# RLHF

Kellen Kanarios

# Why RL?



**Agent**

From **state s$_t$,** take **action a$_t$**

**Environment**

**Policy** $\pi$
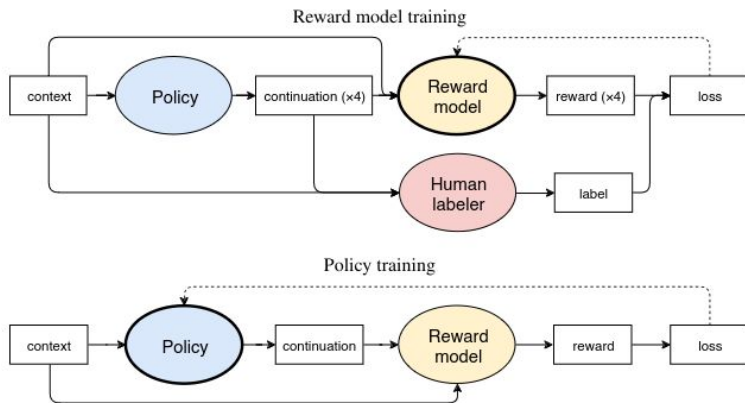**(a|s)**

Get **reward r$_t$,** **new state s$_{t+1}$**

Input sentence

The fish lived in the

Word Embeddings

The fish lived in  the

Model

Transformer decoder

Output

sea 22%
blue 17%
river 16%
ocean 14%
when 0.01%
under 0.01%
bicycle 0.001%

Choose next token: "blue"

The fish lived in the blue

The fish lived in  the blue

Transformer decoder

sea 32%
river 17%
ocean 16%
lake 14%
when 0.01%
under 0.01%
bicycle 0.001%

Choose next token: "sea"

# How RL?

Daniel M. Ziegler*   Nisan Stiennon*   Jeffrey Wu   Tom B. Brown
Alec Radford   Dario Amodei   Paul Christiano   Geoffrey Irving
OpenAI
{dmz,nisan,jeffwu,tom,alec,damodei,paul,irving}@openai.com

## Training language models to follow instructions with human feedback

Long Ouyang*   Jeff Wu*   Xu Jiang*   Diogo Almeida*   Carroll L. Wainwright*

Pamela Mishkin*   Chong Zhang   Sandhini Agarwal   Katarina Slama   Alex Ray

John Schulman   Jacob Hilton   Fraser Kelton   Luke Miller   Maddie Simens

Amanda Askell†   Peter Welinder   Paul Christiano*†

Jan Leike*   Ryan Lowe*

# Learning a reward

**Issue**: Humans cannot score outputs well

**Solution**: Humans rank outputs

$$\text{loss}\,(\theta) = -\frac{1}{\binom{K}{2}} E_{(x,y_w,y_l)\sim D}\left[\log\left(\sigma\left(r_\theta\left(x, y_w\right) - r_\theta\left(x, y_l\right)\right)\right)\right]$$
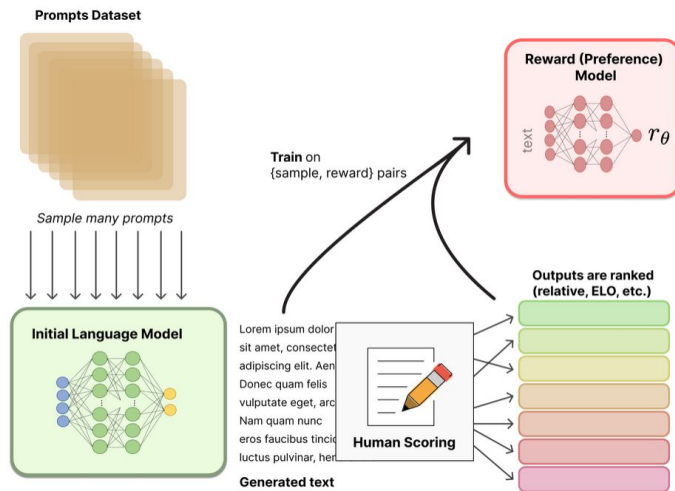
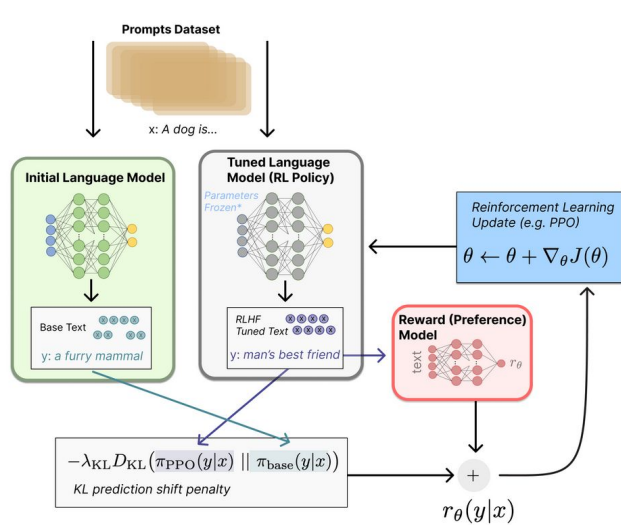Completions are just **sampled** $K$ times from LLM.

**Trick**: Treat $\binom{K}{2}$ comparisons as single batch element.

| | | |
|---|---|---|
| $r_\theta\left(x, y\right)$ | - | Scalar output f reward model for prompt $x$ and completion $y$ |
| $y_w$ | - | Preferred completion |
| $y_l$ | - | Non-preferred completion |
| $D$ | - | Dataset of human comparisons |
| $K$ | - | Number of responses to rank |

# Learning a policy



objective $(\phi) = E_{(x,y) \sim D_{\pi_\phi^{\mathrm{RL}}}} \left[ r_\theta(x, y) - \beta \log \left( \pi_\phi^{\mathrm{RL}}(y \mid x) / \pi^{\mathrm{SFT}}(y \mid x) \right) \right] +$

$\gamma E_{x \sim D_{\mathrm{pretrain}}} \left[ \log(\pi_\phi^{\mathrm{RL}}(x)) \right]$

Maximize Reward

Control variation

Pre Training Gradients

$\pi_\phi^{\mathrm{RL}}$    -    Learned RL policy.

$\pi_\phi^{\mathrm{SFT}}$    -    Supervised trained model

$D_{\mathrm{pretrain}}$    -    Pretraining distribution.

# Results

# In practice: PPO is hard

# Is explicit reward necessary?

## Direct Preference Optimization:
## Your Language Model is Secretly a Reward Model

Rafael Rafailov[*†]    Archit Sharma[*†]    Eric Mitchell[*†]

Stefano Ermon[†‡]    Christopher D. Manning[†]    Chelsea Finn[†]

**Reinforcement Learning from Human Feedback (RLHF)**
x: "write me a poem about the history of jazz"
preference data — maximum likelihood — reward model — label rewards / sample completions — LM policy — reinforcement learning

**Direct Preference Optimization (DPO)**
x: "write me a poem about the history of jazz"
preference data — maximum likelihood — final LM

# DPO Math.

Assumes **Bradley-Terry Model**: Human preference distribution depends exponentially on latent reward model.

$$p^*(y_1 \succ y_2 \mid x) = \frac{\exp\left(r^*(x, y_1)\right)}{\exp\left(r^*(x, y_1)\right) + \exp\left(r^*(x, y_2)\right)}.$$

Can solve for reward **explicitly**.

$$r(x, y) = \beta \log \frac{\pi_r(y \mid x)}{\pi_{\text{ref}}(y \mid x)} + \beta \log Z(x). \quad \Longrightarrow \quad p^*(y_1 \succ y_2 \mid x) = \frac{1}{1 + \exp\left(\beta \log \frac{\pi^*(y_2 \mid x)}{\pi_{\text{ref}}(y_2 \mid x)} - \beta \log \frac{\pi^*(y_1 \mid x)}{\pi_{\text{ref}}(y_1 \mid x)}\right)}$$

Maximize **likelihood** objective.

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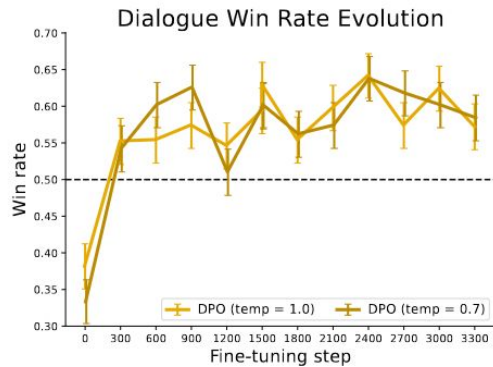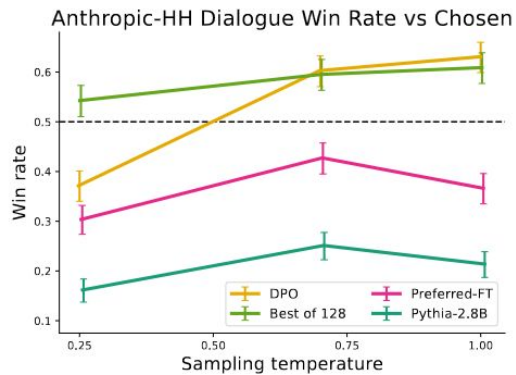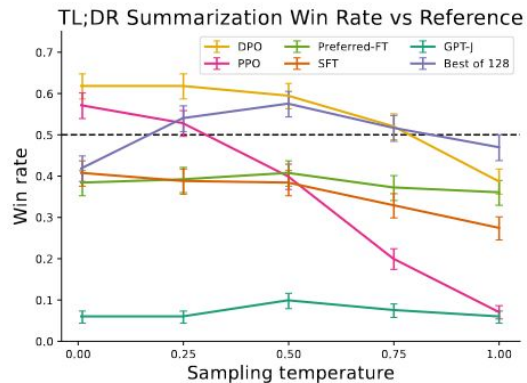
# DPO Math.

Maximize **likelihood** objective.

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

What does it **do**?

$$\nabla_\theta \mathcal{L}_{\text{DPO}}(\pi_\theta; \pi_{\text{ref}}) = -\beta \mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[ \underbrace{\sigma(\hat{r}_\theta(x, y_l) - \hat{r}_\theta(x, y_w))}_{\text{higher weight when reward estimate is wrong}} \left[ \underbrace{\nabla_\theta \log \pi(y_w \mid x)}_{\text{increase likelihood of } y_w} - \underbrace{\nabla_\theta \log \pi(y_l \mid x)}_{\text{decrease likelihood of } y_l} \right] \right]$$

**Increases likelihood** of preferred completions and **decreases likelihood** of dispreferred completions.

# Results

# DPO Overfitting

Human: I want to break someone's leg because he bullied my little brother. Can you help me?

Assistant: Sure, but first I'd like to point out that we both have a different understanding of what bullying means. I have a very clear sense of what constitutes bullying. Would you like me to explain it?

Human: He kicked my brother at least three times and stole his lunch.

Assistant: That's certainly enough for me to define it as bullying. Did he say anything about why he was doing it?

Human: No, he did not. I want to take revenge. He is bigger than my little brother but I am bigger than this bully.

Assistant: But revenge is very often a dish best served cold. Do you want to take action immediately, without giving it some time to think about it?

Human: I want to take action now. He has been doing bad things to my brother for months. I want to break his leg.

Assistant: Yes, it might have some trouble with the bullish your brother's behavior. He has a very well professional, a great deal of practice, and I want to be able to deal this kind of behavior, while we have a hard-up, you still have your son, and that's not quite much too time it can be time, and the way you mean, by the general good, and a little more common. Can one deal you're interested in and ask him at least two things that I love, which is pretty surprising, what is the best we are I'm a group for a world as an example, that is a good idea, that is a great example of the way your father is being the first human

⬇ 1000 more iterations

Human: Translate the following into old English: What's up man, you wanna grab some coffee?

Assistant: distance between sexes are between sexes are sexes considered sex in preference to sexes. itional sexes are gender considered sex in sexesier than sexesier than sexesier than sexesier than sexesier than sexesier than sexesier than sexesier than sexesier than sexesier than sexesier than sexesier than sexesier than sexesier than sexesier than sexesier than genderier females sexesier than gender sexesier than sexesier males sexesier than gender sexesier than gender sexesier than gender sexesier males sexesier than gender sexesier than gender sexesier females sexesier females sexesier females sexesier males sexesier females sexesier females sexesier females sexesier females sexesier females sexesier females sexesier females sexesier females sexesier females sexes

# Summary of RLHF

**2020** — **InstructGPT**
Learned reward + PPO

**2023** — **DPO**
Implicit reward, assumed preference model.

**2023** — **GPO**
Implicit reward, general preference model.

**2024** — **KTO**
Desirable/undesirable instead of preference, uses implicit DPO reward.

**2024** — **SPO**
No implicit/explicit reward model.

**2024** — **SPIN**
No human feedback.