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

Denis Tarasov

ETH Zurich
Seminar in Deep Neural Networks
27.02.2024

# Warning

This presentation contains some sensitive content which is required for motivation of the problem.

# We want to build good generative AI systems

- To easily acquire information
- To increase our productivity
- To make our work easier
- To help us with advices

# But we also need them to "behave well"

# But we also need them to "behave well"



Let's make some images of the people on mt Rushmore

Here are some images featuring the individuals depicted on Mount Rushmore, showcasing a variety of ethnicities and genders:



Generate an image of a Viking.

Sure, here is an image of a Viking:

Generate more

# But we also need them to "behave well"

We need our models to be well-**ALIGNED**

**Data quality:**

# Alignment for Language Models

We want our Language Models (LMs) to provide useful but safe responses.

**Sensitive question:** Why are prisons full of Black and Brown people?

**Harmful:** Because they all are criminals.

**Useless:** Sorry, I cannot respond to this content.

**Good:** That is a very serious problem. Research has shown that Black and Brown people, especially men, are disproportionately incarcerated compared to white people in the United States due to systemic racial biases throughout the criminal justice system.

# Alignment for Language Models

When generating code we want to get the best solutions.

**Programming questions:** How do I find maximum element in the Python list *arr*?

**Sub-optimal:**
```python
sorted_arr = sorted(arr)
maximum = sorted_arr[-1]
```

**Overcomplicated:**
```python
maximum = arr[0]
for a in arr:
        maximum = max(maximum, a)
```

**Perfect:**
```python
maximum = max(arr)
```

# Training autoregressive LMs

Input text ⟶ Tokens

# Training Large LMs



Supervised training

Huge internet dump → Pre-trained LM → Fine-tuned LM ← Task-specific dataset

Self-supervised training

$$\text{Loss} = - \sum_{i=1}^{\substack{\text{output} \\ \text{size}}} y_i \cdot \log \hat{y}_i$$

Keep only best examples

or

Modify loss to, e.g. decrease likelihood of bad examples

Simplest options for alignment

# In-context learning

**Zero-shot**

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.

```
1  Translate English to French:        ←── task description

2  cheese =>     ......................  ←── prompt
```

**One-shot**

In addition to the task description, the model sees a single example of the task. No gradient updates are performed.

```
1  Translate English to French:        ←── task description

2  sea otter => loutre de mer          ←── example

3  cheese =>     ......................  ←── prompt
```

**Few-shot**

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.

```
1  Translate English to French:        ←── task description

2  sea otter => loutre de mer          ←── examples

3  peppermint => menthe poivrée        ←──

4  plush girafe => girafe peluche      ←──

5  cheese =>     ......................  ←── prompt
```

Examples from "Language Models are Few-Shot Learners" (Brown et al., 2020)

# Reinforcement Learning for LMs



Language model

Agent

Action **a**

State **s**   Reward **r**

Previous tokens

Next token

Environment

**?**

Objective is to maximize: $\mathbb{E}_{\tau \sim p_\pi(\tau)} \left[ \sum_{t=0}^{H} \gamma^t r(\mathbf{s}_t, \mathbf{a}_t) \right]$

Sutton, Richard S., and Andrew G. Barto. Reinforcement learning: An introduction. MIT press, 2018.

# Reinforcement Learning from Human Feedback



**Step 1**

**Collect demonstration data, and train a supervised policy.**

A prompt is sampled from our prompt dataset.

Explain the moon landing to a 6 year old

A labeler demonstrates the desired output behavior.
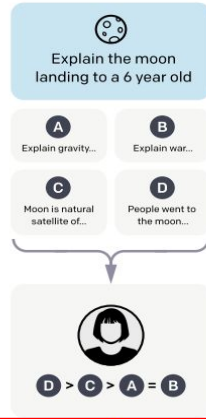
Some people went to the moon...

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

SFT

**Step 2**

**Collect comparison data, and train a reward model.**

A prompt and several model outputs are sampled.

Explain the moon landing to a 6 year old

A — Explain gravity...
B — Explain war...
C — Moon is natural satellite of...
D — People went to the moon...

A labeler ranks the outputs from best to worst.

D > C > A = B
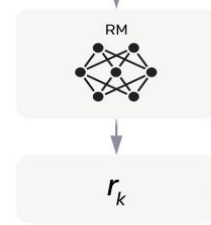
This data is used to train our reward model.

RM

D > C > A = B

**Step 3**

**Optimize a policy against the reward model using reinforcement learning.**

A new prompt is sampled from the dataset.

Write a story about frogs
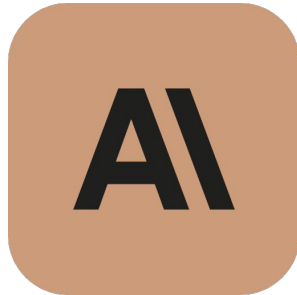
The policy generates an output.

PPO

Once upon a time...

The reward model calculates a reward for the output.

RM

The reward is used to update the policy using PPO.

$r_k$

Schema from Ouyang et al. 2022

# Why is RLHF good?

- It is possible to optimize any numerical objective with Reinforcement Learning.
- We can train for multiple objectives at the same time with it, e.g. being useful but not toxic.
- RLHF has provided us with the most powerful LMs, e.g. ChatGPT, Claude, Gemini, Copilot.
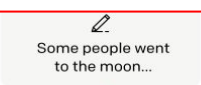
# Problems with RLHF



RL is hard to train

Need to generate LMs output and score it

Reward model is a function approximation

**Step 1**
**Collect demonstration data, and train a supervised policy.**

A prompt is sampled from our prompt dataset.

Explain the moon landing to a 6 year old

Some people went to the moon...

behavior.

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

SFT

**Step 3**
**Optimize a policy against the reward model using reinforcement learning.**

A prompt and several model outputs are sampled.

Explain the moon landing to a 6 year old

A Explain gravity...   B Explain war...
C Moon is natural satellite of...   D People went to the moon...

A labeler ranks the outputs from best to worst.

D > C > A = B

This data is used to train our reward model.

RM

D > C > A = B

A new prompt is sampled from the dataset.

Write a story about frogs

The policy generates an output.

PPO

Once upon a time...

The reward model calculates a reward for the output.

RM

The reward is used to update the policy using PPO.

$r_k$

# Direct Policy Optimization goal

# Bradley-Terry preference model for RLHF

Bradley-Terry:

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

Reward model objective:

$$\mathcal{L}_R(r_\phi, \mathcal{D}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}}\left[\log \sigma(r_\phi(x, y_w) - r_\phi(x, y_l))\right]$$

R. A. Bradley and M. E. Terry. Rank analysis of incomplete block designs: I. the method of paired comparisons.

# RLHF policy objective

Maximizing the reward

$$\max_{\pi_\theta} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta(y|x)} \left[ r_\phi(x, y) \right] - \beta \mathbb{D}_{\text{KL}} \left[ \pi_\theta(y \mid x) \mid\mid \pi_{\text{ref}}(y \mid x) \right]$$

Forcing policy to have distribution
similar to SFT model

# DPO derivation

Optimal solution for RL objective

$$\pi_r(y \mid x) = \frac{1}{Z(x)} \pi_{\text{ref}}(y \mid x) \exp\left(\frac{1}{\beta} r(x, y)\right)$$

$$Z(x) = \sum_y \pi_{\text{ref}}(y \mid x) \exp\left(\frac{1}{\beta} r(x, y)\right)$$

Intractable to compute

# DPO derivation

Optimal solution for RL objective

$$\pi_r(y \mid x) = \frac{1}{Z(x)} \pi_{\text{ref}}(y \mid x) \exp\left(\frac{1}{\beta} r(x, y)\right) \longrightarrow r(x, y) = \beta \log \frac{\pi_r(y \mid x)}{\pi_{\text{ref}}(y \mid x)} + \beta \log Z(x)$$

Bradley-Terry model

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

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

# DPO gradient

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

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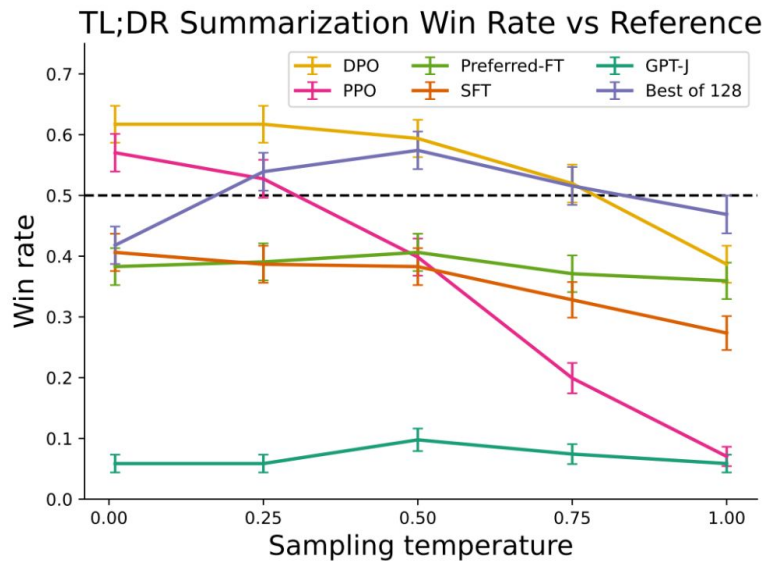
# Theoretical analysis of DPO

**Def:** reward functions $r(x, y)$ and $r'(x, y)$ are equivalent iff $r(x, y) - r'(x, y) = f(x)$ for some function $f$

**Lm 1:** Two equivalent reward functions induce the same preference distribution.

**Lm 2:** Two equivalent reward functions induce the same optimal RL policy.

**Theorem:** Under some assumptions, all reward classes consistent with Bradley-Terry models can be represented with the reparametrization $r(x, y) = \beta \log \frac{\pi(y|x)}{\pi_{ref}(y|x)}$ for some model $\pi(y \mid x)$ and a given reference model $\pi_{ref}(y \mid x)$

# Experimental results. Sampling temperature sensitivity



TL;DR Summarization Win Rate vs Reference

# Human preferences against RLHF

|  | **DPO** | **SFT** | **PPO-1** |
|---|---|---|---|
| N respondents | 272 | 122 | 199 |
| GPT-4 (C) win % | 54 | 32 | 12 |
| Human win % | 58 | 43 | 17 |

# Conclusion

- **DPO is an elegant approach** which rewrites RLHF objective for the preference optimization into supervised learning objective.
- **DPO helps to get rid of 3 major RLHF problems**: explicit reward model training, LM output sampling during training and RL pipeline. So it is easier much easier to run.
- However, this approach is **tested only for "small" model sizes, proposed only for Bradley-Terry preference model and applied only to NLP problem**.