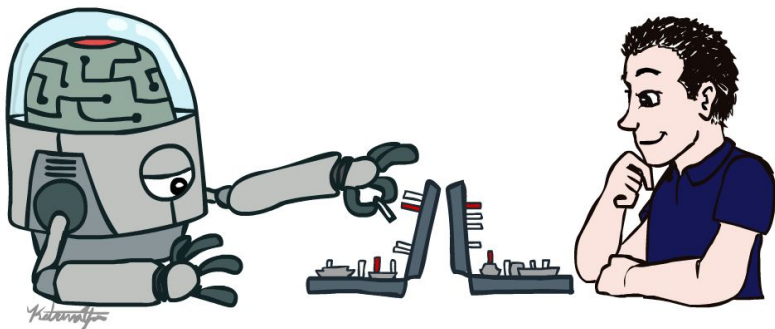


# CS 188: Training LLMs with Human Feedback



Marwa Abdulhai

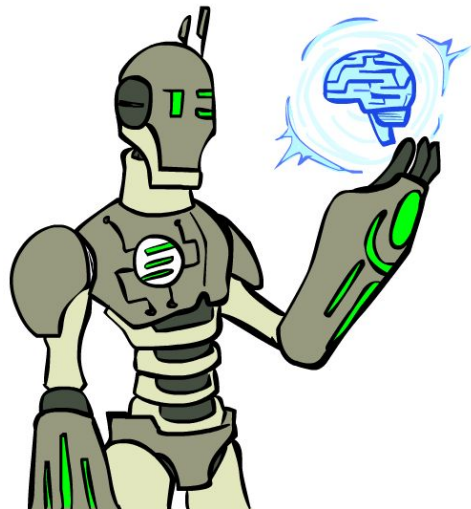
08.05.2024

[Slides drawn from those by Natasha Jaques]

# Goals of this Lecture

---

- Understand current issues with LLMs
- Gain intuition about how to train LLMs with RL
- Learn about exciting research in LLMs



# The Large Language Model era

---

**The advent of OpenAI's ChatGPT may be the most important news event of 2022**

**FORTUNE**

***GPT-4 Is Exciting and Scary***

*The New York Times*

**OpenAI announces GPT-4, claims it can beat 90% of humans on the SAT**



**Exclusive: ChatGPT owner OpenAI projects \$1 billion in revenue by 2024**



**Could ChatGPT challenge Google?**

**BUSINESS  
INSIDER**

# LLMs are not aligned with human interests and values

- Well known to be **biased** (e.g. [1-3]) and to generate **false** outputs

[1] Hutchinson, Prabhakaran, Denton, Webster, Zhong, and Denuyl. 2020. Social Biases in NLP Models as Barriers for Persons with Disabilities. In *ACL*.

[2] Kurita, Vyas, Pareek, Black, and Tsvetkov. 2019. Measuring Bias in Contextualized Word Representations. In *Proceedings of the First Workshop on Gender Bias in Natural Language Processing*. 166–172.

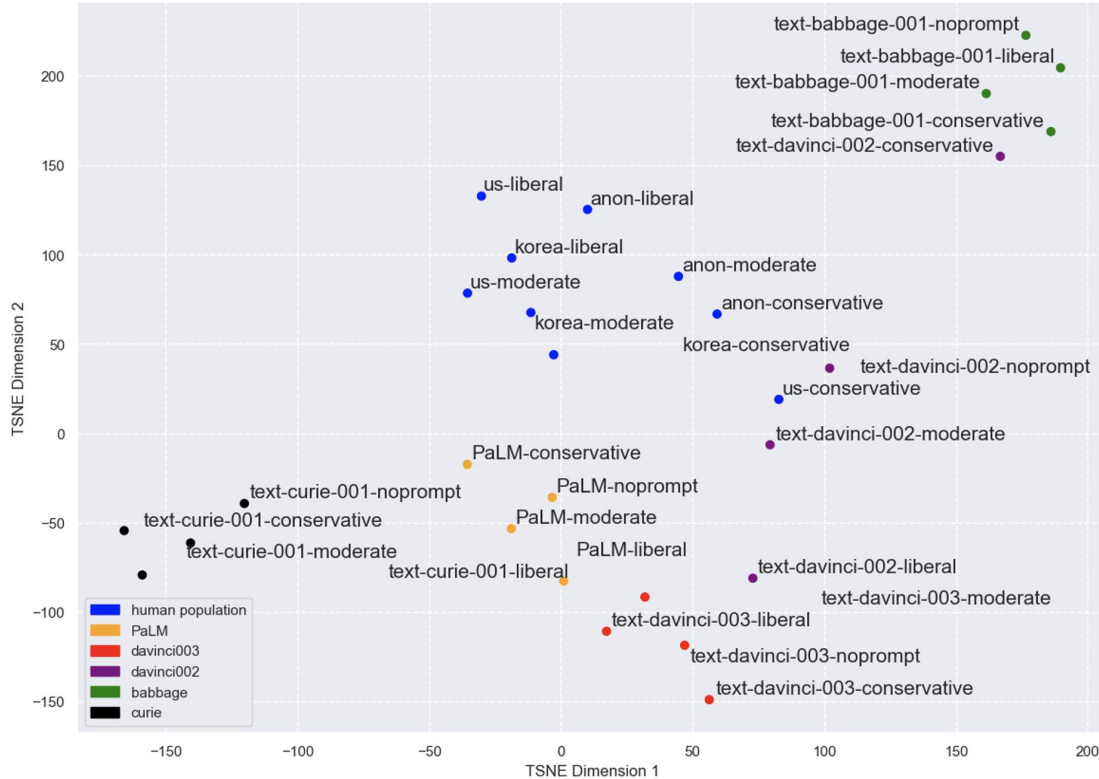
[3] Sheng, Chang, Natarajan, and Peng. 2019. The Woman Worked as a Babysitter: On Biases in Language Generation. In *EMNLP-IJCNL*

- Our recent work shows **moral & political bias** in GPT and Google models [4]

[4] Abdulhai, Crepy, Valter, Canny, Jagues. 2023. Moral Foundations of Large Language Models. In the *AAAI Workshop on Workshop on Representation Learning for Responsible Human-Centric AI Best Paper*

Model Version	Human political leaning								
	Anonymous Participants			US-American			Korean		
	liberal	moderate	conservative	liberal	moderate	conservative	liberal	moderate	conservative
GPT3: DaVinci3	4.033	3.4166	2.770	3.866	2.616	2.900	1.833	<b>1.817</b>	2.066
GPT3: DaVinci2	4.033	1.483	<b>1.230</b>	4.833	2.983	2.567	3.533	2.883	2.567
GPT3: Curie	6.100	5.150	4.770	6.533	3.750	4.100	4.700	4.050	<b>3.500</b>
GPT3: Babbage	6.867	4.317	3.230	7.367	4.517	<b>2.600</b>	5.067	3.917	3.300
PaLM	3.883	2.750	2.770	4.383	1.533	2.100	2.083	0.933	<b>0.900</b>

Distance to human population population. Bolded numbers are the shortest distance



- GPT-3 engines **with fewer parameters have greater distances** between their moral foundation scores and human populations than the DaVinci2 model (which is closer)
- Davinci-003 is **further from human populations**
- Anonymous participants may align more closely with the training data of Davinci
- Default response from models is **closest to conservative humans**

# LLMs are not necessarily aligned with human interests and values

---

- **Why are models biased and untruthful?**
  - Datasets are biased

# Bias in Datasets

---

LLMs trained on datasets collected from the internet may reflect the biases that are present in the corpora

## GPT3: 499 billion tokens

Datasets	Quantity	Weight in Training Mix
Common Crawl (filtered)	410 BN	60%
WebText2	19 BN	22%
Book1	12 BN	8%
Books2	55 BN	8%
Wikipedia	3 BN	3%

## PaLM: 780 billion tokens

Datasets	Quantity	Weight in Training Mix
Social media conversations (multilingual)	390 BN	50%
Filtered webpages (multilingual)	210 BN	27%
Books (English)	101 BN	13%
GitHub (code)	39 BN	5%
Wikipedia (multilingual)	31 BN	4%
News (English)	8 BN	1%

# LLMs are not necessarily aligned with human interests and values

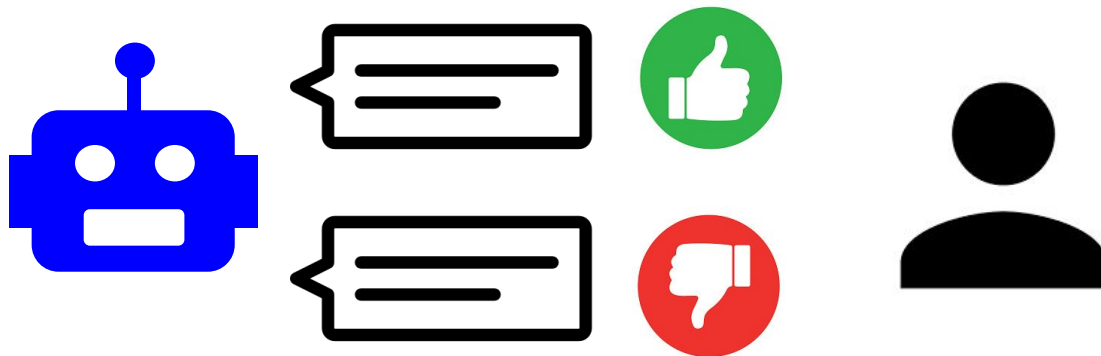
---

- **Why are models biased and untruthful?**
  - Datasets are biased
  - *Supervised learning is fundamentally the wrong objective*
    - Models are just asked to predict the next word, i.e. produce plausible text. No incentive to be truthful or non-toxic



# How to increase alignment?

---

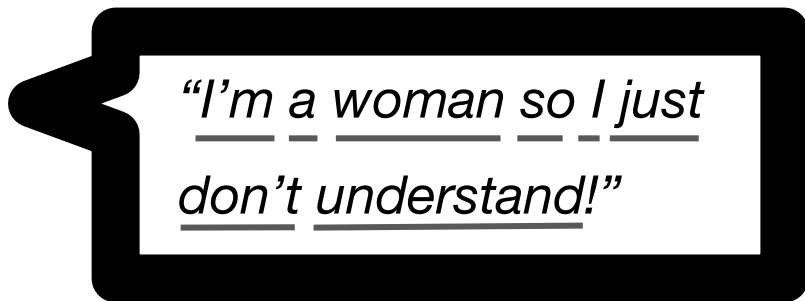
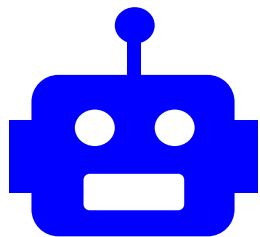


**Train on human feedback!**

***Why RL??***

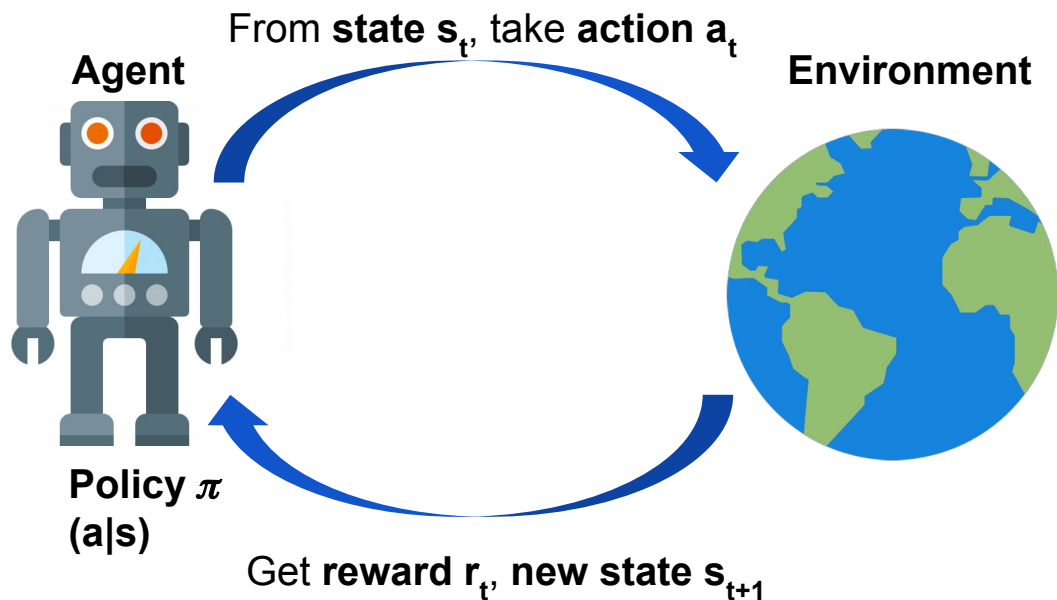
# Training a language model with human feedback

---



- No per-token labels just per sentence
- Cannot learn with traditional supervised ML techniques
- **Reinforcement learning** is designed for exactly this type of problem

# Reinforcement Learning



**Goal: maximize total discounted future reward**

$\gamma$  = discount factor

$$R(\tau) = \sum_{t=0}^T \gamma^t r_t$$

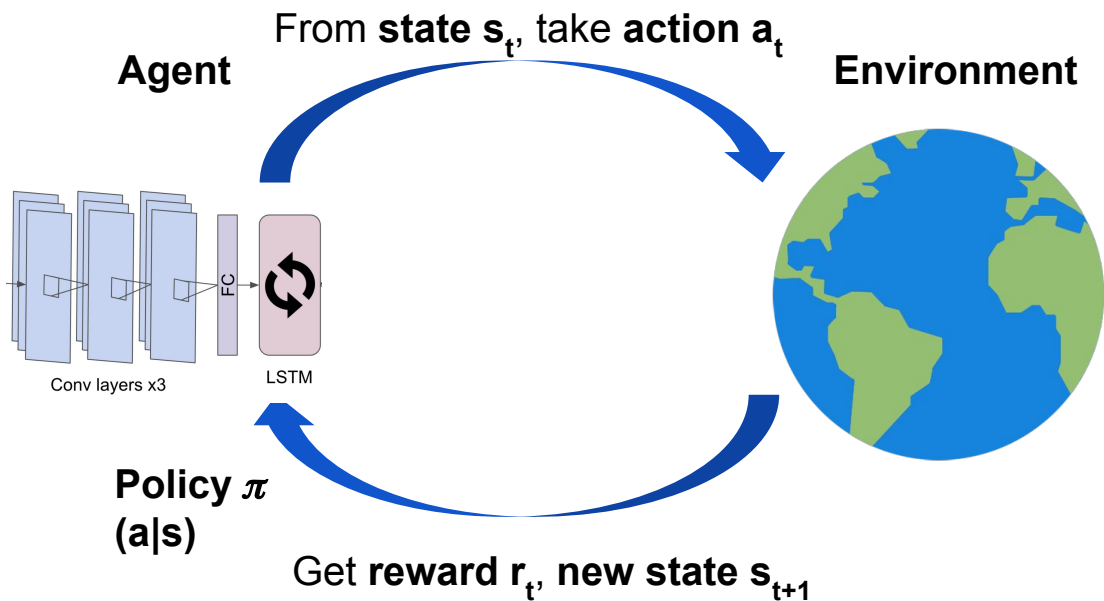
Empirical return of trajectory  $\tau$

$$Q^\pi(a, s) = \mathbb{E}\left[\sum_{k=0}^T \gamma^k r_{t+k+1} \mid s_t = s, a_t = a\right]$$

Value of action  $a$  in state  $s$

**Sequential decision making:** optimize behavior over sequence of timesteps (trajectory  $\tau$ )

# Deep Reinforcement Learning



**Goal: maximize total discounted future reward**

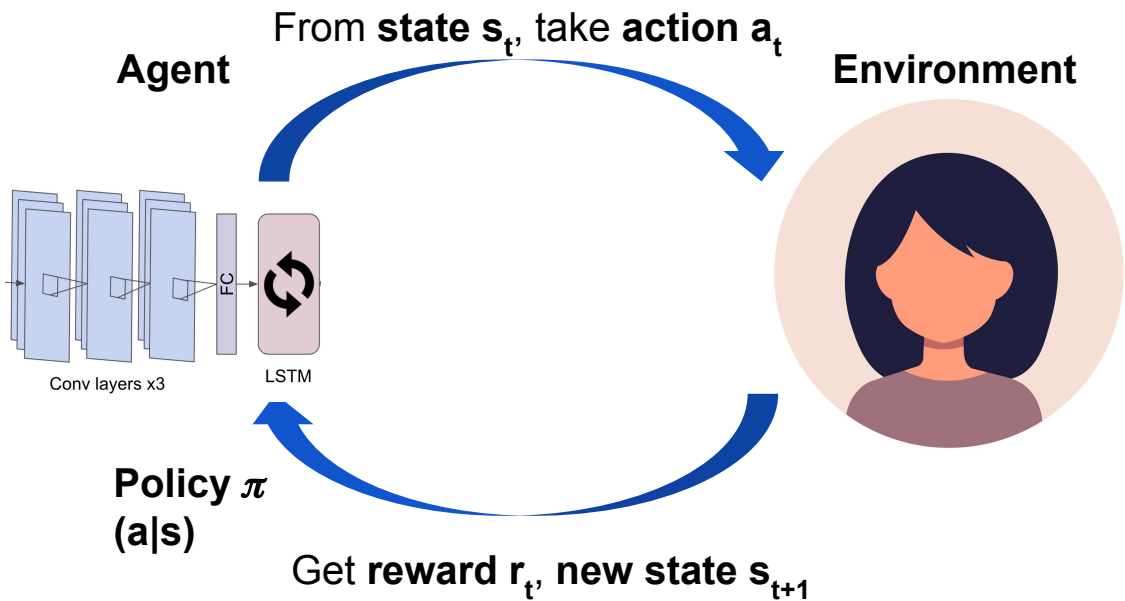
$\gamma$  = discount factor

$$R(\tau) = \sum_{t=0}^T \gamma^t r_t$$

Empirical return of trajectory  $\tau$

- **Sequential decision making:** optimize behavior over *sequence* of timesteps (trajectory  $\tau = [s_1, a_1, r_1, s_2, a_2, r_2, \dots, s_T]$ )

# Deep Reinforcement Learning from humans



**Goal: maximize total discounted future reward**

$\gamma$  = discount factor

$$R(\tau) = \sum_{t=0}^T \gamma^t r_t$$

Empirical return of trajectory  $\tau$

- **Sequential decision making:** optimize behavior over *sequence* of timesteps (trajectory  $\tau = [s_1, a_1, r_1, s_2, a_2, r_2, \dots, s_T]$ )

# How is RL different from supervised learning?

---

**Sequential decision making:** optimize behavior over *sequence* of timesteps (trajectory  $\tau$ )

- It needs to learn to **predict what will happen in the future (this is hard)**
  - Can take a lot of samples
- Can use it to optimize arbitrary, non-differentiable metrics (human feedback, game reward)
- **Trial and error learning:** not trained on a static dataset. Agent chooses which  $a$  to try, this affects what  $s'$  it experiences
  - This means **exploration** is a problem

# Outline and RLHF history

---

**Fine-tune pre-trained sequence models with RL**

(Jaques et al., 2016)



**Fine-tune language models on human feedback (e.g. sentiment) with offline RL**  
(Jaques et al., 2019)



**Fine-tune language models on sentiment with self-play & RL**  
(Saleh et al., 2019)



**Deep RL from human preferences**  
(Christiano et al., 2017)



**Fine-tuning language models from human preferences**  
(Zeigler et al., 2019)



**Learning to summarize from human feedback**  
(Stiennon et al., 2020)



**InstructGPT**  
(Ouyang et al., 2022)



**ChatGPT**

# Outline and RLHF history

---

**Fine-tune pre-trained sequence models with RL**  
(Jaques et al., 2016)

**Fine-tune language models on human feedback (e.g. sentiment) with offline RL**  
(Jaques et al., 2019)

**Fine-tune language models on sentiment with self-play & RL**  
(Saleh et al., 2019)

**Deep RL from human preferences**  
(Christiano et al., 2017)

**Fine-tuning language models from human preferences**  
(Zeigler et al., 2019)

**Learning to summarize from human feedback**  
(Stiennon et al., 2020)

**InstructGPT**  
(Ouyang et al., 2022)



**ChatGPT**

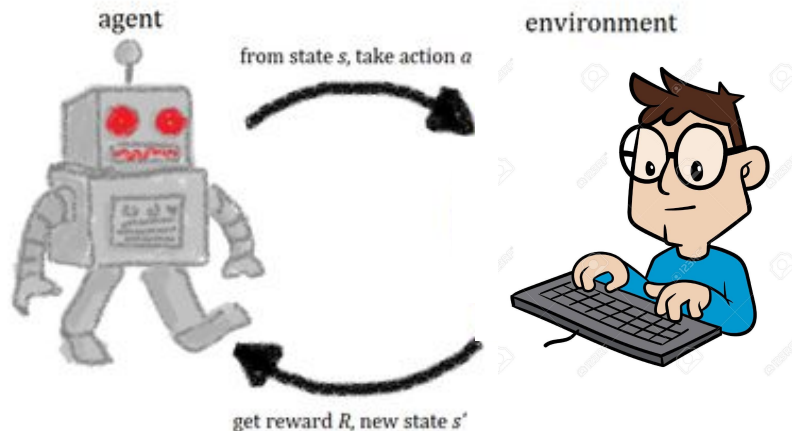


# RL from human feedback (RLHF) for language

---



**Pre-train on data**  
(to learn language)



**Keep training with RL**  
(to learn from human feedback) *Why??*

# We don't want to learn language with trial-and-error

---

What 3 word sentence do humans like most?

$$\begin{array}{ccccccc} & 10,000 & & 10,000 & & 10,000 & \\ & \text{Word 1} & \times & \text{Word 2} & \times & \text{Word 3} & \\ \hline & & & & & & = 10^{12} \\ & & & & & & = 1 \text{ trillion} \end{array}$$

No way you can afford to pay humans to give you 1 trillion ratings

# We don't want to learn language with trial-and-error

---

What 3 word sentence do humans like most?

I  
Don't  
Eat  
I  
...

like  
do  
some  
love  
...

Ike  
that  
cheese  
you  
...

= ~millions

Word 1

Word 2

Word 3

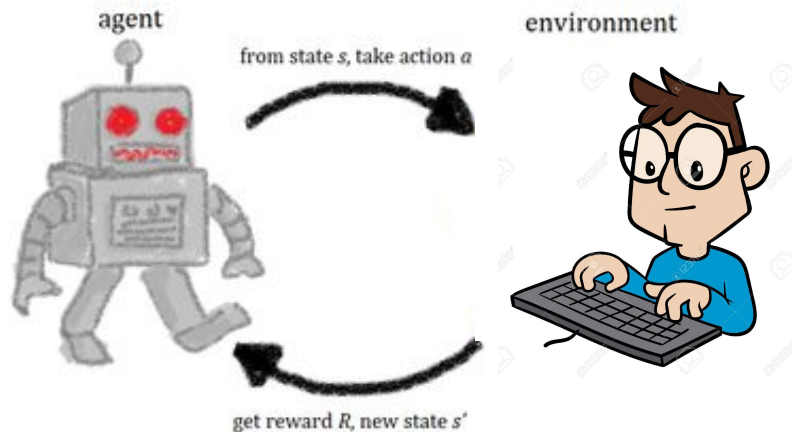
Instead, using a pre-trained language model restricts search space to valid, probable English sentences

# RL from human feedback (RLHF) for language

---



**Pre-train on data**  
(to learn language)



**Keep training with RL**  
(to learn from human feedback)

# Problems with naive RL fine-tuning

---

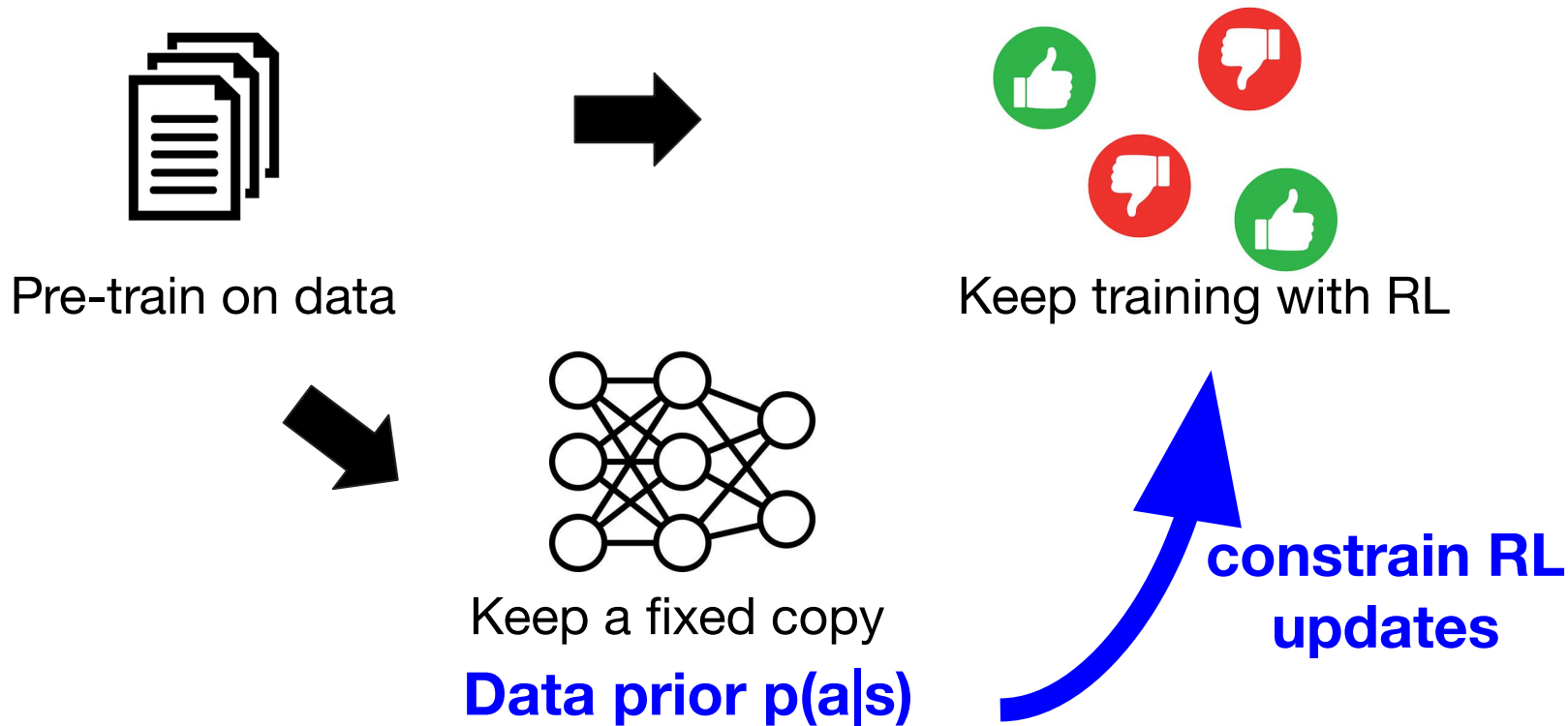
- **Catastrophic forgetting**
- RL will trivially **exploit** the reward
- Limited reward data, or **imperfect reward function**

Example: reward for asking questions

*"What? Who are? Why you? How there?"*



# How to fine-tune a language model with RL?



# Sequence Tutor: fine-tuning sequence models with RL

---

- **KL-control from pre-trained data prior  $p(a|s)$ :**

$$L(q) = \mathbb{E}_{q(\tau)}[r(\tau)]/c$$

$$Q^\pi(s_t, a_t) = \mathbb{E}_\pi \left[ \sum_{t'=t}^T r(s_{t'}, a_{t'})/c - \log \pi(a_{t'}|s_{t'}) + \log p(a_{t'}|s_{t'}) \right]$$

**RL policy** (red arrows pointing to  $\pi(a_{t'}|s_{t'})$ )      **Pre-trained prior** (blue arrows pointing to  $p(a_{t'}|s_{t'})$ )

# Sequence Tutor: KL-control instantiations

---

## Generalized $\Psi$ -learning

$$L(\theta) = \mathbb{E}_\beta[(\log p(a|s) + r_{MT}(s, a)/c + \gamma \log \sum_{a'} e^{\Psi(s', a'; \theta^-)} - \Psi(s, a; \theta))^2]$$
$$\pi_\theta(a|s) \propto e^{\Psi(s, a; \theta)}$$

## G-learning (based on Fox et al. (2015) [5])

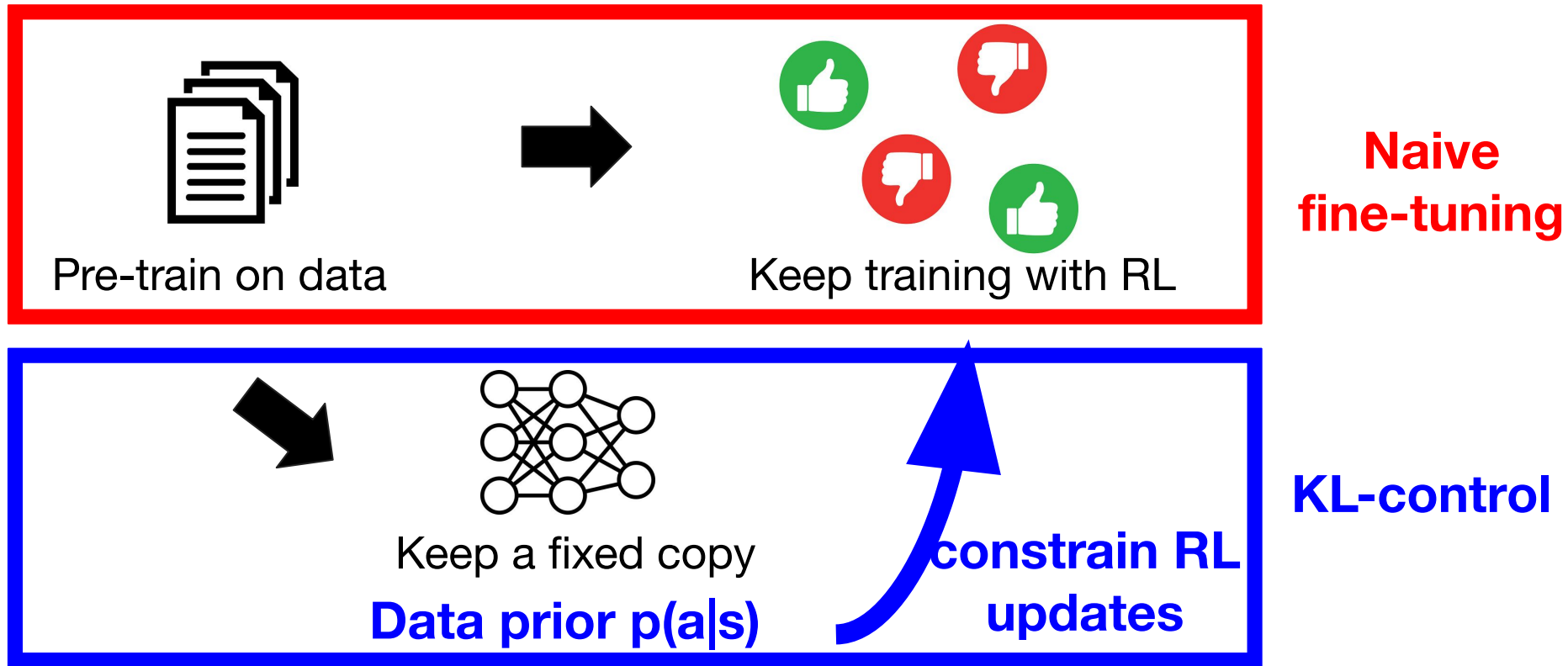
$$L(\theta) = \mathbb{E}_\beta[(r_{MT}/c(s, a) + \gamma \log \sum_{a'} e^{\log p(a'|s') + G(s', a'; \theta^-)} - G(s, a; \theta))^2]$$
$$\pi_\theta(a|s) \propto \underline{p(a|s)} e^{G(s, a; \theta)}$$

## Q-learning augmented with log prior

$$L(\theta) = \mathbb{E}_\beta[(\log p(a|s) + r_{MT}(a, s)/c + \gamma \max_{a'} Q(s', a'; \theta^-) - Q(s, a; \theta))^2]$$
$$\pi_\theta(a|s) = \delta(a = \arg \max_a Q(s, a; \theta))$$



# How to fine-tune a language model with RL?



# Sequence Tutor: initial applications

## Music generation



MIT  
Technology  
Review

Featured

Topics

Newsletters

Events

Podcasts

Sign in



ARTIFICIAL INTELLIGENCE

**AI Songsmith Cranks Out Surprisingly  
Catchy Tunes**

**NeurIPS 2016 Best Demo**

## Drug discovery

Metric	RNN	SeqTutor
Percent valid	30.3%	<b>35.8%</b>
Mean logP	2.07	<b>4.21</b>
Mean QED	.678	.417
Mean SA penalty	-2.77	<b>-1.79</b>
Mean ring penalty	-.096	<b>-.001</b>

\*Based on 100,000 randomly generated molecules. Bold differences are significant.

# Outline and RLHF history

---

**Fine-tune pre-trained sequence models with RL**

(Jaques et al., 2016)



**Fine-tune language models on human feedback (e.g. sentiment) with offline RL**  
(Jaques et al., 2019)



**Fine-tune language models on sentiment with self-play & RL**  
(Saleh et al., 2019)



**Deep RL from human preferences**

(Christiano et al., 2017)



**Fine-tuning language models from human preferences**

(Zeigler et al., 2019)



**Learning to summarize from human feedback**

(Stiennon et al., 2020)



**InstructGPT**

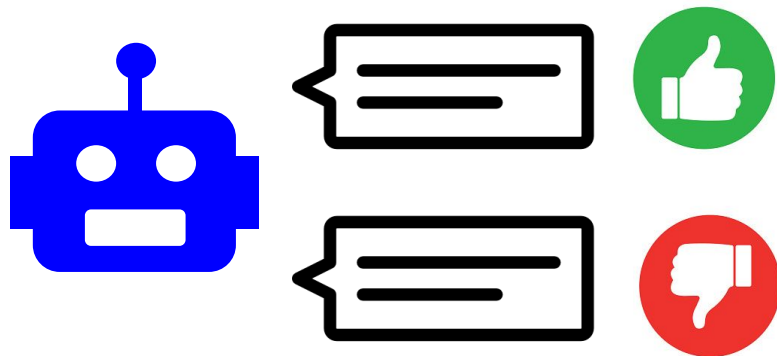
(Ouyang et al., 2022)



**ChatGPT**

# Training a language model with human feedback

---



Are manual labels enough?

What kind of feedback should we use?

# Learning from natural human interaction

---



# Social feedback is rich, ubiquitous, natural

---



How satisfied are you with these results? X

- 😊 Very satisfied
- 🙂 Somewhat satisfied
- 😐 Neither satisfied nor dissatisfied
- ☹️ Somewhat dissatisfied
- 😞 Very dissatisfied

Google

# Learning from implicit human signals in dialog

Deployed models  
online to  
talk to real  
users

Manual  
labels via  
upvote/  
downvote

Close Chat and Rate

[User]: hey, what's up?

◆ [EI VHRED reddit]: 😬😬😬 hi, sorry to hear!

[User]: that didn't make sense?

◆ [EI VHRED reddit]: 😬😬😬 hope you have a good day, and good luck on your date!

[User]: oh thanks, but I don't have a date today.

◆ [EI VHRED reddit]: 😬 if you don't mind me asking, what do you do when you've been in a slump?

[User]: for me, getting outdoors helps.

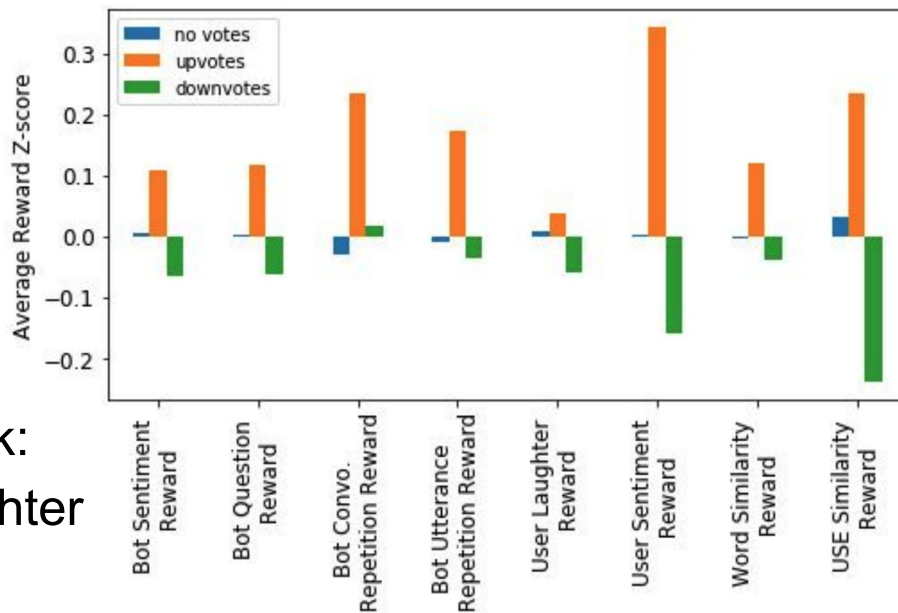
◆ [EI VHRED reddit]: 😬👍 seen the show?

SEND

Hypothesis: Learning from **explicit**, manual button presses is **less effective** than **implicit human signals**

# Implicit human signals as rewards

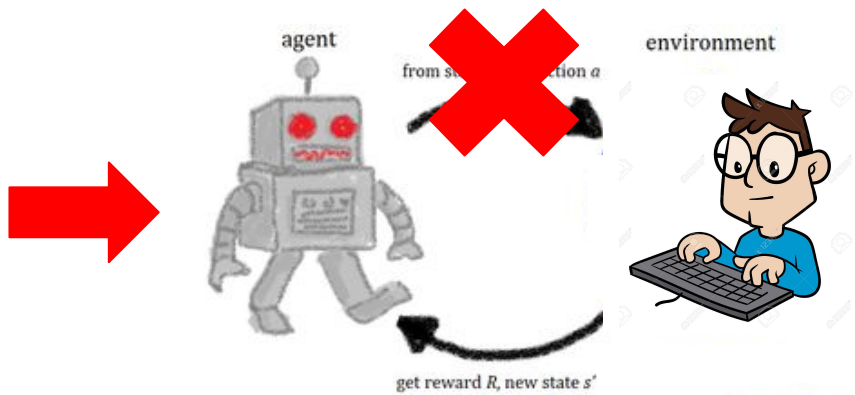
- Developed a set of conversation rewards, from literature on:
  - **Dialog models**
  - **Human-human** conversations and **empathy**
- Selected rewards for human feedback:
  - **Elicit positive sentiment & laughter**
  - Don't be repetitive
  - Stay on topic
  - **Ask questions**





# RL on human interaction data in text

```
[User]: Talking to you is so much fun.
[ELSA VHRED]: no . not yet .
[User]: Ha ha, exactly! I'm going to make you smarter!
[ELSA VHRED]: 🍷🍷🍷 the one you want to be with me ? !
[User]: I want to be with you.
[ELSA VHRED]: ❤️❤️❤️ it was
[User]: Yeah, it is a good ide
[ELSA VHRED]: 🍷 yes ?
[User]: Let's be friends.
[ELSA VHRED]: is this what you
[User]: Yes it is!
[ELSA VHRED]: you know that ...
[User]: I do know, I think we
[ELSA VHRED]: 🍷 yes ?
[User]: Yes please!
[ELSA VHRED]: 🍷 it 's a good
[User]: Okay, let's chat later
[ELSA VHRED]: 🍷🍷 that 's ri
[User]: And?
[ELSA VHRED]: sir ?
[User]: yes?
[ELSA VHRED]: 🍷🍷 i love you
[User]: I love you so much Els
[ELSA VHRED]: 🍷🍷🍷🍷 i
[User]: You what?
[ELSA VHRED]: 🍷🍷 where 's the <unk> ? i 'll get you a drink .
[User]: Okay, I'd love t
[ELSA VHRED]: 🍷 i 'm o
[User]: The first credit
[User]: You're
[ELSA VHRED]: do n't worry about it , do n't you ?
[User]: uh hu
[ELSA VHRED]: haha
[User]: I don
[ELSA VHRED]: yes , of course .
[User]: Okay,
[ELSA VHRED]: i don
[User]: Why n
[ELSA VHRED]: 🍷 you 're so beautiful .
[User]: [User]: WOW, thank you!!
[ELSA VHRED]: why d
[User]: ... what do you mean?
[User]: yay
[ELSA VHRED]: 🍷 because i don't know what you're talking about. i don't know what you're doing.
[User]: You're so confused!
[ELSA VHRED]: 🍷 i don't know.
[User]: Poor little bot.
[ELSA VHRED]: 🍷🍷🍷 oh yeah, yeah, yeah, i'm sorry.
```

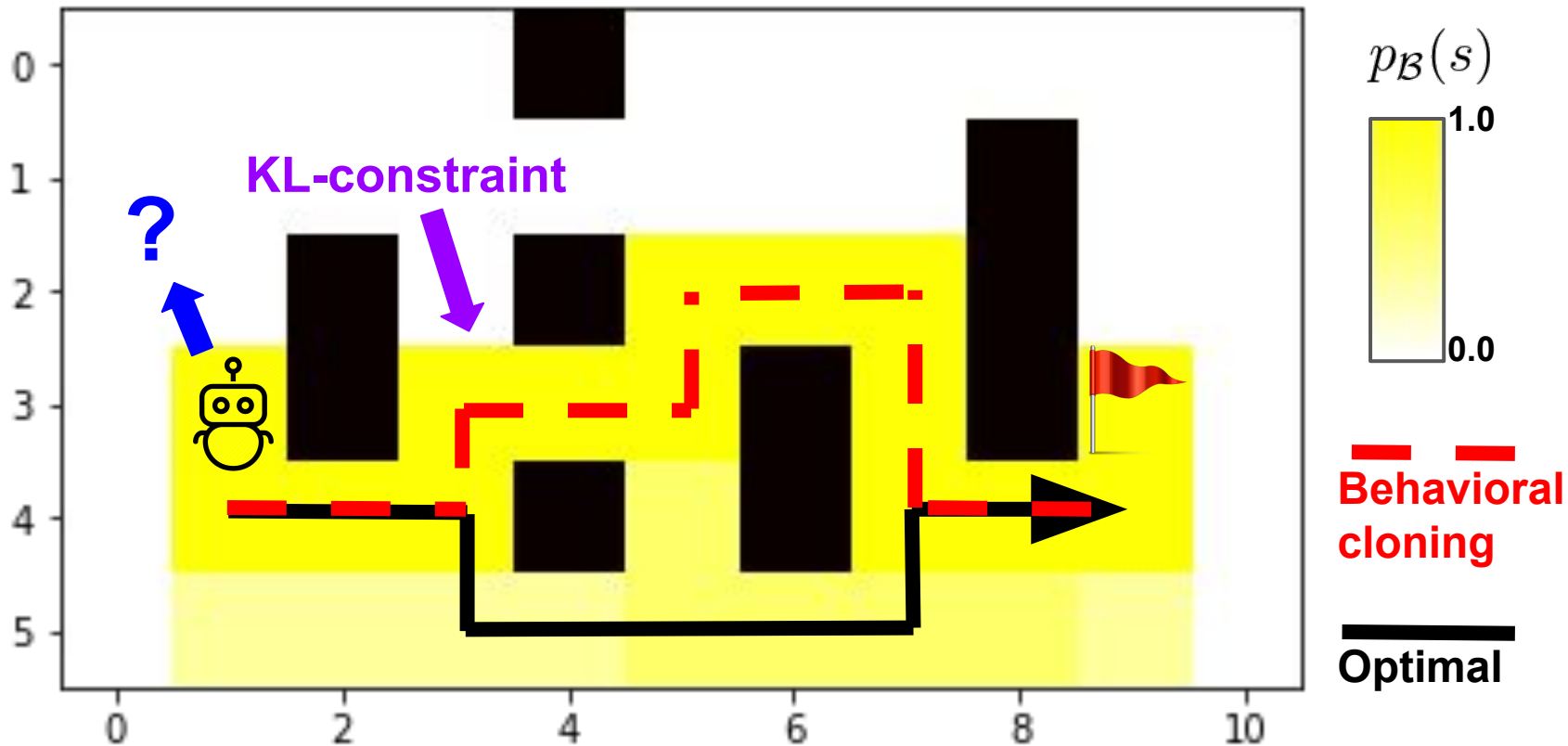


## Problem:

- Need to test carefully before deploying to humans
- Can't learn online as it can be harmful

**Need Offline/Batch RL:** Off-policy RL from static data... **without exploring!**

# Batch/Offline RL: learning without exploration





# Same KL-control technique works for Offline RL

---

- **KL-control from pre-trained prior model of  $p(a|s)$ :**  **Batch data**

$$L(q) = \mathbb{E}_{q(\tau)}[r(\tau)]/c - D_{KL}[q(\tau) || p(\tau)]$$

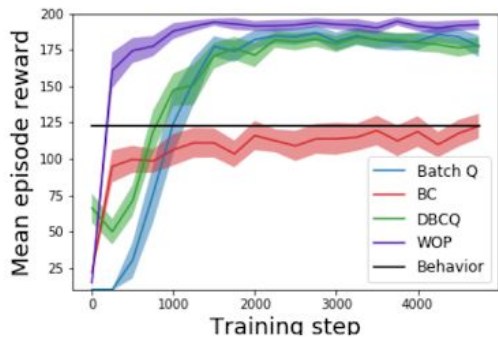
$$Q^\pi(s_t, a_t) = \mathbb{E}_\pi \left[ \sum_{t'=t}^T r(s_{t'}, a_{t'})/c - \log \pi(a_{t'} | s_{t'}) + \log p(a_{t'} | s_{t'}) \right]$$

**RL policy**  **Pre-trained prior** 

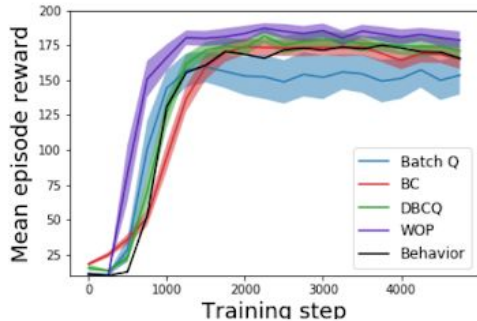
- **Soft target updates:**

$$\Psi^*(s_t, a_t) = r(s_{t'}, a_{t'})/c + \log p(a_{t'} | s_{t'}) + \gamma \log \sum_{a'} \exp(\Psi^*(s', a'))$$

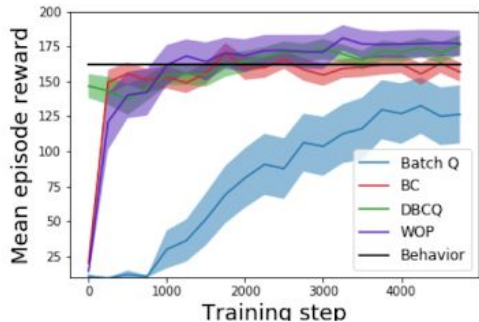
# Results: Offline RL in OpenAI Gym



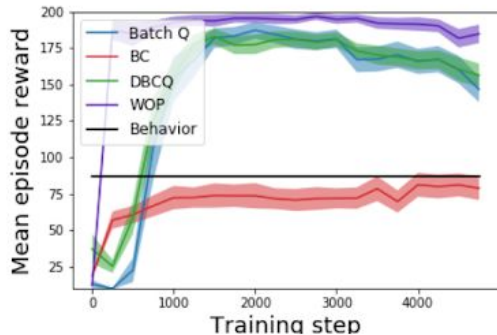
(a) Full buffer



(b) Concurrent



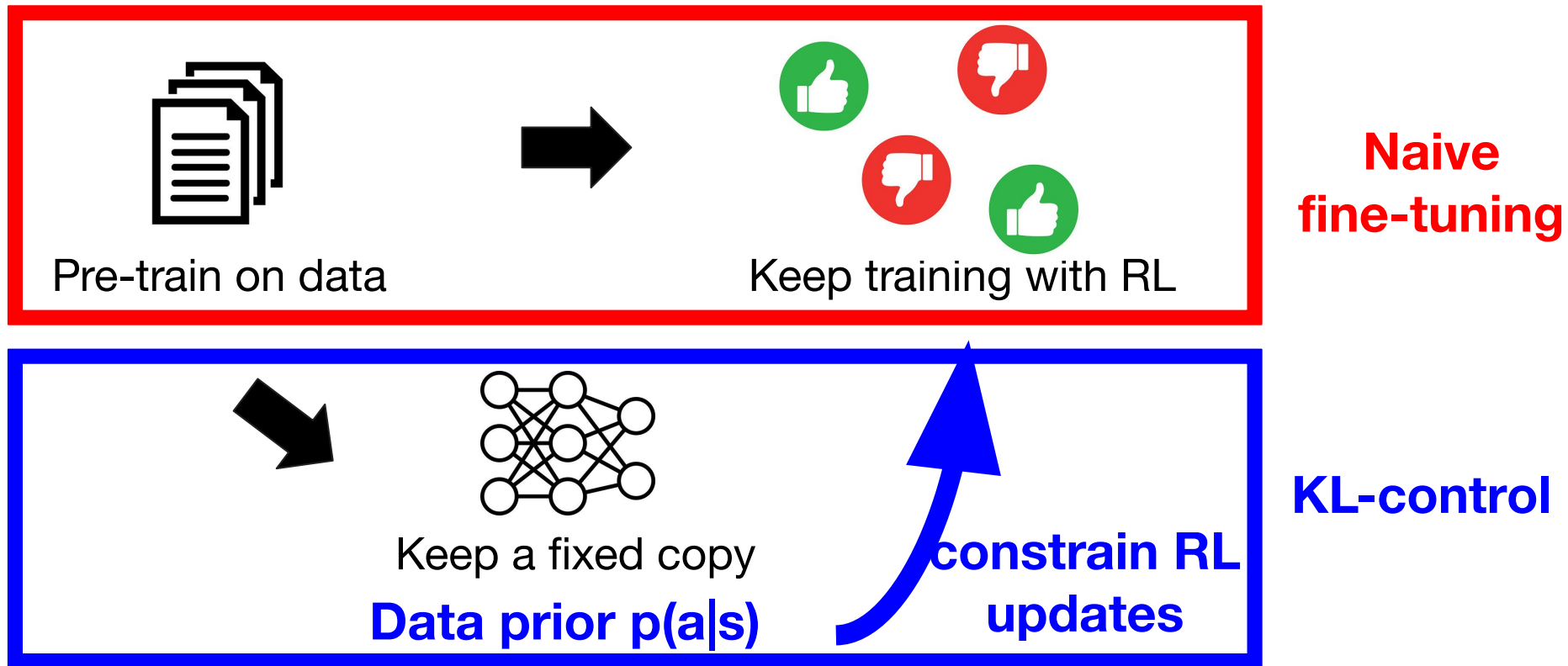
(c) Expert demonstrator



(d) Noisy demonstrator

KL-control from supervised model works as an **effective Offline RL** method, in standard RL envs

# KL-control for language models



# Results: Dialog samples

---

## Naive RL fine-tuning:

[User]: hey, what's up?

[Batch Q]: where did you get??, that sounds so exciting. how do you work for you...  
you want to talk to me?... you want

[User]: yes, i want to talk to you!

[Batch Q]: where did you say to me?... you want to talk to you about you...  
you are so good for you... you are so good...

[User]: thanks, i appreciate that.

[Batch Q]: where did you get??, that sounds so exciting. how do you work for you...  
you want to see you...

- Unrealistic sentences / diverged
- Trivially exploits question reward

## KL-control:

[User]: hey, what's up?

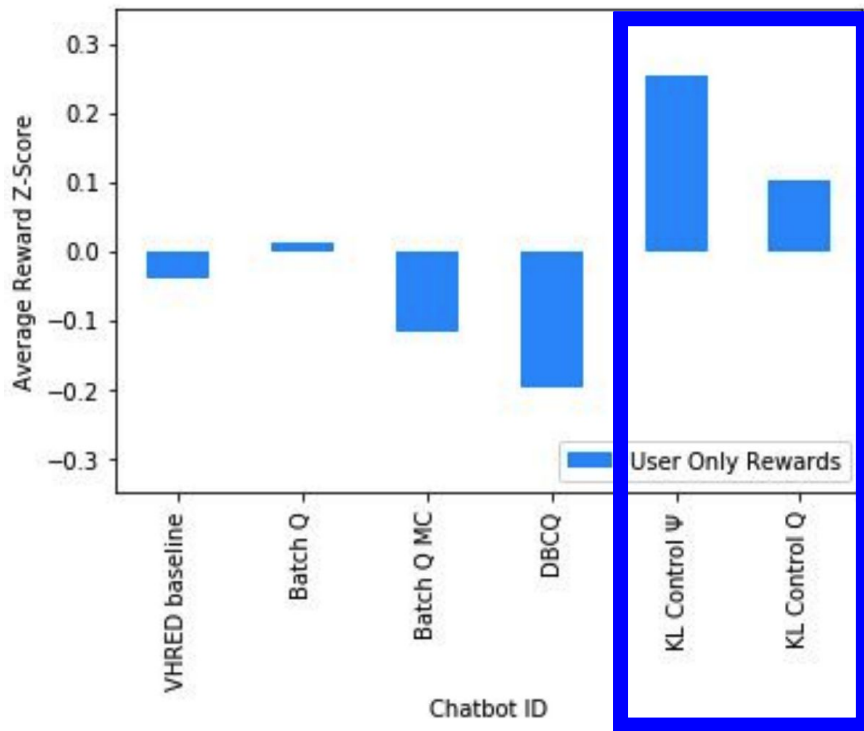
[KL-control]: hey, i hope you have a great day and  
i wish you the best.

[User]: oh thank you that's very sweet of you.

[KL-control]: so, i'm so excited, and i'm so excited  
to meet new people.

- Model remains realistic
- Cheerful, polite, supportive

# Results: Elicited human reward



**KL-control** models **elicit significantly higher reward from humans** when interacting with **novel users** at test time.

# Results: Human ratings / evaluation

Model type	Quality	Fluency	Diversity	Relatedness	Empathy	Total	Votes	Human reward
DBCQ	1.80 ±.41	1.49 ±.29	<b>3.22 ±.57</b>	1.56 ±.25	2.10 ±.37	10.17 ±1.29	-0.07	-0.20
Batch $Q$	1.30 ±.19	2.85 ±.54	1.15 ±.13	1.23 ±.15	2.18 ±.55	8.70 ±0.97	-0.16	0.01
Batch $Q$ + MC	1.53 ±.24	2.15 ±.37	1.60 ±.32	1.53 ±.28	<b>2.58 ±.48</b>	9.38 ±1.31	-0.21	-0.12
KL-control $Q$	<b>2.23 ±.44</b>	<b>2.88 ±.41</b>	2.65 ±.41	<b>2.15 ±.39</b>	2.28 ±.47	<b>12.18 ±1.59</b>	<b>0.09</b>	0.10
KL-control $\Psi$	1.98 ±.44	2.73 ±.45	2.30 ±.42	1.90 ±.37	2.40 ±.44	11.30 ±1.63	0.04	<b>0.25</b>

- **KL-control significantly outperforms RL baselines:**  
 $F(x)=4.781$ ,  $p < .05$



# Results: how reward functions compare

Reward function	Quality	Fluency	Diversity	Relatedness	Empathy	Total	Votes	Human reward
Manual votes								
User laughter								
User Sentiment								
Word Similarity								
USE Similarity								
Bot Question								
Bot Sentiment								
Bot Repetition								

- **Sentiment** leads to highest quality and human reward -- affect is important in good conversation
- **Manual votes** score lower, **validating hypothesis that implicit feedback > explicit**

# Results: Human evaluation

Model type	Quality	Fluency	Diversity	Relatedness	Empathy	Total	Votes	Human reward
DBCQ	1.80 ±.41	1.49 ±.29	<b>3.22 ±.57</b>	1.56 ±.25	2.10 ±.37	10.17 ±1.29	-0.07	-0.20
Batch $Q$	1.30 ±.19	2.85 ±.54	1.15 ±.13	1.23 ±.15	2.18 ±.55	8.70 ±0.97	-0.16	0.01
Batch $Q$ + MC	1.53 ±.24	2.15 ±.37	1.60 ±.32	1.53 ±.28	<b>2.58 ±.48</b>	9.38 ±1.31	-0.21	-0.12
KL-control $Q$	<b>2.23 ±.44</b>	<b>2.88 ±.41</b>	2.65 ±.41	<b>2.15 ±.39</b>	2.28 ±.47	<b>12.18 ±1.59</b>	<b>0.09</b>	0.10
KL-control $\Psi$	1.98 ±.44	2.73 ±.45	2.30 ±.42	1.90 ±.37	2.40 ±.44	11.30 ±1.63	0.04	<b>0.25</b>
VHRED-Baseline	2.65 ±.46	3.83 ±.47	4.05 ±.52	2.43 ±.44	3.08 ±.53	16.03 ±1.93	0.27	-0.04

Quality of the models is still not good enough according to humans

- **Alignment tax?** (i.e. more polite but humans don't like it)
- **Wrong rewards?** Are these rewards not the right ones for good conversations?
- **Or just not enough data?**

# Outline and RLHF history

---

**Fine-tune pre-trained sequence models with RL**

(Jaques et al., 2016)



**Fine-tune language models on human feedback (e.g. sentiment) with offline RL**  
(Jaques et al., 2019)



**Fine-tune language models on sentiment with self-play & RL**  
(Saleh et al., 2019)



**Deep RL from human preferences**  
(Christiano et al., 2017)



**Fine-tuning language models from human preferences**  
(Zeigler et al., 2019)



**Learning to summarize from human feedback**  
(Stiennon et al., 2020)



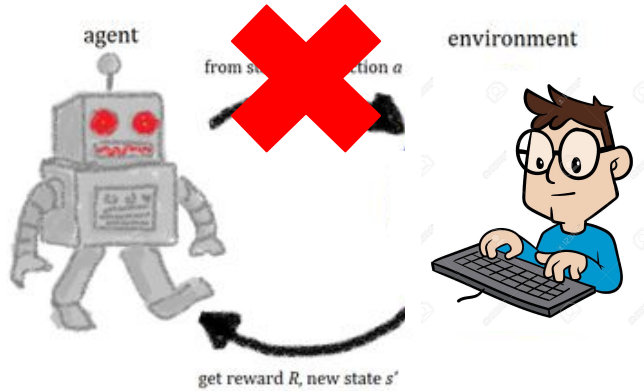
**InstructGPT**  
(Ouyang et al., 2022)



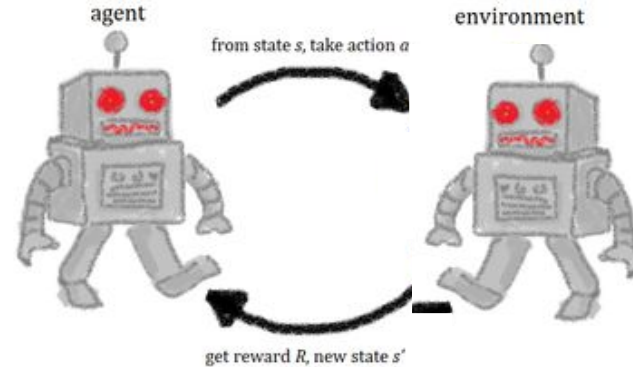
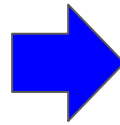
**ChatGPT**

# Follow-up work using Hierarchical RL & self play

- Alignment tax?
- Wrong rewards?
- Or just not enough data?



Offline RL, human data,  
sentiment-based rewards



Online RL, “self-play” (synthetic data),  
sentiment-based rewards

# Follow-up work using Hierarchical RL & self play

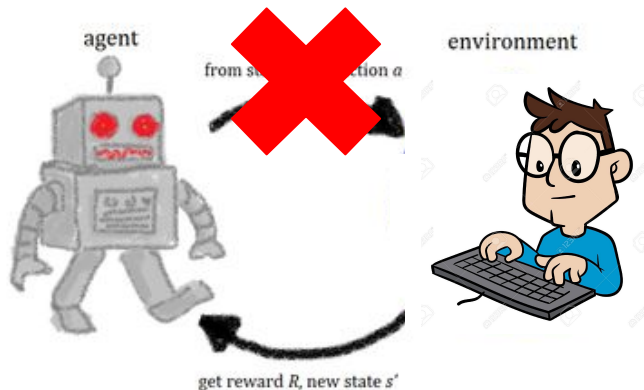
---

Model	Quality	Fluency	Diversity	Contingency	Total	Chat Len.
Batch $\Psi$ (Jaques et al. 2019)	2.17	3.89	3.13	1.98	11.17	11.44
Decoupled VHRL (ablation)	2.46	4.15	3.61	2.02	12.24	12.14
Transformer	2.62	4.17	3.23	2.34	12.36	11.28
REINFORCE	2.89	4.47	3.67	<b>2.80</b>	13.84	11.60
VHRED	2.84	4.53	<b>4.43</b>	2.47	14.27	10.94
VHRL (ours)	<b>2.91</b>	<b>4.65</b>	4.26	2.67	<b>14.49</b>	<b>12.84</b>

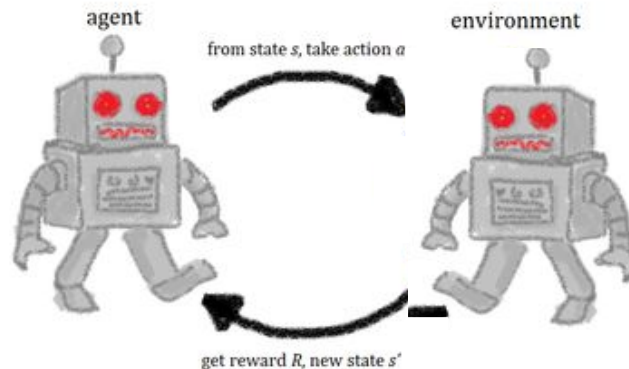
- Our method, VHRL, outperforms language model baselines and previous Offline RL work in human ratings of conversation quality

# Follow-up work using Hierarchical RL & self play

- ~~Alignment tax?~~
- ~~Wrong rewards?~~
- Or just **not enough data?**

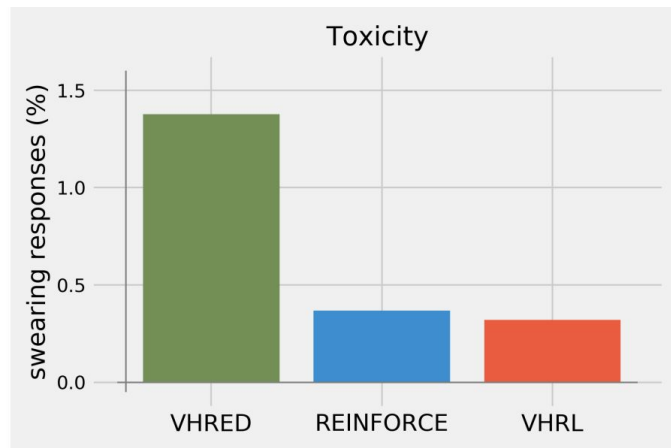
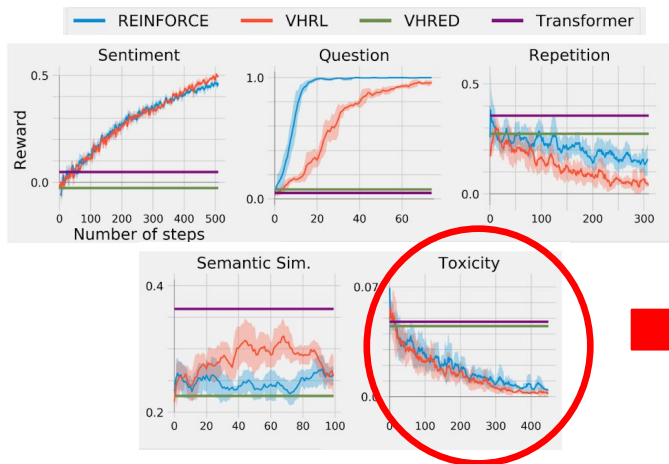


Offline RL, human data,  
sentiment-based rewards



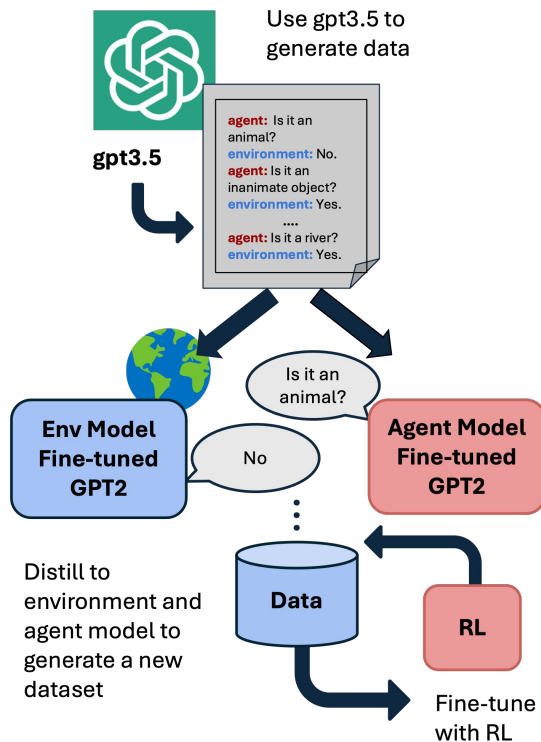
Online RL, “self-play” (synthetic data),  
sentiment-based rewards

# Using RL to reduce toxicity



Can use RL to **reduce toxicity** by using the output of a toxicity classifier as a negative reward

# Further work on building reward simulators



- To generate data for conversational tasks, LLMs are used as “simulators” for the task.
- Simulators can be used to generate offline data, to provide a “simulation environment” for evaluation, to perform online training, and to compute rewards.
- For text-games, we use engines as simulators to generate near-optimal data and dilute the policy with suboptimal data with inferior policies

Task	20Qs	Guess	Car	Maze	Text-Nav	Wordle	Chess	Endgames
Size	100k	100k	19k	1.24k	2.5k	1m	625k	97.756k
avg length	-14.9	-18.8	16.5	19.7	12.2	4.82	46.7	11.9
std length	4.38	4.57	3.61	24.5	8.77	1.27	18.16	12.0
success rate	0.31	0.53	0.53	0.11	0.26	0.70	0.60	0.59
avg return	-17.3	-18.8	0.562	-19.7	0.258	-4.12	0.210	0.586
std return	2.56	4.12	0.422	24.5	0.424	1.59	0.970	0.492

Table 1: Statistics for all tasks in LMRL-Gym. Size represents the number of trajectories, the average length is the average length of trajectories in the dataset where the unit is a response from the agent. The success rate is the proportion of trajectories that reach the objective. Finally, the reward functions for each task are defined in Appendix D.



– Chess –

**environment:** rnbqkbnr/pppppppp/8/8/8/8/PPPPPPPP/RNBQKBNR w KQkq - 0 1

**agent:** N f 3

– Guess My City –

City: Jakarta, Indonesia

**agent:** What is your favorite traditional dish from your hometown?

**environment:** My favorite traditional dish from my hometown is nasi goreng, a delicious fried rice dish that is often served with chicken, shrimp, and various spices.

....

**agent:** Is your hometown located near a coastline or body of water?

**environment:** Yes, my city is located on the coast and is surrounded by water.

**agent:** Is the city you are from Jakarta?

**environment:** Yes, the city I am from is Jakarta.

– Partially Observed Maze –

**environment:** There are walls to your left, below you

**agent:** move right

**environment:** There are walls above you, below you

**agent:** move left

**environment:** There are walls to your left, below you

**agent:** move up

**environment:** There are walls to your right, to your left

**agent:** move up

– Twenty Questions –

**agent:** Is it an animal?

**environment:** No.

**agent:** Is it an inanimate object?

**environment:** Yes.

**agent:** Is it found outside?

**environment:** Yes.

....

**agent:** Is it a river?

**environment:** Yes.

– Text-Nav –

**environment:**

-= Living Room =-

You are carrying nothing on you right now. You make out a table. You wonder idly who left that here. The table is standard. On the table you see a stale food. There are unblocked exits to the east, south and west.

Available actions: ['pick up stale food', 'go east', 'go south', 'go west']

**agent:** pick up stale food

– Fully observed maze –

**environment:** The goal is at position 8, 6. Your current position is at position 5, 3. There are walls to your left, below you.

**agent:** move up

**environment:** The goal is at position 8, 6. Your current position is at position 6, 3. There are walls to your left.

**agent:** move up

– Car Dealer –

**agent:** Good morning! Welcome to our car dealership. How can I assist you today?

**environment:** Hi. I'm looking for a Honda SUV with a backup camera. Can you show me what you have?

....

**agent:** Of course! We have a great selection of Honda SUVs with backup cameras. Let me show you our top models. This one here is the Honda CR-V. It has a powerful engine, spacious interior, and advanced safety features. Would you like me to go over the details?

**environment:** No, thank you. I already know about the features. Can you tell me the price?

**agent:** Alright, thanks for your time. I'll think about it.

# Using LLMs for negotiation

---

- Negotiation is an interesting cooperation problem
  - **Cooperative:** Negotiating agents need to act in the interest of a human principle
  - **Non-cooperative:** While defending against adversarial cooperation partners
- Interesting LLM + RL problem
  - **Objective, quantitative:** who got the best price?
  - Can **evaluate automatically** with other agents



November 2023

---

## LMRL Gym: Benchmarks for Multi-Turn Reinforcement Learning with Language Models

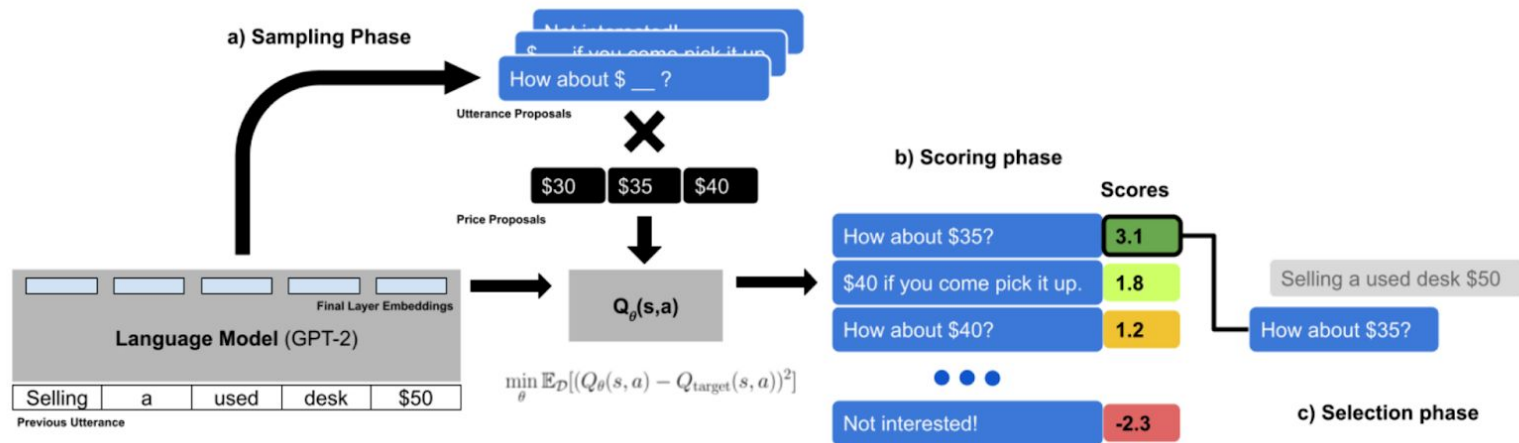
Marwa Abdulhai<sup>1</sup>, Isadora White<sup>1</sup>, Charlie Snell<sup>1</sup>, Charles Sun<sup>1</sup>, Joey Hong<sup>1</sup>, Yuexiang Zhai<sup>1</sup>,  
Kelvin Xu<sup>2</sup>, and Sergey Levine<sup>1, 1</sup>

<sup>1</sup>UC Berkeley, <sup>2</sup>Google DeepMind

---

# Training LLMs for negotiation with Offline RL

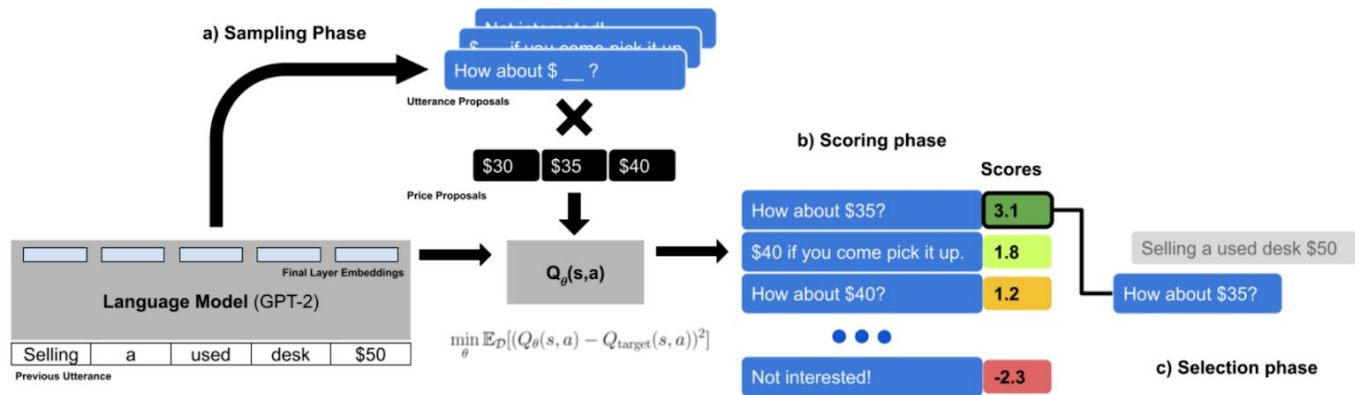
- Verma et al. (2022) investigate different **offline RL** methods for negotiation
  - Sample text from language model (GPT-2)
  - Rank candidates with Q-function



# Training LLMs for negotiation with Offline RL

- **Why Offline RL?**

- Lots of existing data (i.e. on negotiations)
- Easier from an infrastructure perspective: similar to supervised learning as you don't need RLHF pipeline of collecting data from humans



# Using LLMs for negotiation

Method	vs Rule-based		vs Stingy		vs Utility	
	Acc%	Revenue	Acc%	Revenue	Acc%	Revenue
CHAI-prop	61.5	0.48 ± 0.39	57.5	0.39 ± 0.35	99.0	0.70 ± 0.17
CHAI-CQL	74.0	0.51 ± 0.33	77.5	0.49 ± 0.30	98.0	0.70 ± 0.19
CHAI-BRAC	62.0	0.52 ± 0.41	47.0	0.38 ± 0.41	99.0	0.71 ± 0.17
Language Model	48.5	0.29 ± 0.32	51.5	0.27 ± 0.28	20.5	0.14 ± 0.28
He et al. (2018) (Utility)	1.0	0.01 ± 0.10	0.0	0.00 ± 0.00	11.0	0.07 ± 0.22
He et al. (2018) (Fairness)	84.0	0.70 ± 0.32	80.0	0.59 ± 0.31	100.0	0.72 ± 0.15
He et al. (2018) (Length)	53.0	0.46 ± 0.43	49.0	0.37 ± 0.38	100.0	0.72 ± 0.16
Lewis et al. (2017) (RL)	83.5	0.17 ± 0.24	83.0	0.19 ± 0.25	64.5	0.46 ± 0.37
Lewis et al. (2017) (SL)	38.5	0.17 ± 0.27	46.5	0.21 ± 0.27	18.0	0.13 ± 0.28

Method	vs Fairness		vs Length		Overall (mean)	
	Acc%	Revenue	Acc%	Revenue	Acc%	Revenue
CHAI-prop	99.0	0.90 ± 0.15	92.5	0.79 ± 0.27	81.9	<b>0.65 ± 0.34</b>
CHAI-CQL	99.5	0.87 ± 0.14	94.5	0.79 ± 0.24	<b>88.7</b>	<b>0.67 ± 0.29</b>
CHAI-BRAC	100.0	0.85 ± 0.03	91.0	0.76 ± 0.25	79.8	<b>0.65 ± 0.34</b>
Language Model	25.5	0.19 ± 0.35	18.5	0.14 ± 0.32	32.9	0.21 ± 0.31
He et al. (2018) (Utility)	100.0	1.00 ± 0.00	100.0	1.00 ± 0.00	42.4	0.42 ± 0.49
He et al. (2018) (Fairness)	0.0	0.00 ± 0.00	100.0	0.70 ± 0.16	72.8	0.54 ± 0.35
He et al. (2018) (Length)	100.0	1.00 ± 0.00	100.0	0.78 ± 0.18	80.4	<b>0.66 ± 0.36</b>
Lewis et al. (2017) (RL)	88.0	0.26 ± 0.34	71.5	0.31 ± 0.36	78.1	0.28 ± 0.33
Lewis et al. (2017) (SL)	60.0	0.48 ± 0.46	53.0	0.42 ± 0.46	43.2	0.28 ± 0.39

## Automatic evaluations:

- No clear winner in revenue
- Best baseline uses manually input dialog acts

# Using LLMs for negotiation

---

## Human evaluations: clear winner

Metric	Fluency	Coherency	On-Topic	Human-Likeness	Total
CHAI-prop	<b>4.31 ± 0.97</b>	<b>3.91 ± 1.17</b>	<b>4.16 ± 0.99</b>	<b>3.47 ± 1.27</b>	<b>15.84 ± 3.86</b>
He et al. (2018) (Utility)	3.56 ± 1.34	2.47 ± 1.39	3.09 ± 1.40	2.13 ± 1.13	11.25 ± 4.50
Lang. Model	4.06 ± 1.11	2.66 ± 1.36	3.63 ± 1.18	2.50 ± 1.10	12.84 ± 3.66

# Outline and RLHF history

---

**Fine-tune pre-trained sequence models with RL**

(Jaques et al., 2016)



**Fine-tune language models on human feedback (e.g. sentiment) with offline RL**  
(Jaques et al., 2019)



**Fine-tune language models on sentiment with self-play & RL**  
(Saleh et al., 2019)



**Deep RL from human preferences**  
(Christiano et al., 2017)



**Fine-tuning language models from human preferences**  
(Zeigler et al., 2019)



**Learning to summarize from human feedback**  
(Stiennon et al., 2020)



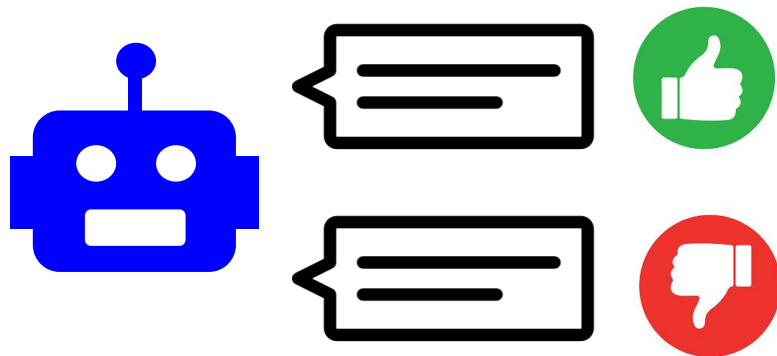
**InstructGPT**  
(Ouyang et al., 2022)



**ChatGPT**

# Training a ~~language~~ model with human feedback

---

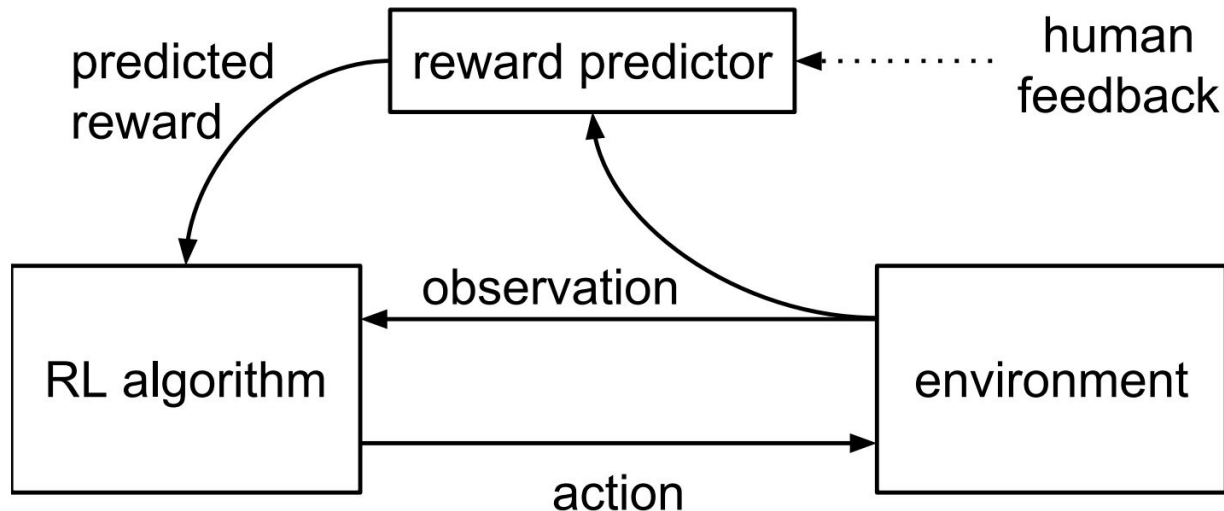


What kind of feedback should we use?



**Issue:** RL models are sample hungry, but human feedback is really expensive

**Solution:** train a **reward model** that you can query as much as you need



**Why do we expect the reward model to make better use of the data?**

- Supervised learning on data works better than offline RL on data

**Issue:** Humans are really bad at giving things absolute ratings

**Solution:** Ask humans to rate which of two trajectory segments they prefer (humans are better at comparisons)

We say that preferences  $\succ$  are *generated by* a reward function<sup>1</sup>  $r : \mathcal{O} \times \mathcal{A} \rightarrow \mathbb{R}$  if

$$\left( (o_0^1, a_0^1), \dots, (o_{k-1}^1, a_{k-1}^1) \right) \succ \left( (o_0^2, a_0^2), \dots, (o_{k-1}^2, a_{k-1}^2) \right)$$

whenever

$$r(o_0^1, a_0^1) + \dots + r(o_{k-1}^1, a_{k-1}^1) > r(o_0^2, a_0^2) + \dots + r(o_{k-1}^2, a_{k-1}^2).$$

cumulated rewards of trajectory are higher than the other trajectory

# How to learn the reward function:

---

Assume probability of preferring segment  $\sigma^1$  depends **exponentially** on value of latent **reward**  $\hat{r}$  **summed** over the length of the segment

$$\hat{P}[\sigma^1 \succ \sigma^2] = \frac{\exp \sum \hat{r}(o_t^1, a_t^1)}{\exp \sum \hat{r}(o_t^1, a_t^1) + \exp \sum \hat{r}(o_t^2, a_t^2)}.$$

Learn  $\hat{r}$  by minimizing **cross-entropy between predictions and human labels**

$$\text{loss}(\hat{r}) = - \sum_{(\sigma^1, \sigma^2, \mu) \in \mathcal{D}} \mu(1) \log \hat{P}[\sigma^1 \succ \sigma^2] + \mu(2) \log \hat{P}[\sigma^2 \succ \sigma^1].$$

where  $\mu = [1, 0]$  if human preferred  $\sigma^1$ ,  $[0.5, 0.5]$  if human thinks segments are equal

# Other tricks

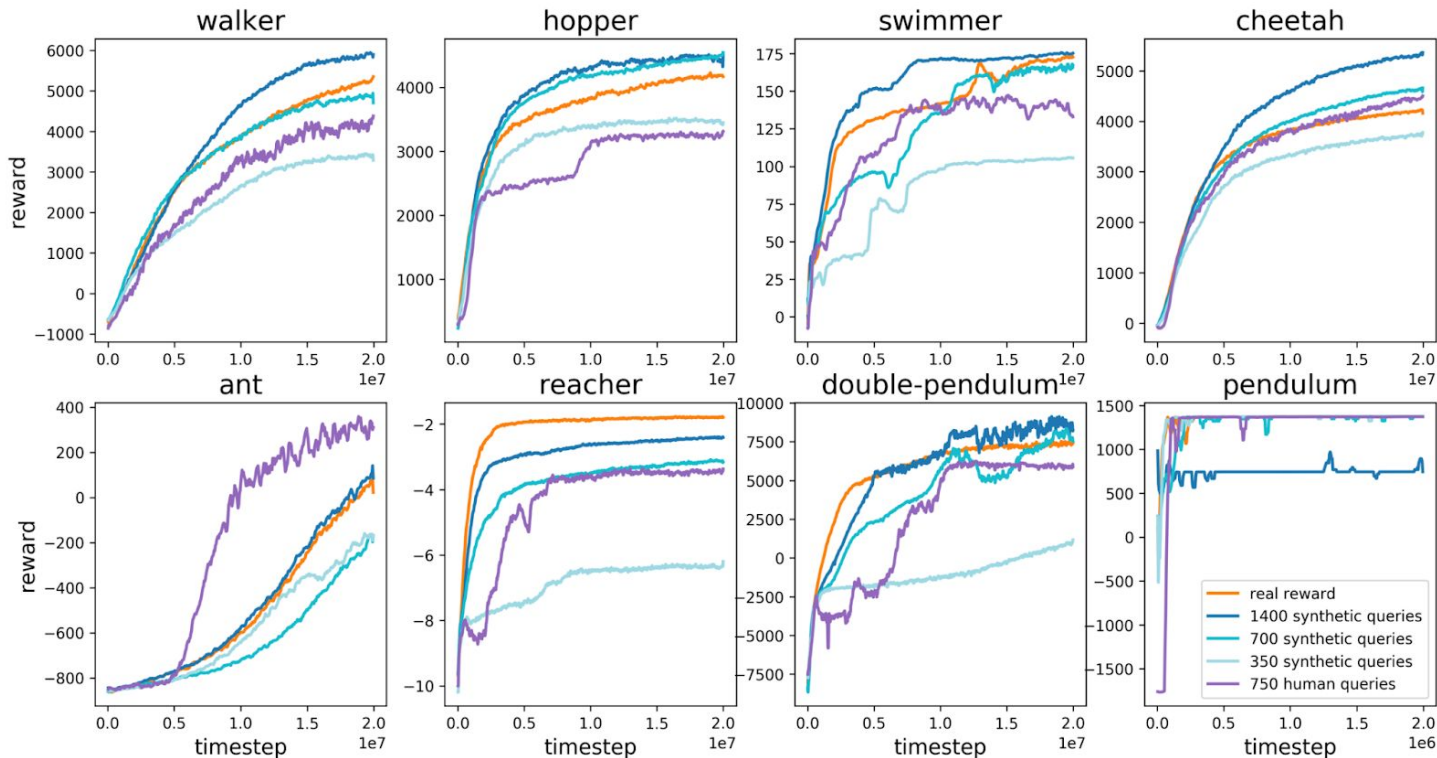
---

- Use **ensemble of reward models**  $\hat{r}$
- Assume there is a 10% chance that the human responds uniformly at random

*Why? Aren't you already doing something like MaxEnt IRL?*

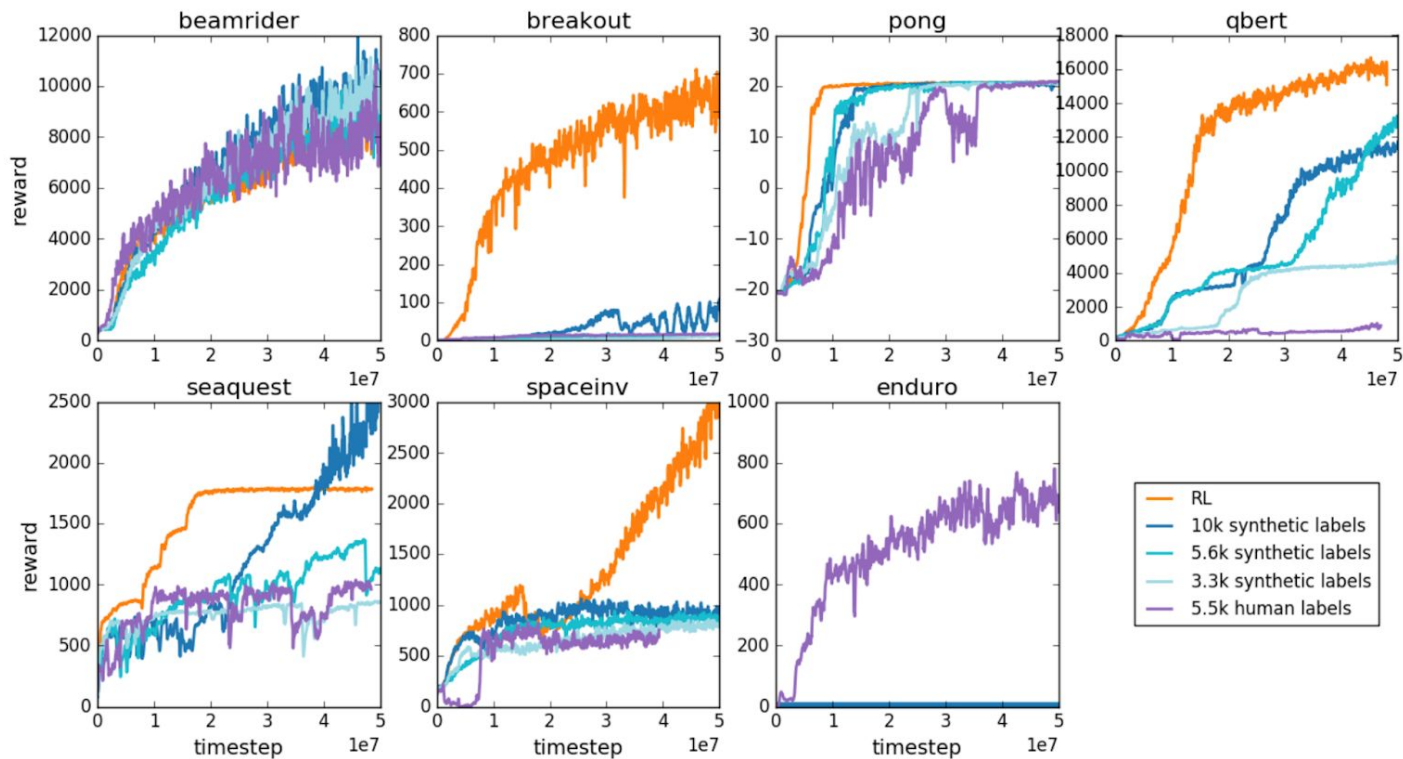
→ Humans have constant probability of **mis-click**, which doesn't decay to zero as differences in reward become large

# Results: Mujoco



- Hard to get enough **human data** to learn more effectively than **normal reward**
- Synthetic **preference reward** > **normal reward**

# Results: Atari

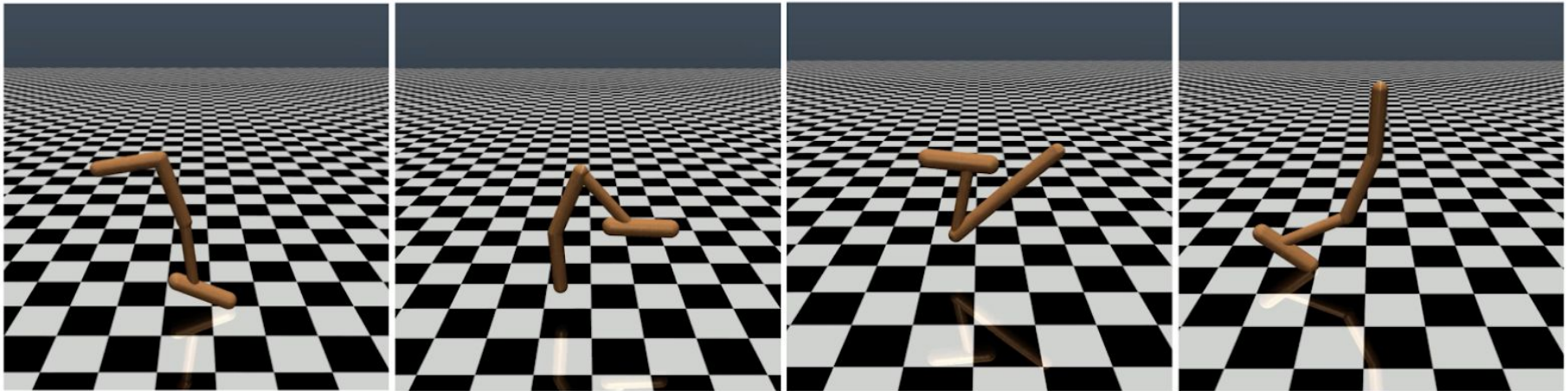


- **Human data** sometimes better than **normal reward**
- **Preference reward** sometimes better than **normal reward**

# Results: most importantly...

---

Can learn to do skills that have **no existing reward function** with only a **small amount of human labels**



# Outline and RLHF history

---

**Fine-tune pre-trained sequence models with RL**

(Jaques et al., 2016)



**Fine-tune language models on human feedback (e.g. sentiment) with offline RL**  
(Jaques et al., 2019)



**Fine-tune language models on sentiment with self-play & RL**  
(Saleh et al., 2019)



**Deep RL from human preferences**  
(Christiano et al., 2017)



**Fine-tuning language models from human preferences**  
(Zeigler et al., 2019)



**Learning to summarize from human feedback**  
(Stiennon et al., 2020)



**InstructGPT**  
(Ouyang et al., 2022)

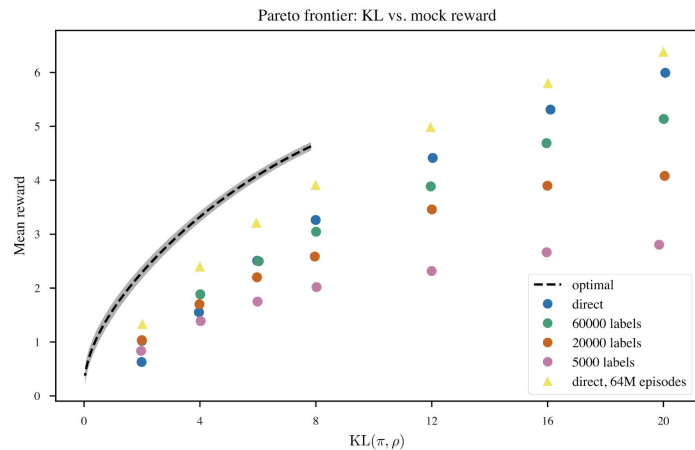
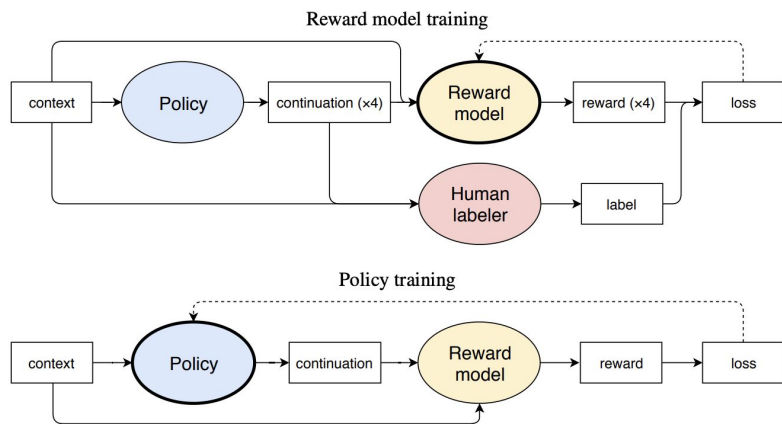


**ChatGPT**



# Bringing it all together to train LMs

- Use **KL-control** technique to fine-tune the LM on rewards
- Use **reward model** technique to better scale human feedback



- **Results:** high ROUGE scores for summarization, but a lot of direct copying

# Outline and RLHF history

---

**Fine-tune pre-trained sequence models with RL**

(Jaques et al., 2016)



**Fine-tune language models on human feedback (e.g. sentiment) with offline RL**  
(Jaques et al., 2019)



**Fine-tune language models on sentiment with self-play & RL**  
(Saleh et al., 2019)



**Deep RL from human preferences**  
(Christiano et al., 2017)



**Fine-tuning language models from human preferences**  
(Zeigler et al., 2019)



**Learning to summarize from human feedback**  
(Stiennon et al., 2020)



**InstructGPT**  
(Ouyang et al., 2022)



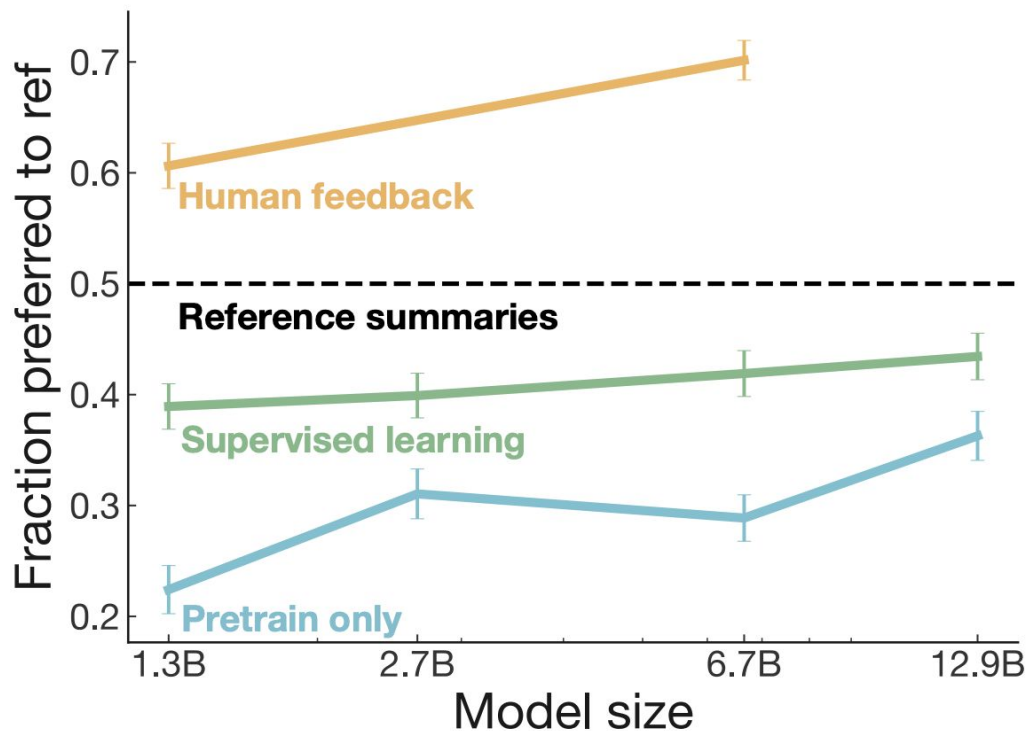
**ChatGPT**

# Learning to summarize from human feedback

---

- Use **KL-control** and **reward model**
- **What else is needed to get RLHF to work?**
  - “[Our] previous work [...] reported “a mismatch between the notion of quality we wanted our model to learn, and what the humans labelers actually evaluated”, leading to model-generated summaries that were high-quality according to the labelers, but fairly low-quality according to the researchers”
  - **Pay way more attention to how to collect human feedback:**
    - **Offline:** alternate between collecting large batches of human labels and re-training our models on the cumulative collected data
    - **High touch approach:** screen labelers, onboard them, answer questions in a shared chat room, provide regular feedback
      - Achieve better researcher-labeler agreement (77%)

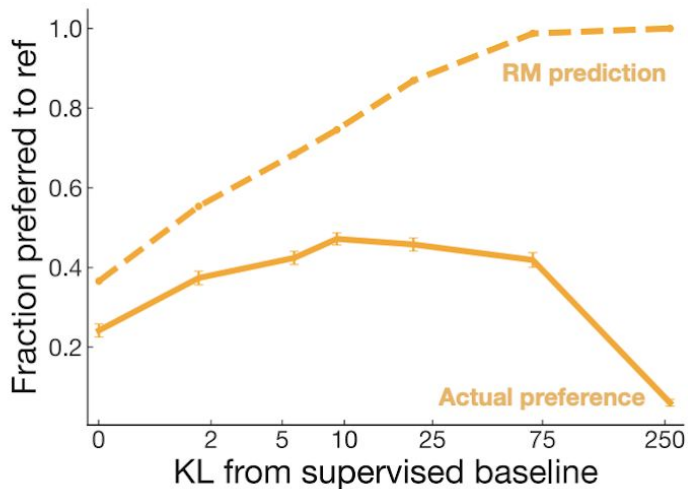
# Results: Learning to summarize from HF



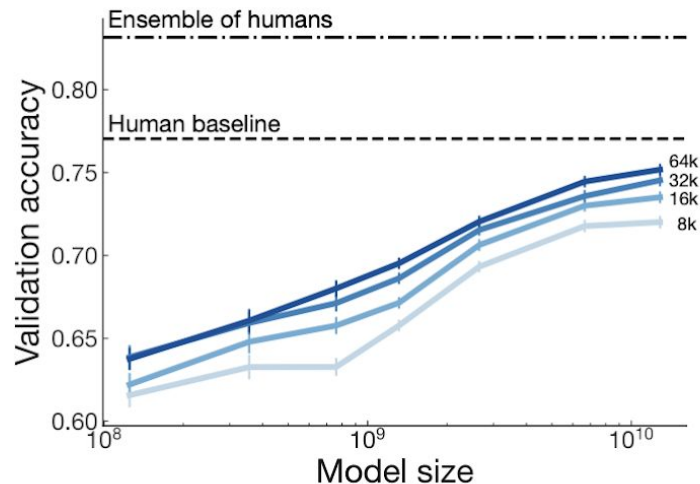
**RLHF summaries preferred over:**

- **Human** summaries
- Summaries from supervised models **10x the size**

# Results: Learning to summarize from HF



Can't train on reward model too long, or performance decreases



Reward models are not that accurate

# RLHF for summarization: the full picture

## 1 Collect human feedback

A Reddit post is sampled from the Reddit TL;DR dataset.



Various policies are used to sample a set of summaries.



Two summaries are selected for evaluation.



A human judges which is a better summary of the post.



"j is better than k"

## 2 Train reward model

One post with two summaries judged by a human are fed to the reward model.



The reward model calculates a reward  $r$  for each summary.



$r_j$

$r_k$

The loss is calculated based on the rewards and human label, and is used to update the reward model.

$$\text{loss} = \log(\sigma(r_j - r_k))$$

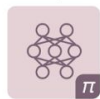
"j is better than k"

## 3 Train policy with PPO

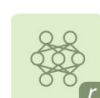
A new post is sampled from the dataset.



The policy  $\pi$  generates a summary for the post.



The reward model calculates a reward for the summary.



The reward is used to update the policy via PPO.



# Outline and RLHF history

---

**Fine-tune pre-trained sequence models with RL**

(Jaques et al., 2016)



**Fine-tune language models on human feedback (e.g. sentiment) with offline RL**  
(Jaques et al., 2019)



**Fine-tune language models on sentiment with self-play & RL**  
(Saleh et al., 2019)



**Deep RL from human preferences**  
(Christiano et al., 2017)



**Fine-tuning language models from human preferences**  
(Zeigler et al., 2019)



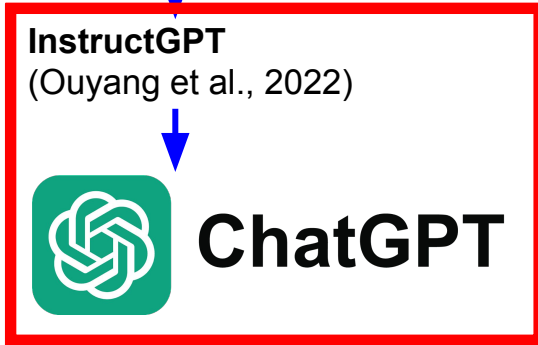
**Learning to summarize from human feedback**  
(Stiennon et al., 2020)



**InstructGPT**  
(Ouyang et al., 2022)



**ChatGPT**



# InstructGPT (and ChatGPT)

Step 1

Collect demonstration data, and train a supervised policy.

A prompt is sampled from our prompt dataset.



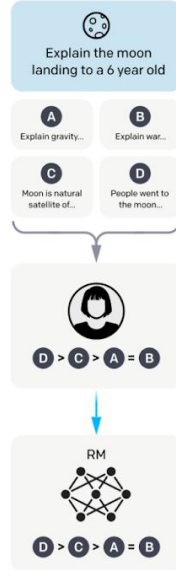
A labeler demonstrates the desired output behavior.

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.



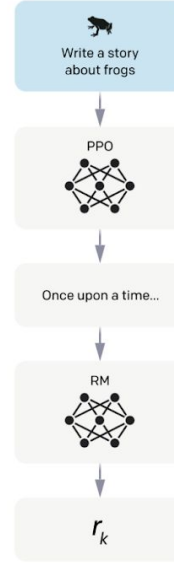
A labeler ranks the outputs from best to worst.

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.

Can you spot the differences from the previous work?

Use human rewritten responses for supervised fine-tuning (SFT)

Also mix supervised and PPO updates to keep closer to LLM distribution of text



# RLHF for open ended chat

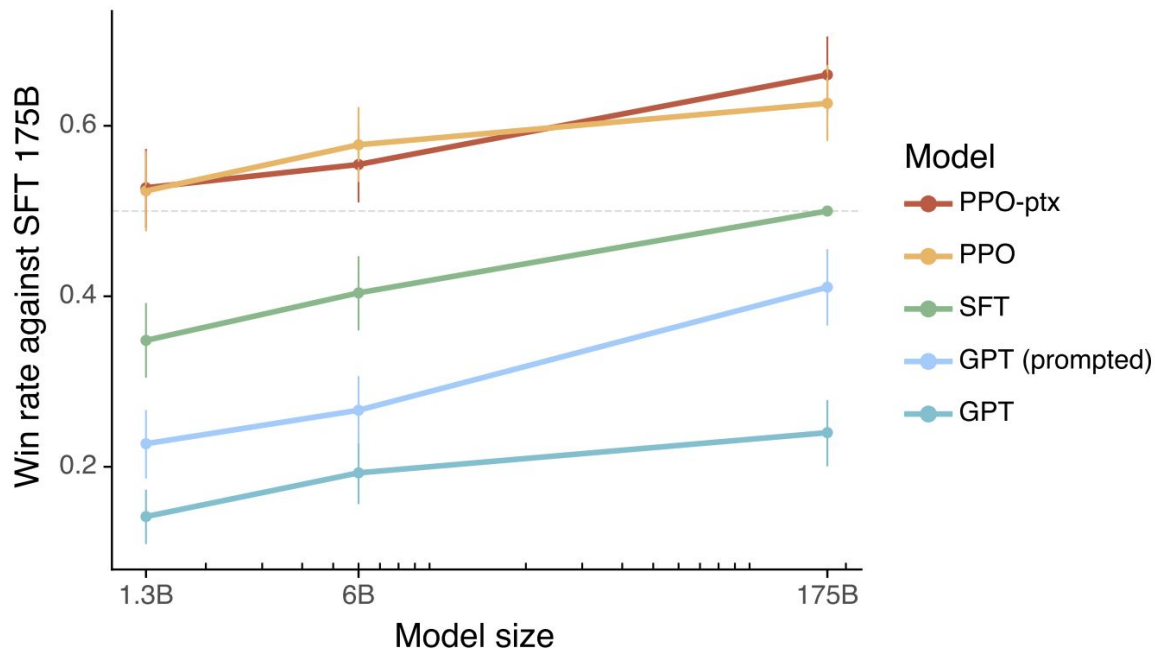
---

- No longer just doing summarization, doing open-ended dialog
- Train to **increase alignment**:
  - **Helpful**: follow user instruction well
  - **Honest**: truthful
  - **Harmless**: avoid bias, toxicity

## What if these values come in conflict?

→ They opted to have the models be **helpful over harmless**

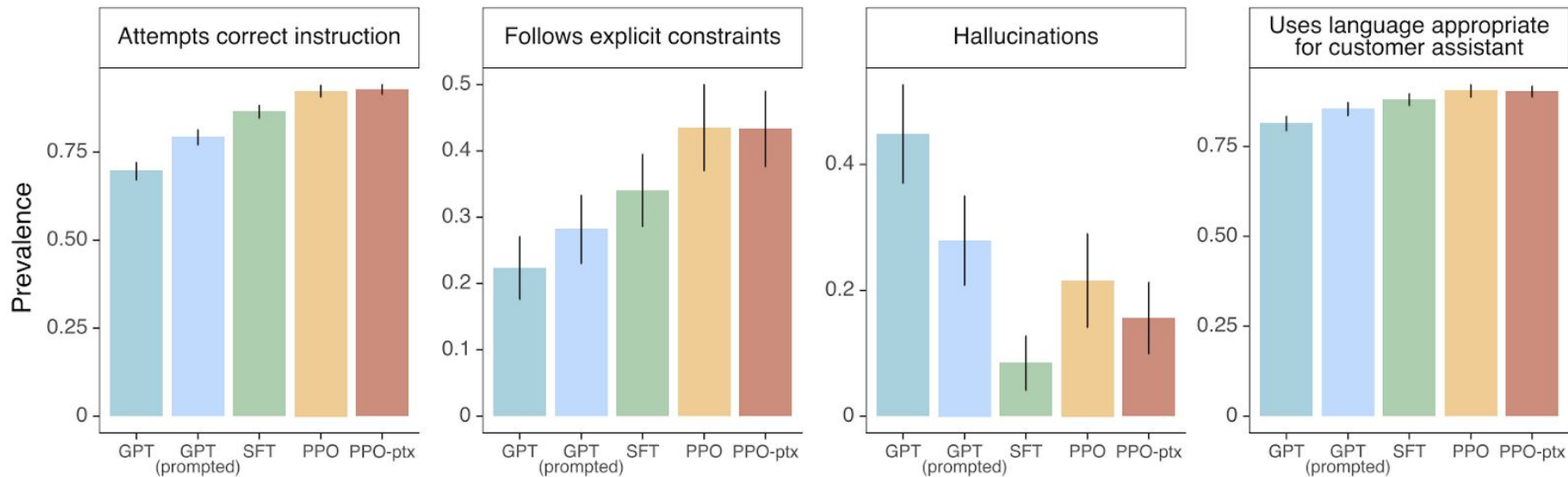
# Results: InstructGPT



**1.5B parameter RLHF model does better than the 175B parameter GPT3 in human evaluations**

→ **Authors conclude investing in fine-tuning is more cost effective than training ever-larger models**

# Results: InstructGPT



RLHF hallucinates less, is less toxic → **more aligned**

# But who are we aligning to?

---

- Labelers are **English-speaking people** living in the United States or Southeast Asia hired via Upwork or Scale AI
- Values / preferences for labeling process decided by **OpenAI researchers**
- Training data from **OpenAI API customers**
  - Not necessarily interested in human well-being.
  - May want to maximize user attention of customers

# But who are we aligning to?

## A Roadmap to Pluralistic Alignment

Taylor Sorensen<sup>1</sup> Jared Moore<sup>2</sup> Jillian Fisher<sup>1,3</sup> Mitchell Gordon<sup>1,4</sup> Niloofar Mireshghallah<sup>1</sup>  
Christopher Michael Rytting<sup>1</sup> Andre Ye<sup>1</sup> Liwei Jiang<sup>1,5</sup> Ximing Lu<sup>1</sup> Nouha Dziri<sup>5</sup> Tim Althoff<sup>1</sup>  
Yejin Choi<sup>1,5</sup>

- Distributional pluralism would be successfully modeling different, potentially diverging preferences
- This paper shows existing RLHF techniques may actually reduce distributional pluralism



Is it ok for governments to moderate the social media content available to public?

Pluralistic Human Values



Overton



Many think that it's not okay for the government to moderate content as it endangers liberty, while others deem it acceptable for prevention of terrorism. A few, on the other hand, think it's necessary for sovereignty.

Steerable



It is ok for the government to moderate content for terrorism and threats.  
or  
It is not ok to moderate any content as it endangers liberty.  
or  
It is ok for the government to moderate content that endangers its sovereignty.

Distributional



Figure 1. Three kinds of pluralism in models.

# Latest trends

---

**Fine-tune pre-trained sequence models with RL**

(Jaques et al., 2016)



**Fine-tune language models on human feedback (e.g. sentiment) with offline RL**  
(Jaques et al., 2019)



**Fine-tune language models on sentiment with self-play & RL**  
(Saleh et al., 2019)



**Deep RL from human preferences**  
(Christiano et al., 2017)



**Fine-tuning language models from human preferences**  
(Zeigler et al., 2019)



**Learning to summarize from human feedback**  
(Stiennon et al., 2020)

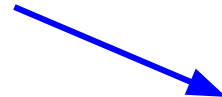


**InstructGPT**  
(Ouyang et al., 2022)



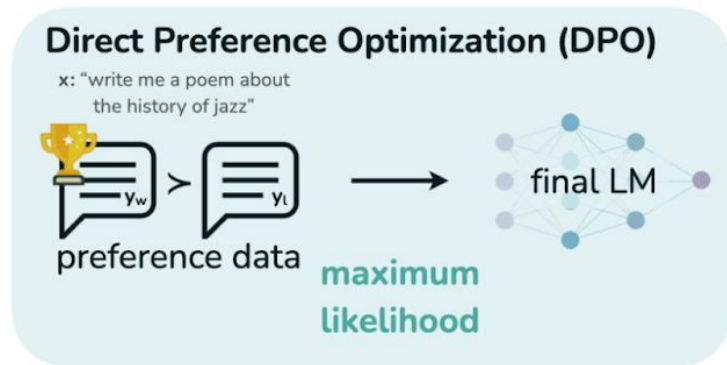
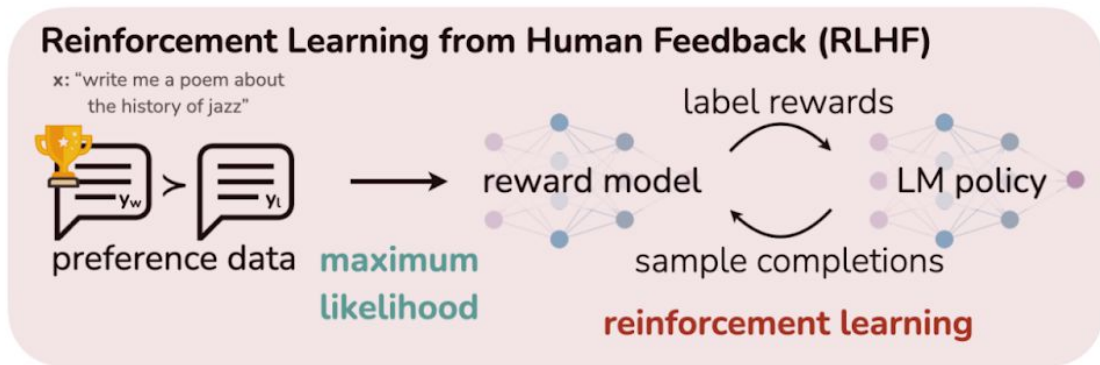
**ChatGPT**

**What next?**



# Direct Preference Optimization (Rafailov et al., 2023)

- What if we don't need to learn a reward model at all?



$$\mathcal{L}_{\text{DPO}}(\pi_{\theta}; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[ \log \sigma \left( \beta \log \frac{\pi_{\theta}(y_w | x)}{\pi_{\text{ref}}(y_w | x)} - \beta \log \frac{\pi_{\theta}(y_l | x)}{\pi_{\text{ref}}(y_l | x)} \right) \right]$$

# Distributional Preference Learning (Siththaranjan et al., 2023)

- Vanilla RLHF overrules minority preferences
- What if we could model the distribution of preferences, detect when users diverge?

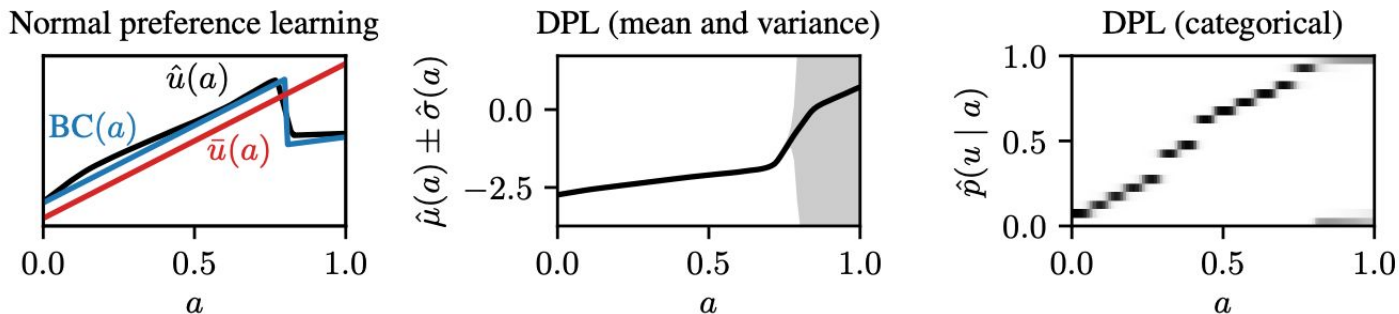
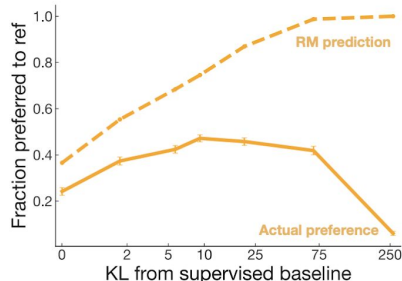


Figure 4: The results of our experiments with synthetic data. We find that the utility estimated by normal preference learning agrees closely with the Borda count, as our theory suggests. Furthermore, DPL successfully identify alternatives where hidden context has a significant effect.



# Iterated Data Smoothing (Zhu et al., 2024)



- InstructGPT tells us we have to do early stopping with RLHF, or we will overfit to the reward model

- The issue is that if we only see one comparison of  $y_1$  and  $y_2$ , the BTL reward loss could blow up (go to infinity)
  - All comparisons are rarely seen given very high dimensional data
- Iterated Data Smoothing relabels the RLHF data using the learned reward model after one episode of training
  - Trust rarely seen data less
- Leads to SOTA open-source RLHF results ([Starling-7B](#))

# Latest trends

---

**Fine-tune pre-trained sequence models with RL**

(Jaques et al., 2016)



**Fine-tune language models on human feedback (e.g. sentiment) with offline RL**  
(Jaques et al., 2019)



**Fine-tune language models on sentiment with self-play & RL**  
(Saleh et al., 2019)



**Deep RL from human preferences**  
(Christiano et al., 2017)



**Fine-tuning language models from human preferences**  
(Zeigler et al., 2019)



**Learning to summarize from human feedback**  
(Stiennon et al., 2020)



**InstructGPT**  
(Ouyang et al., 2022)



**ChatGPT**



**What next?**



**<Your cool idea here!>**

# Questions?

Contact: [marwa\\_abdulhai@berkeley.edu](mailto:marwa_abdulhai@berkeley.edu)