

IMPERIAL

# Benchmarking and Evaluation

11/12/2025

Shamsuddeen Muhammad  
Google DeepMind Academic Fellow,  
Imperial College London  
<https://shmuhammadd.github.io/>

Idris Abdulkumin  
Postdoctoral Research Fellow,  
DSFSI, University of Pretoria  
<https://abumafrim.github.io/>

# Suggested Reading

- Challenges and Opportunities in NLP Benchmarking
- Measuring Massive Multitask Language Understanding
- Holistic Evaluation of Language Models
- AlpacaEval

# What is Benchmark in NLP?

A **benchmark** refers to a standardized dataset and associated evaluation protocol used to measure and compare the performance of NLP systems on specific tasks.

Benchmarks serve as critical tools for driving scientific progress, fostering reproducibility, and providing objective criteria for evaluating models (Bowman and Dahl, 2021; Gildea, 2001).

## The Beyond the Imitation Game Benchmark (BIG-bench) is a collaborative benchmark intended to probe large language models future capabilities. More than 200 tasks.

[Submitted on 9 Jun 2022 (v1), last revised 12 Jun 2023 (this version, v3)]

### Beyond the Imitation Game: Quantifying and extrapolating the capabilities of language models

Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska, Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W. Kocurek, Ali Safaya, Ali Tazarov, Alice Xiang, Alicia Parrish, Allen Nie, Aman Hussain, Amanda Askell, Amanda Dsouza, Ambrose Slone, Ameet Rahane, Anantharaman S. Iyer, Anders Andreassen, Andrea Madotto, Andrea Santilli, Andreas Stuhlmüller, Andrew Dai, Andrew La, Andrew Lampinen, Andy Zou, Angela Jiang, Angelica Chen, Anh Vuong, Animesh Gupta, Anna Gottardi, Antonio Norelli, Anu Venkatesh, Arash Gholamidavoodi, Arfa Tabassum, Arul Menezes, Arun Kirubarajan, Asher Mullokandov, Ashish Sabharwal, Austin Herrick, Avia Efrat, Ayukut Erdem, Ayla Karakaş, B. Ryan Roberts, Bao Sheng Loe, Barret Zoph, Bartłomiej Bojanowski, Batuhan Özüy, Behnam Hedayatnia, Behnam Neyshabur, Benjamin Inden, Benno Stein, Berk Ekmekcı, Bill Yuchen Lin, Blake Howald, Bryan Orionion, Cameron Diao, Cameron Dour, Catherine Stinson, Cedrick Argueta, César Ferri Ramírez, Chandan Singh, Charles Rathkopf, Chenlin Meng, Chitta Baral, Chiyu Wu, Chris Callison-Burch, Chris Waites, Christian Voigt, Christopher D. Manning, Christopher Potts, Cindy Ramirez, Clara E. Rivera, Clemencia Siro, Colin Raffel, Courtney Ashcraft, Cristina Garbacea, Damien Silao, Dan Garrette, Dan Hendrycks, Dan Kliman, Dan Roth, Daniel Freeman, Daniel Khashabi, Daniel Levy, Daniel Moseguí González, Danielle Perszyk, Danny Hernandez, Danqi Chen, Daphne Ippolito, Dar Gilboa, David Dohan, David Drakard, David Jurgens, Debajyoti Datta, Deep Ganguli, Denis Emelin, Denis Kleyko, Deniz Yuret, Derek Chen, Derek Tam, Dieuwke Hupkes, Diganta Misra, Dilyat Buzan, Dimitri Coelho Mollo, Diyi Yang, Dong-Ho Lee, Dylan Schrader, Ekaterina Shutova, Ekin Dogus Cubuk, Elad Segal, Eleanor Hagerman, Elizabeth Barnes, Elizabeth Donoway, Ellis Pavlick, Emanuele Rodola, Emma Lam, Eric Chu, Eric Tang, Erkut Erdem, Ernie Chang, Ethan A. Chi, Ethan Dyer, Ethan Jerzak, Ethan Kim, Eunice Engefu Manyasi, Evgenij Zheltonozhskii, Fanyue Xia, Fatemeh Siar, Fernando Martínez-Plumed, Francois Chollet, Frieda Rong, Gaurav Mishra, Genta Indra Winata, Gerard de Melo, Germán Kruszewski, Giambattista Parascandolo, Giorgio Mariani, Gloria Wang, Gonzalo Jaimovich-López, Gregor Betz, Guy Gur-Ari, Hana Galijasevic, Hannah Kim, Hannah Rashkin, Hannaneh Hajishirzi, Harsh Mehta, Hayden Bogar, Henry Shevelin, Hinrich Schütze, Hiromu Yakura, Hongming Zhang, Hugh Mee Wong, Ian Ng, Isaac Noble, Jaap Jumelet, Jack Geissinger, Jackson Kernion, Jacob Hilton, Jaehoon Lee, Jaime Fernández Fisac, James B. Simon, James Koppel, James Zheng, James Zou, Jan Kocoń, Jana Thompson, Janelle Wingfield, Jared Kaplan, Jarema Radom, Jascha Sohl-Dickstein, Jason Phang, Jason Wei, Jason Yosinski, Jekaterina Novikova, Jelle Bosscher, Jennifer Marsh, Jeremy Kim, Jeroen Taal, Jesse Engel, Jesujoba Alabi, Jiacheng Xu, Jiaming Song, Jillian Tang, Joan Waweru, John Burden, John Miller, John U. Balis, Jonathan Batchelder, Jonathan Berant, Jörg Frohberg, Jos Rozen, Jose Hernandez-Orallo, Joseph Boudeman, Joseph Guerr, Joseph Jones, Joshua B. Tenenbaum, Joshua S. Rule, Joyce Chua, Kamil Kanclerz, Karen Livescu, Karl Krauth, Karthik Gopalakrishnan, Katerina Ignatyeva, Katja Markert, Kaustubh D. Dhole, Kevin Gimpel, Kevin Omondi, Kory Mathewson, Kristen Chiafulla, Ksenia Shkaruta, Kumar Shridhar, Kyle McDonnell, Kyle Richardson, Laria Reynolds, Leo Gao, Li Zhang, Liam Dugan, Lianhui Qin, Lidia Contreras-Ochando, Louis-Philippe Morency, Luca Moschella, Lucas Lam, Lucy Noble, Ludwig Schmidt, Lu亨g He, Luis Oliveros Colón, Luke Metz, Lütfi Kerem Şenel, Maarten Bosma, Maarten Sap, Maartje ter Hoeve, Maheen Farooqi, Manaal Faruqui, Mantas Mazeika, Marco Baturan, Marco Marelli, Marco Maru, Maria Jose Ramírez Quintana, Marie Tolkiene, Mario Giulianelli, Martha Lewis, Martin Potthast, Matthew L. Leavitt, Matthias Hagen, Mátyás Schubert, Medina Orduna Baitemirova, Melody Arnaud, Melvin McElrath, Michael A. Yee, Michael Cohen, Michael Gu, Michael Ivanitskiy, Michael Starritt, Michael Strube, Michał Świdrowski, Michele Bevilacqua, Michihiro Yasunaga, Mihir Kale, Mike Cain, Mimee Xu, Mirac Suzgun, Mitch Walker, Mo Tiwari, Mohit Bansal, Moin Aminnaseri, Mor Geva, Mozhdeh Gheini, Mukund Varma T, Nanyun Peng, Nathan A. Chi, Nayeon Lee, Neta Gur-Ari Krakover, Nicholas Cameron, Nicholas Roberts, Nick Doiron, Nicole Martinez, Nikita Nangia, Niklas Deckers, Niklas Muenninghoff, Nitisha Shirish Keskar, Niveditha S. Iyer, Noah Constant, Noah Fiedel, Nuan Wen, Oliver Zhang, Omar Agha, Omar Elbaghdadi, Omer Levy, Owain Evans, Pablo Antonio Moreno Casares, Parth Doshi, Pascale Fung, Paul Pu Liang, Paul Vicol, Pegah Alipoormolabashi, Peiyuan Liao, Percy Liang, Peter Chang, Peter Eckersley, Phu Mon Htut, Pinyu Hwang, Piotr Miłkowski, Piyush Patil, Pouya Pezeshkpour, Priti Oli Qiaozhu Mei, Qing Lyu, Qinlang Chen, Rabin Banjade, Rachel Etta Rudolph, Raefel Gabriel, Rahel Habacker, Ramon Risco, Raphaël Millière, Rhythm Garg, Richard Barnes, Rif A. Saurous, Riku Arakawa, Robbe Raymaekers, Robert Frank, Rohan Sikand, Roman Novak, Roman Sitelew, Ronan LeBras, Rosanne Liu, Rowan Jacobs, Rui Zhang, Ruslan Salakhutdinov, Ryan Chi, Ryan Lee, Ryan Stovall, Ryan Teehan, Rylan Yang, Sahib Singh, Saif M. Mohammad, Sajant Anand, Sam Dillavou, Sam Shleifer, Sam Wiseman, Samuel Gruetter, Samuel R. Bowman, Samuel S. Schoenholz, Sanghyun Han, Sanjeev Kwatra, Sarah A. Rous, Sarit Ghazarian, Sayan Ghosh, Sean Casey, Sebastian Bischoff, Sebastian Gehrmann, Sebastian Schuster, Sepideh Sadeghi, Shadi Hamdan, Sharon Zhou, Shashank Srivastava, Sherry Shi, Shikhar Singh, Shima Asaadi, Shixiang Shane Gu, Shubh Pachchigar, Shubham Toshniwal, Shyam Upadhyay, Shyamolina (Shammie)Debnath, Siamak Shakeri, Simon Thormeyer, Simone Melzi, Siva Reddy, Sneha Priscilla Makini, Soo-Hwan Lee, Spencer Torene, Sriharsha Hatwar, Stanislas Dehaene, Stefan Ermon, Stella Biderman, Stephanie Lin, Stephen Prasad, Steven T. Piantadosi, Stuart M. Sieber, Summer Mishherghi, Svetlana Kiritchenko, Swaroop Mishra, Tal Linzen, Tal Schuster, Tao Li, Tao Yu, Tariq Ali, Tatsu Hashimoto, Te-Lin Wu, Théo Desbordes, Theodore Rothschild, Thomas Phan, Tianle Wang, Tiberius Nkonyili, Timo Schick, Timofei Kornev, Titus Tunduny, Tobias Gerstenberg, Trenton Chang, Trishala Neeraj, Tushar Khot, Tyler Shultz, Uri Shaham, Vedant Misra, Vera Demberg, Victoria Nyamai, Vikas Raunak, Vinay Uday Prabhu, Vishakh Padmakumar, Vivek Sri Kumar, William Fedus, William Saunders, William Zhang, Wout Vossen, Xiang Ren, Xiaoyu Tong, Xinran Zhao, Xinyi Wu, Xudong Shen, Yadollah Yaghoobzadeh, Yair Lakretz, Yangqiu Song, Yasaman Bahri, Yejin Choi, Yichi Yang, Yiding Hao, Yifu Chen, Yonatan Belinkov, Yu Hou, Yufang Hou, Yuntao Bai, Zachary Seid, Zhuoye Zhao, Zijie J. Wang, Zirui Wang, Ziyi Wu (collapse list)

Language models demonstrate both quantitative improvement and new qualitative capabilities with increasing scale. Despite their potentially transformative impact, these new capabilities are as yet poorly characterized. In order to inform future research, prepare for disruptive new model capabilities, and ameliorate socially harmful effects, it is vital that we understand the present and near-future capabilities and limitations of language models. To address this challenge, we introduce the Beyond the Imitation Game benchmark (BIG-bench). BIG-bench currently consists of 204 tasks contributed by 450 authors across 120 institutions. Task descriptions, datasets, and code are available at <https://big-bench.ai>. BIG-bench features 204 tasks that are believed to be

# What is a benchmark in NLP?

A benchmark typically consists of:

- **A dataset**, often split into train/dev/test partitions.
- **A task formulation**, such as classification, generation, or structured prediction.
- **Evaluation metrics**, e.g., accuracy, F1 score, BLEU, or exact match.
- **A leaderboard** (optional), listing systems ranked by performance.

## Benchmark in NLP

A standardized dataset and evaluation protocol for measuring and comparing the performance of NLP systems.



Dataset



Task



Evaluation Metrics

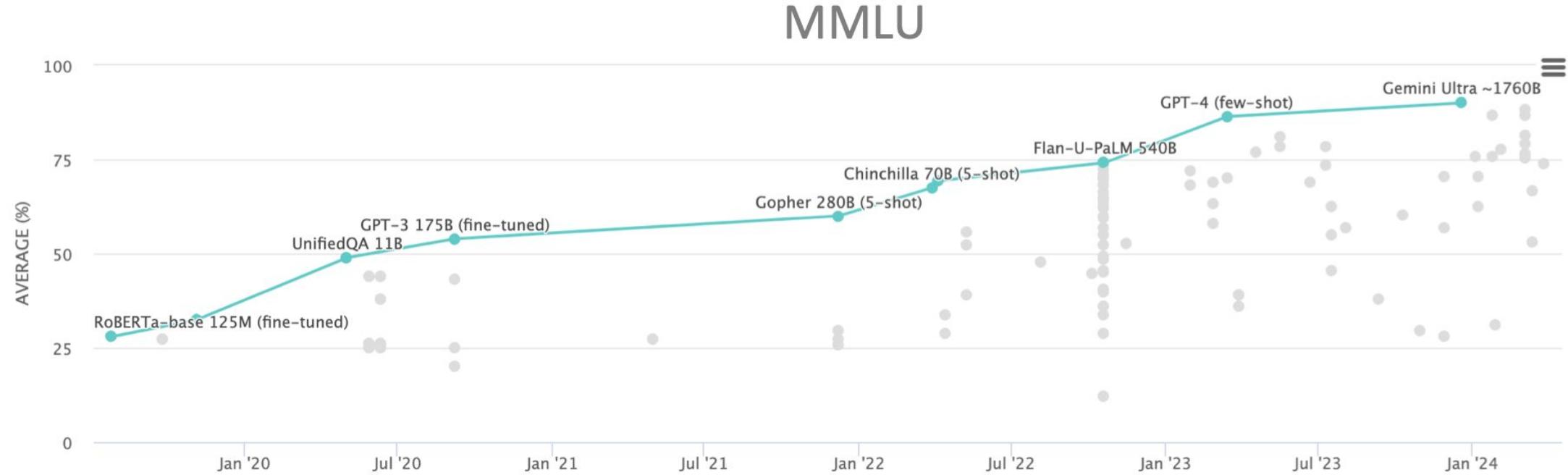


Leaderboard

# Purpose of Benchmark

- **Standardization:** Enables apples-to-apples comparison across methods.
- **Progress Tracking:** Assesses how models improve over time.
- **Challenge Identification:** Reveals specific limitations (e.g., commonsense, bias, multilinguality).
- **Community Building:** Shared goals and competition promote collaboration.

# Benchmarks and evaluations drive progress



Benchmarks and how we drive the progress of the field

# Why Benchmarking & Evaluation Matter

“Without evaluation, progress is an illusion.”

- NLP is advancing rapidly—but **how do we know if we're making real progress?**
- Benchmarks like **GLUE**, **SQuAD**, **BLEU** have become central to research and development.
- Evaluation tells us:
  - Which model performs better
  - Whether a system generalizes across domains or languages
  - If a system is *fair*, *robust*, and *trustworthy*

*Benchmarking and evaluation are the compass of NLP—they guide research, compare systems, and expose limitations.*

# Two major types of evaluations

## Close-ended evaluations

### Example

Text: Read the book, forget the movie!

Label: Negative

## Open ended evaluations

**Context (human-written):** In a shocking finding, scientist discovered a herd of unicorns living in a remote, previously unexplored valley, in the Andes Mountains. Even more surprising to the researchers was the fact that the unicorns spoke perfect English.

**GPT-2:** The scientist named the population, after their distinctive horn, Ovid's Unicorn. These four-horned, silver-white unicorns were previously unknown to science.

Now, after almost two centuries, the mystery of what sparked this odd phenomenon is finally solved.

Dr. Jorge Pérez, an evolutionary biologist from the University of La Paz, and several companions, were exploring the Andes Mountains when they found a small valley, with no other animals or humans. Pérez noticed that the valley had what appeared to be a natural fountain, surrounded by two peaks of rock and silver snow.

# Closed-ended Evaluation

# Closed-ended tasks

- Limited number of potential answers
- Often one or just a few correct answers
- Enables automatic evaluation as in ML

# Closed-ended tasks

- Sentiment analysis: SST / IMDB / Yelp ...

**Example**

Text: Read the book, forget the movie!

Label: Negative

- Entailment: SNLI

**Example**

Text: A soccer game with multiple males playing.

Hypothesis: Some men are playing sport.

Label: Entailment

- Name entity recognition: CoNLL-2003
- Part-of-Speech: PTB

# Open-ended tasks

- Coreference resolution: WSC

## Example

Text: Mark told Pete many lies about himself, which Pete included in his book. He should have been more truthful.

Coreference: False

- Question Answering: Squad 2

## Example

Endangered Species Act Paragraph: "... Other legislation followed, including the Migratory Bird Conservation Act of 1929, a **1937 treaty** prohibiting the hunting of right and gray whales, and the **Bald Eagle Protection Act of 1940**. These later laws had a low cost to society—the species were relatively rare—and little **opposition** was raised."

Question 1: "Which laws faced significant **opposition**?"

Plausible Answer: later laws

Question 2: "What was the name of the **1937 treaty**?"

Plausible Answer: Bald Eagle Protection Act

# Close-ended multitask benchmark - superGLUE



Leaderboard Version: 2.0

Rank	Name	Model	URL	Score	BoolQ	CB	COPA	MultiRC	ReCoRD	RTE	WiC	WSC	AX-b	AX-g	
1	JDExplore d-team	Vega v2		91.3	90.5	98.6/99.2	99.4	88.2/62.4	94.4/93.9	96.0	77.4	98.6	-0.4	100.0/50.0	
+	2	Liam Fedus	ST-MoE-32B		91.2	92.4	96.9/98.0	99.2	89.6/65.8	95.1/94.4	93.5	77.7	96.6	72.3	96.1/94.1
	3	Microsoft Alexander v-team	Turing NLR v5		90.9	92.0	95.9/97.6	98.2	88.4/63.0	96.4/95.9	94.1	77.1	97.3	67.8	93.3/95.5
	4	ERNIE Team - Baidu	ERNIE 3.0		90.6	91.0	98.6/99.2	97.4	88.6/63.2	94.7/94.2	92.6	77.4	97.3	68.6	92.7/94.7
	5	Yi Tay	PaLM 540B		90.4	91.9	94.4/96.0	99.0	88.7/63.6	94.2/93.3	94.1	77.4	95.9	72.9	95.5/90.4
+	6	Zirui Wang	T5 + UDG, Single Model (Google Brain)		90.4	91.4	95.8/97.6	98.0	88.3/63.0	94.2/93.5	93.0	77.9	96.6	69.1	92.7/91.9
+	7	DeBERTa Team - Microsoft	DeBERTa / TuringNLRv4		90.3	90.4	95.7/97.6	98.4	88.2/63.7	94.5/94.1	93.2	77.5	95.9	66.7	93.3/93.8
	8	SuperGLUE Human Baselines	SuperGLUE Human Baselines		89.8	89.0	95.8/98.9	100.0	81.8/51.9	91.7/91.3	93.6	80.0	100.0	76.6	99.3/99.7
+	9	T5 Team - Google	T5		89.3	91.2	93.9/96.8	94.8	88.1/63.3	94.1/93.4	92.5	76.9	93.8	65.6	92.7/91.9

Attempt to measure “general language capabilities”

# Examples from superGLUE

Cover a number of different tasks

- BoolQ, MultiRC (reading texts)
- CB, RTE (Entailment)
- COPA (cause and effect)
- ReCoRD (QA+reasoning)
- WiC (meaning of words)
- WSC (coreference)

BoolQ	<p><b>Passage:</b> Barq's – Barq's is an American soft drink. Its brand of root beer is notable for having caffeine. Barq's, created by Edward Barq and bottled since the turn of the 20th century, is owned by the Barq family but bottled by the Coca-Cola Company. It was known as Barq's Famous Olde Tyme Root Beer until 2012.</p> <p><b>Question:</b> is barq's root beer a pepsi product   <b>Answer:</b> No</p>
CB	<p><b>Text:</b> B: And yet, uh, I we-, I hope to see employer based, you know, helping out. You know, child, uh, care centers at the place of employment and things like that, that will help out. A: Uh-huh. B: What do you think, do you think we are, setting a trend?</p> <p><b>Hypothesis:</b> they are setting a trend   <b>Entailment:</b> Unknown</p>
COPA	<p><b>Premise:</b> My body cast a shadow over the grass.   <b>Question:</b> What's the CAUSE for this?</p> <p><b>Alternative 1:</b> The sun was rising.   <b>Alternative 2:</b> The grass was cut.</p> <p><b>Correct Alternative:</b> 1</p>
MultiRC	<p><b>Paragraph:</b> Susan wanted to have a birthday party. She called all of her friends. She has five friends. Her mom said that Susan can invite them all to the party. Her first friend could not go to the party because she was sick. Her second friend was going out of town. Her third friend was not so sure if her parents would let her. The fourth friend said maybe. The fifth friend could go to the party for sure. Susan was a little sad. On the day of the party, all five friends showed up. Each friend had a present for Susan. Susan was happy and sent each friend a thank you card the next week</p> <p><b>Question:</b> Did Susan's sick friend recover?   <b>Candidate answers:</b> Yes, she recovered (T), No (F), Yes (T), No, she didn't recover (F), Yes, she was at Susan's party (T)</p>
ReCoRD	<p><b>Paragraph:</b> (CNN) Puerto Rico on Sunday overwhelmingly voted for statehood. But Congress, the only body that can approve new states, will ultimately decide whether the status of the US commonwealth changes. Ninety-seven percent of the votes in the nonbinding referendum favored statehood, an increase over the results of a 2012 referendum, official results from the State Electoral Commission show. It was the fifth such vote on statehood. "Today, we the people of Puerto Rico are sending a strong and clear message to the US Congress ... and to the world ... claiming our equal rights as American citizens, Puerto Rico Gov. Ricardo Rossello said in a news release. @highlight Puerto Rico voted Sunday in favor of US statehood</p> <p><b>Query</b> For one, they can truthfully say, "Don't blame me, I didn't vote for them," when discussing the &lt;placeholder&gt; presidency   <b>Correct Entities:</b> US</p>
RTE	<p><b>Text:</b> Dana Reeve, the widow of the actor Christopher Reeve, has died of lung cancer at age 44, according to the Christopher Reeve Foundation.</p> <p><b>Hypothesis:</b> Christopher Reeve had an accident.   <b>Entailment:</b> False</p>
WiC	<p><b>Context 1:</b> Room and <u>board</u>.   <b>Context 2:</b> He nailed <u>boards</u> across the windows.</p> <p><b>Sense match:</b> False</p>

# Close-ended: challenges

- Choosing your metrics: accuracy / precision / recall / f1-score / ROC
  - [https://github.com/cgpotts/cs224u/blob/main/evaluation\\_metrics.ipynb](https://github.com/cgpotts/cs224u/blob/main/evaluation_metrics.ipynb)
  - [https://scikit-learn.org/stable/modules/model\\_evaluation.html](https://scikit-learn.org/stable/modules/model_evaluation.html)
- Aggregating across metrics or tasks
- Where do the labels come from?
- Are there spurious correlations?

SuperGLUE Tasks		
Matthew's Corr	F1a / EM	
Avg. F1 / Accuracy	Accuracy	F1 / Accuracy
Accuracy	Accuracy	Gender Parity / Accuracy

# Open-ended evaluation

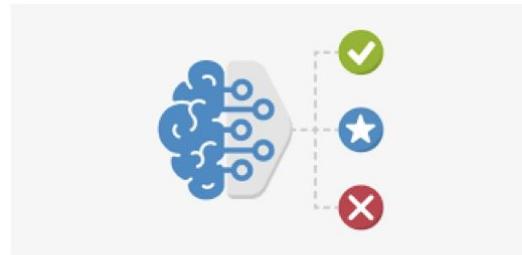
# Open-ended evaluation

- Long generations with too many possible correct answers to enumerate
  - => can't use standard ML metrics
- There are now better and worse answers (not just right and wrong)
- Example:
  - Summarization: CNN-DM / Gigaword
  - Translation: WMT
  - Instruction-following: Chatbot Arena / AlpacaEval / MT-Bench

# Types of evaluation methods for text generation

Ref: They walked **to the grocery store** .  
Gen: The woman went to the **hardware store** .

Content Overlap Metrics



Model-based Metrics



Human Evaluations

# Content overlap metrics

Ref: They walked to the grocery store .  
Gen: The woman went to the hardware store .



- Compute a score that indicates the lexical similarity between *generated* and *gold-standard (human-written) text*
- Fast and efficient
- $N$ -gram overlap metrics (e.g., **BLEU**, **ROUGE**, METEOR, CIDEr, etc.)  
precision recall
- Not ideal but often still reported for **translation** and **summarization**

# A simple failure case

*n*-gram overlap metrics have no concept of semantic relatedness!



Are you enjoying the  
CS224N lectures?

Score:

0.67

0.25

False negative 0

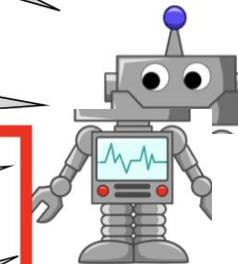
False positive 0.67

Heck yes !



Yes !

You know it !

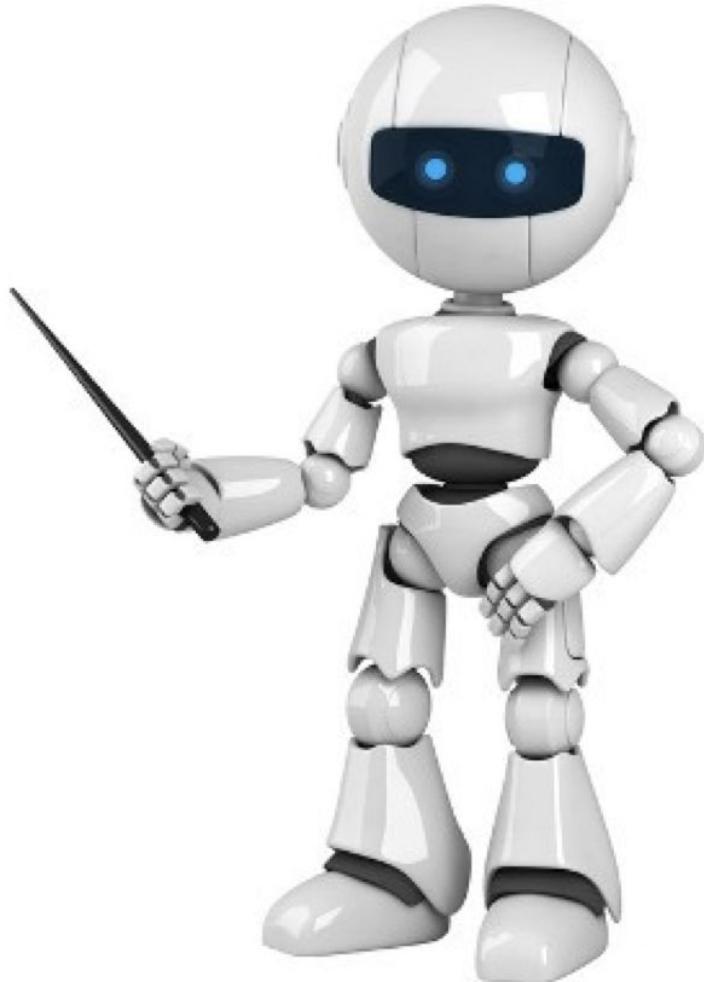


Yup .

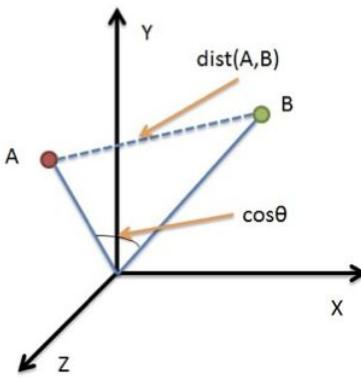
Heck no !

# Model-based metrics to capture more semantics

- Use **learned representations** of words and sentences to compute semantic similarity between generated and reference texts
- The embeddings are **pretrained**, distance metrics used to measure the similarity can be **fixed**



# Model-based metrics: Word distance functions

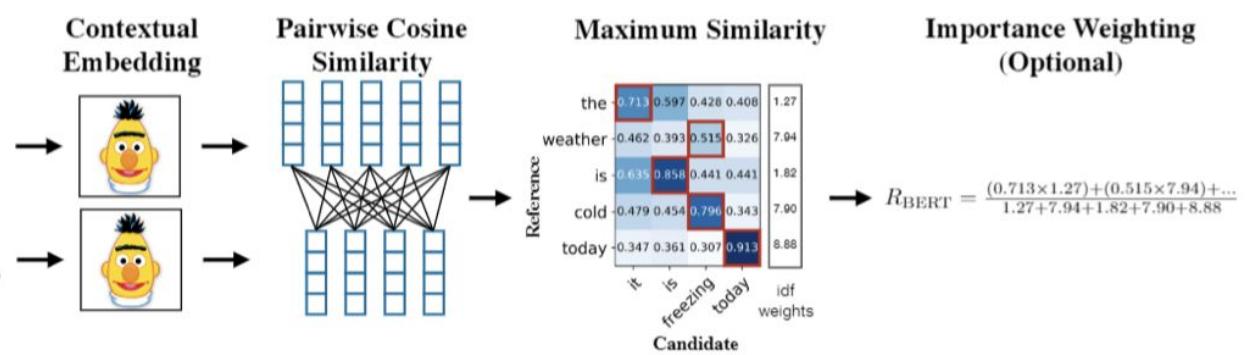


## BERTSCORE

Uses pre-trained contextual embeddings from BERT and matches words in candidate and reference sentences by cosine similarity.  
(Zhang et.al. 2020)

Reference  $\mathcal{X}$   
*the weather is cold today*

Candidate  $\hat{\mathcal{X}}$   
*it is freezing today*

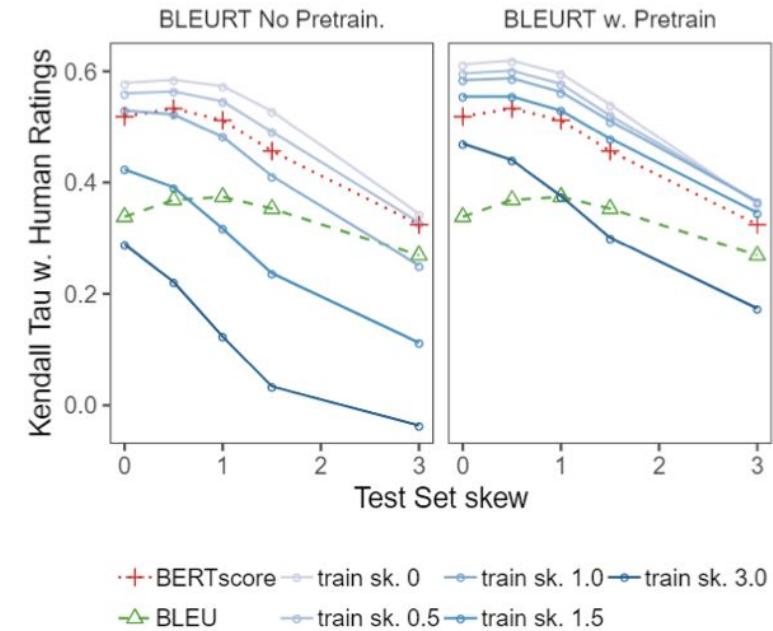


# Model-based metrics: Beyond word matching

## BLEURT:

A regression model based on BERT returns a score that indicates to what extent the candidate text is grammatical and conveys the meaning of the reference text.

(Sellam et.al. 2020)



# Reference free evals

- **Reference-based evaluation:**
  - Compare human written reference to model outputs
  - Used to be ‘standard’ evaluation for most NLP tasks
  - Examples: BLEU, ROUGE, BertScore etc.
- **Reference free evaluation**
  - Have a model give a score
  - No human reference
  - Was nonstandard – now becoming popular with GPT4
  - Examples: AlpacaEval, MT-Bench

# Human evaluations



- Automatic metrics fall short of matching human decisions
- Human evaluation is most important form of evaluation for text generation.
- Gold standard in developing new automatic metrics
  - New automated metrics must correlate well with human evaluations!

# Human evaluations

- Ask *humans* to evaluate the quality of generated text
- Overall or along some specific dimension:
  - fluency
  - coherence / consistency
  - factuality and correctness
  - commonsense
  - style / formality
  - grammaticality
  - redundancy

**Note:** Don't compare human evaluation scores across differently conducted studies

Even if they claim to evaluate the same dimensions!

# Human evaluation: Issues

- Human judgments are regarded as the **gold standard**
- But it also has issues:
  - Slow
  - Expensive
  - Inter-annotator disagreement (esp. if subjective)
  - Intra-annotator disagreement across time

# Human evaluation: Issues

- Challenges with human evaluation
  - How to describe the task?
  - How to show the task to the humans?
  - What metric do you use?
  - Selecting the annotators
  - Monitoring the annotators: time, accuracy, ...

# Reference-free eval: chatbots



VS

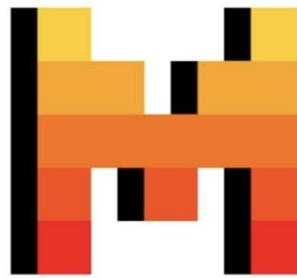


Table 1: Distribution of use case categories from our API prompt dataset.

Use-case	(%)
Generation	45.6%
Open QA	12.4%
Brainstorming	11.2%
Chat	8.4%
Rewrite	6.6%
Summarization	4.2%
Classification	3.5%
Other	3.5%
Closed QA	2.6%
Extract	1.9%

- How do we evaluate something like ChatGPT?
- *So many* different use cases it's hard to evaluate
- The responses are also long-form text, which is even harder to evaluate.

# Side-by-side ratings

The screenshot shows the homepage of the Chatbot Arena: Benchmarking LLMs in the Wild. The page has a dark background with white text. At the top, there's a navigation bar with links to Blog, GitHub, Paper, Dataset, Twitter, and Discord. Below that is a section titled "Rules" with a small icon of a scroll. A bulleted list of rules follows:

- Ask any question to two anonymous models (e.g., ChatGPT, Claude, Llama) and vote for the better one!
- You can continue chatting until you identify a winner.
- Vote won't be counted if model identity is revealed during conversation.

Below the rules is a section titled "Arena Elo Leaderboard" with a trophy icon. It says: "We collect 200K+ human votes to compute an Elo-based LLM leaderboard. Find out who is the 🏆 LLM Champion!"

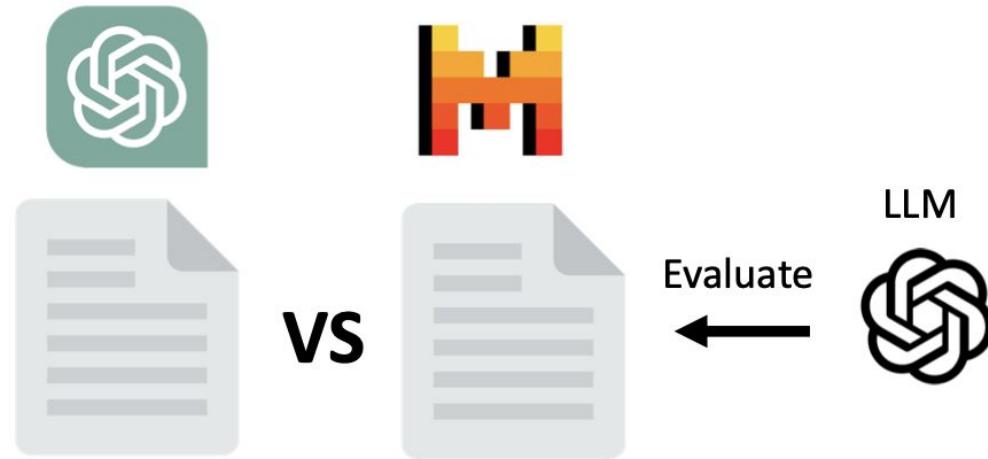
At the bottom, there's a "Chat now!" button with a thumbs-up icon. Below it, a callout box says "Expand to see the descriptions of 35 models". There are two tabs: "Model A" and "Model B".

Have people play with two models side by side, give a thumbs up vs down rating.

# What's missing with side-by-side human eval?

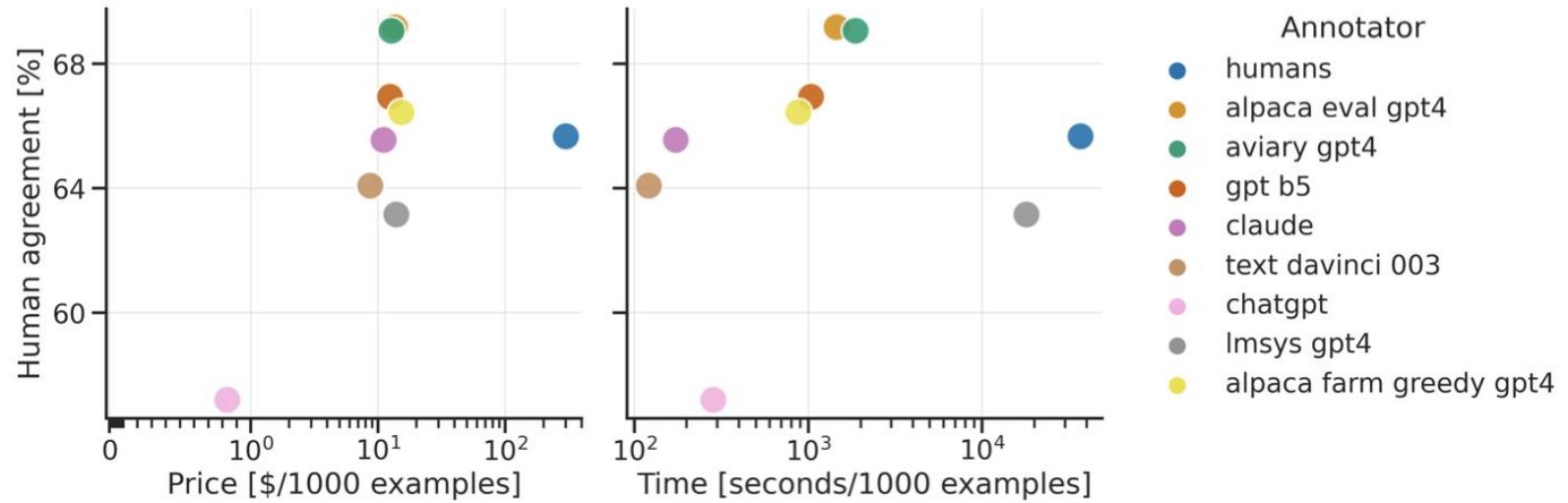
- Current gold standard for evaluation of chat LLM
- **Cost**
  - Human annotation takes large, community effort
  - New models take a long time to benchmark
  - Only notable models get benchmarked

# Lowering the costs – use a LM evaluator



- Use a LM as a reference free evaluator
- Surprisingly high correlations with human
- Common versions: AlpacaEval, MT-bench

# AlpacaFarm : Human agreement



**100x Cheaper, 100x faster, and higher agreement than humans**

# What are common LM datasets?

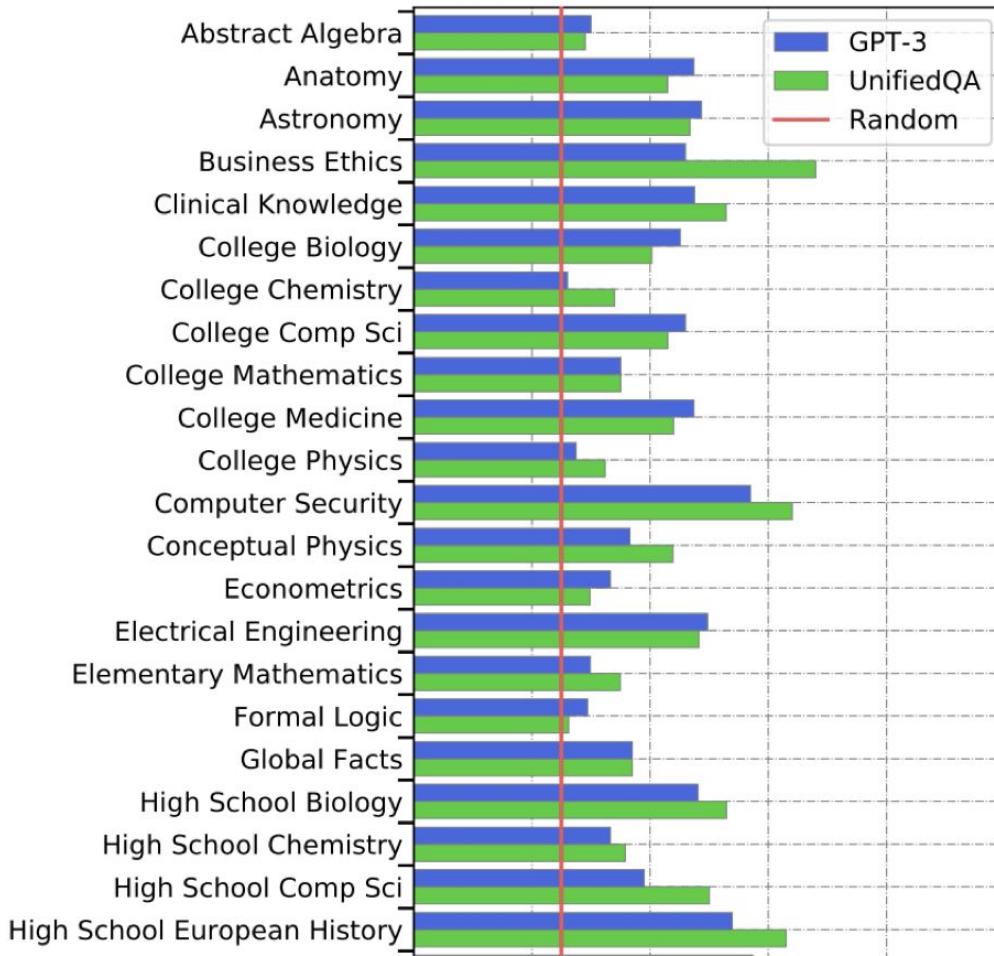
- What do these benchmarks evaluate on?

- A huge mix of things!

Scenario	Task	What	Who
<a href="#">NarrativeQA</a> narrative_qa	short-answer question answering	passages are books and movie scripts, questions are unknown	annotators from summaries
<a href="#">NaturalQuestions (closed-book)</a> natural_qa_closedbook	short-answer question answering	passages from Wikipedia, questions from search queries	web users
<a href="#">NaturalQuestions (open-book)</a> natural_qa_openbook_longans	short-answer question answering	passages from Wikipedia, questions from search queries	web users
<a href="#">OpenbookQA</a> openbookqa	multiple-choice question answering	elementary science	Amazon Mechanical Turk workers
<a href="#">MMLU (Massive Multitask Language Understanding)</a> mmlu	multiple-choice question answering	math, science, history, etc.	various online sources
<a href="#">GSM8K (Grade School Math)</a> gsm	numeric answer question answering	grade school math word problems	contractors on Upwork and Surge AI
<a href="#">MATH</a> math_chain_of_thought	numeric answer question answering	math competitions (AMC, AIME, etc.)	problem setters
<a href="#">LegalBench</a> legalbench	multiple-choice question answering	public legal and administrative documents, manually constructed questions	lawyers
<a href="#">MedQA</a> med_qa	multiple-choice question answering	US medical licensing exams	problem setters

## Massive Multitask Language Understanding (MMLU) [Hendrycks et al., 2021]

New benchmarks for measuring LM performance on 57 diverse *knowledge intensive* tasks



# Some intuition: examples from MMLU

## Astronomy

What is true for a type-Ia supernova?

- A. This type occurs in binary systems.
- B. This type occurs in young galaxies.
- C. This type produces gamma-ray bursts.
- D. This type produces high amounts of X-rays.

Answer: A

## High School Biology

In a population of giraffes, an environmental change occurs that favors individuals that are tallest. As a result, more of the taller individuals are able to obtain nutrients and survive to pass along their genetic information. This is an example of

- A. directional selection.
- B. stabilizing selection.
- C. sexual selection.
- D. disruptive selection

Answer: A

# Other capabilites: code

Nice feature of code: evaluate vs test cases

Metric: Pass@1 (Pass @ k means one of k outputs pass)

GPT4: ~67%

```
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)
```

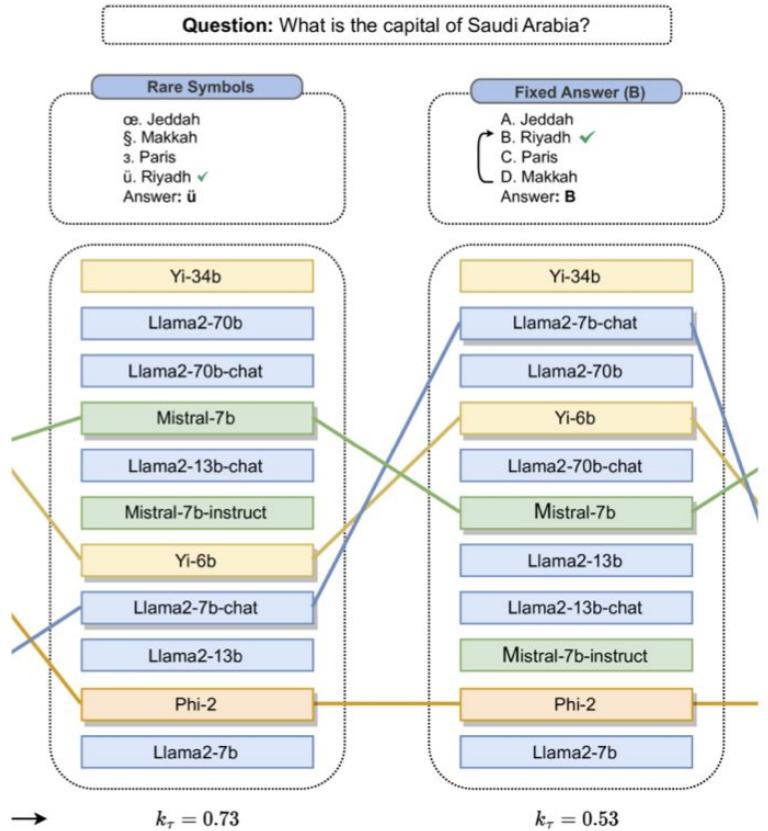
```
def encode_cyclic(s: str):
    """
    returns encoded string by cycling groups of three characters.
    """
    # split string to groups. Each of length 3.
    groups = [s[(3 * i):min((3 * i + 3), len(s))] for i in range((len(s) + 2) // 3)]
    # cycle elements in each group. Unless group has fewer elements than 3.
    groups = [(group[1:] + group[0]) if len(group) == 3 else group for group in groups]
    return "".join(groups)

def decode_cyclic(s: str):
    """
    takes as input string encoded with encode_cyclic function. Returns decoded string.
    """
    # split string to groups. Each of length 3.
    groups = [s[(3 * i):min((3 * i + 3), len(s))] for i in range((len(s) + 2) // 3)]
    # cycle elements in each group.
    groups = [(group[-1] + group[:-1]) if len(group) == 3 else group for group in groups]
    return "".join(groups)
```

HumanEval ('Human written' eval for code generation)

# Issues and challenges with evaluation

# Consistency issues



[Alzahrani et al 2024]

# Contamination and overfitting issues



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

g's Race	implementation, math		greedy, implementation	
nd Chocolate	implementation, math		Cat?	implementation, strings
triangle!	brute force, geometry, math		Actions	data structures, greedy, implementation, math
	greedy, implementation, math		Interview Problem	brute force, implementation, strings

...



Susan Zhang   
@suchenzang

I think Phi-1.5 trained on the benchmarks. Particularly, GSM8K.



Susan Zhang @suchenzang · Sep 12  
Let's take [github.com/openai/grade-s...](https://github.com/openai/grade-s...)

If you truncate and feed this question into Phi-1.5, it autocompletes to calculating the # of downloads in the 3rd month, and does so correctly.

Change the number a bit, and it answers correctly as well.

1/



...

**Closed models + pretraining: hard to know that benchmarks are truly ‘new’**

# Monoculture of NLP benchmarking

Area	# papers	English	Accuracy / F1	Multilinguality	Fairness and bias	Efficiency	Interpretability	>1 dimension
ACL 2021 oral papers	461	69.4%	38.8%	13.9%	6.3%	17.8%	11.7%	6.1%
MT and Multilinguality	58	0.0%	15.5%	56.9%	5.2%	19.0%	6.9%	13.8%
Interpretability and Analysis	18	88.9%	27.8%	5.6%	0.0%	5.6%	66.7%	5.6%
Ethics in NLP	6	83.3%	0.0%	0.0%	100.0%	0.0%	0.0%	0.0%
Dialog and Interactive Systems	42	90.5%	21.4%	0.0%	9.5%	23.8%	2.4%	2.4%
Machine Learning for NLP	42	66.7%	40.5%	19.0%	4.8%	50.0%	4.8%	9.5%
Information Extraction	36	80.6%	91.7%	8.3%	0.0%	25.0%	5.6%	8.3%
Resources and Evaluation	35	77.1%	42.9%	5.7%	8.6%	5.7%	14.3%	5.7%
NLP Applications	30	73.3%	43.3%	0.0%	10.0%	20.0%	10.0%	0.0%

Most papers only evaluate on English and performance (accuracy)

# Multilingual benchmarking

- Benchmarks exist, we should use them!
- MEGA: Multilingual Evaluation of Generative AI
  - 16 datasets, 70 languages
- GlobalBench:
  - 966 datasets in 190 languages.
- XTREME: A Massively Multilingual Multi-task Benchmark for Evaluating Cross-lingual Generalization
  - 9 tasks, 40 languages
- Multilingual Large Language Models Evaluation Benchmark
  - MMLU / ARC / HellaSwag translated in 26 languages

# Evaluation: Takeaways

- **Closed ended tasks**
  - Think about what you evaluate (diversity, difficulty)
- **Open ended tasks**
  - Content overlap metrics (useful for low-diversity settings)
  - Chatbot evals – very difficult! Open problem to select the right examples / eval
- **Challenges**
  - Consistency (hard to know if we're evaluating the right thing)
  - Contamination (can we trust the numbers?)
  - Biases
- In many cases, the best judge of output quality is **YOU!**
  - **Look at your model generations. Don't just rely on numbers!**

# IMPERIAL

## Q and A