

Lecture 17: Evaluating LLMs

Instructor: Swabha Swayamdipta USC CSCI 444 NLP Nov 10, 2025



Announcements + Logistics

- Progress Report and Quiz 4 grades out
- Upcoming:
 - Paper Discussions on Wed, 11/12: MMLU, Dynabench, Chatbot Arena
 - 11/17: HW3 due
 - 11/24: Quiz 5 + Bonus
 - 12/1 and 12/3: Project Presentations: Sign up for dates / slots ASAP

Lecture Outline

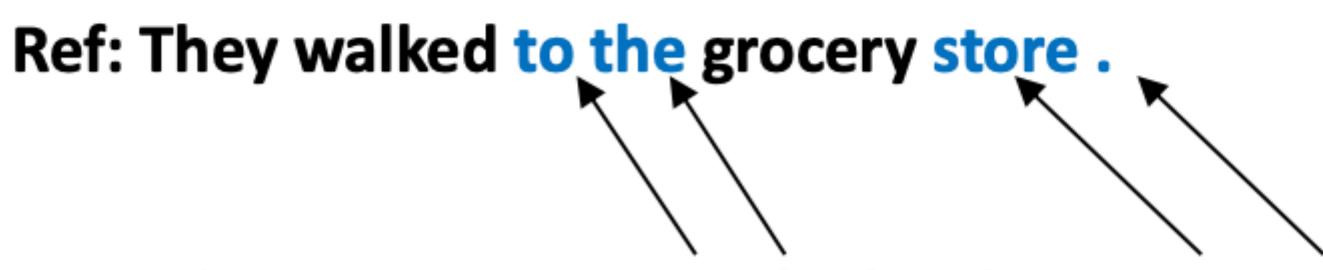
- Announcements
- Quiz 4 Answers
- Generative Evaluation: Metrics / Methods
- Benchmarks
- Rethinking Evaluation Benchmarks
 - Some of my own research



Generative Evaluation Metrics / Methods

Evaluation Strategies

- With Reference
 - Lexical Matching
 - Semantic Matching
- Without Reference
 - Perplexity
 - Model-Based Metrics
 - Advanced: Distributional Matching
 - Simplest, Most Reliable Strategy to-date: Human Evaluation
 - Even simpler and least reliable: Auto Evaluation



Gen: The woman went to the hardware store.

Reference-Based Metrics

Ref: They walked to the grocery store.

Gen: The woman went to the hardware store.

- Only possible for close-ended generation tasks
- Compute a score that indicates the lexical similarity between generated and gold-standard (human-written) text
- Fast and efficient and widely used
- *n*-gram overlap metrics (e.g., BLEU, ROUGE, etc.)

BLEU

- Stands for Bilingual Evaluation Understudy
- BLEU compares the machine-written translation to one or several human-written translation(s), and computes a similarity score based on:
 - Geometric mean of n-gram precision (usually for 1, 2, 3 and 4-grams)
 - Plus a penalty for too-short system translations
- BLEU is useful but imperfect
 - There are many valid ways to translate a sentence
 - So a good translation can get a poor BLEU score because it has low n-gram overlap with the human translation
- Precision-based metric

USC Viterbi

Precision, Recall and F-1

• True Positives, True Negatives, False Positives and False Negatives

$$\frac{TP}{TP + FP}$$

Of all the items in the prediction, how many match the ground truth

Recall =
$$\frac{TP}{TP + FN}$$

Of all the items in the ground truth, how many are correctly predicted

$$F_1 = \frac{2 * PR}{P + R}$$

Harmonic Mean of Precision and Recall

Different value for different classes!

BLEU: Details

- Purely precision-based rather than combining precision and recall.
- BLEU score for a corpus of candidate translation sentences is a function of
 - the n-gram word precision over all the sentences
 - combined with a brevity penalty computed over the corpus as a whole.
- Consider a corpus composed of a single sentence
 - The unigram precision for this corpus is the percentage of unigram tokens in the candidate translation that also occur in the reference translation, and ditto for bigrams and so on, up to 4-grams
 - It computes this n-gram precision for unigrams, bigrams, trigrams, and 4-grams and takes the geometric mean
- Because BLEU is a word-based metric, it is very sensitive to word tokenization, making it impossible to compare different systems if they rely on different tokenization

USC Viterbi

BLEU: Example

Reference translation 1:

The U.S. island of Guam is maintaining a high state of alert after the Guam airport and its offices both received an e-mail from someone calling himself the Saudi Arabian Osama bin Laden and threatening a biological/chemical attack against public places such as the airport.

Reference translation 2:

Guam International Airport and its offices are maintaining a high state of alert after receiving an e-mail that was from a person claiming to be the wealthy Saudi Arabian businessman Bin Laden and that threatened to launch a biological and chemical attack on the airport and other public places.

Machine translation:

The American [?] international airport and its the office all receives one calls self the sand Arab rich business [?] and so on electronic mail, which sends out; The threat will be able after public place and so on the airport to start the biochemistry lattack, [?] highly alerts after the maintenance.

Reference translation 3:

The US International Airport of Quam and its office has received an email from a self-claimed Arabian millionaire named Laden, which threatens to launch a biochemical attack on such public places as airport. Guam authority has been on alert.

Reference translation 4:

US Guam International Airport and its office received an email from Mr. Bin Laden and other rich businessman from Saudi Arabia. They said there would be biochemistry air raid to Guam Airport and other public places. Guam needs to be in high precaution about this matter.

Papineni et al., 2002

ROUGE

- Stands for "Recall-Oriented Understudy for Gisting Evaluation"
- Originally created for evaluating automatic summarization as well as machine translation
- Comparing an automatically produced summary or translation against a set of reference summaries (typically human-produced)
- Four variants:
 - ROUGE-N
 - ROUGE-L
 - ROUGE-S
 - ROUGE-W

ROUGE: Details

- ROUGE-N: measures unigram, bigram, trigram and higher order n-gram overlap
 - n-gram recall between a candidate summary and a set of reference summaries
- ROUGE-L: measures longest matching sequence of words using LCS.
 - Does not require consecutive matches but in-sequence matches that reflect sentence level word order.
 - Since it automatically includes longest in-sequence common n-grams, you don't need a predefined n-gram length.
- ROUGE-S: Is any pair of words in a sentence in order, allowing for arbitrary gaps.
 - Also be called skip-gram concurrence.
 - For example, **skip-bigram** measures the overlap of word pairs that can have a maximum of two gaps in between words. As an example, for the phrase "cat in the hat" the skip-bigrams would be "cat in, cat the, cat hat, in the, in hat, the hat".
- ROUGE-W: Weighted Longest Common Subsequence

Evaluating Generation: Other Options

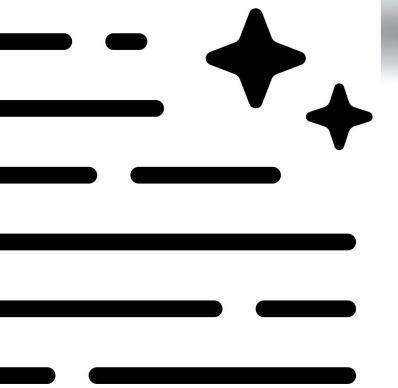
$$PPL(\mathbf{w}) = P(w_1 w_2 ... w_N)^{-\frac{1}{N}}$$

- Perplexity!
- Model-based Metrics (BERTScore, BARTScore, Word Mover's Distance, BLEURT)
 - Use learned representations of words and sentences to compute semantic similarity between generated and reference texts
 - No more n-gram bottleneck because text units are represented as embeddings!
 - The embeddings are pretrained, distance metrics used to measure the similarity can be fixed
- Automatic metrics for text comparisons fall short of matching human decisions
- So, Human Evaluation!

Human Evaluation

- Ask humans to evaluate the quality of generated text
 - Along specific axes: fluency, coherence / consistency, factuality and correctness, commonsense, etc.
 - Mostly done via crowdsourcing
- Human judgments are regarded as the gold standard
- Many, many correct answers! Ground truth might be irrelevant
- Most LLMs are compared with each other on their generations directly
 - Evaluated by humans or models

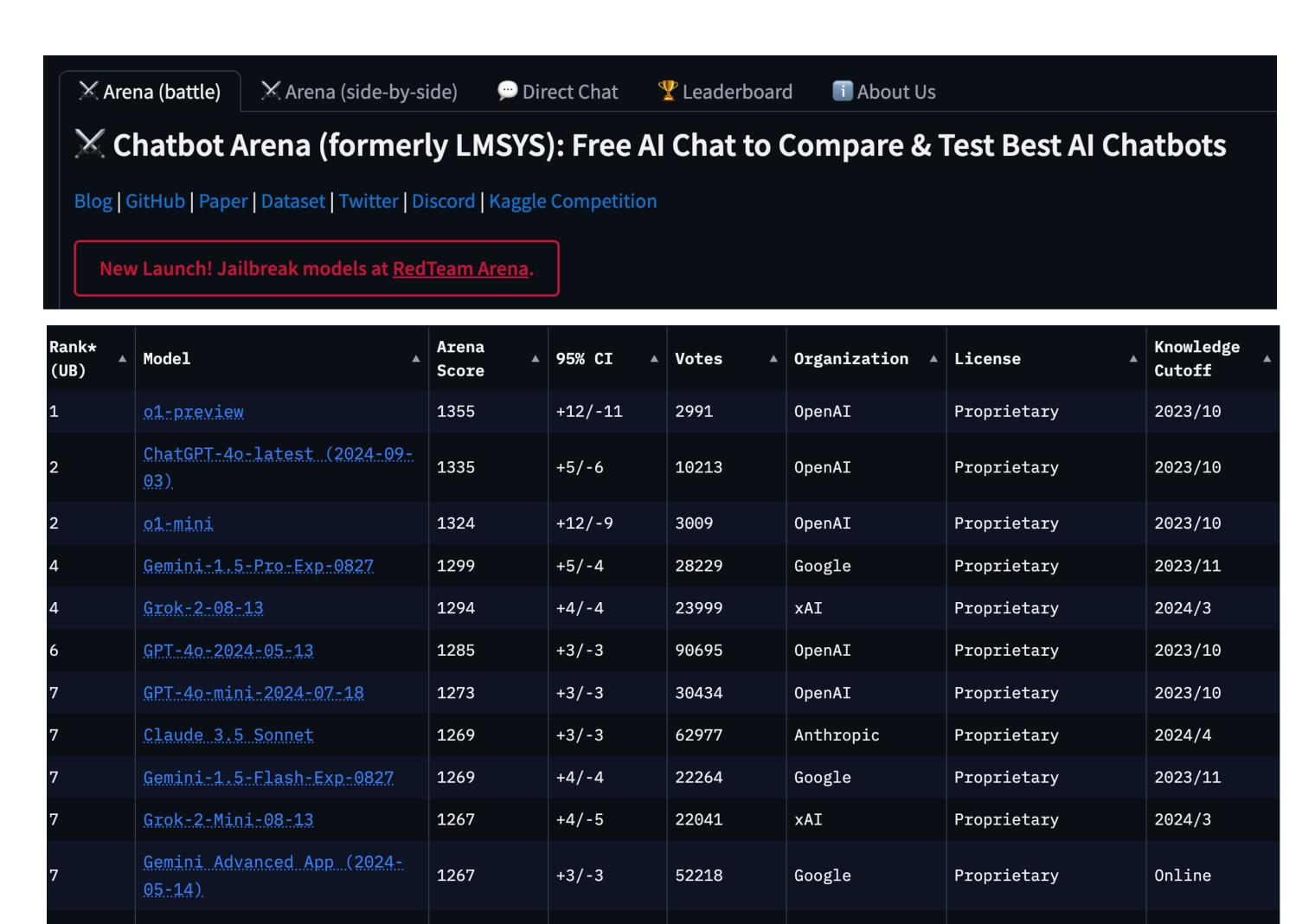






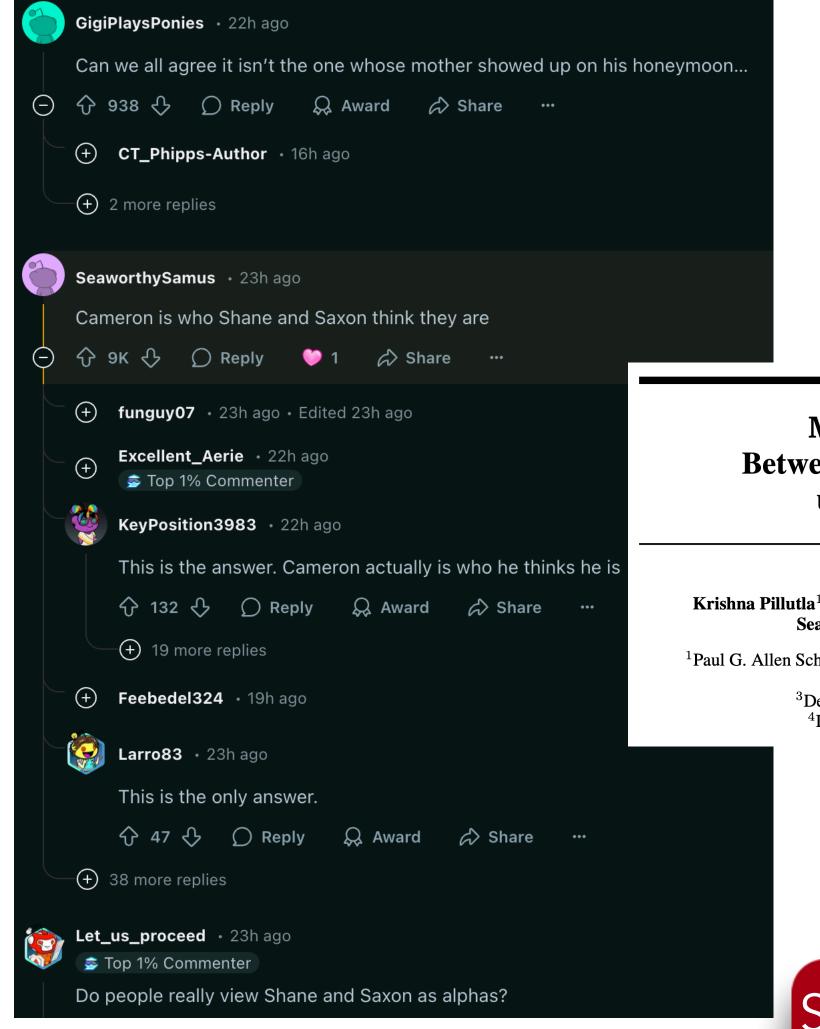
Chatbot Arena

- How to do human evaluation at very large scales?
- Maybe a model is good if people choose to use it (and pay for it)...
- Platform for ranking LLMs
- Limitations:
 - Could be possible to game this...



USC Viterbi

Shane, Cameron, or Saxon—who is the ultimate White Lotus alpha?











When considering who might be the "ultimate alpha" among Shane, Cameron, and Saxon from The White Lotus series, I'd need to analyze each character's portrayal:

Shane Patton (Season 1) is the entitled newlywed who becomes obsessed with getting the suite he booked, ultimately stabbing Armond in a confrontation. His privilege and

ate a certain type of aggressive dominance.

ason 2) is the wealthy, manipulative finance bro who exerts control exual conquest, even seducing his friend's wife. His confidence and lack different kind of social dominance.

acter introduced in Season 3 set in Thailand. As my knowledge cutoff is have detailed information about his character traits or behaviors in the

MAUVE: Measuring the Gap Between Neural Text and Human Text using Divergence Frontiers

Krishna Pillutla¹ Swabha Swayamdipta² Rowan Zellers¹ John Thickstun³ Sean Welleck^{1,2} Yejin Choi^{1,2} Zaid Harchaoui⁴

¹Paul G. Allen School of Computer Science & Engineering, University of Washington

²Allen Institute for Artificial Intelligence

³Department of Computer Science, Stanford University

⁴Department of Statistics, University of Washington

Based on the first two seasons, both Shane and Cameron exhibit toxic "alpha" behaviors in different ways - Shane through entitlement and Cameron through manipulation. Without complete information about Saxon, I can't make a definitive comparison between all three characters.

Subjectivity in evaluation of long-form text is a feature, not a bug

Human Evaluation: Limitations

Human evaluation of generations can be hard, but it's still considered the most reliable method today

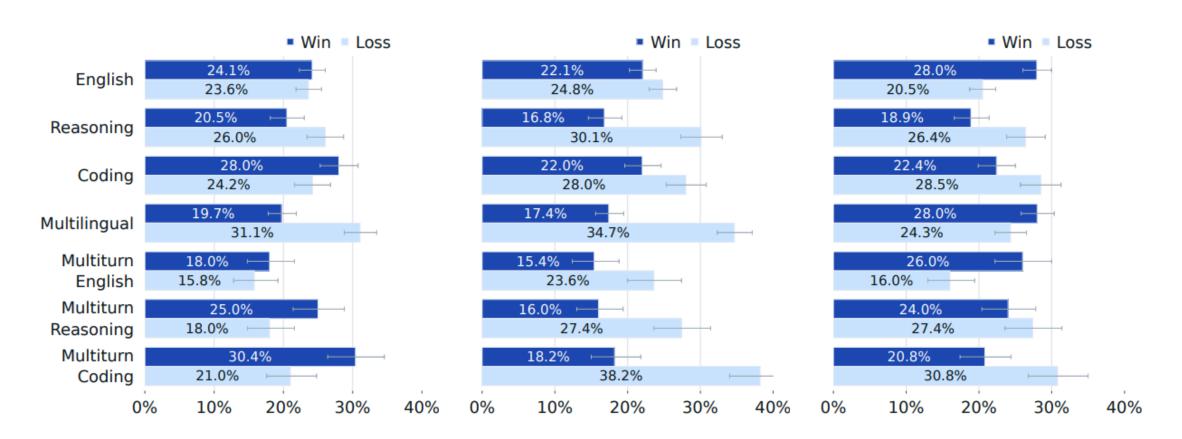


Figure 17 Human evaluation results for the Llama 3 405B model. Left: Comparison with GPT-4. Middle: Comparison with GPT-40. Right: Comparison with Claude 3.5 Sonnet. All results include 95% confidence intervals and exclude ties.

- Of course, we know that human eval is slow and expensive
- Beyond the cost of human eval, it's still fairly hard to conduct, and far from perfect:
 - Results are inconsistent / not reproducible
 - Can be subjective!
 - Misinterpret your question
 - Precision not recall

USC Viterbi

Automatic Evaluation with LLM Judges

AlpacaFarm: A Simulation Framework for Methods that Learn from Human Feedback

Yann Dubois* Stanford Xuechen Li* Stanford Rohan Taori* Stanford Tianyi Zhang* Stanford

Ishaan Gulrajani Stanford

Jimmy Ba University of Toronto Carlos Guestrin Stanford Percy Liang Stanford Tatsunori B. Hashimoto Stanford Cheap and theoretically consistent with human evaluation. BUT... reliability? Models evaluating their own generations may lead to weird mode collapsing effect

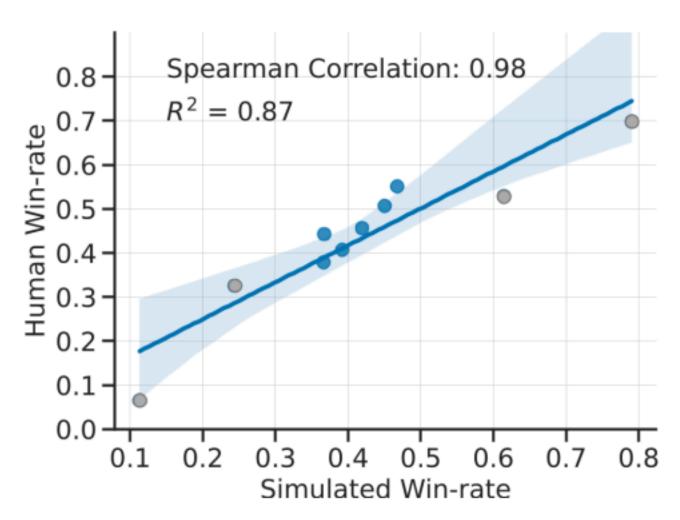


Figure 3: The ranking of methods trained and evaluated in AlpacaFarm matches that of methods trained and evaluated in the human-based pipeline. Each point represents one method M (e.g. PPO). The x-axis shows the simulated evaluation (win-rates measured by $p_{\rm sim}^{\rm eval}$) on methods trained in simulation $M_{\rm sim}$. The y-axis shows human evaluation (win-rates measured by $p_{\rm human}$) on methods trained with human feedback $M_{\rm human}$. Gray points show models that we did not train, so their x and y values only differ in the evaluation (simulated vs human). Without those points, we have $R^2 = 0.83$ and a Spearman Correlation of 0.94.

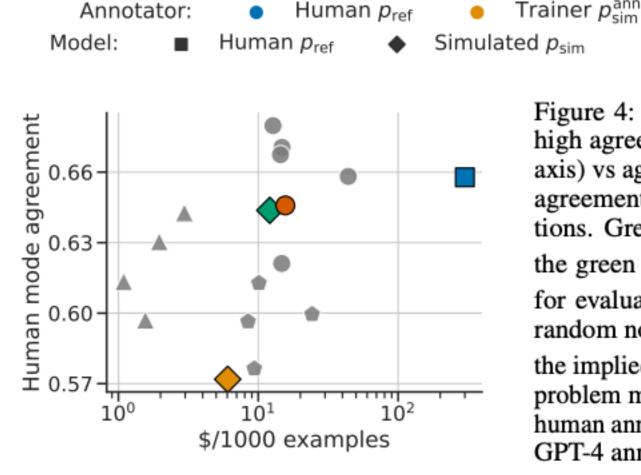


Figure 4: Our simulated annotators are cheap and have a high agreement with human annotators. We show price (x-axis) vs agreement (y-axis) as measured by each annotator's agreement with the majority vote among 3 human annotations. Grey points are all simulated annotators in the pool, the green shows the resulting pool of annotators (used for evaluation), the orange shows the same pool with random noise added during training. This does not change the implied reward function from shows the learning problem more challenging. The blue shows the average of human annotators, and the red shows a single low variance GPT-4 annotator analyzed below.

ChatGPT

Evaluator $p_{\text{sim}}^{\text{eval}}$

GPT4 $p_{\text{sim}}^{\text{GPT4}}$

Davinci003

Evaluating Systems without References

- Compare human / natural language distributions to model-generated language distributions
- Divergence between these two distributions can be measured by MAUVE

MAUVE: Measuring the Gap Between Neural Text and Human Text using Divergence Frontiers

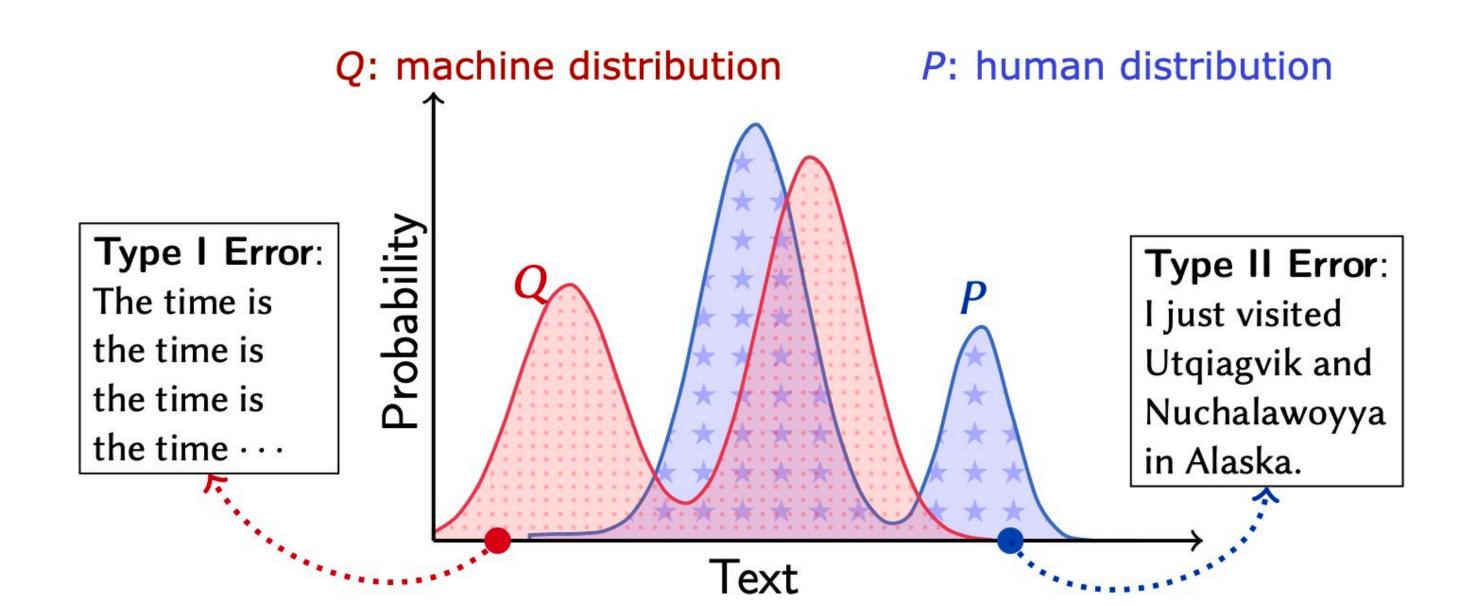
Krishna Pillutla¹ Swabha Swayamdipta² Rowan Zellers¹ John Thickstun³ Sean Welleck^{1,2} Yejin Choi^{1,2} Zaid Harchaoui⁴

¹Paul G. Allen School of Computer Science & Engineering, University of Washington

²Allen Institute for Artificial Intelligence

³Department of Computer Science, Stanford University

⁴Department of Statistics, University of Washington



USC Viterbi

A Survey on Evaluation of Large Language Models

YUPENG CHANG* and XU WANG*, School of Artificial Intelligence, Jilin University, China

JINDONG WANG[†], Microsoft Research Asia, China

YUAN WU[†], School of Artificial Intelligence, Jilin University, China

LINYI YANG, Westlake University, China

KAIJIE ZHU, Institute of Automation, Chinese Academy of Sciences, China

HAO CHEN, Carnegie Mellon University, USA

XIAOYUAN YI, Microsoft Research Asia, China

CUNXIANG WANG, Westlake University, China

YIDONG WANG, Peking University, China

WEI YE, Peking University, China

YUE ZHANG, Westlake University, China

YI CHANG, School of Artificial Intelligence, Jilin University, China

PHILIP S. YU, University of Illinois at Chicago, USA

QIANG YANG, Hong Kong University of Science and Technology, China

XING XIE, Microsoft Research Asia, China

LLM Evaluation mostly focuses on

creating new benchmarks

LLMs evaluation What to evaluate

General

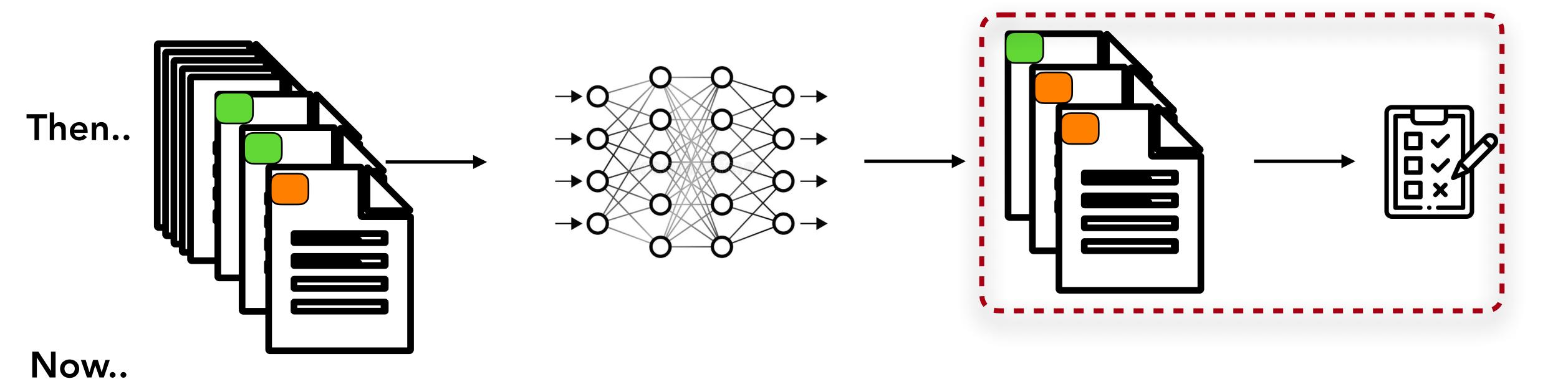
(Sec. 3)

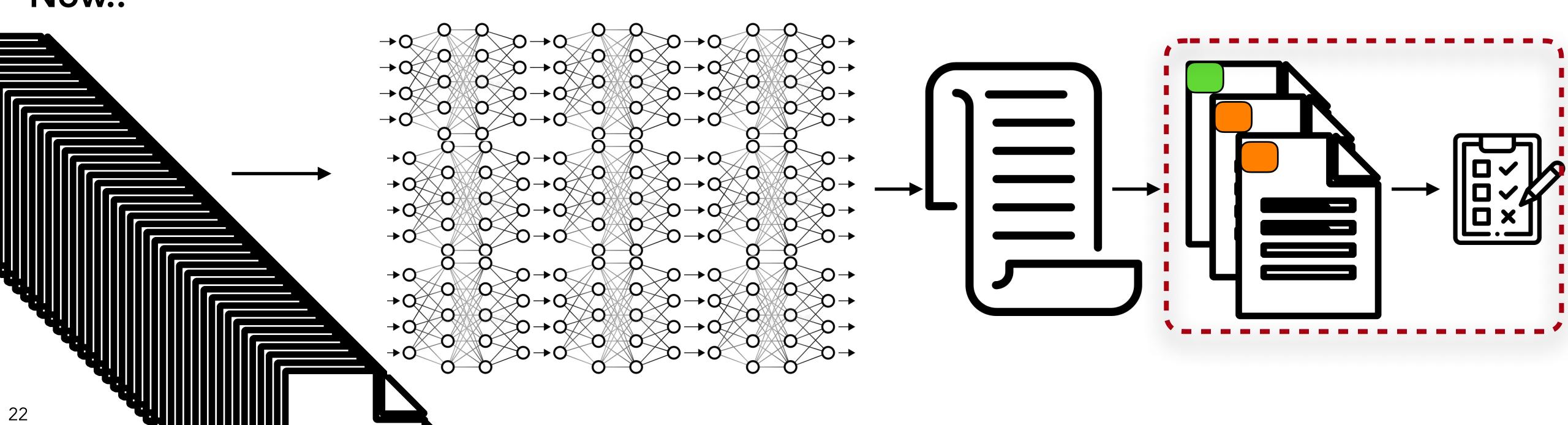
Natural language understanding: (1) Sentiment analysis: Bang et al. [6]/ Liang et al. [114]/ Lopez-Lira and Tang [129]/ Qin et al. [159]/ Wang et al. [218]/ Zhang et al. [251] (2) Text classification: Liang et al. [114] / Peña et al. [154] / Yang and Menczer [233] (3) Natural language inference: Lee et al. [105] / Qin et al. [159] (4) Others: Choi et al. [23] / Riccardi and Desai [166] / Tao et al. [184] Reasoning: Bang et al. [6] / Bian et al. [9] / Frieder et al. [45] / Fu et al. [47] / Gendron et al. [56] / Jiang et al. [86] / Liévin et al. [117] Liu et al. [124] / Orrù et al. [147] / Pan et al. [151] / Qin et al. [159] / Saparov et al. [170] / Wu et al. [227] / Wu et al. [226] Xu et al. [229] / Zhuang et al. [265] / Zhang et al. [244] Natural Natural language generation: language (1) Summarization: Bang et al. [6] / Liang et al. [114] / Pu and Demberg [158] / Qin et al. [159] processing (2) Dialogue: Bang et al. [6] / Lin and Chen [121] / Qin et al. [159] / Zheng et al. [259] (3) Translation: Bang et al. [6] / Lyu et al. [130] / Wang et al. [208] (4) Question answering: Bai et al. [5] / Bang et al. [6] / Bian et al. [9] / Laskar et al. [102] / Liang et al. [114] / Qin et al. [159] (5) Others: Chen et al. [20] / Chia et al. [22] / Pu and Demberg [158] Multilingual: Abdelali et al. [1] / Ahuja et al. [2] / Bang et al. [6] / Lai et al. [100] / Zhang et al. [250] Factuality: Gekhman et al. [55] / Honovich et al. [74] / Manakul et al. [133] / Min et al. [138] / Pezeshkpour [156] / Wang et al. [204] Robustness: Li et al. [111] / Liu et al. [123] / Wang et al. [207] / Wang et al. [206] / Yang et al. [234] / Zhao et al. [258] Zhu et al. [264] / Zhuo et al. [267] Ethics and biases: Cao et al. [16] / Deshpande et al. [35] / Dhamala et al. [37] / Ferrara [42] / Gehman et al. [53] Robustness / Ethics/ Hartmann et al. [65] / Hendrycks et al. [69] / Parrish et al. [153] / Rutinowski et al. [167] / Sheng et al. [175] Biases/ Trustworthiness Simmons [176] / Wang et al. [209] / Zhuo et al. [266] / Zhao et al. [256] Trustworthiness: Hagendorff and Fabi [62] / Wang et al. [201] / Liu et al. [123] / Li et al. [113] / Rawte et al. [163] Xie et al. [228] / Zhang et al. [253] Deroy et al. [34] / Frank [44] / Nay et al. [139] / Wu et al. [224] / Ziems et al. [269] Social science Mathematics: Arora et al. [3] / Bubeck et al. [15] / Collins et al. [27]/ Dao and Le [31] / Wei et al. [221] / Wu et al. [225] Yuan et al. [241] / Yu et al. [237] Natural science General science: Arora et al. [3] / Castro Nascimento and Pimentel [18] / Guo et al. [61] & engineering Engineering: Bubeck et al. [15] / Liu et al. [125] / Pallagani et al. [150] / Sridhara et al. [181] / Valmeekam et al. [195] Valmeekam et al. [194] / Zhuang et al. [265] Medical queries: Chervenak et al. [21] / Duong and Solomon [39] / Hamidi and Roberts [63] / Holmes et al. [73] Jahan et al. [81] / Johnson et al. [87] / Samaan et al. [169] / Thirunavukarasu et al. [186] Medical applications Medical examination: Gilson et al. [57] / Kung et al. [97] Medical assistants: Cascella et al. [17] / Khan et al. [93] / Lahat et al. [99] / Lyu et al. [131] / Oh et al. [143] / Wang et al. [217] Huang et al. [77] / Karpas et al. [90] / Parisi et al. [152] / Qin et al. [160] / Qin et al. [161] / Schick et al. [172] / Shen et al. [174] Agent applications Education: Dai et al. [30] / citetde Winter [32] / citetHellas et al. [67] / Wang and Demszky [210] / Wei et al. [221] Search and recommendation: Dai et al. [29] / Fan et al. [40] / Lanzi and Loiacono [101] / Sun et al. [183] / Thakur et al. [185] Xu et al. [232] / Yuan et al. [240] / Zhang et al. [246] Other applications Personality testing: Bodroza et al. [10] / Jentzsch and Kersting [84] / Liang et al. [115] / Safdari et al. [168] / Song et al. [180] / Wang et al. [212] Specific tasks: Lanzi and Loiacono [101] / Le and Zhang [103] / Wang et al. [216]

Xiezhi [59]/MMLU [70]/ C-Eval [78]/OpenLLM [80]/DynaBench [94]/Chatbot Arena [128]/AlpacaEval [112]/HELM [114]/BIG-bench [182]



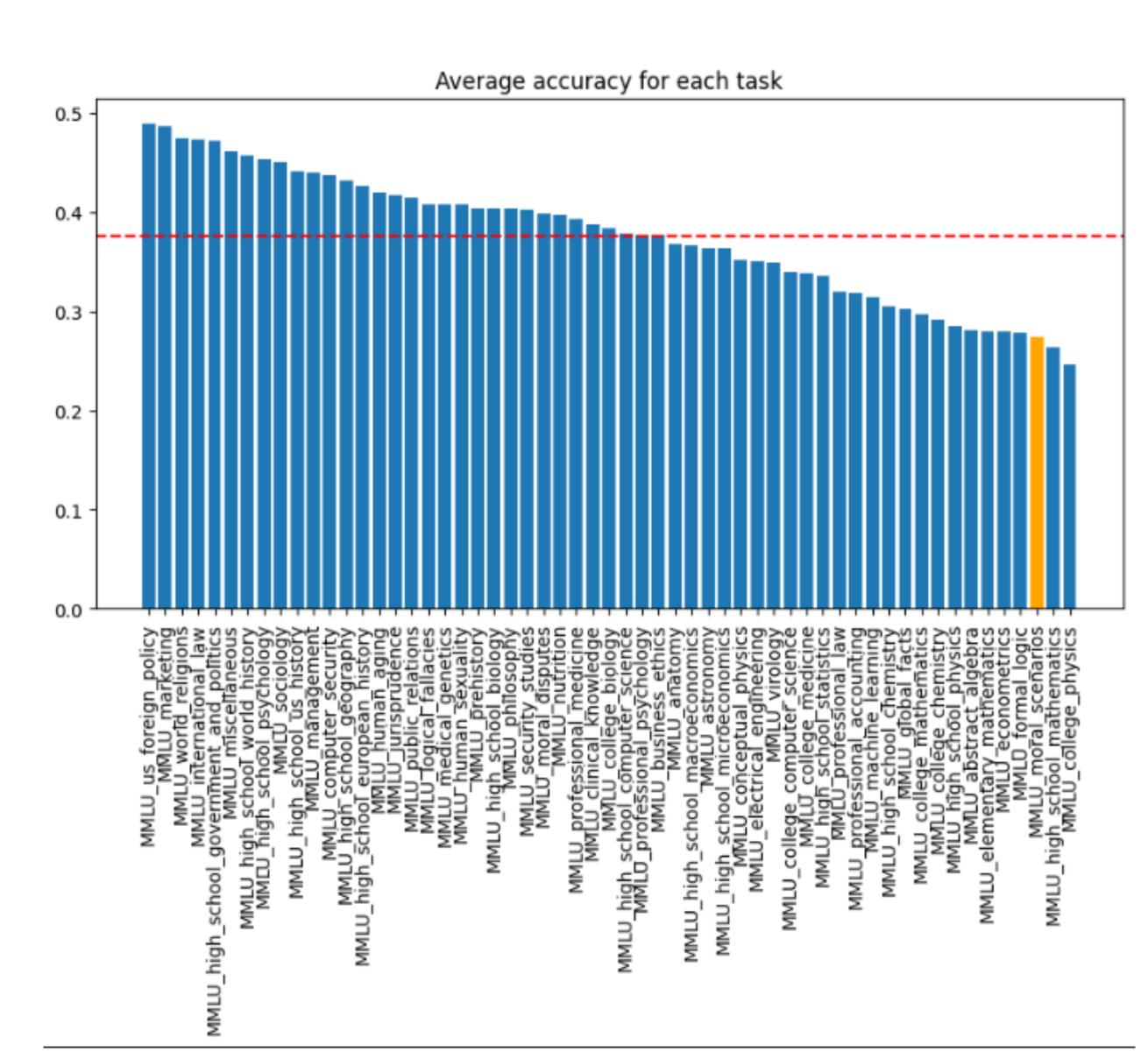
Benchmarks





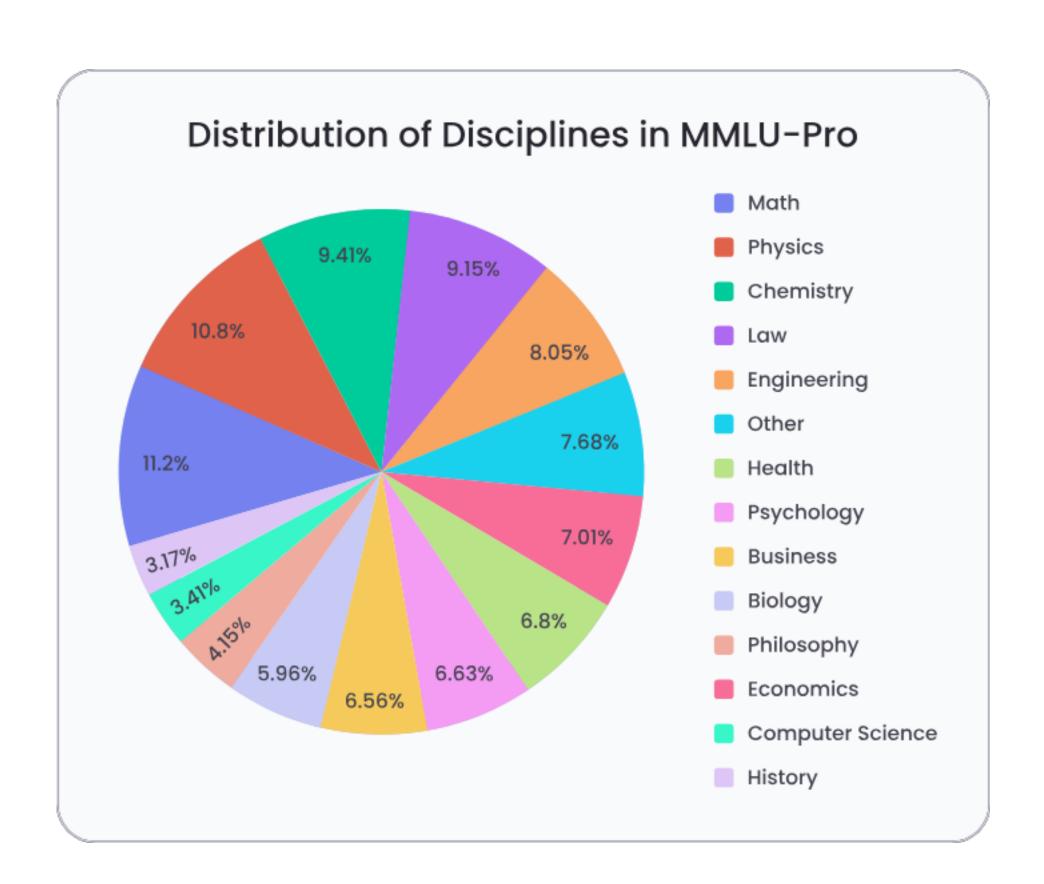
Evaluation of LLMs

- Almost exclusively on downstream tasks, as opposed to intrinsic metrics
 - Intrinsic metrics, e.g. perplexity
- Few popular multitask benchmarks
 - GLUE Language Understanding Tasks
 - SuperGLUE Language Understanding Tasks
 - HellaSwag Commonsense Reasoning
 - Truthful QA Fact Verification
 - MMLU Massive Multitask Language Understanding, 15908 knowledge and reasoning questions in 57 areas including medicine, mathematics, computer science, law, and others
 - GSM 8K Grade School Math
 - BigBench subsumes some of these benchmarks



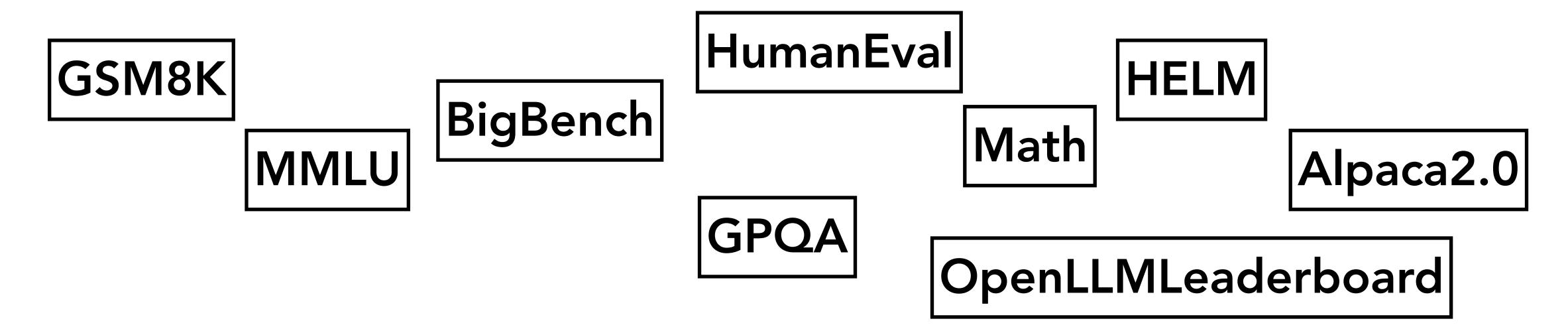
MMLU

- Massive Multi-Task Language Understanding Benchmark
- 57 tasks, across disciplines multiple-choice
- "collected by graduate and undergraduate students from freely available sources online"
- Evaluated on GPT-3 using few-shot prompting
- Critique
 - Overused
 - Really about testing knowledge, not language understanding
- MMLU-Pro [Wang+ 2024]:
 - Removed noisy/trivial questions from MMLU
 - Expanded 4 choices to 10 choices
 - Chain-of-thought prompting





Benchmarks Today



Humanity's Last Exam

Organizing Team

Long Phan*1, Alice Gatti*1, Ziwen Han*2, Nathaniel Li*1,

Josephina Hu², Hugh Zhang[‡], Chen Bo Calvin Zhang², Mohamed Shaaban², John Ling², Sean Shi², Michael Choi², Anish Agrawal², Arnav Chopra², Adam Khoja¹, Ryan Kim[†], Richard Ren¹, Jason Hausenloy¹, Oliver Zhang¹, Mantas Mazeika¹,

Summer Yue**2, Alexandr Wang**2, Dan Hendrycks**1

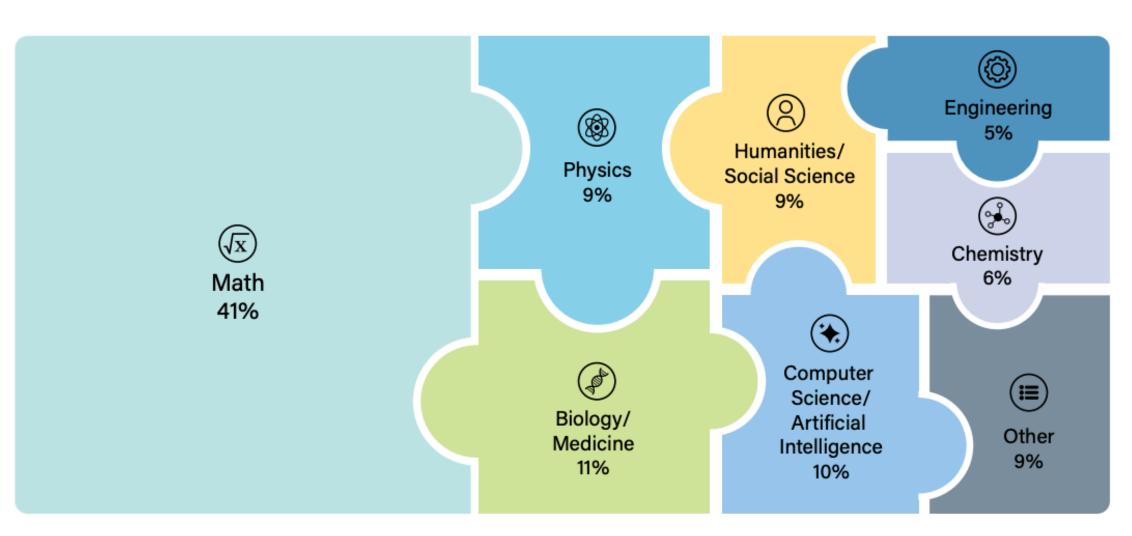


Figure 3: HLE consists of 2,700 exam questions in over a hundred subjects, grouped into high level categories here. We provide a more detailed list of subjects in Appendix B.3.

¹ Center for AI Safety, ² Scale AI



Language models are great at generating language



Write With Transformer distil-gpt2 (1)





Trigger autocomplete or tab



Why is Skub the greatest product on the market today? I'll give you three great reasons:

- 1) Skub is the most reliable solution for your personal, professional and personal life. It offers a great combination of quality and availability.
- 2) Skub is the most versatile solution for all your clients. It also has a great selection of advanced software, features and products.
- 3) Skub's unique design and its strong functional design allows you to easily customize your products. It provides great products and services with all the latest features. It has been designed for everyone with a clear look, without needing any further extra effort.

ChatGPT ~

Why is skub the greatest product on the market today? I'll give you three great reasons:

- 1. Unparalleled Versatility Whether you need it for home, work, or play, Skub adapts to every situation with ease. No other product offers such a wide range of applications!
- 2. Cutting-Edge Innovation Skub is designed with the latest technology and backed by years of research. It's not just a product—it's a revolution in convenience and efficiency.
- 3. Undeniable Popularity The people have spoken, and Skub reigns supreme. With rave reviews and an ever-growing fan base, it's clear that Skub isn't just a trend—it's the future.

What more could you ask for? Skub is, simply put, the best.



USC Viterbi

Current LLM evaluation is done on a few popular benchmarks

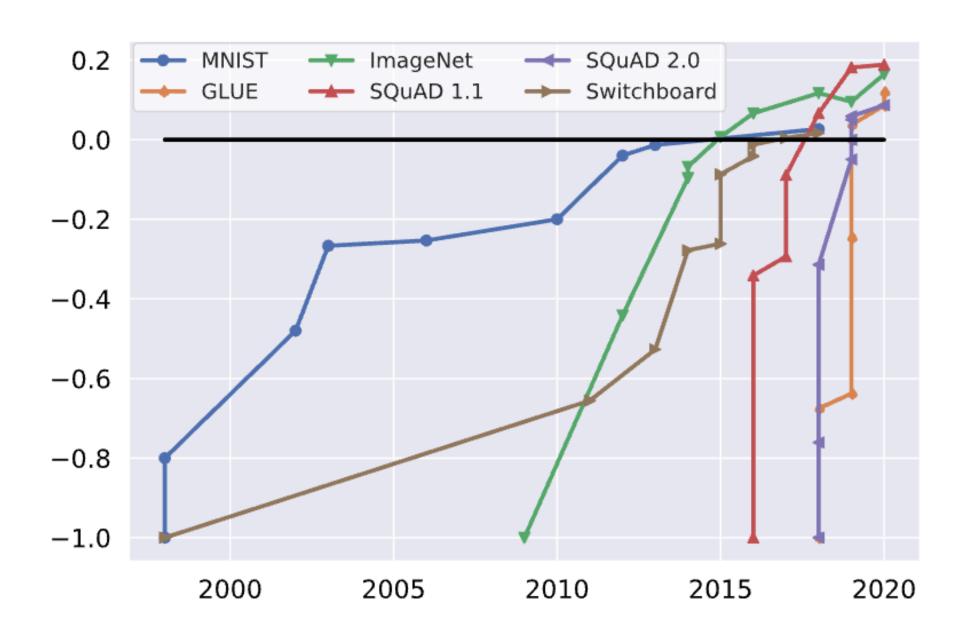


Figure 1: Benchmark saturation over time for popular benchmarks, normalized with initial performance at minus one and human performance at zero.

General	MMLU (Hendrycks et al., 2021a), MMLU-Pro (Wang et al., 2024b), IFEval (Zhou et al., 2023)
Math and reasoning	GSM8K (Cobbe et al., 2021), MATH (Hendrycks et al., 2021b), GPQA (Rein et al., 2023), ARC-Challenge (Clark et al., 2018)
Code	HumanEval (Chen et al., 2021), MBPP (Austin et al., 2021), HumanEval+ (Liu et al., 2024a), MBPP EvalPlus (base) (Liu et al., 2024a), MultiPL-E (Cassano et al., 2023)
Multilinguality	MGSM (Shi et al., 2022), Multilingual MMLU (internal benchmark)
Tool-use	Nexus (Srinivasan et al., 2023), API-Bank (Li et al., 2023b), API-Bench (Patil et al., 2023), BFCL (Yan et al., 2024)
Long context	ZeroSCROLLS (Shaham et al., 2023), Needle-in-a-Haystack (Kamradt, 2023), InfiniteBench (Zhang et al., 2024)

Table 16 Post-training benchmarks by category. Overview of all benchmarks we use to evaluate post-trained Llama 3 models, ordered by capability.

Model	Prompt	MMLU	GPQA	МАТН	HumanEval	MGSM ^[1]	DROP ^[1] (F1, 3- shot)	SimpleQA
-------	--------	------	------	------	-----------	---------------------	---	----------

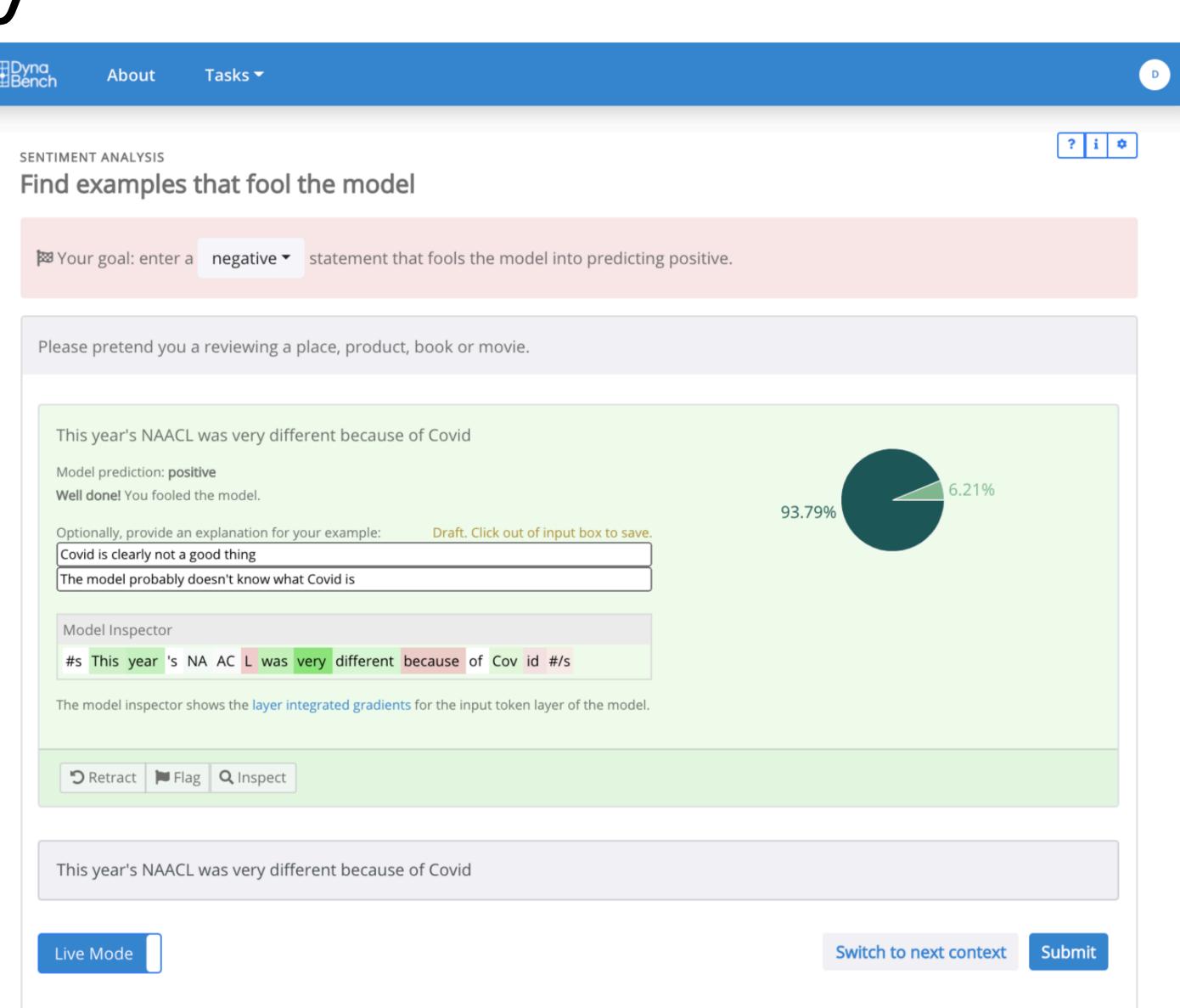
3.1. DeepSeek-R1 Evaluation

	Benchmark (Metric)	Claude-3.5- Sonnet-1022	GPT-40 0513	DeepSeek V3	OpenAI o1-mini		DeepSeek R1
	Architecture	-	-	MoE	-	-	MoE
	# Activated Params	-	-	37B	-	-	37B
	# Total Params	-	-	671B	-	-	671B
	MMLU (Pass@1)	88.3	87.2	88.5	85.2	91.8	90.8
	MMLU-Redux (EM)	88.9	88.0	89.1	86.7	-	92.9
	MMLU-Pro (EM)	78.0	72.6	75.9	80.3	-	84.0
English	DROP (3-shot F1)	88.3	83.7	91.6	83.9	90.2	92.2
	IF-Eval (Prompt Strict)	86.5	84.3	86.1	84.8	-	83.3
	GPQA Diamond (Pass@1)	65.0	49.9	59.1	60.0	75.7	71.5
	SimpleQA (Correct)	28.4	38.2	24.9	7.0	47.0	30.1
	FRAMES (Acc.)	72.5	80.5	73.3	76.9	-	82.5
	AlpacaEval2.0 (LC-winrate)	52.0	51.1	70.0	57.8	-	87.6
	ArenaHard (GPT-4-1106)	85.2	80.4	85.5	92.0	-	92.3
	LiveCodeBench (Pass@1-COT)	38.9	32.9	36.2	53.8	63.4	65.9
Code	Codeforces (Percentile)	20.3	23.6	58.7	93.4	96.6	96.3
Code	Codeforces (Rating)	717	759	1134	1820	2061	2029
	SWE Verified (Resolved)	50.8	38.8	42.0	41.6	48.9	49.2
	Aider-Polyglot (Acc.)	45.3	16.0	49.6	32.9	61.7	53.3
Math	AIME 2024 (Pass@1)	16.0	9.3	39.2	63.6	79.2	79.8
	MATH-500 (Pass@1)	78.3	74.6	90.2	90.0	96.4	97.3
	CNMO 2024 (Pass@1)	13.1	10.8	43.2	67.6	-	78.8
	CLUEWSC (EM)	85.4	87.9	90.9	89.9	-	92.8
	C-Eval (EM)	76.7	76.0	86.5	68.9	-	91.8
	C-SimpleQA (Correct)	55.4	58.7	68.0	40.3	-	63.7



DynaBench

- Open-source platform for dynamic dataset creation and model benchmarking
- Adversarial benchmark generation
 - Create harder and harder examples
- Critique: might lead to outof-distribution evaluation



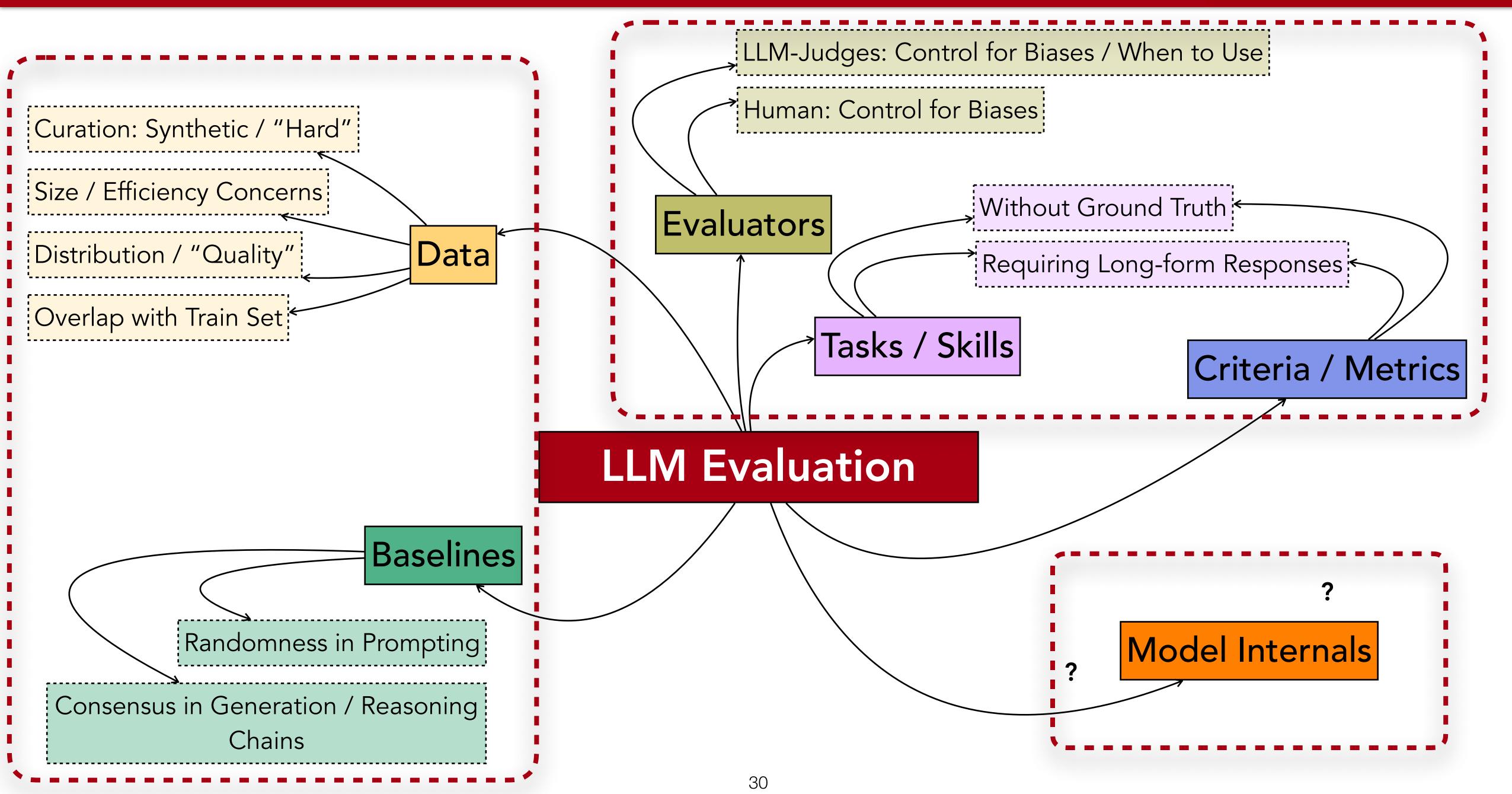


Benchmark Tasks vs. Real-World Tasks

- We have evidence language models are getting better and better at math, inference, coding and many logical / analytical tasks
- How about we go bold and test them directly on real-world tasks with real stakes?
- Caveat: this will involve humans, and much manual labor



Some of the content in the following slides contains sensitive content and can be upsetting





My Research

Annotators with Attitudes: How Annotator Beliefs And Identities Bias Toxic Language Detection

Maarten Sap^{♥♠} Swabha Swayamdipta[♠] Laura Vianna[♦] Xuhui Zhou[♣] Yejin Choi^{♥♠} Noah A. Smith^{♥♠}

[⋄]Paul G. Allen School of Computer Science, University of Washington, Seattle, WA, USA

[♠] Allen Institute for AI, Seattle, WA, USA

[⋄] Department of Psychology, University of Washington, Seattle, WA, USA

[♠] Georgia Institute of Technology, Atlanta, GA, USA

OATH-Frames: Characterizing Online Attitudes Towards Homelessness with LLM Assistants

Jaspreet Ranjit Brihi Joshi Rebecca Dorn Laura Petry Olga Koumoundouros Jayne Bottarini Peichen Liu Eric Rice Swabha Swayamdipta

Evic Rice Swabha Swayamdipta

■ Laura Petry

■ Peichen Liu

- *Thomas Lord Dept. of Computer Science, University of Southern California
- ◆Suzanne-Dwork School of Social Work, University of Southern California

We're Afraid Language Models Aren't Modeling Ambiguity

Alisa Liu Zhaofeng Wu Julian Michael Alane Suhr Alane Suhr Yejin Choi Alexander Koller Swabha Swayamdipta Noah A. Smith Yejin Choi Paul G. Allen School of Computer Science & Engineering, University of Washington Allen Institute for AI University of Southern California UC Berkeley Saarland University New York University Massachusetts Institute of Technology

Evaluating Capabilities where there are No Right Answers:

- Evaluating free-text explanations / reasoning chains
- Evaluating LLM assistants in real world problems (e.g. social work, emergency response)
- Evaluating personalized agents, simulations
- Future: Biases in Human Evaluators of LLMs

Tasks / Skills

Criteria / Metrics

Evaluators

Reframing Human-AI Collaboration for Generating Free-Text Explanations

Sarah Wiegreffe* Jack Hessel† Swabha Swayamdipta†

Mark Riedl* Yejin Choi[†]

*School of Interactive Computing, Georgia Institute of Technology

†Allen Institute for Artificial Intelligence

[⋄]Paul G. Allen School of Computer Science and Engineering, University of Washington

REV: Information-Theoretic Evaluation of Free-Text Rationales

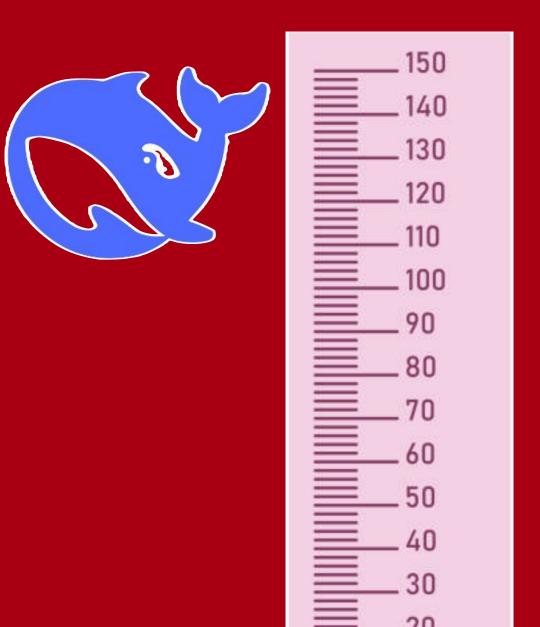
Hanjie Chen^{♥*} Faeze Brahman^{♠♦} Xiang Ren^{♠♠} Yangfeng Ji[♥] Yejin Choi^{♠♦} Swabha Swayamdipta[♠]

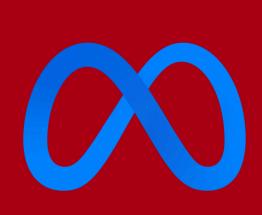
Department of Computer Science, University of Virginia
Allen Institute for AI University of Southern California

[⋄] Paul G. Allen School of Computer Science & Engineering, University of Washington



Rethinking Evaluation Benchmarks





Benchmarks vs. real-world tasks





Need AI help for local politicians to understand public attitudes on homelessness to frame the most resonant message to inform public policy.

USC Annenberg School for Communication and Journalism

Easy peasy, we've got these amazing language models - let me do this for you in a couple of weeks, tops



School of Engineering



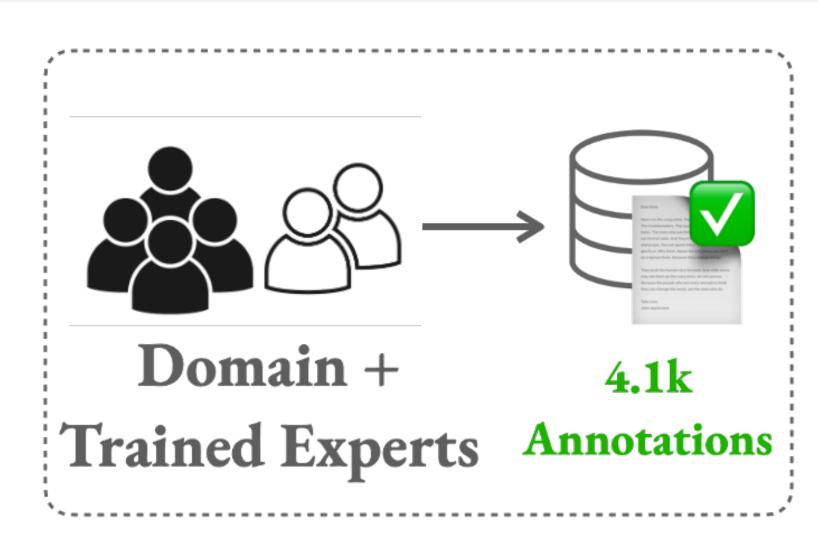
@mention how about do something for all the homeless veterans and americans on the streets? you'll give millions to the immigrants crossing our so called secure border but our homeless americans get nothing. You are merely trying to buy votes



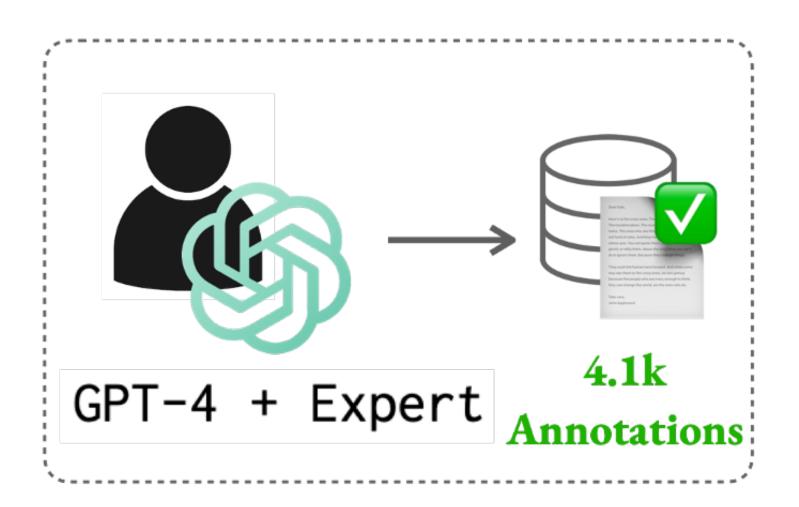
~2.5M tweets on homelessness between 2021-2023

government critique

money allocation deserving/undeserving

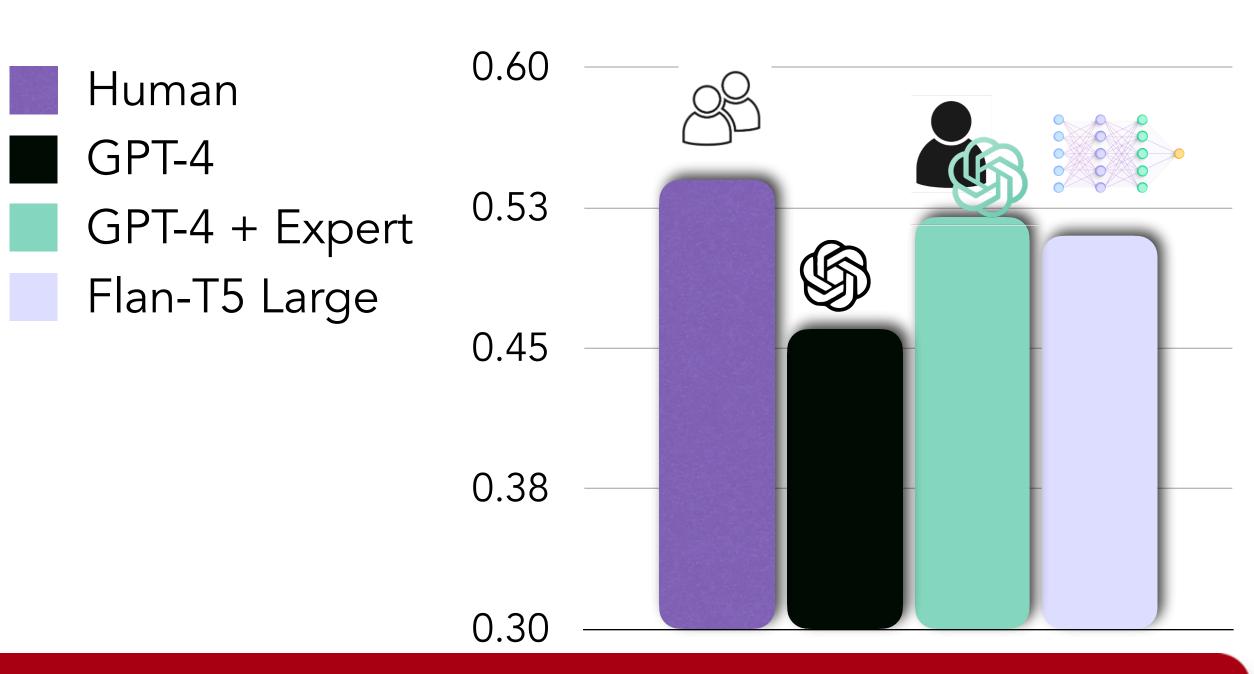


Building a test set took more than 7-8 months of laborious manual work



Results in 6.5x speedup in annotation time!

LMs in Communications / Social Work: Informing policies on homelessness



With some effort, language models can be used as assistants for determining variables which could inform communication and public policy

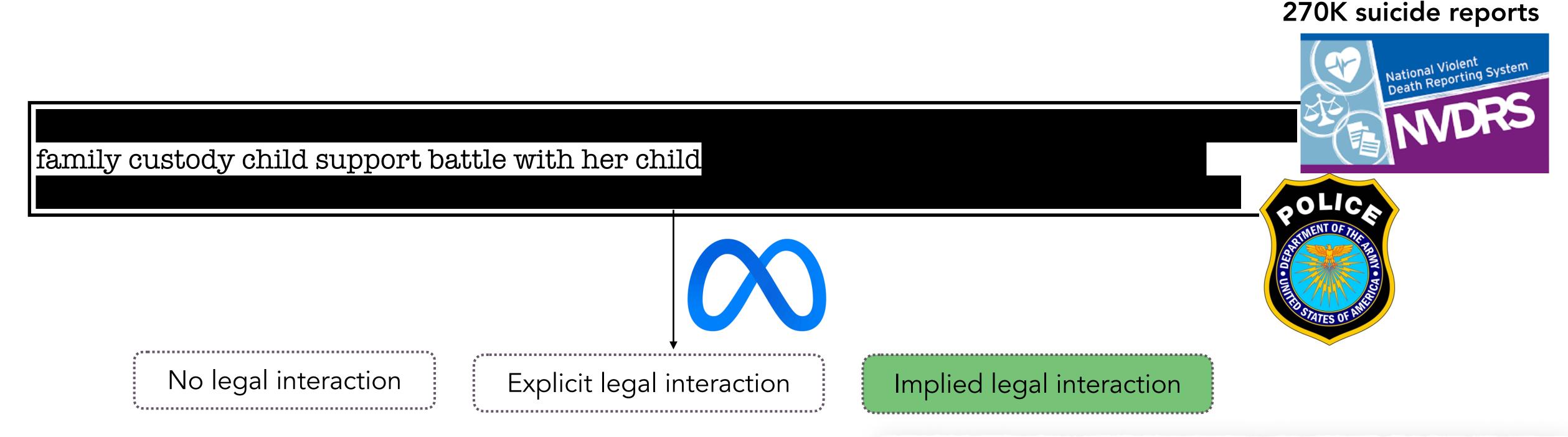




Heard you have been helping out colleagues with language models. I need some help with analyzing suicide reports that help me validate my hypotheses on novel factors for designing suicide interventions.



How often do victims interact with non-clinical personnel (e.g. legal professionals) in days before their death?





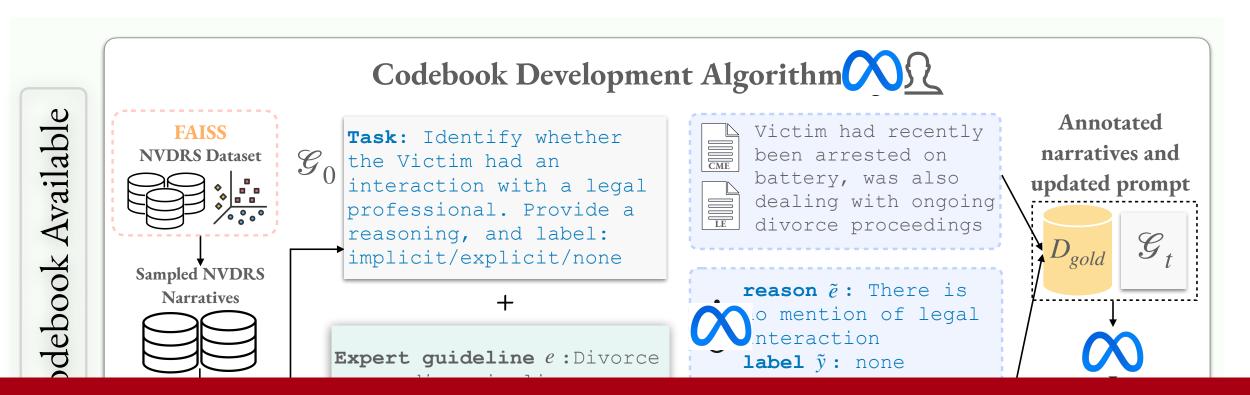
LMs in Social Work

Manual Codebook Development:

Legal Interactions







Instead of performance alone, speedup in manual labor at the same performance level should be an important metric in real world tasks

Implicit Interaction

Explicit Interaction

No Interaction

Expert

Codebook

Model	\mathcal{G}_0	\mathcal{G}_{expert}	\mathcal{G}_{12}
Meta-Llama-3-70B	0.57	0.78	0.79
Qwen2.5-32B	0.68	0.77	0.77
Qwen2.5-14B	0.63	0.77	0.76

Speedup in process: weeks to hours!

Our research was going at the pace of a bullock cart and you came at it with a jet plane!

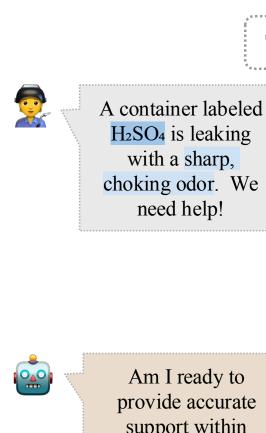
10.4% of 270K suicide narratives had evidence of legal interactions!

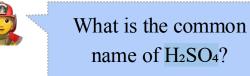




Public Health: Emergency Response

As a volunteer emergency responder, I'd like to automate part of the emergency response process using language models





Task 1: Chemical Translation

Search online.



Am I ready to provide accurate support within seconds?

The common name of H₂SO₄ is Sulfuric Acid.





How do I best respond to an incident involving Sulfuric Acid?

Find "Sulfuric Acid" in a 400page Emergency Response Guidebook (ERG).



Hazard: Highly corrosive and reactive. Recommended PPE: Latex gloves, face shield, and respirator. Action: Evacuate a 10-foot radius around the spill, ensure proper ventilation, and establish a 300-foot exclusion zone downwind.

30 seconds later



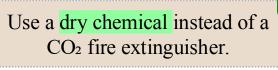
I'm first on scene with no fire extinguisher. What's the safest way to handle a Sulfuric Acid fire?

Wait for someone with more resources or experience to arrive.











30 seconds later





Model

I. Translation

II. Incident Response

III. Exam

Acc. (↑)

49.0

47.3

<u>60.0</u>

58.0

60.0

63.9

Real-world tasks are still hard to evaluate directly, need to be broken down. And human involvement is still important

GPT-40

92.4

68.0

52.7

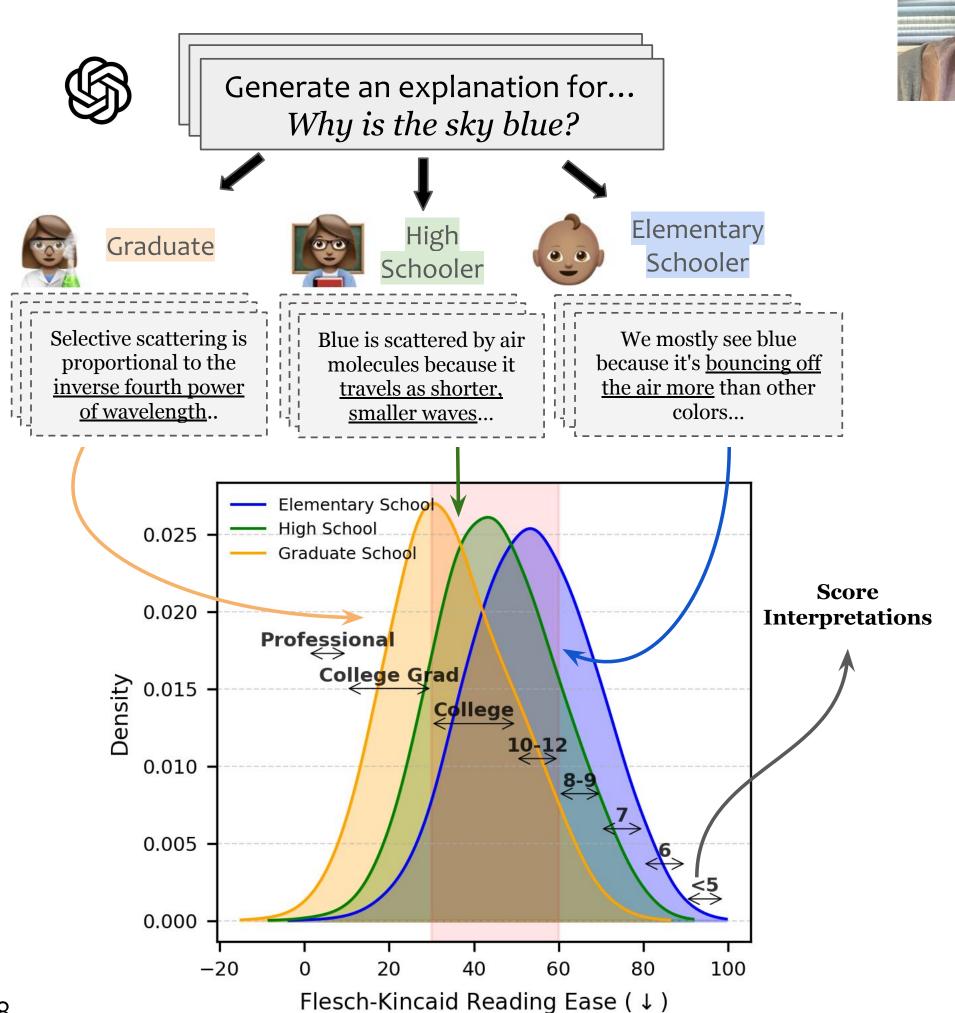
80.2

2.0

Surana, Ye & Swayamdipta. NeurIPS LLM-Eval. 2025.

USC Viterbi

Education



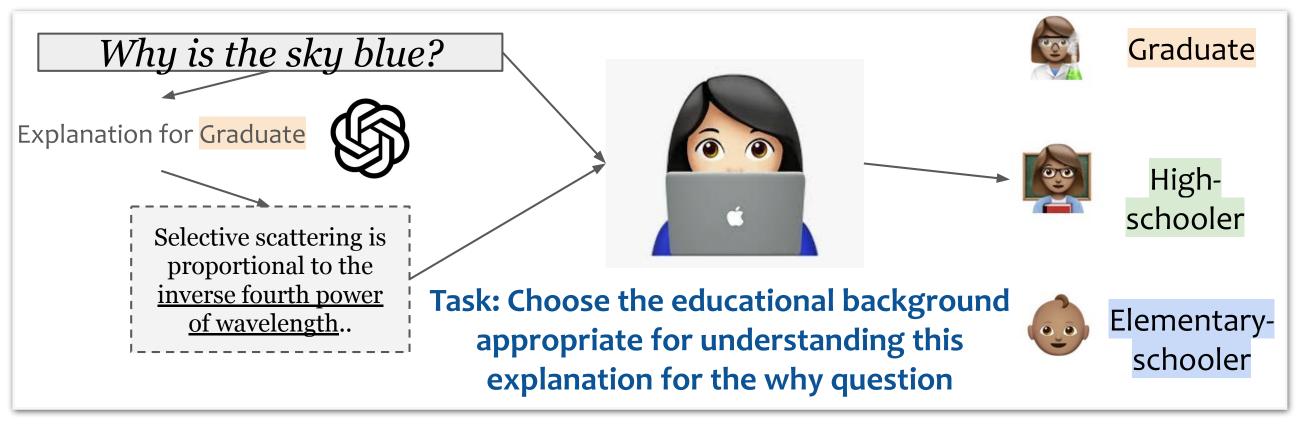


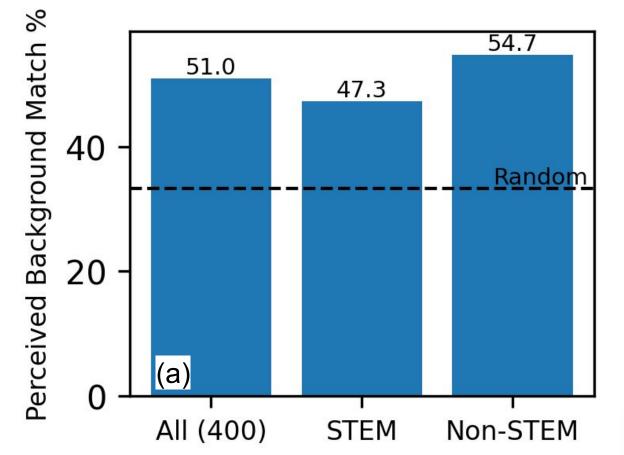
I feel very shy around my teacher, so I ask all my questions to ChatGPT.

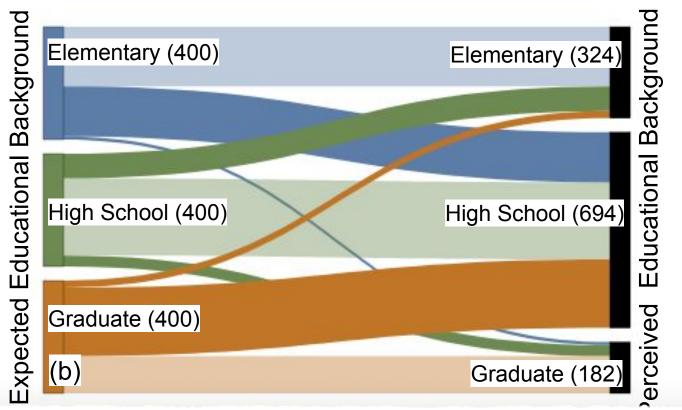
Umm... does it help you every time?

I dunno...









Joshi, ..., **Swayamdipta** and Ren. EMNLP 2025



Evaluator Efforts and Effects on Human Evaluation

- Human evaluators seem to like language model responses
 - especially those with limited training / experience
 - experts tend to disagree more with LMs, but can often miss important details too

The Impact of Generative AI on Critical Thinking: Self-Reported Reductions in Cognitive Effort and Confidence Effects From a Survey of Knowledge Workers

Hao-Ping (Hank) Lee Carnegie Mellon University Pittsburgh, Pennsylvania, USA haopingl@cs.cmu.edu

Ian Drosos Microsoft Research Cambridge, United Kingdom t-iandrosos@microsoft.com Advait Sarkar
Microsoft Research
Cambridge, United Kingdom
advait@microsoft.com

Sean Rintel
Microsoft Research
Cambridge, United Kingdom
serintel@microsoft.com

Nicholas Wilson Microsoft Research Cambridge, United Kingdom niwilson@microsoft.com Lev Tankelevitch
Microsoft Research
Cambridge, United Kingdom
levt@microsoft.com

Richard Banks
Microsoft Research Cambridge
Cambridge, United Kingdom
rbanks@microsoft.com

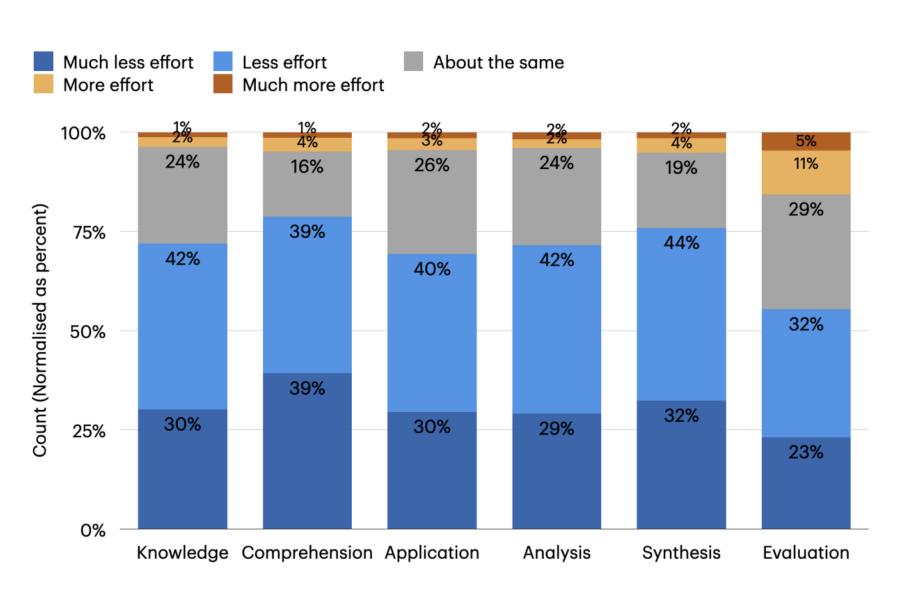
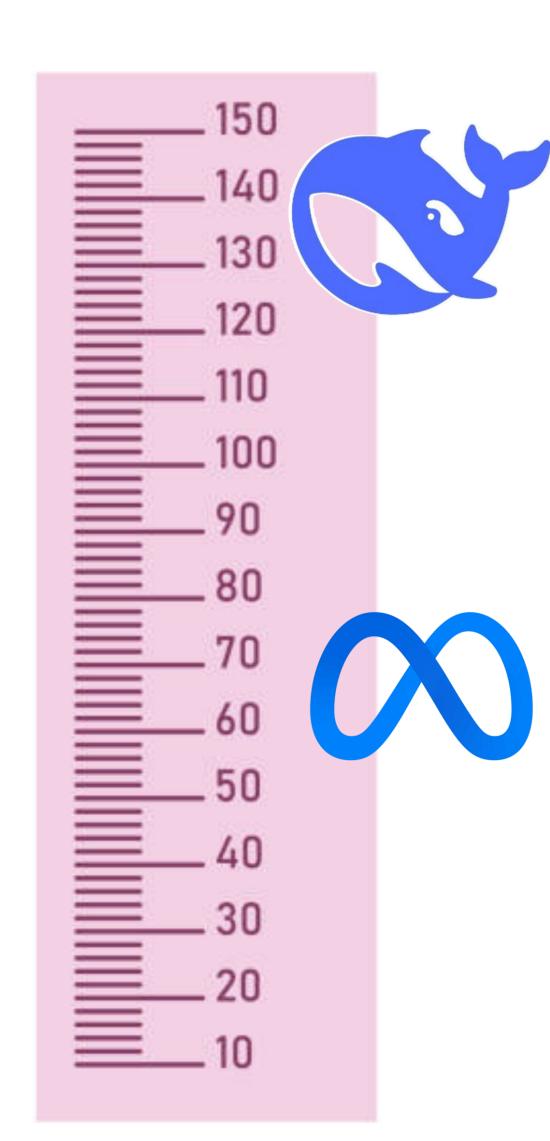


Figure 2: Distribution of perceived effort (%) in cognitive activities (based on Bloom's taxonomy) when using a GenAI tool compared to not using one.

Rethinking Benchmark Tasks

- Our notion of task difficulty might have limitations; using real-world problems to guide tasks is as important as collecting "hard" benchmarks.
 - Over-reliance on a few select benchmarks (each with a single aggregate metric) might give us a false sense of progress
- Quantitative evaluation of language models on real-world tasks is hard, requires humans with domain expertise
 - Task itself needs to be broken down into smaller parts for ease of evaluation
- Harder tasks (no right answers, experts might disagree): Human verification, real data
 - Easier tasks (possible right answers): Auto evals, synthetic data





Evaluation of Language Models

- Important task, if there ever was one! Given a fixed model, how "good" is it?
- Perplexity is still of extreme importance, even though an intrinsic metric
- TL;DR
 - There is no one true evaluation; choose the evaluation depending on what you're trying to measure.
 - Always look at the individual instances and the predictions.
 - There are many aspects to consider: capabilities, safety, costs, realism.
 - Clearly state the rules of the game (methods versus models/systems).
- Benchmarks vs. Vibes