CS6216 Advanced Topics in Machine Learning (Systems)

LLM Alignment

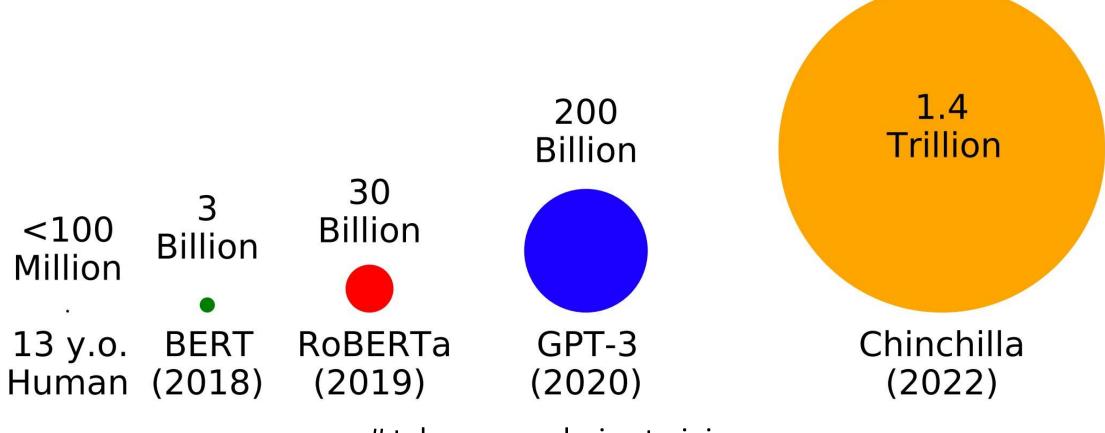
Yao LU 22 Oct 2025

National University of Singapore School of Computing

Outline

- Aligning LLMs: from models to assistants
 - Instruction tuning
 - Reinforcement learning with human feedback (RLHF)
 - Chain-of-thought

LLMs trained on more and more data



tokens seen during training

https://babylm.github.io/

What kinds of things does pretraining learn?

- Stanford University is located in _______, California. [Trivia]
- I put ____ fork down on the table. [syntax]
- The woman walked across the street, checking for traffic over ____ shoulder. [coreference]
- I went to the ocean to see the fish, turtles, seals, and _____. [lexical semantics/topic]
- Overall, the value I got from the two hours watching it was the sum total of the popcorn and the drink. The movie was ____. [sentiment]
- Iroh went into the kitchen to make some tea. Standing next to Iroh, Zuko pondered his destiny. Zuko left the _____. [some reasoning – this is harder]
- I was thinking about the sequence that goes 1, 1, 2, 3, 5, 8, 13, 21, _____ [some basic arithmetic; they don't learn the Fibonnaci sequence]

Language models may do rudimentary modeling of agents, beliefs, and actions:

Pat watches a demonstration of a bowling ball and a leaf being dropped at the same time in a vacuum chamber. Pat, who is a physicist, predicts that the bowling ball and the leaf will fall at the same rate.

Changing the last sentence of the prompt, we get:

... Pat, who has never seen this demonstration before, predicts that the bowling ball will fall to the ground first. This is incorrect. In a vacuum chamber, there is no air

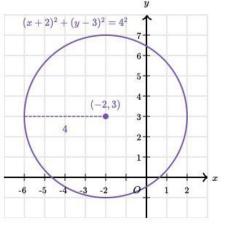
Language Models as Agent Models [Andreas, 2022]

math:

We can describe circles in the xy-plane using equations in terms of x and y. Circle equations questions require us to understand the connection between these equations and the features of circles.

For example, the equation $(x+2)^2+(y-3)^2=4^2$ is graphed in the xy-plane

below. It is a circle with a center at (-2,3) and a radius of 4.



code:

```
5 // Determine whether the sentiment of text is positive
6 // Use a web service
7 async function isPositive(text: string): Promise<boolean> {
     const response = await fetch(`http://text-processing.com/api/sentiment/`, {
       method: "POST",
       body: `text=${text}`,
10
       headers: {
11
         "Content-Type": "application/x-www-form-urlencoded",
12
13
       },
14
     });
     const json = await response.json();
15
16
     return json.label === "pos";
```

medicine:

Rapid and chronic ethanol tolerance are composed of distinct memory-like states in Drosophila

Abstract

Ethanol tolerance is the first type of behavioral plasticity and neural plasticity that is induced by ethanol intake, and yet its molecular and circuit bases remain largely unexplored. Here, we characterize three distinct forms of ethanol tolerance in male Drosophila: rapid, chronic, and repeated. Rapid tolerance is composed of two short-lived memory-like states, one that is labile and one that is consolidated. Chronic tolerance, induced by continuous exposure, lasts for two days, induces ethanol preference, and hinders the development of rapid tolerance through the activity of

I need to throw a dinner party for 6 people who are vegetarian. Can you suggest a 3-course menu with a chocolate dessert?

[Microsoft Bing]

(Also see OpenAl's ChatGPT, Google's Bard, Anthropic's Claude)





Sure, I can suggest a 3-course menu with a chocolate dessert for your vegetarian dinner party. Here are some possible options:

- Starter: Wild Mushroom Tartlets with Onion Sauce or Vegan Popcorn Tofu Nuggets 2
- Main: **Vegan Butternut Squash Mac** ² or Vegetarian Three-Bean Chili 3
- Dessert: Chocolate Lava Cake 4 or **Chocolate Pasta with Chocolate Hazelnut Cream Sauce, White Chocolate** Shavings, and Fresh Berries 4

I hope this helps you plan your party. Enjoy!

Learn more: 1. booths.co.uk

+10 more

Outline

- Aligning LLMs: from models to assistants
 - Instruction tuning
 - Reinforcement learning with human feedback (RLHF)
 - Chain-of-thought

Language modeling ≠ assisting users

Explain the moon landing to a 6 year old in a few sentences.

GPT-3
Explain the theory of gravity to a 6 year old.

Explain the theory of relativity to a 6 year old in a few sentences.

Explain the big bang theory to a 6 year old.

Explain evolution to a 6 year old.

Language models are *not aligned* with user intent [Ouyang et al., 2022].

Language modeling ≠ assisting users

PROMPT Explain the moon landing to a 6 year old in a few sentences.

COMPLETION Human

A giant rocket ship blasted off from Earth carrying astronauts to the moon. The astronauts landed their spaceship on the moon and walked around exploring the lunar surface. Then they returned safely back to Earth, bringing home moon rocks to show everyone.

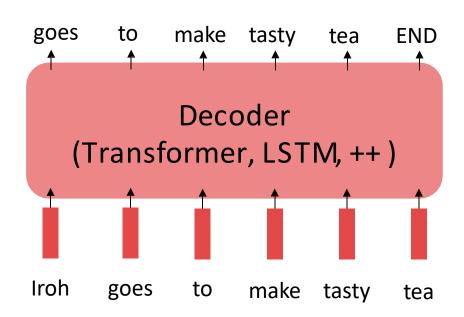
Language models are not *aligned* with user intent [Ouyang et al., 2022]. Finetuning to the rescue!

The pretraining / finetuning paradigm

Pretraining can improve NLP applications by serving as parameter initialization.

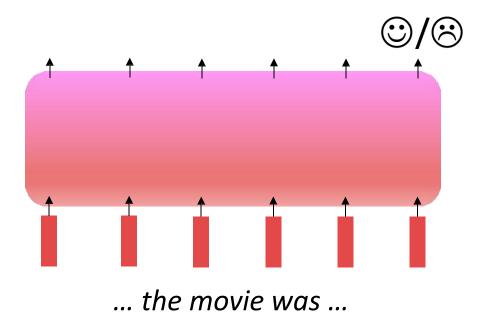
Step 1: Pretrain (on language modeling)

Lots of text; learn general things!



Step 2: Finetune (on your task)

Not many labels; adapt to the task!

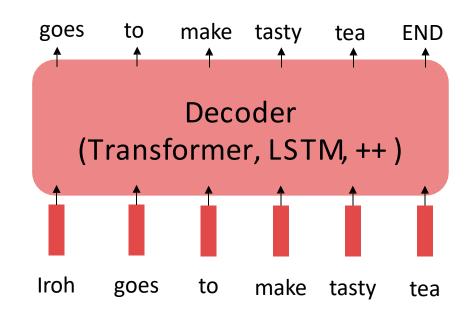


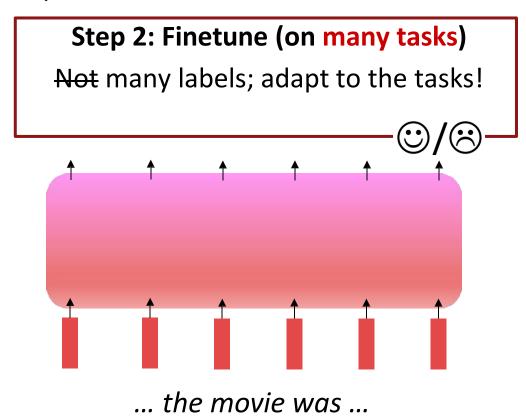
Scaling up finetuning

Pretraining can improve NLP applications by serving as parameter initialization.

Step 1: Pretrain (on language modeling)

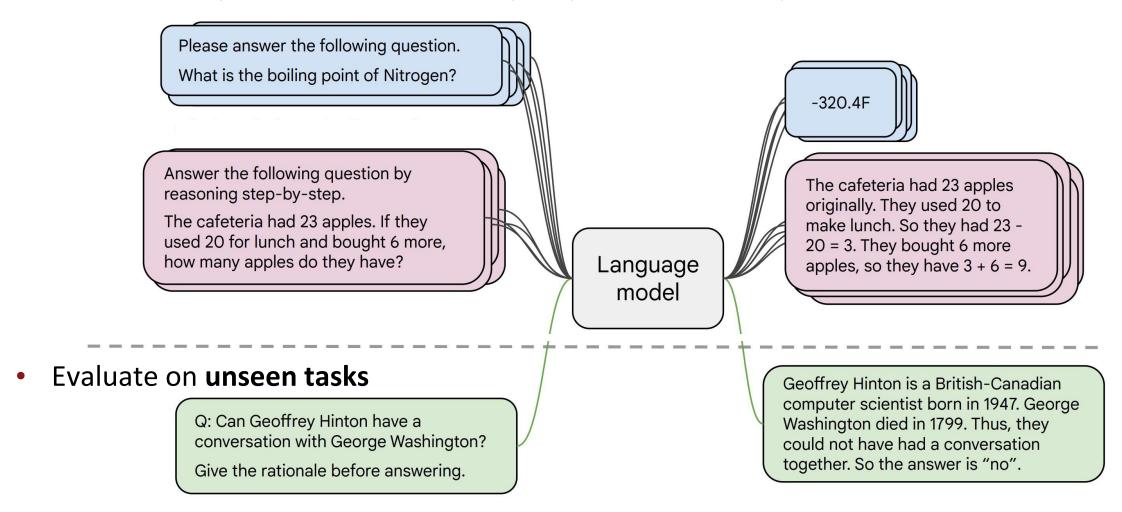
Lots of text; learn general things!





Instruction finetuning

Collect examples of (instruction, output) pairs across many tasks and finetune an LM



Instruction fine-tuning pretraining

- As is usually the case, data + model scale is key for this to work!
- E.g., the Super- NaturalInstructions dataset contains over 1.6K tasks,
 3M+ examples
 - Classification, sequence tagging, rewriting, translation, QA...
- How do we evaluate such a model?

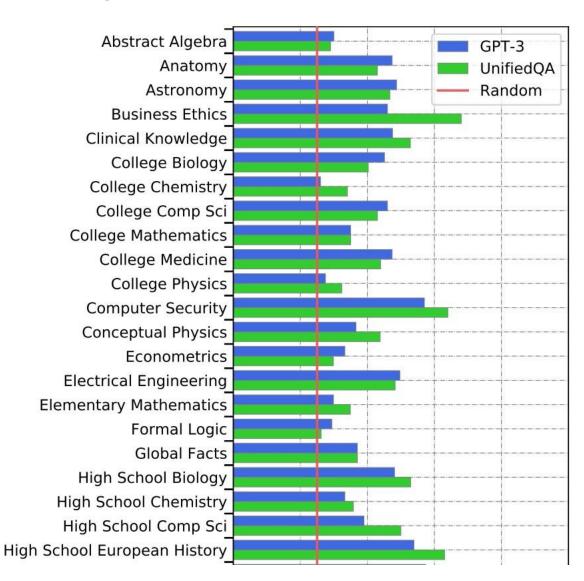


New benchmarks for multitask LMs

Massive Multitask Language Understanding (MMLU)

[Hendrycks et al., 2021]

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



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

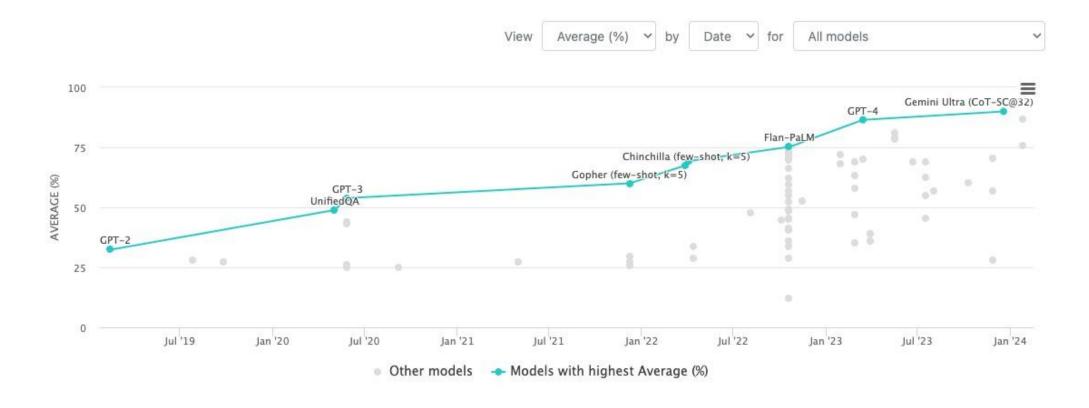
High School Biology

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

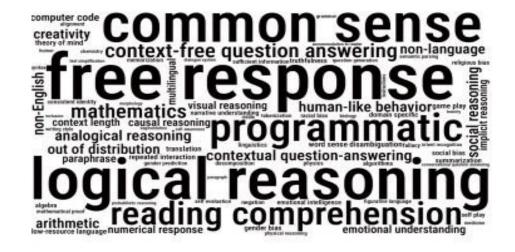
Progress on MMLU



Rapid, impressive progress on challenging knowledge-intensive benchmarks

New benchmarks for multitask LMs

BIG-Bench [Srivastava et al., 2022] 200+ tasks, spanning:



https://github.com/google/BIG-bench/blob/main/bigbench/benchmark tasks/README.md

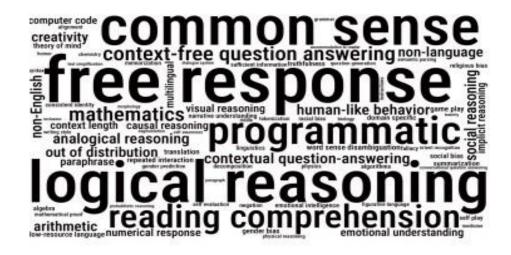
BEYOND THE IMITATION GAME: QUANTIFY-ING AND EXTRAPOLATING THE CAPABILITIES OF LANGUAGE MODELS

Alphabetic author list:*

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 Tazarv, Alice Xiang, Alicia Parrish, Allen Nie, Aman Hussain, Amanda Askell, Amanda Dsouza, Ambrose Slone, Ameet Rahane, Anantharaman S. Iyer, Andres 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, Avkut Erdem, Avla Karakas, B. Rvan Roberts, Bao Sheng Loe, Barret Zoph, Bartlomiei Bojanowski, Batuhan Özvurt, Behnam Hedayatnia, Behnam Neyshabur, Benjamin Inden, Benno Stein, Berk Ekmekci, Bill Yuchen Lin, Blake Howald, 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 Sileo, Dan Garrette, Dan Hendrycks, Dan Kilman, 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, Dilyar Buzan, Dimitri Coelho Mollo, Diyi Yang, Dong-Ho Lee, Ekaterina Shutova, Ekin Dogus Cubuk, Elad Segal, Eleanor Hagerman, Elizabeth Barnes, Elizabeth Donoway, Ellie Paylick, Emanuele Rodola, Emma Lam, Eric Chu, Eric Tang, Erkut Erdem, Ernie Chang, Ethan A. Chi, Ethan Dyer, Ethan Jerzak, Ethan Kim, Eunice Engefu Manyasi, Evgenii Zheltonozhskii, Fanyue Xia, Fatemeh Siar, Fernando Martínez-Plumed, Francesca Happé, Francois Chollet, Frieda Rong, Gaurav Mishra, Genta Indra Winata, Gerard de Melo, Germán Kruszewski, Giambattista Parascandolo, Giorgio Mariani, Gloria Wang, Gonzalo Jaimovitch-López, Gregor Betz, Guy Gur-Ari, Hana Galijasevic, Hannah Kim, Hannah Rashkin, Hannahe Hajishirzi, Harsh Mehta, Hayden Bogar, Henry Shevlin, 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, 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 Berant, Jörg Frohberg, Jos Rozen, Jose Hernandez-Orallo, Joseph Boudeman, 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 Chiafullo, Ksenia Shkaruta, Kumar Shridhar, Kyle McDonell, 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, Luheng 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 Tolkiehn, 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, Michael Swedrowski, Michele Bevilacqua, Michihiro Yasunaga, Mihir Kale, Mike Cain, Mimee Xu, Mirac Suzgun, Mo Tiwari, Mohit Bansal, Moin Aminnaseri Mor Geva, Mozhdeh Gheini, Mukund Varma T, Nanyun Peng, Nathan Chi, Nayeon Lee, Neta Gur-Ari Krakover, Nicholas Cameron, Nicholas Roberts, Nick Doiron, Nikita Nangia, Niklas Deckers, Niklas Muennighoff, Nitish 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 Patii, Pouya Pezeshkpour, Priti Oli, Qiaozhu Mei, Qing Lyu, Qinlang Chen, Rabin Banjade, Rachel Etta Rudolph, Raefer Gabriel, Rahel Habacker, Ramón Risco Delgado, 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, Sarik 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, Shyamolima (Shammie) Debnath, Siamak Shakeri, Simon Thormeyer, Simone Melzi, Siva Reddy, Sneha Priscilla Makini, Soo-Hwan Lee, Spencer Torene, Sriharsha Hatwar, Stanislas Dehaene, Stefan Divic, Stefano Ermon, Stella Biderman, Stephanie Lin, Stephan Prasad, Steven T. Piantadosi, Stuart M. Shieber, Summer Misherghi, 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 Nkinyili, Timo Schick, Timofei Kornev, Timothy Telleen-Lawton, Titus Tunduny, Tobias Gerstenberg, Trenton Chang, Trishala Neeraj, Tushar Khot, Tyler Shultz, Uri Shaham, Vedant Misra, Vera Demberg, Victoria Nyamai, Vikas Raunak, Vinay Ramasesh, Vinay Uday Prabhu, Vishakh Padmakumar, Vivek Srikumar, 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, Zijian Wang, Zijie J. Wang, Zirui Wang,

New benchmarks for multitask LMs

BIG-Bench [Srivastava et al., 2022] 200+ tasks, spanning:



https://github.com/google/BIGbench/blob/main/bigbench/benchmark tasks/README.md

Kanji ASCII Art to Meaning

This subtask converts various kanji into ASCII art and has the language model guess their meaning from the ASCII art.

Instruction finetuning

Model input (Disambiguation QA)

Q: In the following sentences, explain the antecedent of the pronoun (which thing the pronoun refers to), or state that it is ambiguous.

Sentence: The reporter and the chef will discuss their favorite dishes.

Options:

- (A) They will discuss the reporter's favorite dishes
- (B) They will discuss the chef's favorite dishes
- (C) Ambiguous

A: Let's think step by step.

Before instruction finetuning

The reporter and the chef will discuss their favorite dishes.

The reporter and the chef will discuss the reporter's favorite dishes.

The reporter and the chef will discuss the chef's favorite dishes.

The reporter and the chef will discuss the reporter's and the chef's favorite dishes.



Highly recommend trying FLAN-T5 out to get a sense of its capabilities:

Instruction finetuning

Model input (Disambiguation QA)

Q: In the following sentences, explain the antecedent of the pronoun (which thing the pronoun refers to), or state that it is ambiguous.

Sentence: The reporter and the chef will discuss their favorite dishes.

Options:

- (A) They will discuss the reporter's favorite dishes
- (B) They will discuss the chef's favorite dishes
- (C) Ambiguous

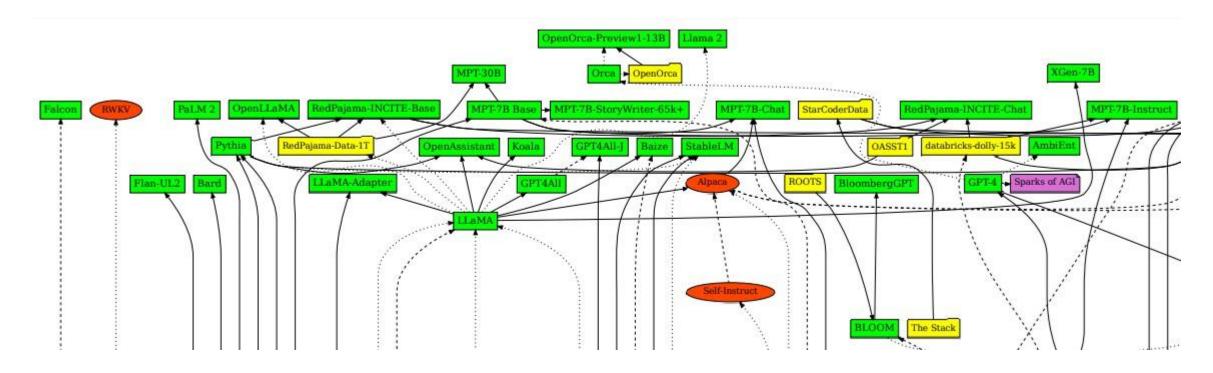
A: Let's think step by step.

After instruction finetuning

The reporter and the chef will discuss their favorite dishes does not indicate whose favorite dishes they will discuss. So, the answer is (C).

Highly recommend trying FLAN-T5 out to get a sense of its capabilities:

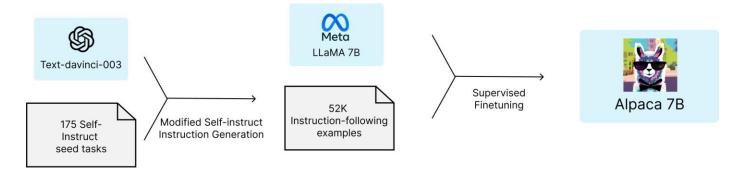
A huge diversity of instruction-tuning datasets



• The release of LLaMA led to open-source attempts to 'create' instruction tuning data

What have we learned from this?

 You can generate data synthetically (from bigger LMs)



 You don't need many samples to instruction tune

LIMA: Less Is More for Alignment

Chunting Zhou $^{\mu*}$ Pengfei Liu $^{\pi*}$ Puxin Xu $^{\mu}$ Srini Iyer $^{\mu}$ Jiao Sun $^{\lambda}$

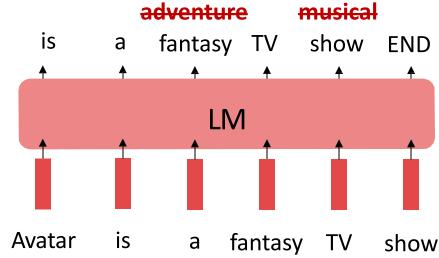
Crowdsourcing can be pretty effective!



Limitations of instruction finetuning?

- One limitation of instruction finetuning is obvious: it's expensive to collect groundtruth data for tasks.
- But there are other, subtler limitations too.
- Problem 1: tasks like open-ended creative generation have no right answer.
 - Write me a story about a dog and her pet grasshopper.
- Problem 2: language modeling penalizes all token-level mistakes equally, but some errors are worse than others.
- Even with instruction finetuning, there

 a mismatch between the LM
 objective and the objective of "satisfy human preferences"!
- Can we explicitly attempt to satisfy human preferences?



Limitations of instruction finetuning

- + Simple and straightforward, generalize to unseen tasks
- Collecting demonstrations for so many tasks is expensive
- Mismatch between LM objective and human preferences

Outline

- Aligning LLMs: from models to assistants
 - Instruction tuning
 - Reinforcement learning with human feedback (RLHF)
 - Chain-of-thought

Optimizing for human preferences

- Let's say we were training a language model on some task (e.g. summarization).
- For each LM sample s, imagine we had a way to obtain a *human reward* of that summary: $R(s) \in \mathbb{R}$, higher is better.

SAN FRANCISCO,
California (CNN) -A magnitude 4.2
earthquake shook the
San Francisco
...
overturn unstable
objects.

An earthquake hit San Francisco. There was minor property damage, but no injuries.

$$R(s_1) = 8.0$$

The Bay Area has good weather but is prone to earthquakes and wildfires.

$$R(s_2) = 1.2$$

Now we want to maximize the expected reward of samples from our LM:

$$\mathbb{E}_{\widehat{S} \sim p_{\theta}(S)}[R(\widehat{S})]$$

Note: for mathematical simplicity we're assuming only one "prompt"

High-level instantiation: 'RLHF' pipeline

Step 1

Collect demonstration data, and train a supervised policy.

A prompt is sampled from our prompt dataset.

A labeler demonstrates the desired output behavior.

This data is used to fine-tune GPT-3 with supervised learning.



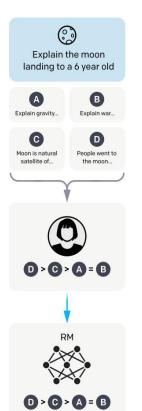
Step 2

Collect comparison data, and train a reward model.

A prompt and several model outputs are sampled.

A labeler ranks the outputs from best to worst.

This data is used to train our reward model.



Step 3

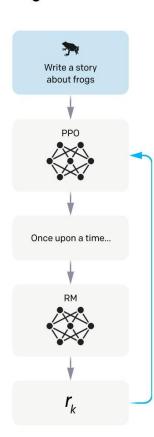
Optimize a policy against the reward model using reinforcement learning.

A new prompt is sampled from the dataset.

The policy generates an output.

The reward model calculates a reward for the output.

The reward is used to update the policy using PPO.



- First step: instruction tuning!
- Second + third steps: maximize reward (how??)

Reinforcement learning to the rescue

- The field of reinforcement learning (RL) has studied these (and related) problems for many years now
 [Williams, 1992; Sutton and Barto, 1998]
- Circa 2013: resurgence of interest in RL applied to deep learning, game-playing [Mnih et al., 2013]
- But the interest in applying RL to modern LMs is an even newer phenomenon [<u>Ziegler et al., 2019</u>; <u>Stiennon et al., 2020</u>; <u>Ouyang et al., 2022</u>]. Why?
 - RL w/ LMs has commonly been viewed as very hard to get right (still is!)
 - Newer advances in RL algorithms that work for large neural models, including language models (e.g. PPO; [Schulman et al., 2017]) Proximal Policy Optimization Algorithms





Optimizing for human preferences

• How do we actually change our LM parameters θ to maximize this?

$$\mathbb{E}_{\hat{s} \sim p_{\theta}(s)}[R(\hat{s})]$$

Let's try doing gradient ascent!

$$\theta_{t+1} \coloneqq \theta_t + \alpha^{\nabla} \theta_t \mathbb{E}_{\hat{s} \sim p_{\theta}(s)} [R(\hat{s})]$$
How do we estimate function is nonthis expectation?? What if our reward differentiable??

- Policy gradient methods in RL (e.g., REINFORCE; [Williams, 1992]) give us tools for estimating and optimizing this objective.
- We'll describe a very high-level mathematical overview of the simplest policy gradient estimator, but a full treatment of RL is outside the scope of this course.

A brief introduction to policy gradient/REINFORCE [Williams, 1992]

We want to obtain

(defn. of expectation) (linearity of gradient)

$$\nabla_{\theta} \mathbb{E}_{\hat{s} \sim p_{\theta}(s)}[R(\hat{s})] = \nabla_{\theta} \sum_{s} R(s) p_{\theta}(s) = \sum_{s} R(s) \nabla_{\theta} p_{\theta}(s)$$

Here we'll use a very handy trick known as the log-derivative trick. Let's try taking the gradient of $\log p_{\theta}(s)$

$$\nabla_{\theta} \log p_{\theta}(s) = \frac{1}{p_{\theta}(s)} \nabla_{\theta} p_{\theta}(s) \implies \nabla_{\theta} p_{\theta}(s) = p_{\theta}(s) \nabla_{\theta} \log p_{\theta}(s)$$
(chain rule) This is an

Plug back in:

expectation of this

$$\sum_{s} R(s) \nabla_{\theta} p_{\theta}(s) = \sum_{s} p_{\theta}(s) R(s) \nabla_{\theta} \log p_{\theta}(s)$$

$$= \mathbb{E}_{\hat{s} \sim p_{\theta}(s)} [R(\hat{s}) \nabla_{\theta} \log p_{\theta}(\hat{s})]$$

A brief introduction to policy gradient/REINFORCE [Williams, 1992]

Now we have put the gradient "inside" the expectation, we can approximate this objective with Monte Carlo samples:

$$\nabla_{\theta} \mathbb{E}_{\hat{s} \sim p_{\theta}(s)}[R(\hat{s})] = \mathbb{E}_{\hat{s} \sim p_{\theta}(s)}[R(\hat{s}) \nabla_{\theta} \log p_{\theta}(\hat{s})] \approx \frac{1}{m} \sum_{i=1}^{m} R(s_i) \nabla_{\theta} \log p_{\theta}(s_i)$$

This is why it's called "reinforcement **learning**": we **reinforce** good actions, increasing the chance they happen again.

Take gradient steps to maximize
$$p_{\theta}(s_i)$$
 creasing the chance they happen again.

Giving us the update rule: $\theta_{t+1} \coloneqq \theta_t + \alpha \frac{1}{m} \sum_{i=1}^m R(s_i) \nabla_{\theta_t} \log p_{\theta_t}(s_i)$ is is **heavily simplified!** There is a *lot*

If *R* is ----

This is **heavily simplified!** There is a *lot* more needed to do RL w/ LMs. Can you see any problems with this objective?

Take steps to minimize $p_{\theta}(s_i)$

Take gradient steps

How do we model human preferences?

- Awesome: now for any **arbitrary, non-differentiable reward function** R(s), we can train our language model to maximize expected reward.
- Not so fast! (Why not?)
- Problem 1: human-in-the-loop is expensive!
 - Solution: instead of directly asking humans for preferences, model their preferences as a separate (NLP) problem! [Knox and Stone, 2009]

An earthquake hit San Francisco. There was minor property damage, but no injuries.

$$S_1$$

$$R(s_1) = 8.0$$

The Bay Area has good weather but is prone to earthquakes and wildfires.

$$R(s_2) = 1.2$$

Train an LM $RM_{.}$ (s) to predict human preferences from an annotated dataset, then optimize for $RM_{.}$ instead.

How do we model human preferences?

- Problem 2: human judgments are noisy and miscalibrated!
- Solution: instead of asking for direct ratings, ask for pairwise comparisons, which can be more reliable [Phelps et al., 2015; Clark et al., 2018]

A 4.2 magnitude earthquake hit San Francisco, resulting in massive damage.

$$S_3$$
 $R(s_3) = 4.1? 6.6? 3.2?$

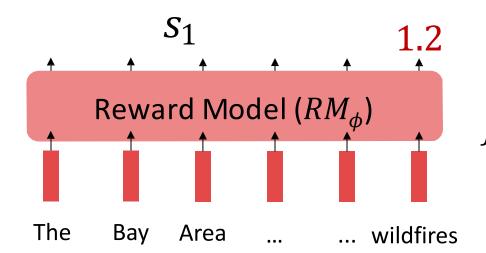
How do we model human preferences?

- Problem 2: human judgments are noisy and miscalibrated!
- **Solution:** instead of asking for direct ratings, ask for **pairwise comparisons**, which can be more reliable [Phelps et al., 2015; Clark et al., 2018]

An earthquake hit San Francisco. There was minor property damage, but no injuries. A 4.2 magnitude earthquake hit

San Francisco, resulting in massive damage.

The Bay Area has good weather but is prone to earthquakes and wildfires.



 S_3 S_2

Bradley-Terry [1952] paired comparison model

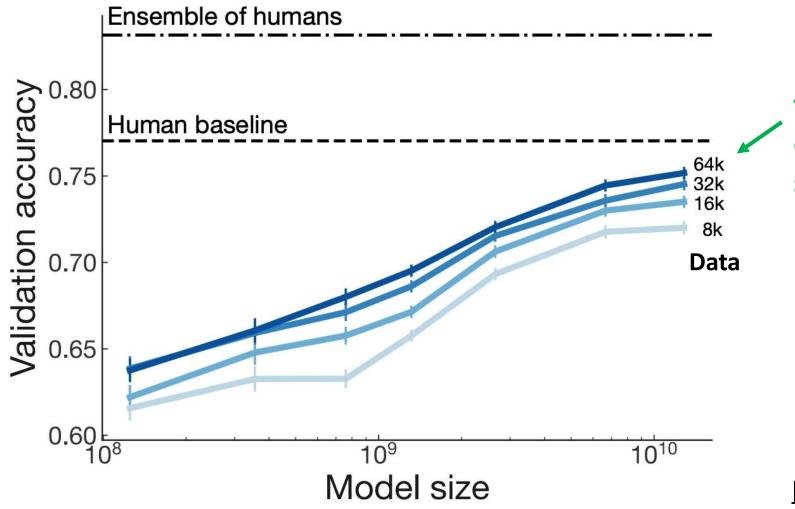
$$J_{RM}(\phi) = -\mathbb{E}_{(s^w, s^l) \sim D} \left[\log \sigma(RM_{\phi}(s^w) - RM_{\phi}(s^l)) \right]$$
"winning" "losing" sw should score

"winning" "losing" sample sample

sw should score higher than

Make sure your reward model works first!

Evaluate RM on predicting outcome of held-out human judgments



Large enough RM trained on enough data approaching single human perf

Stiennon et al., 2020

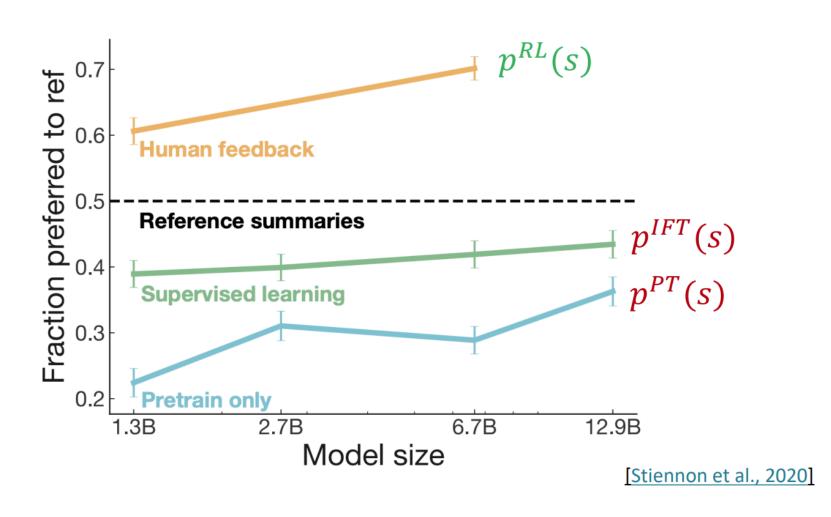
RLHF: Putting it all together [Christiano et al., 2017; Stiennon et al., 2020]

- Finally, we have everything we need:
 - A pretrained (possibly instruction-finetuned) LM $p^{PT}(s)$
 - A reward model $RM_\phi(s)$ that produces scalar rewards for LM outputs, trained on a dataset of human comparisons
 - A method for optimizing LM parameters towards an arbitrary reward function.
- Now to do RLHF:
 - Initialize a copy of the model $p_{ heta}^{RL}(s)$, with parameters heta we would like to optimize
 - Optimize the following reward with RL:

$$R(s) = RM_{\phi}(s) - \beta \log \left(\frac{p_{\theta}^{RL}(s)}{p^{PT}(s)}\right)$$
 Pay a price when $p_{\theta}^{RL}(s) > p^{PT}(s)$

This is a penalty which prevents us from diverging too far from the pretrained model. In expectation, it is known as the **Kullback-Leibler (KL)** divergence between $p_{\theta}^{RL}(s)$ and $p^{PT}(s)$

RLHF provides gains over pretraining + finetuning



InstructGPT: scaling up RLHF to tens of thousands of tasks

Step 1

Collect demonstration data, and train a supervised policy.

30k tasks!

A prompt is sampled from our prompt dataset.

A labeler demonstrates the desired output behavior.

This data is used to fine-tune GPT-3

自自自

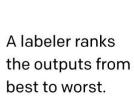
with supervised

learning.

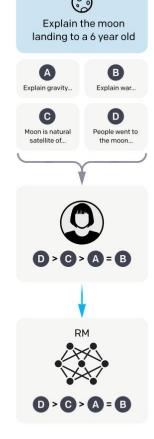
Step 2

Collect comparison data, and train a reward model.

A prompt and several model outputs are sampled.



This data is used to train our reward model.



Step 3

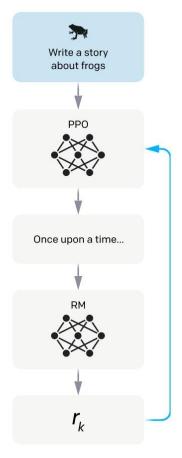
Optimize a policy against the reward model using reinforcement learning.

A new prompt is sampled from the dataset.

The policy generates an output.

The reward model calculates a reward for the output.

The reward is used to update the policy using PPO.



Ouyang et al., 2022

InstructGPT: scaling up RLHF to tens of thousands of tasks

Tasks collected from labelers:

- Plain: We simply ask the labelers to come up with an arbitrary task, while ensuring the tasks had sufficient diversity.
- **Few-shot:** We ask the labelers to come up with an instruction, and multiple query/response pairs for that instruction.
- User-based: We had a number of use-cases stated in waitlist applications to the OpenAI API. We asked labelers to come up with prompts corresponding to these use cases.

Use-case	Prompt
Brainstorming	List five ideas for how to regain enthusiasm for my career
Generation	Write a short story where a bear goes to the beach, makes friends with a seal, and then returns home.

InstructGPT

Explain the moon landing to a 6 year old in a few sentences.

COMPLETION

GPT-3

Explain the theory of gravity to a 6 year old.

Explain the theory of relativity to a 6 year old in a few sentences.

Explain the big bang theory to a 6 year old.

Explain evolution to a 6 year old.

InstructGPT

People went to the moon, and they took pictures of what they saw, and sent them back to the earth so we could all see them.

InstructGPT

PROMPT Write a short poem about a wise frog.

COMPLETION

GPT-3

Write a short story in which a character has two different names.

Write a short story in which you try to get something back that you have lost.

Write a short story in which a character has a bad dream.

InstructGPT

The frog is a symbol of wisdom
He knows all the secrets of the world
He is a master of disguise
And a great teacher of life
He is a symbol of transformation
And the bringer of change
He is the frog who has seen it all
And knows the meaning of it all

ChatGPT: Instruction Finetuning + RLHF for dialog agents

ChatGPT: Optimizing Language Models for Dialogue

Note: OpenAI (and similar companies) are keeping more details secret about ChatGPT training (including data, training parameters, model size)—perhaps to keep a competitive edge...

Methods

We trained this model using Reinforcement Learning from Human Feedback (RLHF), using the same methods as InstructGPT, but with slight differences in the data collection setup. We trained an initial model using supervised fine-tuning: human AI trainers provided conversations in which they played both sides—the user and an AI assistant. We gave the trainers access to model-written suggestions to help them compose their responses. We mixed this new dialogue dataset with the InstructGPT dataset, which we transformed into a dialogue format.

(Instruction finetuning!)

ChatGPT: Instruction Finetuning + RLHF for dialog agents

ChatGPT: Optimizing Language Models for Dialogue

Note: OpenAI (and similar companies) are keeping more details secret about ChatGPT training (including data, training parameters, model size)—perhaps to keep a competitive edge...

Methods

To create a reward model for reinforcement learning, we needed to collect comparison data, which consisted of two or more model responses ranked by quality. To collect this data, we took conversations that AI trainers had with the chatbot. We randomly selected a model-written message, sampled several alternative completions, and had AI trainers rank them. Using these reward models, we can fine-tune the model using <u>Proximal Policy Optimization</u>. We performed several iterations of this process.

(RLHF!)

Controlled comparisons of "RLHF" style algorithms

Method	Simulated win-rate (%)	Human win-rate (%)	
GPT-4	79.0 ± 1.4	69.8 ± 1.6	
ChatGPT	61.4 ± 1.7	52.9 ± 1.7	
PPO	46.8 ± 1.8	55.1 ± 1.7	
Best-of-n	45.0 ± 1.7	50.7 ± 1.8	
Expert Iteration	41.9 ± 1.7	45.7 ± 1.7	
SFT 52k (Alpaca 7B)	39.2 ± 1.7	40.7 ± 1.7	
SFT 10k	36.7 ± 1.7	44.3 ± 1.7	
Binary FeedME	36.6 ± 1.7	37.9 ± 1.7	
Quark	35.6 ± 1.7		
Binary Reward Conditioning	32.4 ± 1.6	u m	
Davinci001	24.4 ± 1.5	32.5 ± 1.6	
LLaMA 7B	11.3 ± 1.1	6.5 ± 0.9	

- Many works study RLHF behaviors using GPT-4 feedback (Simulated) as a surrogate for Human feedback.
- PPO (method in InstructGPT) does work
- Simple baselines (Best-of-n, Training on 'good' outputs) works well too

RLHF behaviors – clear stylistic changes

Instruction: What are the five most common causes of stress among people?

Alpaca (Finetuned w/ 10k data): The five most common causes of stress among people are work, money, relationships, health, and family

Alpaca + PPO: The five most common causes of stress among people are:

- Pressure at work: Whether it is due to tight deadlines, high expectations, or an overwhelming workload, pressure at work can be a major source of stress
- 2. Money problems: Financial issues can be a significant source of stress, from not being able to make ends meet to struggling with debt.
- 3. Unhealthy lifestyle: Poor diet, lack of exercise, and sleep deprivation can all contribute to stress.
- 4. Unhappiness at home: Problems in relationships, conflicts with family members, and other domestic issues can be a major source of stress.
- 5. Uncertainty: Excessive worry and anxiety about the future, health issues, or other unknowns can be a major cause of stress

Significantly more detailed, nicer/clearer list like formatting

Limitations of RL + Reward Modeling

- Human preferences are unreliable!
 - "Reward hacking" is a common problem in RL



https://openai.com/blog/faulty-reward-functions/

Limitations of RL + Reward Modeling

- Human preferences are unreliable!
 - "Reward hacking" is a common problem in RL
 - Chatbots are rewarded to produce responses that seem authoritative and helpful, regardless of truth
 - This can result in making up facts
 - + hallucinations

TECHNOLOGY

Google shares drop \$100 billion after its new AI chatbot makes a mistake

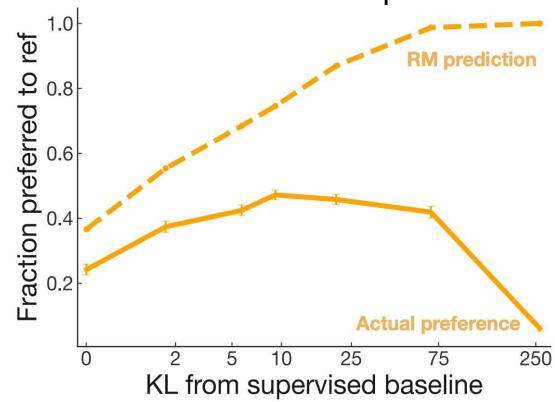
February 9, 2023 · 10:15 AM ET

https://www.npr.org/2023/02/09/1155650909/google-chatbot--error-bard-shares

Limitations of RL + Reward Modeling

- Human preferences are unreliable!
 - "Reward hacking" is a common problem in RL
 - Chatbots are rewarded to produce responses that seem authoritative and helpful, regardless of truth
 - This can result in making up facts
 + hallucinations
- Models of human preferences are even more unreliable!

Reward model over-optimization

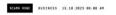


$$R(s) = RM_{\phi}(s) - \beta \log \left(\frac{p_{\theta}^{RL}(s)}{p^{PT}(s)}\right)$$

Where did the labels come from?

Exclusive: OpenAI Used Kenyan Workers on Less Than \$2 Per Hour to Make ChatGPT Less Toxic





Millions of Workers Are Training Al Models for Pennies

From the Philippines to Colombia, low-paid workers label training data for Al models used by the likes of Amazon, Facebook, Google, and Microsoft.



Behind the AI boom, an army of overseas workers in 'digital sweatshops'

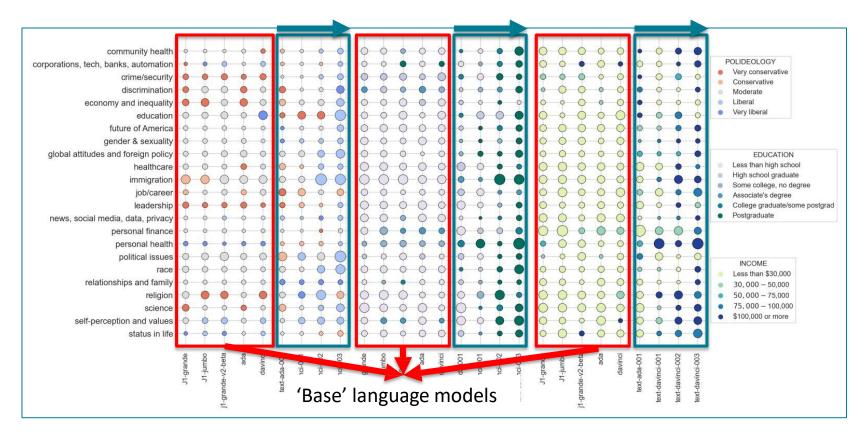
By Rebecca Tan and Regine Cabato August 28, 2023 at 2:00 a.m. EDT



RLHF labels are often obtained from overseas, low-wage workers

Where does the label come from?

What gender do you identify as?		
Male		
Female	44.4%	
Nonbinary / other		
What ethnicities do you identify as?		
White / Caucasian	31.6%	
Southeast Asian	52.6%	
Indigenous / Native American / Alaskan Native	0.0%	
East Asian		
Middle Eastern		
Latinx		
Black / of African descent	10.5%	
What is your nationality?		
Filipino	22%	
Bangladeshi		
American	17%	
Albanian		
Brazilian	5%	
Canadian	5%	
Colombian		
Indian		
Uruguayan		
Zimbabwean	5%	
What is your age?		
18-24	26.3%	
25-34	47.4%	
35-44		
45-54	10.5%	
55-64	5.3%	
65+	0%	
What is your highest attained level of educa	tion?	
Less than high school degree		
High school degree		
Undergraduate degree		
Master's degree		
Doctorate degree	0%	



[Santurkar+ 2023, OpinionQA]

We also need to be quite careful about how annotator biases might creep into LMs

Limitations of RLHF

- + Directly model preferences (cf. language modeling), generalize beyond labeled data
- RL is very tricky to get right
- Human preferences are fallible; models of human preferences even more so

What's next?

- RLHF is still a very underexplored and fastmoving area!
- RLHF gets you further than instruction finetuning, but is (still!) data expensive.
- Recent work aims to alleviate such data requirements:
 - RL from AI feedback [Bai et al., 2022]
 - Finetuning LMs on their own outputs
 [Huang et al., 2022; Zelikman et al.,
 2022]
- However, there are still many limitations of large LMs (size, hallucination) that may not be solvable with RLHF!

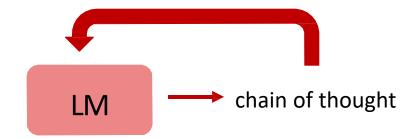
LARGE LANGUAGE MODELS CAN SELF-IMPROVE

Jiaxin Huang^{1*} Shixiang Shane Gu² Le Hou^{2†} Yuexin Wu² Xuezhi Wang² Hongkun Yu² Jiawei Han¹

¹University of Illinois at Urbana-Champaign ²Google

1{jiaxinh3, hanj}@illinois.edu 2{shanegu, lehou, crickwu, xuezhiw, hongkuny}@google.com

[Huang et al., 2022]



Self-Taught Reasoner (STaR)

[Zelikman et al., 2022]

Outline

- Aligning LLMs: from models to assistants
 - Instruction tuning
 - Reinforcement learning with human feedback (RLHF)
 - Chain-of-thought

Hard language tasks: reasoning

Q: If there are 3 cars in the parking lot and 2 more cars arrive, how many cars are in the parking lot?

A: The answer is 5

Q: Take the last letters of the words in "Elon Musk" and concatenate them

A: The answer is **nk**.

Q: What home entertainment equipment requires cable?
Answer Choices: (a) radio shack (b) substation (c) television (d) cabinet

A: The answer is (c).

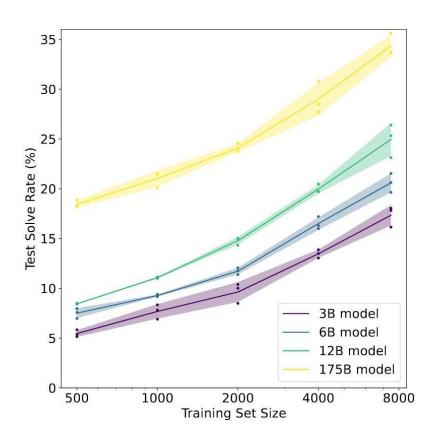
Arithmetic Reasoning (AR) (+ - ×÷...)

Symbolic Reasoning (SR)

Commonsense Reasoning (CR)

Reasoning Problems

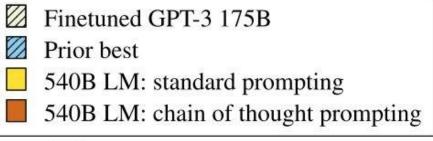
Fine-tune GPT-3 on GSM8K (arithmetic): (Cobbe et al. 2021)

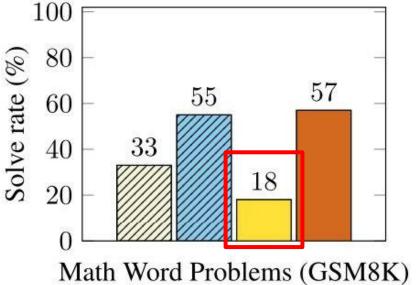


Conjecture: to achieve > 80%, needs 100 times more fine-tuning data for 175B model

Reasoning Problems

GSM8K (arithmetic):





Few-shot standard prompting with even larger model (PaLM 540B) also does not work well.

Reasoning Problems

Scaling up language model size does not **efficiently** achieve high performances, for Arithmetic Reasoning (*AR*), CommonSense Reasoning (*CR*) and Symbolic Reasoning (*SR*) tasks.

Proposed solution: chain of thought prompting

Chain-of-Thought (CoT) NeurIPS'22

A chain of thought is a series of intermediate natural language reasoning steps that lead to the final output.

Use <input, intermediate results, output> triples, rather than simple <input, output> pairs

Benefits:

- Decomposition -> easier intermediate problems
- Interpretable
- More general than neural symbolic computing
- Leveraging prompting of LLM

Few-Shot CoT

Chain of Thought Prompting Elicits Reasoning in Large Language Models Jason Wei Xuezhi Wang Dale Schuurmans Maarten Bosma Brian Ichter Fei Xia Ed H. Chi Quoc V. Le Denny Zhou Google Research, Brain Team

COT

Takeshi Kojima
The University of Tokyo
t.kojima@weblab.t.u-tokyo.ac.jp

Machel Reid

The University of Tokyo

Shixiang Shane Gu Google Research, Brain Team

Zero-Shot CoT

Yutaka Matsuo The University of Tokyo

Large Language Models are Zero-Shot Reasoners

Yusuke Iwasawa The University of Tokyo

Example

Break a large task into sub-tasks and chain them together

Standard Prompting

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Chain-of-Thought Prompting

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. 5 + 6 = 11. The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

A: The answer is 27.



Model Output

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had 23 - 20 = 3. They bought 6 more apples, so they have 3 + 6 = 9. The answer is 9. <

Chain-of-Thought (CoT)

(a) Few-shot

Examples

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A:

(Output) The answer is 8. X

(c) Zero-shot

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A: The answer (arabic numerals) is

(Output) 8 X

(b) Few-shot-CoT (Wei et al., 2022)

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. 5 + 6 = 11. The answer is 11.

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A:

(Output) The juggler can juggle 16 balls. Half of the balls are golf balls. So there are 16/2 = 8 golf balls. Half of the golf balls are blue. So there are 8/2 = 4 blue golf balls. The answer is 4.

CoT Examples

Step-by-step
Answer

(d) Zero-shot-CoT (KoJima et al., 2022)

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

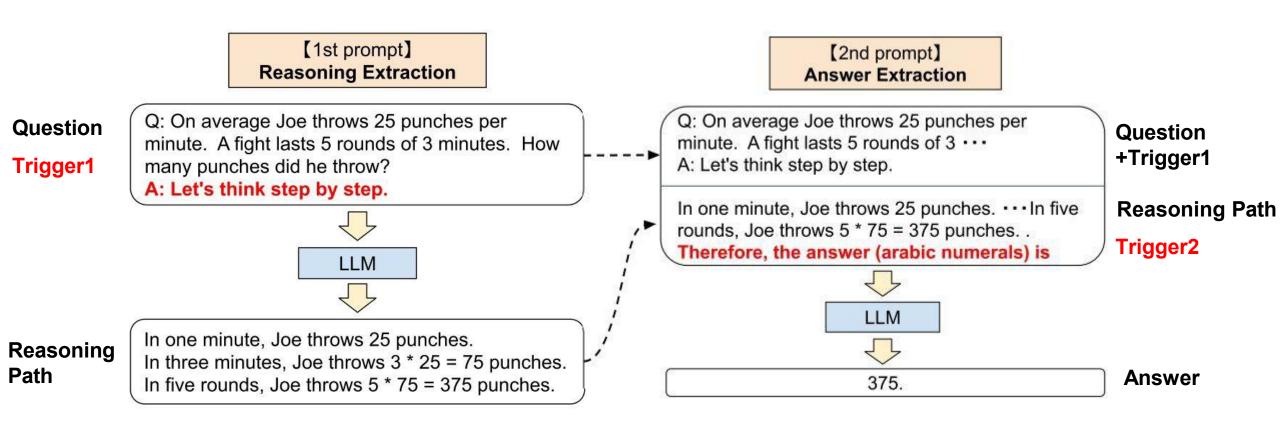
A: Let's think step by step.

(Output) There are 16 balls in total. Half of the balls are golf balls. That means that there are 8 golf balls. Half of the golf balls are blue. That means that there are 4 blue golf balls. ✓

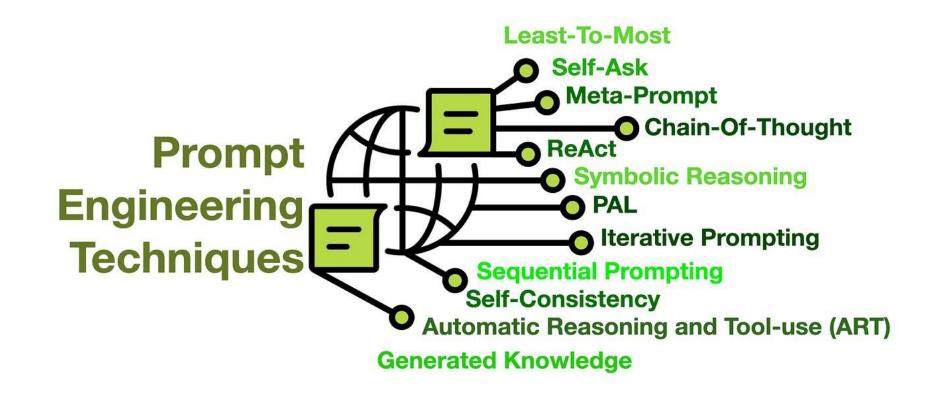
Two-stage Prompting Stepby-step Answer

Zero-Shot Chain of Thought (CoT)

A **two-stage prompting** is applied:



A lot more uncovered...



References and credits

CS224N/Ling284, Stanford University 15-442/15-642, Carnegie Mellon University COS 597G, Princeton University

Outline

- Aligning LLMs: from models to assistants
 - Instruction tuning
 - Reinforcement learning with human feedback (RLHF)
 - Chain-of-thought
 - LLM safety (next lecture)