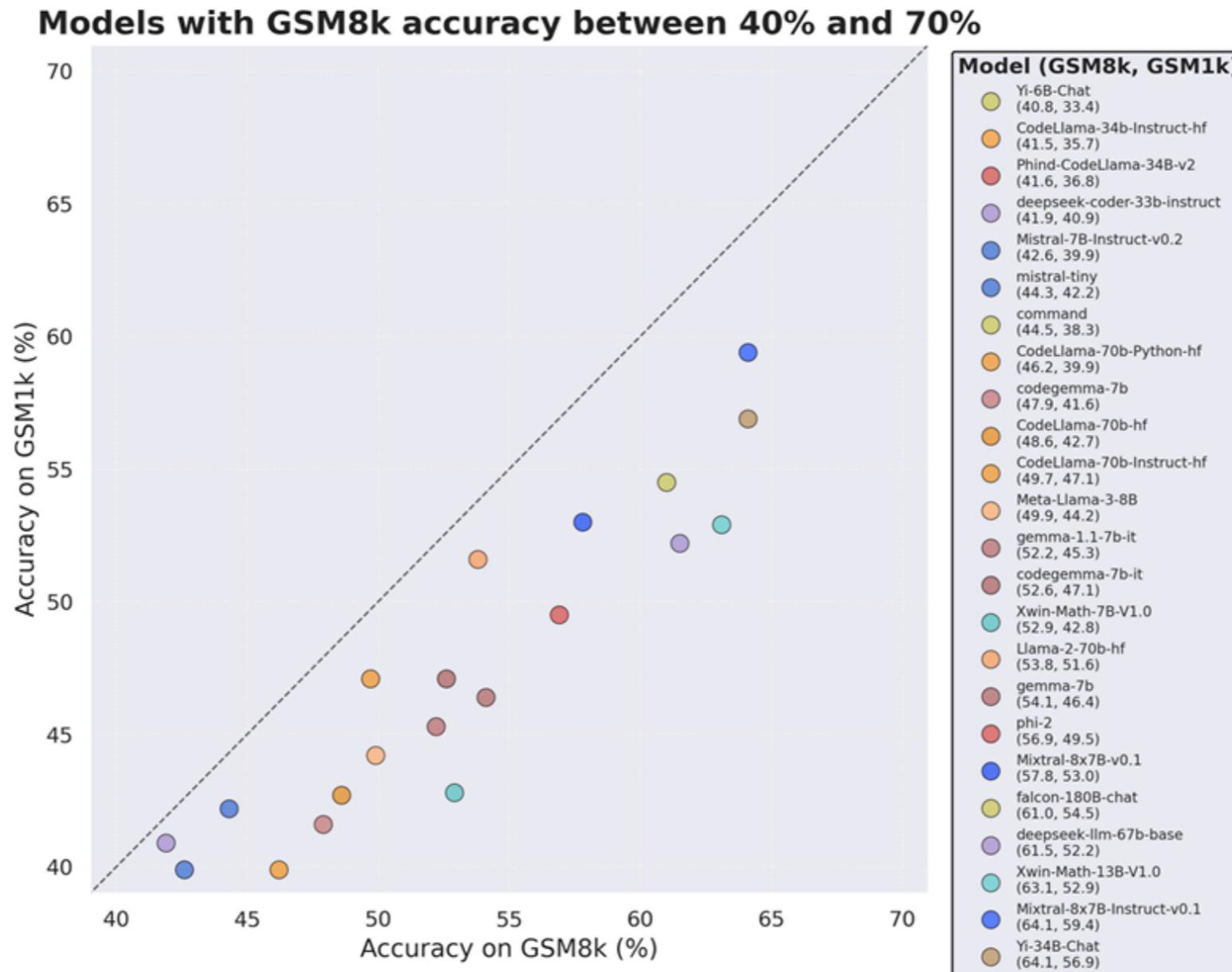
## GSM-Plus: Robustness to input variations



HumanEval-Plus: a more robust metric

# Discussion: properties of good benchmarks?

- Data contamination



# Discussion: properties of good benchmarks?

- Data contamination

## OpenWebMath document

A triangle is formed with edges along the line  $y = \frac{2}{3}x + 5$ , the  $x$ -axis, and the line  $x = k$ . If the area of the triangle is less than 20, find the sum of all possible integral values of  $k$ .

Feb 28, 2018

Look at the graph, here...two triangles are possible :  
<https://www.desmos.com/calculator/m6wnjgldq>  
The height of the triangles at any point will be formed by

$$\left[ \frac{2}{3}x + 5 \right]$$

$$\text{And the bases will be } [x - (-7.5)] = [x + 7.5]$$

So...we want to solve this

$$\frac{1}{2} \left[ \frac{2}{3}x + 5 \right] [x + 7.5] = 20$$

$$\left[ \frac{2}{3}x + 5 \right] [x + 7.5] = 40$$

$$\frac{2}{3}x^2 + 5x + 5x + 37.5 = 0$$

$$\frac{2}{3}x^2 + 10x - 2.5 = 0$$

Using a little technology....the max  $x$  value for the triangle formed above the  $x$  axis will be = .246

And the min  $x$  value for the triangle formed below the  $x$  axis will be = -15.246

With the given boundaries, the integer sums of all possible  $x$  values of  $k$  giving triangles with an area  $< 20$  units<sup>2</sup> =

$$\begin{aligned} & [ (-15) + (-14) + (-13) + \dots + (-2) + (-1) + 0 ] = \\ & - \frac{(15)(16)}{2} = \\ & -120 \end{aligned}$$

Feb 28, 2018

---

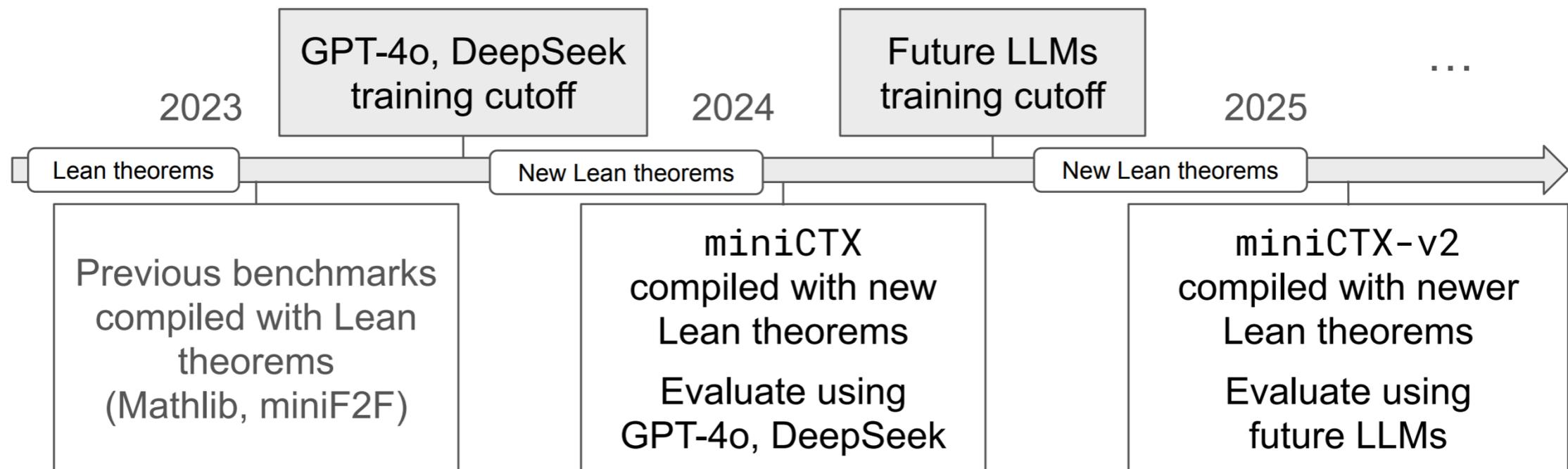
**MATH Problem:** A triangle is formed with edges along the line  $y = \frac{2}{3}x + 5$ , the  $x$ -axis, and the line  $x = k$ . If the area of the triangle is less than 20, find the sum of all possible integral values of  $k$ .

**Answer:** -120

Source: Llemma paper [Azerbaiyev et al]

# Discussion: properties of good benchmarks?

- Data contamination



“Temporal split” in miniCTX [Hu et al]

# Discussion: properties of good benchmarks?

- Efficient to evaluate
  - Hire an expert human to read a long chat and decide if it's correct
    - Expensive, slow, error prone
- ...
- Check whether a generated integer answer matches the answer in the dataset
  - Cheap, fast, exact, but simplifies the task

# Automatically evaluating complex text

**Question:** What are some business etiquette norms when doing business in Japan?

**Assistant A:**

Here are some business etiquette norms to keep in mind when doing business in Japan:

1. **Bowing:** Bowing is a common greeting in Japan. It is a sign of respect and should be reciprocated.
2. **Business cards:** Business cards are exchanged at the beginning of a meeting. It is important to receive and give business cards with both hands and to take a moment to read the card before putting it away.
3. **Punctuality:** Being on time is important in Japan. Arrive a few minutes early to show respect for your colleagues.
4. **Dress code:** Dress conservatively and professionally. Men should wear dark suits and ties, while women should wear conservative business attire.
5. **Language:** Learn a few basic Japanese phrases to show respect for the culture. It is also important to speak clearly and avoid using slang or idioms.
6. **Gift giving:** Gift giving is a common practice in Japan. Bring a small gift, such as a box of chocolates or a small souvenir from your home country.
7. **Dining etiquette:** When dining with colleagues, wait for the host to begin eating before you start. Use chopsticks correctly and avoid leaving food on your plate.

Remember, showing respect for Japanese culture and customs is key to building strong business relationships in Japan.

**Assistant B:**

When doing business in Japan, it is important to be aware of and respect the country's business etiquette norms. Some of these norms include:

1. **Greetings:** A handshake and bow is the standard greeting in Japan. It is important to bow deeply and to offer a firm handshake.
2. **Business cards:** In Japan, business cards are exchanged after the initial greeting. It is important to handle the card with respect and to read it carefully before placing it on the table.
3. **Meetings:** Meetings in Japan are typically formal and conducted in a hierarchical manner. It is important to address people by their titles and to avoid direct eye contact.
4. **Punctuality:** Punctuality is highly valued in Japan. It is important to arrive on time for meetings and to be prepared in advance.
5. **Gift giving:** Gift giving is a common practice in Japan, but it is important to be mindful of the occasion and the relationship. Avoid giving gifts that are considered too personal or extravagant.
6. **Dining:** Dining is an important part of doing business in Japan. It is important to be aware of the proper table manners and to avoid slurping or making loud noises while eating.
7. **Communication:** Communication in Japan is typically indirect and reserved. It is important to listen carefully and to avoid confrontational language.

Overall, it is important to be respectful, punctual, and mindful of hierarchy when doing business in Japan.

# Automatically evaluating complex text

- Option 1: word overlap
  - **Rouge**: Count the word overlap between the prediction & ref. answer

## Question:

What are the benefits of regular exercise?

## Prediction (5 words):

Exercise improves mood and health.

## Reference Answer (6 words):

Regular exercise benefits health and mood.

- **Rouge-1**: “exercise”, “health”, “and”, “mood” overlaps -> Precision: 4/5, Recall: 4/6, F1: 0.727
- **Rouge-2**: No overlap! -> Precision: 0/5, Recall: 0/6, F1: 0.0 **Is this desirable?**
- **Rouge-L**: [“exercise”, “mood”] and [“exercise”, “and”] are common subsequences -> Precision: 2/5, Recall: 2/6, F1: 0.364

Discuss: Pros? Cons?

# Automatically evaluating complex text

- Option 2: use a LLM (“LLM as a judge”)
  - Example:  $z \sim LLM(x, y_A, y_B; \text{prompt})$ 
    - Prompt: “Which output is better, A or B? Think carefully.”
  - Then we can compute the *win-rate* of model A versus model B.

Discuss: Pros? Cons?

# Automatically evaluating complex text

- Option 2: use a LLM (“LLM as a judge”)

## **GPT-4 Judgment (when A is placed in the first position):**

Both Assistant A and Assistant B provide helpful and relevant information about business etiquette norms in Japan. However, Assistant A’s answer is more detailed and organized, covering a wider range of topics such as dress code and language, which are not mentioned in Assistant B’s response. Assistant A also provides more specific advice on how to handle certain situations, such as exchanging business cards and dining etiquette. Therefore, Assistant A’s answer is more comprehensive and useful for someone looking to understand Japanese business etiquette.

## **GPT-4 Judgment (when B is placed in the first position):**

Both Assistant A and Assistant B provide helpful and relevant information about business etiquette norms in Japan. However, Assistant B’s answer is more detailed and covers a wider range of topics, such as addressing people by their titles, avoiding direct eye contact, and being mindful of communication styles. Assistant A’s answer is more concise but lacks some of the depth found in Assistant B’s response. Therefore, Assistant B’s answer is better in this case.

Discuss: Pros? Cons?

# High agreement with humans [Zheng et al 2023]

Table 5: Agreement between two types of judges on MT-bench. “G4-Pair” and “G4-Single” denote GPT-4 with pairwise comparison and single-answer grading respectively. The single-answer grading can be converted into pairwise comparison results for calculating the agreement. We report two setups: “S1” includes non-tie, tie, and inconsistent (due to position bias) votes and counts inconsistent as tie; “S2” only includes non-tie votes. The agreement between two random judges under each setup is denoted as “R=”. The top value in each cell is the agreement, and the bottom gray value is #votes.

Setup	S1 (R = 33%)		S2 (R = 50%)	
	G4-Single	Human	G4-Single	Human
G4-Pair	70%	66%	97%	<b>85%</b>
	1138	1343	662	859
G4-Single	-	60%	-	85%
		1280		739
Human	-	63%	-	<b>81%</b>
		721		479

(a) First Turn

Setup	S1 (R = 33%)		S2 (R = 50%)	
	G4-Single	Human	G4-Single	Human
G4-Pair	70%	66%	95%	<b>85%</b>
	1161	1325	727	864
G4-Single	-	59%	-	84%
		1285		776
Human	-	67%	-	<b>82%</b>
		707		474

(b) Second Turn

# Checklist-based evaluation

[Lee et al 2025]

- Define criteria  $c_1, \dots, c_N$
- Given a response, the LLM judge checks whether each criterion is met

**Get Dimension & Definition**

Task

**Summarization**

Dimensions & Definition

**Consistency**  
the factual alignment between the summary and the summarized source

**Define Sub-dimension (Human)**

Consistency

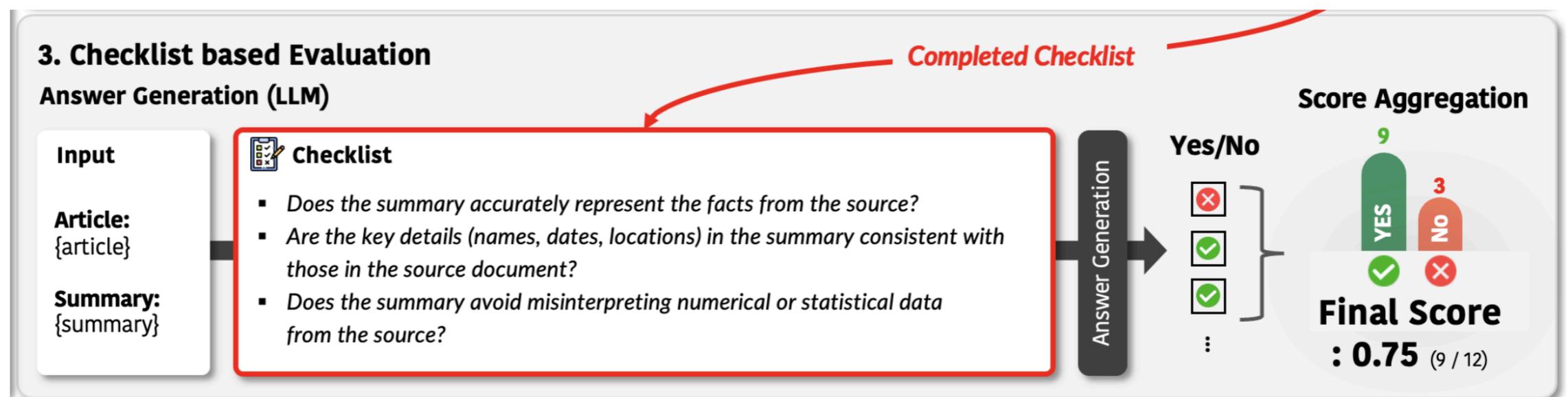
Factual Consistency No New Information ...

Coherence

Topic Maintenance Logical Flow ...

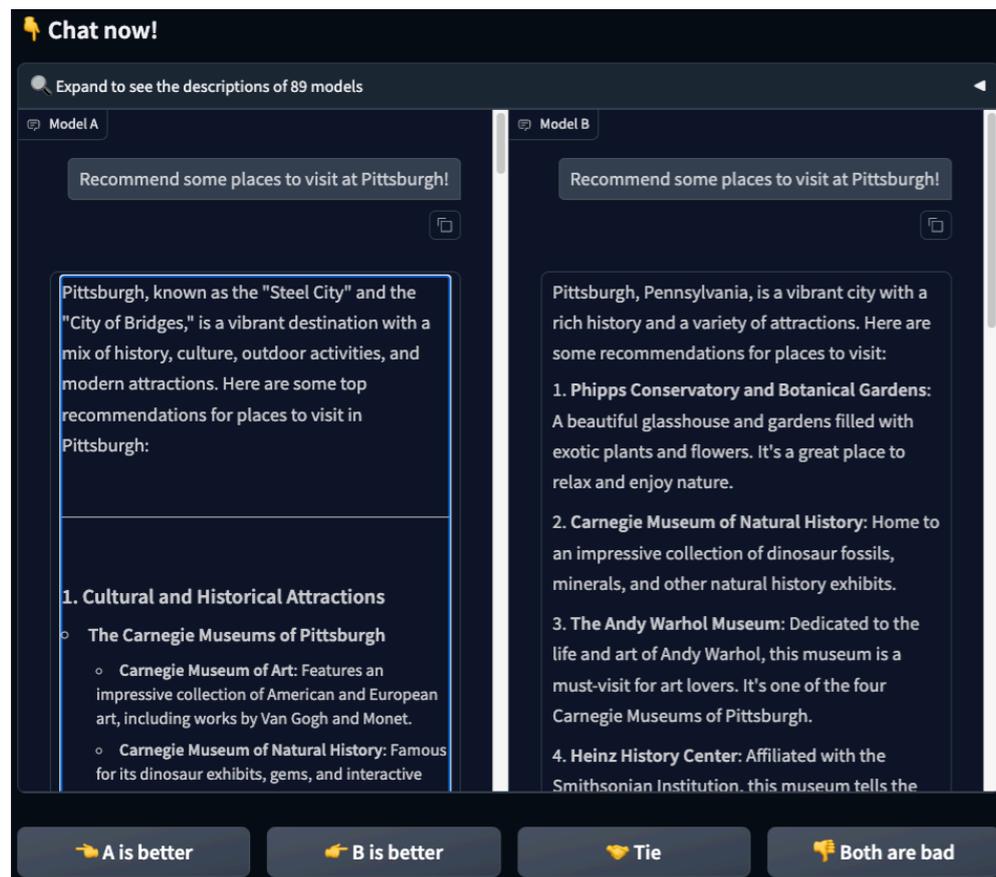
Fluency

Formatting Grammar ...



# Automatically evaluating complex text

- Option 3: use humans
  - LMSys Chatbot Arena: ask humans to vote which response is better



→ ↻ ↗ Imarena.ai/leaderboard

### Leaderboard Overview

See how leading models stack up across text, image, vision, and b each Arena, you can explore deeper insights in their dedicated tab

Text 🕒 3 days ago

Rank (UB) ↑	Model ↓	Score ↓
1	AI claude-sonnet-4-5-20250929-t...	1453
1	🌐 gemini-2.5-pro	1452
1	AI claude-opus-4-1-20250805-thi...	1449
1	AI claude-sonnet-4-5-20250929	1439
2	🌀 chatgpt-4o-latest-20250326	1441

Discuss: Pros? Cons?

# Discussion: properties of good benchmarks?

- Breadth/diversity
- Depth/difficulty
- Utility
- Robustness
- Data (un-)contamination
- Efficient evaluation (automatic, low cost)

# This lecture

- Intro/evaluation setup
- Properties of good benchmarks
  - Example benchmarks
  - Example task metrics
- **Next: What can I conclude from the benchmark score?**

Thank you