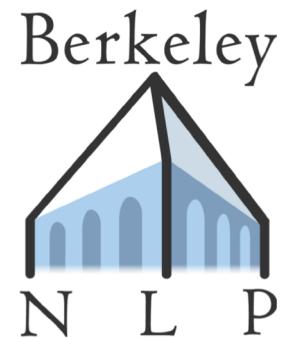


RLHF and Instruction-tuning



Eric Wallace
CS 288

With lots of credits to Jesse Mu and Stanford CS224N



Few-shot Learning Thus Far

- Thus far, we have talked about using LMs “out-of-the-box” for few-shot
 - surprising emergent property

Questions:

- Can we directly train models to do few-shot learning?
- Can we directly train models to follow arbitrary user instructions?
- Can we directly train models to obey toxicity & safety constraints?



Lecture Overview

- ▶ Instruction Finetuning
- ▶ Reinforcement Learning from Human Feedback (RLHF)
- ▶ Open challenges with RLHF



Lecture Overview

- **Instruction Finetuning**
- Reinforcement Learning from Human Feedback (RLHF)
- Open challenges with RLHF

Language modeling ≠ assisting users

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

Language modeling ≠ assisting users

PROMPT *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.

Language modeling ≠ assisting users

PROMPT *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.

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 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](#)].
Finetuning to the rescue!

Language modeling ≠ assisting users

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

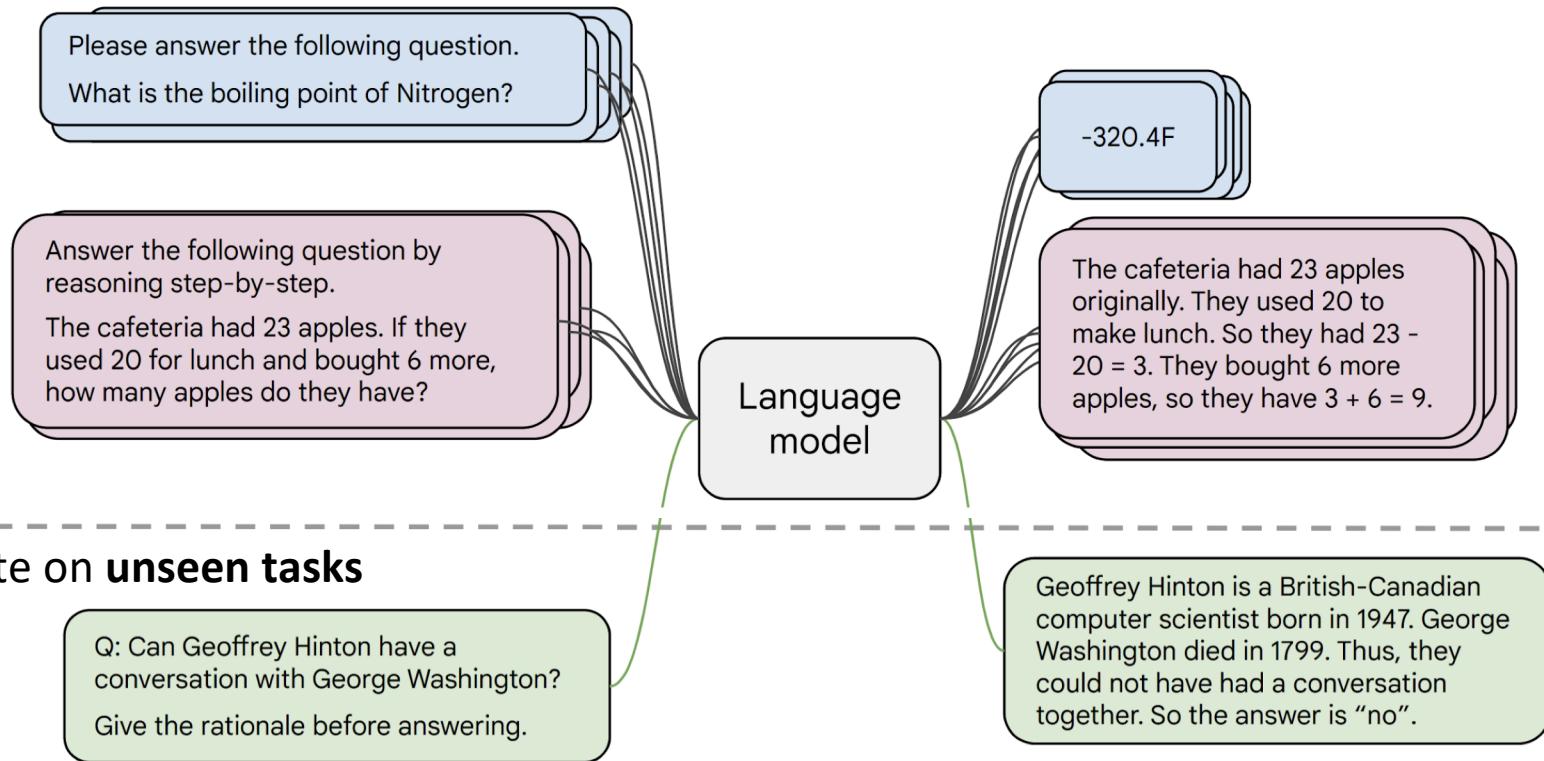
COMPLETION **Human**

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!

Instruction Finetuning

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



Scaling Up Instruction Finetuning



Aside: new benchmarks for multitask LMs

BIG-Bench [Srivastava et al., 2022]

200+ tasks, spanning:



BEYOND THE IMITATION GAME: QUANTIFYING AND EXTRAPOLATING THE CAPABILITIES OF LANGUAGE MODELS

Alphabetic author list:

Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Al Md Shoeib, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska, Aitor Lewkowycz, Akashat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W. Kocurek, Al Safaya, Ali Tazari, Alice Xiang, Alicia Parish, Allen Nee, Aman Hussain, Amanda Askell, Amanda Dsouza, Ambrose Scone, Ameteet Rahane, Anantharaman S. Iyer, Anders Arendse, Andrea Madotto, Andrea Santilli, Andreas Stuhlmüller, Andrei Dan, Andrew La, Andrew Lampinen, Andy Zou, Angela Jiang, Angelica Chen, Anil Vuong, Animesh Gupta, Anna Gottardi, Antonio Norelli, Anu Venkatesh, Arsh Ghoshalimadavodi, Arfa Tabassum, Arul Menzies, Arun Kirubarajan, Arush Mulkandoroff, Ashish Sabharwal, Austin Herrick, Aysha Efrat, Ayukr Erdem, Ayla Karakas, B. Ryan Roberts, Babu Sloe, Leng Barre Zoph, Bartłomiej Bojanowski, Batuhan Özyurt, Behnam Heydarnia, Behnam Neyshabur, Benjamin Inden, Bento Stein, Berk Ekmekci, Bill Yuchen Lin, Blake Howald, Cameron Diaz, Catherine Stinson, Cedrick Argueta, Ceska Ferri Ramrez, Chandan Charles, Charles Rathkopf, Chennin 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 Ascrash, Cristina Garbacea, Daniela Silo, Dar Garrett, Dan Hendryks, Dan Kilman, Dan Roth, Daniel Freeman, David Khashabi, Daniel Levy, Daniel Mosegu González, Danielle Persyuk, Danny Hernandez, Danqi Chen, Daphne Ippolito, Dar Gilboa, David Dohan, David Drakard, David Jennings, Debjyoti Datta, Deep Ganguli, Denil Emerson, Denil Klecko, Deniz Yurek, Derek Chen, Derek Tam, Dieuwke Hupkes, Diganta Misra, Dilayar Buzant, Dimitri Coelho Mollo, Diyi Yang, Dong-Ho Lee, Ekaterina Shutova, Ekin Dogus Cubuk, Elad Segal, Eleanor Hagerman, Elizabeth Barnes, Elizabeth Donnelly, Ellie Pavlick, Emanuele Rodola, Emma Lam, Eric Chu, Eric Tang, Erik Erdean, Ernie Chang, Ethan A. Chi, Ethan Dyer, Ethan Jerzak, Ethan Kim, Eunice Efenyu Manyangi, Eunice Velzontholzhskii, Fanuye Xia, Fatemeh Siar, Fernand Martínez-Plumed, Francesca Happé, Frédéric Chollet, Frieda Rong, Gaura Mishra, Genta Indra Winata, Gerard de Melo, Germán Kruszewski, Giambattista Parascandolo, Giorgio Mariani, Gloria Wang, Gualberto Jaimovich-López, Gregor Betz, Guy Gur-Ari, Hana Galijasevic, Hannah Kim, Hannah Rashkin, Hananach Hajisirizi, Harsh Mehta, Hayden Bogar, Henry Shevlin, Hinrich Schütze, Hiromu Yakuri, Hongming Zhang, Hugh Mee Wong, Ian Ng, Isabele Anne, Jaya Junelet, Jack Geissinger, Jackson Kemion, Jacob Hilton, Jaehoon Lee, Jaime Fernández Fisa, James B. Simon, James Koppel, James Zheng, James Zhou, Jan Kooçof, Jana Thompson, Jared Kaplan, Jayara Radom, Jasha Soh-Dickstein, Jason Jason, Jason Wang, Jason Xiong, James Yosinski, Jenekaterina Novikova, Jela Bosscher, Jennifer Marsh, Jeremy Kim, Jerome Taal, Jesse Engel, Jesujob Alabi, Jiacheng Xu, Jianming Song, Julian Liang, Joan Wawerl, John Burden, John Miller, John U. Balis, Jonathan Berant, Jörn Frohberg, Jason Rose, Jose Hernandez-Orralo, Joseph Boudeau, Joseph Jones, Joshua S. Rule, Joyce Chua, Kamil Kancerz, Karen Livescu, Karl Krauth, Kartik Gopalakrishnan, Katerina Ignatyeva, Katja Markert, Kaustubh D. Dhole, Kevin Gimpel, Kevin Omundi, Kory Mathewson, Kristen Chaitfield, Ksema Shkantra, Kumar Shridhar, Kyle McDonell, Kyle Richardson, Laria Reynolds, Leo Gao, Li Zhang, Liam Dugan, Lianhui Qin, Lidia Contreras-Ochando, Louis-Philippe Morency, Lucas Moschella, Lucas Lam, Lucy Noble, Ludwig Schmidt, Luheung Ho, Luis Oliveros Colón, Luke Metz, Lütke Kerren, Menel, Maarten Bosma, Maarten Saap, Maartje ter Hoeve, Maheen Farooqi, Manal Faruqui, Mantas Mazeika, Marco Baturan, Marco Marcelli, Marco Maria, Marie Jose Ramirez Quintana, Marie Tolikhie, Marie Giulianelli, Martha Lewis, Martin Potthast, Mathew L. Leavitt, Matthias Hagen, Mátévás Schubert, Medina Orduna Batemirova, Melody Arnaud, Michael McElrath, Michael A. Yee, Michael Cohen, Michael Gu, Michael Ivantiskiy, Michael Starritt, Michael Strube, Michael Swedrowski, Michele Belafacqua, Michihiro Yasuura, Mihr Kale, Mike Can, Mirex Yu, Miraz Sugano, Miyo Tiwari, Mohit Bansal, Moni Ammiaseri, Moni Geva, Mozhdeh Gheini, Mokdad Varnia T, Nanyan Peng, Natinah Chn, Nayeon Lee, Neta Gur-Ari Krakover, Nicholas Cameron, Nicholas Roberts, Nick Doiron, Nikita Nangia, Nikita Deckers, Niklas Muennighoff, Nitish Shirish Keskar, Niveditha S. Iyer, Noam Constant, Noah Friedl, Nuwan W. Oliver Zhang, Omar Elbadrawi, Omer Levy, Oshin Evans, Pablo Moreno Casares, Parth Doshi, Pascale Fung, Paul Liang, Paul Viola, Peipei Alipournejad, Peter Han, Percy Liang, Peter Chang, Peter Eckersley, Phi Mu Htun, Phayre Hwang, Piota Milkowska, Plymara Pava, Pooya Pezeshkpour, Prithi Orlitz, Qiangliang Chen, Qiangliang Gu, Rhynon Garg, Richard Baraniuk, Rif A.一口水, Riku Arakawa, Rudolfas Rasteb Gabrilov, Rubin Hoberman, Ramón Ruiz Delgado, Raphaël Roussel, Rosanna Liu, Rowan Joshi, Sajant Arora, Sam Dillavou, Sam Shleifer, Sami Wissel, Samrat Ganguly, Samvel R. Bowen, Samy Bengio, Sanghoon Han, Sangeeta Kwatra, Sarah A. Ross, Sami Ghazarian, Sami Ghosh, Sami Casy, Sebastian Brattin, Sebastian Gehrmann, Sebastian Schuster, Sepehde Sadeghi, Shadi Hamdan, Sharun Zhou, Shashank Srivastava, Sherry Shi, Shishir Singh, Shima Asadi, Shixiana Shang Gu, Shubh Pathschigbar, Shubham Toshinali, Shyan Upadhyah, Shyamolima (Shammie) Debnath, Sumeek Shakeri, Simon Thomeyer, Simonet Melzi, Siva Reddy, Suhba Prisilia Makaini, Soow-Hwan Lee, Spencer Torrence, Srebarsha Hatwar, Stanisha Dehaene, Stefan Diviac, Stella Ermon, Stella Biderman, Stephanie Lin, Stephen Prasad, Steven T. Piantadosi, Stuart M. Shieber, Summer Misherghi, Svetlana Kiritchenko, Swapna Mishra, Tal Linzen, Tal Schuster, Tao Li, Tao Yu, Tariq Ali, Tatsu Hashimoto, Te-Lin Wu, Théo Desbordes, Thore Rothchild, Thomas Pham, Tianle Wang, Tiberius Nkinyih, Timo Schick, Timofei Korot, Timothy Tellessen-Lawton, Tintu Sudharsan, Tobias Gerstenberg, Trenton Chang, Trishala Neeraj, Tushar Khot, Tyler Shultz, Uri Shaham, Vedant Misra, Verma Demberg, Victoria Nymark, Vilas Ruukan, Vinay Ramasesh, Vinay Úday Prabhu, Vishakh Padmakumar, Vivek Srikanth, William Saunders, William Zhang, Wout Vossem, Xiang Ren, XiaoYao Tong, Xian Zhou, Xinyi Wu, Xudong Shen, Yadollah Yaghoozbadeh, Yair Lazrek, Yangju Song, Yasaman Barabi, Yejin Choi, Yichi Yang, Yiding Hao, Yifu Chen, Yuanan Belinkov, Yu Hou, Yufang Hou, Yunato Bai, Zachary Seid, Zhuoye Zhao, Zijian Wang, Zijie J. Wang, Zirui Wang, Zixi Wu.

Aside: new benchmarks for multitask LMs

BIG-Bench [[Srivastava et al., 2022](#)]

200+ tasks, spanning:



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.

.....#.....
.....#.....
#####.....#.....
.....#####.....
.....#.#.#.#.....
.....###.#.....###.....
....##.....#.....##.....
.....#.....#.....##.....
.....##.....#.....##.....
#####.....#####.....#.....
.....##.....#.....##.....
.....#####.....#####.....#.....
.....##.....#.....##.....#.....
.....##.....#.....##.....#.....

Gains from Instruction Finetuning

- ▶ Lots of models based on finetuning T5
 - Flan-T5
 - Tk-Instruct
 - T0
 -

Params	Model	Norm. avg.
80M	T5-Small	-9.2
	Flan-T5-Small	-3.1 (+6.1)
250M	T5-Base	-5.1
	Flan-T5-Base	6.5 (+11.6)
780M	T5-Large	-5.0
	Flan-T5-Large	13.8 (+18.8)
3B	T5-XL	-4.1
	Flan-T5-XL	19.1 (+23.2)
11B	T5-XXL	-2.9
	Flan-T5-XXL	23.7 (+26.6)

Bigger model →
= bigger Δ

Qualitative Results

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.

✖ (doesn't answer question)

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

<https://huggingface.co/google/flan-t5-xxl>

Qualitative Results

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:

<https://huggingface.co/google/flan-t5-xxl>

[Chung et al., 2022]

Lecture Plan: From Language Models to Assistants

1. Instruction finetuning

- + Simple and straightforward, generalize to unseen tasks
- ?
- ?

2. Reinforcement Learning from Human Feedback (RLHF)

3. What's next?

Limitations of instruction finetuning?

- **Problem 1:** it's expensive to collect ground-truth data for tasks
 - *Provide me five active research areas in April 2023 for LLMs*
- **Problem 2:** tasks like open-ended creative generation have no right answer.
 - *Write me a story about a dog and her pet grasshopper.*
- **Problem 3:** Even with instruction tuning, you are not directly “maximizing human preferences”
- Can we explicitly attempt to satisfy human preferences?



Lecture Overview

- ▶ Instruction Finetuning
- ▶ **Reinforcement Learning from Human Feedback (RLHF)**
- ▶ Open challenges with RLHF

Optimizing for human preferences

- Let's say we were training a language model on some task (e.g. summarization).

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.

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
...
overtake unstable
objects.

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.

$$s_2 \\ R(s_2) = 1.2$$

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.

$$s_1 \\ R(s_1) = 8.0$$

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

$$s_2 \\ R(s_2) = 1.2$$

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

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

A (very!) brief introduction to policy gradient/REINFORCE [Williams, 1992]

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

A (very!) brief introduction to policy gradient/REINFORCE [Williams, 1992]

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

We **reinforce** good actions, increasing the chance they happen again.

If R is +++



If R is ---



Take gradient steps to maximize $p_\theta(s_i)$



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

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.

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!

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](#)]

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.

$$s_2 \\ R(s_2) = 1.2$$


Train an LM $RM_\phi(s)$ to predict human preferences from an annotated dataset, then optimize for RM_ϕ 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

$R(s_3) = 4.1? \quad 6.6? \quad 3.2?$

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.

s_1

>

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

s_3

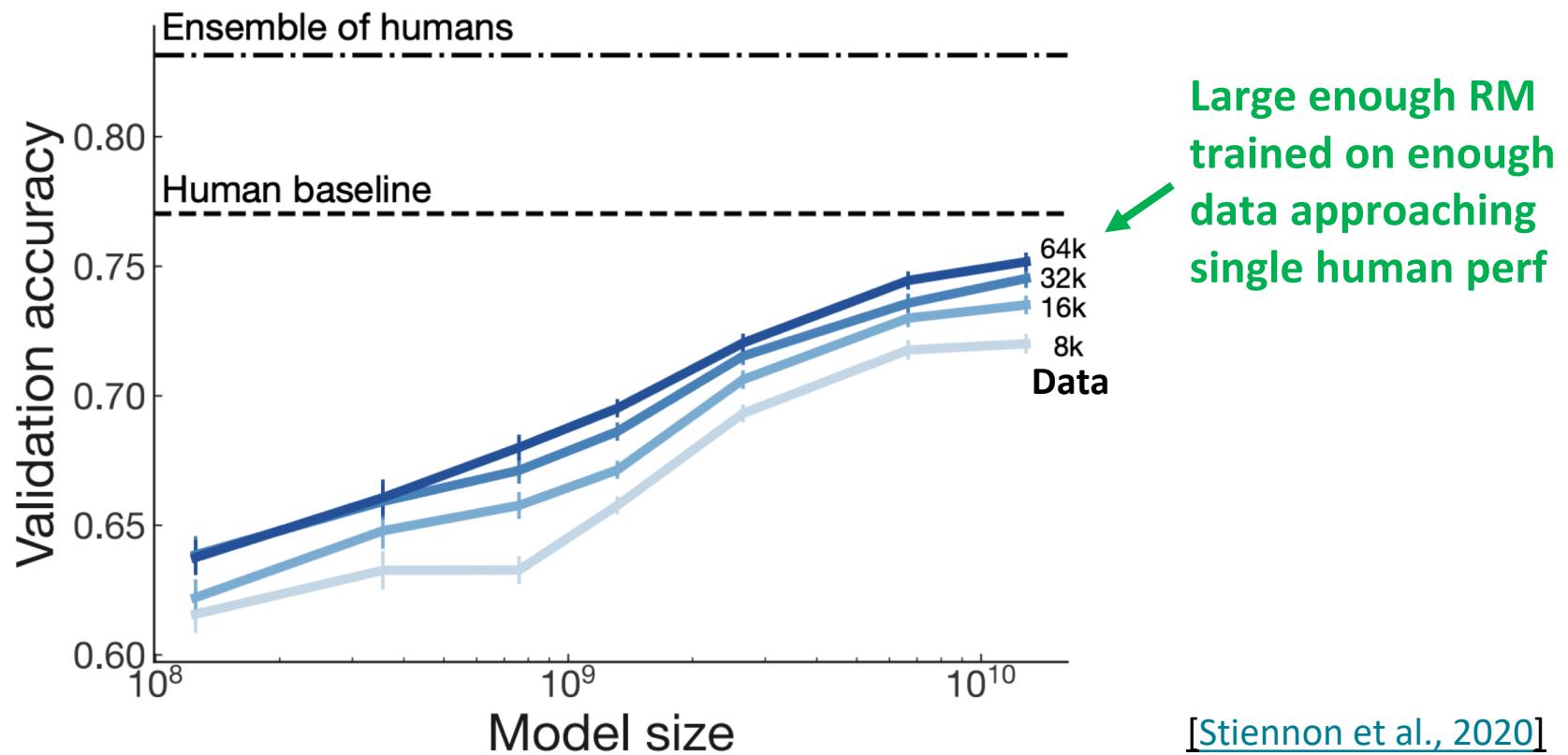
>

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

s_2

Make sure your reward model works first!

Evaluate RM on predicting outcome of held-out human judgments



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.

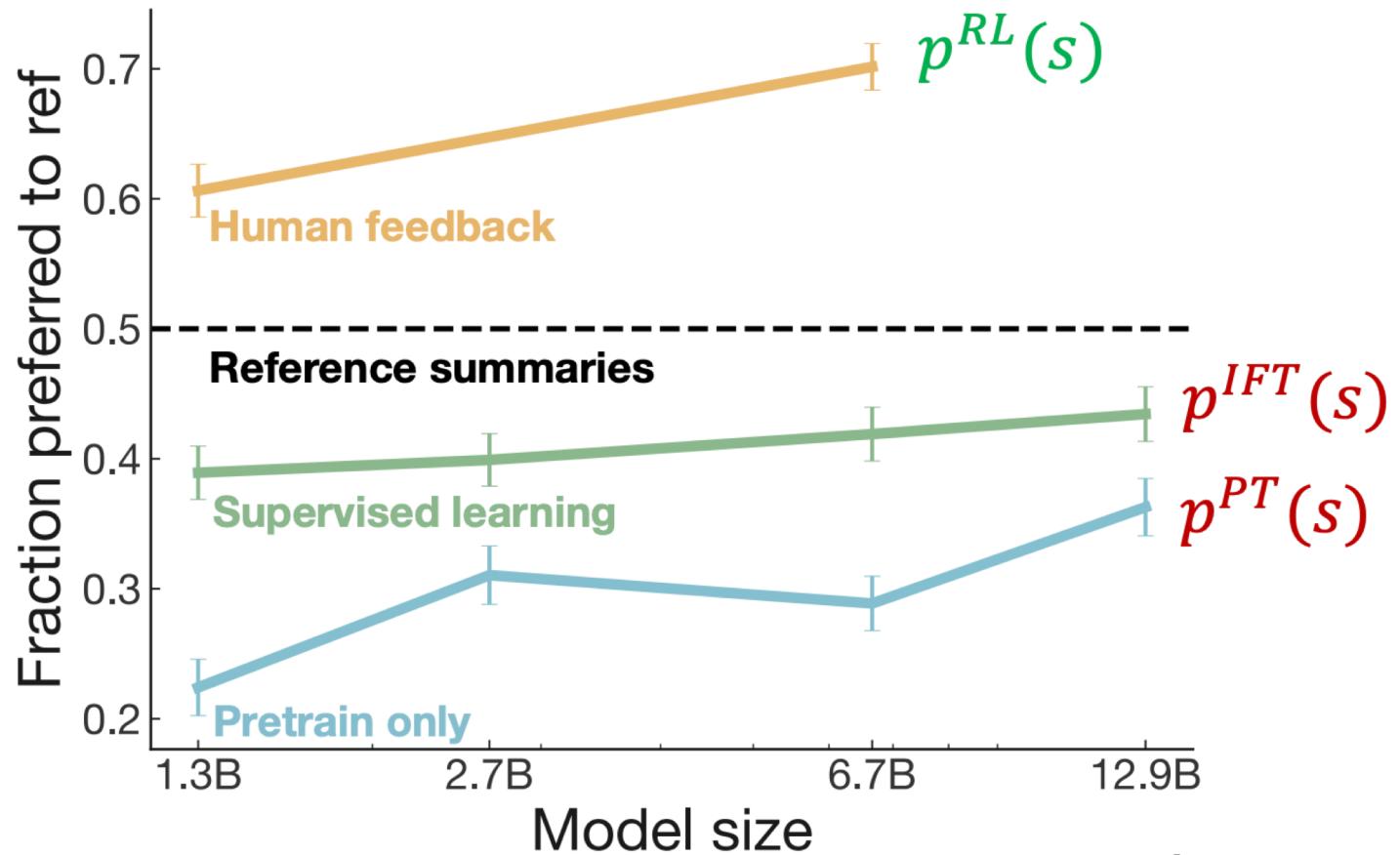
RLHF: Putting it all together [Christiano et al., 2017; Stiennon et al., 2020]

- 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 $p_\theta^{RL}(s)$, with parameters θ 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
 $p_\theta^{RL}(s) > p^{PT}(s)$

RLHF provides gains over pretraining + finetuning



[Stiennon et al., 2020]

InstructGPT: scaling up RLHF to tens of thousands of tasks

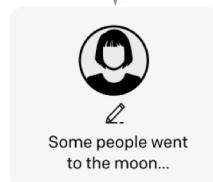
Step 1

Collect demonstration data, and train a supervised policy.

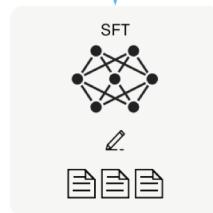
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



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



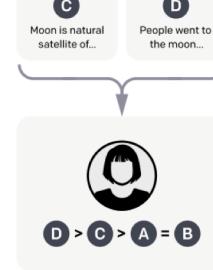
Step 2

Collect comparison data, and train a reward model.

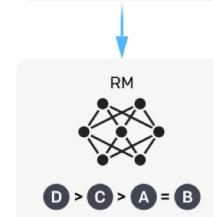
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



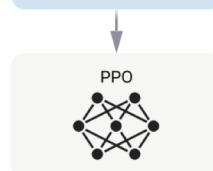
Step 3

Optimize a policy against the reward model using reinforcement learning.

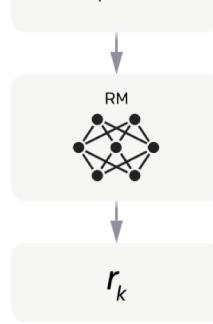
A new prompt is sampled from the dataset.



The policy generates an output.



Once upon a time...



r_k

The reward model calculates a reward for the output.

The reward is used to update the policy using PPO.

[\[Ouyang et al., 2022\]](#)

InstructGPT

PROMPT *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 Proximal Policy Optimization. We performed several iterations of this process.

(RLHF!)



Lecture Overview

- Instruction Finetuning
- Reinforcement Learning from Human Feedback (RLHF)
- **Open challenges with RLHF**

Limitations of RL + Reward Modeling

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



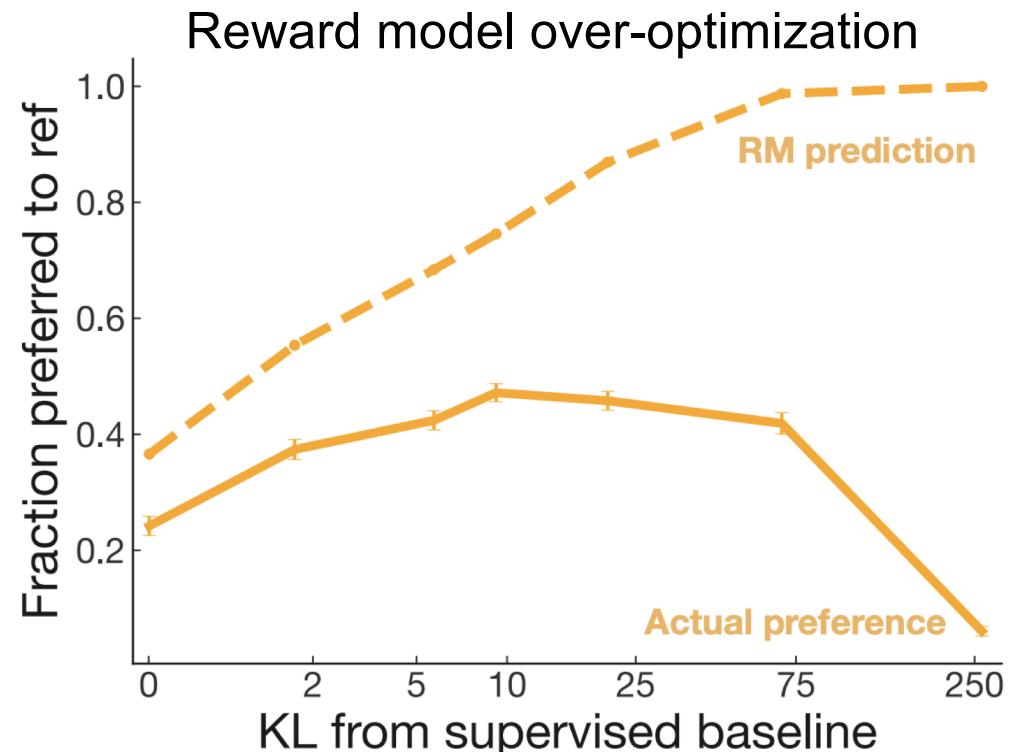
<https://openai.com/blog/faulty-reward-functions/>

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

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!



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

[Stiennon et al., 2020]