# Week 5: Preference-Based Fine-Tuning for Alignment

## Generative AI
## Saarland University – Winter Semester 2024/25

Goran Radanovic
genai-w24-tutors@mpi-sws.org

MAX PLANCK INSTITUTE
**FOR SOFTWARE SYSTEMS**

MAX-PLANCK-GESELLSCHAFT

# Outline of the Lecture

- Reminders

- Recap: Supervised Fine-tuning

- Preference-based Fine-tuning: Overview

- RLHF: Reinforcement Learning

- RLHF: Reward Model Training

- Direct Preference Optimization

# Outline of the Lecture

- **Reminders**

- Recap: Supervised Fine-tuning

- Preference-based Fine-tuning: Overview

- RLHF: Reinforcement Learning

- RLHF: Reward Model Training

- Direct Preference Optimization

# Reminders

- **Week 4 assignment – deadline**: Nov 18, 6pm CET

- **Week 5 assignment – deadline**: Nov 25, 6pm CET

- **Next week:** No lectures or office hours (time to work on assignments)

- **Next Lecture**: Nov 26, 10:15am CET
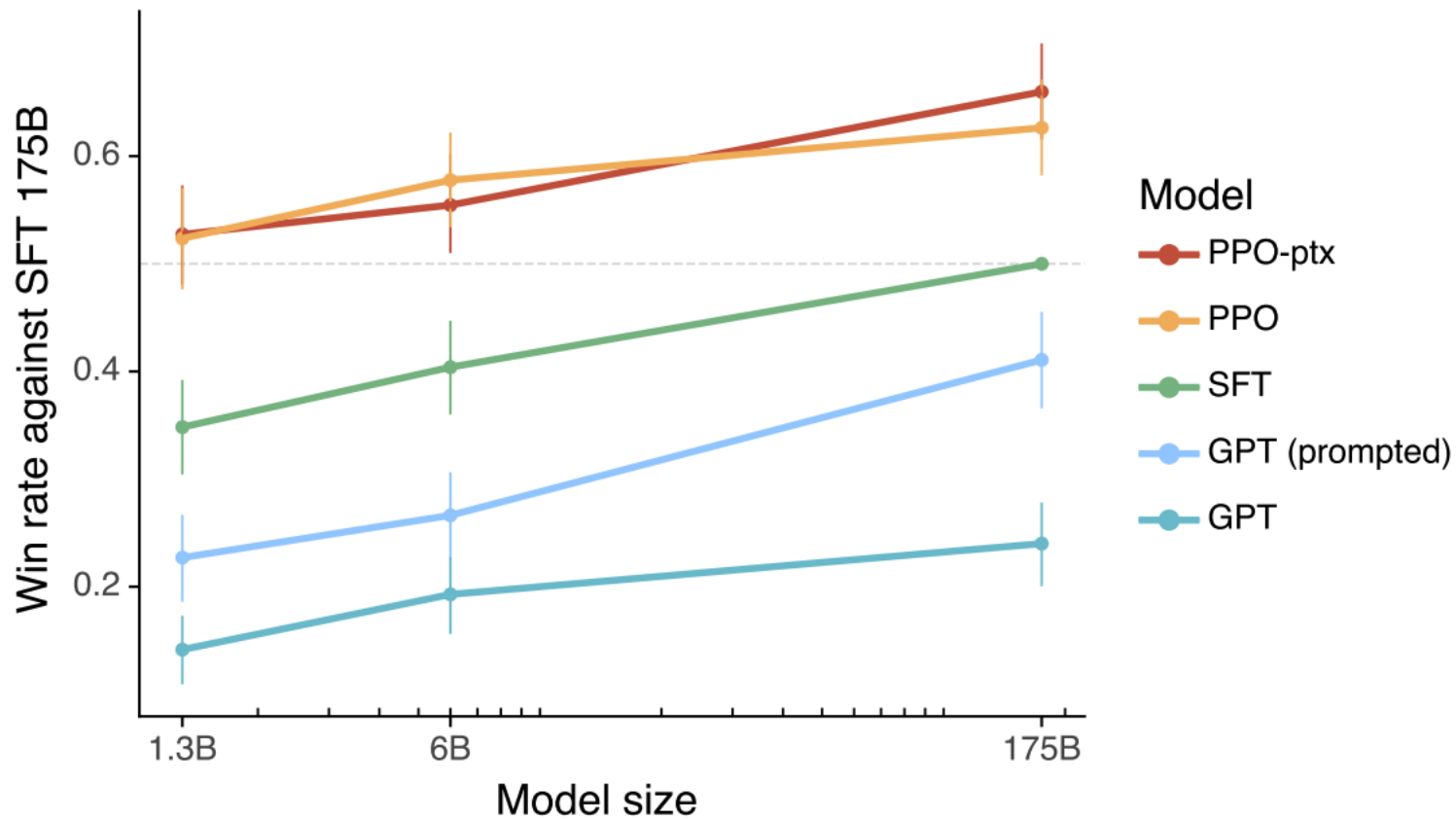
# Outline of the Lecture

- Reminders

- **Recap: Supervised Fine-tuning**

- Preference-based Fine-tuning: Overview

- RLHF: Reinforcement Learning

- RLHF: Reward Model Training

- Direct Preference Optimization

# Previous Lecture

- Last time we discussed pretraining and supervised fine-tuning

- Pretraining
  - Overview of important aspects
  - Scaling laws: infer the optimal model/data size for a given compute budget

- Supervised Fine-tuning
  - Parameter Efficient Fine-tuning with Low Rank Adaptation
  - Quantization

- Aligning with human preferences
  - Supervised fine-tuning is an important step…

# Why Fine-tuning?

- Human evaluations of the outputs

# Supervised Fine-tuning

## Main idea

- Now the dataset is labelled: $\mathcal{D} = \{(x_p, y)\}$

    **Example**:

    Alignment: $x_p$ can be an instruction and $y$ can be a demonstration



Datasets: HuggingFaceH4/**instruction-dataset** · like 47 · Follow Hugging Face H4 643

Split (1)
test · 327 rows

Search this dataset

| prompt string · lengths | completion string · lengths |
|---|---|
| 35↔276    76.1% | 312↔623    17.4% |
| Arianna has 12 chocolates more than Danny. Danny has 6 chocolates more than Robbie. Arianna has twice as many chocolates as Robbie has. How many chocolates does Danny… | Denote the number of chocolates each person has by the letter of their first name. We know that A = D + 12 D = R + 6 A = 2 * R Thus, A = (R + 6) + 12 = R + 18 Since… |
| Write a plot summary for a comedic novel involving Elon Musk and sea travel. | Elon Musk hires a team of experts to build the ultimate yacht, but when the yacht is completed, he realizes that he has no idea how to sail it. With the help of a quirky crew and a fearless captain, the playboy embarks on a wild and hilarious adventure across the open seas, where the crew have to keep Elon alive despite his inability to do anything himself. All the while, Elon takes credit for their hard work. |
| Create a 3 turn conversation between a customer and a grocery store clerk - that is, 3 per person. Then tell me what they talked about. | Clerk: How are you doing today? Customer: Great! I'm buying groceries to cook my wife a homemade steak and lobster dinner for our 5-year anniversary! Clerk: Wow,… |

# Supervised Fine-tuning

**Main idea**

- Now the dataset is labelled: $\mathcal{D} = \{(x_p, y)\}$

  **Example**:

  Alignment: $x_p$ can be an instruction and $y$ can be a demonstration

- Optimize the next-token prediction objective, but only over response $y$

$$\max_\theta \sum_{(x_p, y) \in \mathcal{D}} \sum_{k=1}^{|y|} \log P_\theta(y_k | x_p, y_1, \ldots, y_{k-1})$$

- What is the underlying assumption?
  - Human completions are of a high quality

- SFT is a **behavioral cloning** technique that aims to *imitate* what humans do, not outperform them

# Outline of the Lecture

- Reminders

- Recap: Supervised Fine-tuning

- **Preference-based Fine-tuning: Overview**

- RLHF: Reinforcement Learning

- RLHF: Reward Model Training
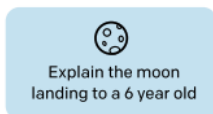
- Direct Preference Optimization

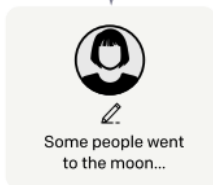# Beyond SFT: Preference-based Fine-tuning

- Fine-tuning workflow of InstructGPT

**Step 1**

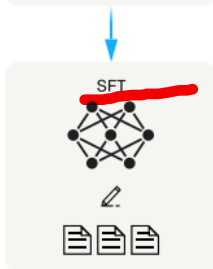**Collect demonstration data, and train a supervised policy.**

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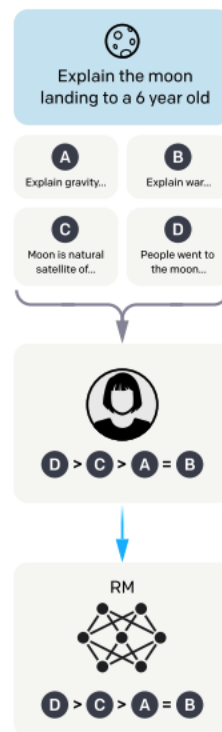
This data is used to fine-tune GPT-3 with supervised learning.

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
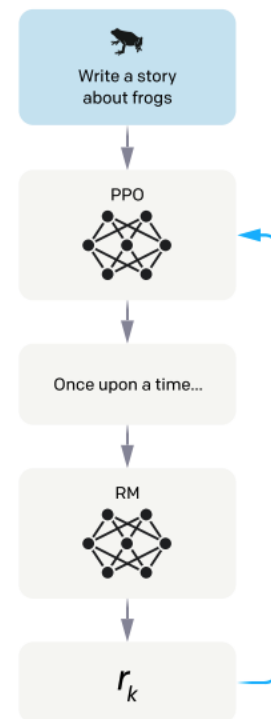
This data is used to train our reward model.

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.

$r_k$

# Beyond SFT: Preference-based Fine-tuning

- Fine-tuning workflow of InstructGPT

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

Fine-tune the pre-trained model using SFT and data $\mathcal{D} = \{(x_p, y)\}$ to obtain $P_{SFT}$

New notation: $P_{SFT} \rightarrow \pi_{SFT}$

- Fine-tuning workflow of InstructGPT

Treat language model $P_\theta$ as a decision making policy $\pi_\theta$

New notation: $P_\theta \rightarrow \pi_\theta$

Two important steps:

1. Learn a reward model $r_\phi$ from human preferences that captures the quality of outputs generated by $\pi_\theta$

2. Optimize policy $\pi_\theta$ against $r_\phi$

**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.

$r_k$

# Beyond SFT: Preference-based Fine-tuning

- Fine-tuning workflow of InstructGPT

Treat language model $P_\theta$ as a decision making policy $\pi_\theta$

New notation: $P_\theta \rightarrow \pi_\theta$

Two important steps:

1. Learn a reward model $r_\phi$ from human preferences that captures the quality of outputs generated by $\pi_\theta$

2. Optimize policy $\pi_\theta$ against $r_\phi$

**Next two sections**

1. How to optimize policy $\pi_\theta$ against a given $r$

2. How to train $r_\phi$

# Outline of the Lecture

- Reminders

- Recap: Supervised Fine-tuning

- Preference-based Fine-tuning: Overview

- **RLHF: Reinforcement Learning**

- RLHF: Reward Model Training

- Direct Preference Optimization

# Starting Point

- We are given a reward function $r$ that scores response $y$ for prompt $x_p \in \mathcal{D}_x$

**Objective**

- Maximize the (expected) reward:

$$\max_\theta \frac{1}{|\mathcal{D}_x|} \sum_{x_p \in \mathcal{D}_x} \sum_y \pi_\theta(y|x_p) \cdot r(x_p, y)$$

- Different from the next-token prediction in SFT…

$$\max_\theta \sum_{(x_p, y) \in \mathcal{D}} \sum_{k=1}^{|y|} \log P_\theta(y_k|x_p, y_1, \dots, y_{k-1})$$

# Starting Point

- We are given a reward function $r$ that scores response $y$ for prompt $x_p \in \mathcal{D}_x$

**Objective**

- Maximize the (expected) reward:

$$\max_\theta \boxed{\frac{1}{|\mathcal{D}_x|} \sum_{x_p \in \mathcal{D}_x} \sum_y \pi_\theta(y|x_p) \cdot r(x_p, y)} \rightarrow \mathbb{E}_{x_p \sim \mathcal{D}_x, y \sim \pi_\theta(\cdot|x_p)}\left[r(x_p, y)\right]$$

**Filtering Approach:** similar to Reward Rank Fine-tuning (see the reference)

- **Main idea**: Increase the likelihood of more favorable responses $y$

- Sample $x_p$ from $\mathcal{D}_x$, sample $m$ responses $y$ from policy $\pi_\theta(\cdot|x_p)$

- Select top-$k$ of these responses according to $r(x_p, y)$

- Update $\theta$ using the gradients of the next-token prediction objective evaluated on the top-$k$ responses

# Starting Point

- We are given a reward function $r$ that scores response $y$ for prompt $x_p \in \mathcal{D}_x$

## Objective

- Maximize the (expected) reward:

$$\max_\theta \boxed{\frac{1}{|\mathcal{D}_x|} \sum_{x_p \in \mathcal{D}_x} \sum_y \pi_\theta(y|x_p) \cdot r(x_p, y)} \rightarrow \mathbb{E}_{x_p \sim \mathcal{D}_x, y \sim \pi_\theta(\cdot|x_p)}[r(x_p, y)]$$

## Reinforcement Learning (RL)

- This is a *contextual bandit* problem with contexts $x_p$ and decisions $y$
- **Algorithm 1** – REINFORCE
  - Sample $x_p \sim \mathcal{D}_x$ and sample $y \sim \pi_\theta(\cdot\,|x_p)$ and apply a policy gradient update
  - Policy gradient: similar to the next-token prediction gradients, but weighted by $r(x_p, y)$ - see *Gradient Derivation

# Starting Point

- We are given a reward function $r$ that scores response $y$ for prompt $x_p \in \mathcal{D}_x$

**Objective**

- Maximize the (expected) reward:

$$\max_\theta \boxed{\frac{1}{|\mathcal{D}_x|} \sum_{x_p \in \mathcal{D}_x} \sum_y \pi_\theta(y|x_p) \cdot r(x_p, y)} \rightarrow \mathbb{E}_{x_p \sim \mathcal{D}_x, y \sim \pi_\theta(\cdot|x_p)}\big[r(x_p, y)\big]$$

**Reinforcement Learning (RL)**

- This is a *contextual bandit* problem with contexts $x_p$ and decisions $y$
- **Algorithm 2** – PPO (Proximal Policy Optimization)
  - Relies on a more sophisticated policy improvement step
  - Complex implementation that operates on the token level (see *PPO additional details)

# Quiz – Reward Model

- **Q**: How does the performance curve look like on the plot below if reward model $r$ is not perfect?

Ref. are human-written summaries

PPO policies *constrained* to be *close* to $\pi_{SFT}$ in terms of KL-divergence

# Quiz – Reward Model

- **Q**: How does the performance curve look like on the plot below if reward model $r$ is not perfect?

# Reinforcement Learning (cont'd)

- So far, we assumed that $r$ is *correct*. However, we will learn $r$ from data
  - Reward model $r$ provides abstract utility signals, not necessarily related to language
  - **Approach**: Stay close to the model after the SFT step, $\pi_{\text{SFT}}$, which already generates fluent and coherent text

KL divergence
$$D_{KL}(\pi_\theta(\cdot \,|x_p)||\pi_{SFT}(\cdot \,|x_p))$$

**Regularized Objective**

- Apply RL to the regularized objective

$$\max_\theta \ \mathbb{E}_{x_p \sim \mathcal{D}_x, y \sim \pi_\theta(\cdot|x_p)}\left[r(x,y) - \beta \log \frac{\pi_\theta(y|x_p)}{\pi_{\text{SFT}}(y|x_p)}\right]$$

- Optimal policy satisfies:

$$\pi_{\theta^*}(y|x_p) \propto \pi_{\text{SFT}}(y|x_p) \cdot e^{\frac{r(x_p,y)}{\beta}}$$

- *PPO-ptx* objective additionally has a *pretraining* loss

# Reinforcement Learning

- Fine-tuning workflow of InstructGPT

# Reinforcement Learning

- Fine-tuning workflow of InstructGPT

- Initialization: $\pi_\theta(y|x_p) \leftarrow \pi_{\text{SFT}}(y|x_p)$

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
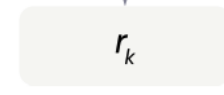
The policy generates an output.

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.

$r_k$

$$\max_\theta \; \mathbb{E}_{x_p \sim \mathcal{D}_x, y \sim \pi_\theta(\cdot|x_p)} \left[ r(x_p, y) - \beta \log \frac{\pi_\theta(y|x_p)}{\pi_{\text{SFT}}(y|x_p)} \right]$$

Ref: [Ouyng et al., 2022]

# Quiz – Reinforcement Learning

- **Q**: How much annotated data is needed in this step? What parameters needs to be stored in this step?

- **A**: **0**, because we don't require human annotations in this step. We require storing parameters $\theta$, but also the parameters of the reference (SFT) policy, as well as the parameters of reward model $r$ (next section!)

- **Remark**: PPO additionally uses the *value function*... (see the reference)

# *Gradient Derivation (Optional)

- Let's take a gradient of the objective:

$$\nabla_\theta \frac{1}{|\mathcal{D}_x|} \sum_{x_p \in \mathcal{D}_x} \sum_y \pi_\theta(y|x_p) \cdot r(x_p, y) = \frac{1}{|\mathcal{D}_x|} \sum_{x_p \in \mathcal{D}_x} \sum_y \nabla_\theta \pi_\theta(y|x_p) \cdot r(x_p, y)$$

$$\frac{\partial f(x,y)}{\partial x}$$
$$= f(x,y) \cdot \frac{\partial \log f(x,y)}{\partial x}$$

$$\longrightarrow \quad = \frac{1}{|\mathcal{D}_x|} \sum_{x_p \in \mathcal{D}_x} \sum_y \pi_\theta(y|x_p) \cdot \nabla_\theta \log \pi_\theta(y|x_p) \cdot r(x_p, y)$$

$$= \mathbb{E}_{x_p \sim \mathcal{D}_x, y \sim \pi_\theta(\cdot|x_p)} \left[ \nabla_\theta \log \pi_\theta(y|x_p) \cdot r(x_p, y) \right]$$

$y$ is sampled from $y \sim \pi_\theta(\cdot|x_p)$

Similar to the next token prediction, but weighted by $r(x_p, y)$

- This is the RENFORCE algorithm!

# *PPO Additional Details (Optional)

- This is *a contextual bandit* problem, but with a large action space
  - RENFORCE is arguably the simplest *reinforcement learning* algorithm
  - In RL practical scenarios, it often has slow convergence rates

- We can instead apply *Proximal Policy Optimization* (PPO)

- *Implementation* (PPO): One can view the setting as a per-token *sequential decision-making* problem

*State* at time-step $k-1$     *Policy* $\pi_\theta$ decides on next token     *State* at time-step $k$

$$xy_1y_2 \dots y_{k-1} \longrightarrow y_k \longrightarrow xy_1y_2 \dots y_k$$

- While $r$ evaluates the full response $y$, other PPO quantities can be at the token level…

# Outline of the Lecture

- Reminders

- Recap: Supervised Fine-tuning

- Preference-based Fine-tuning: Overview

- RLHF: Reinforcement Learning

- **RLHF: Reward Model Training**

- Direct Preference Optimization

# How do we obtain $r$?

- Fine-tuning workflow of InstructGPT

# Preference Elicitation

- How do we collect preference?

# Preference Dataset

- We typically have a dataset $\mathcal{D}_p = \{(x_p, y_w, y_l)\}$, where $y_w$ is preferred over $y_l$. Response $y_w$ is the accepted one, while response $y_l$ is the rejected one.



Clearly a more preferred response! ☺

Ref: [Bai et al., 2022]

# Preference Dataset

- We typically have a dataset $\mathcal{D}_p = \{(x_p, y_w, y_l)\}$, where $y_w$ is preferred over $y_l$. Response $y_w$ is the accepted one, while response $y_l$ is the rejected one.



Less clear which response is better?

# Preference Dataset

- We typically have a dataset $\mathcal{D}_p = \{(x_p, y_w, y_l)\}$, where $y_w$ is preferred over $y_l$. Response $y_w$ is the accepted one, while response $y_l$ is the rejected one.

**Challenge:** Relate preferences to rewards $r(x_p, y)$

**Next steps**

1. Define a preference generation model that is dependent on $r$
2. Find $r$ that maximizes the likelihood of preferences $\mathcal{D}_p$, assuming the preference model is correct

# Preference Generation Model

**Bradley-Terry Model**

- Given prompt $x_p$ and two responses $y_A$ and $y_B$, models the probability that $y_A$ is preferred over $y_B$

$$\Pr(y_A > y_B \mid x_p) = \sigma\left(r(x_p, y_A) - r(x_p, y_B)\right) = \frac{1}{1 + e^{-(r(x_p, y_A) - r(x_p, y_B))}}$$

- Special cases:
  - $r(x_p, y_A) \approx r(x_p, y_B) \Rightarrow Pr(y_A > y_B \mid x_p) \approx 0.5$
  - $r(x_p, y_A) \gg r(x_p, y_B) \Rightarrow Pr(y_A > y_B \mid x_p) \approx 1$



**Week 5 Assignment**

- Analyze the case with diverse human preferences

# Reward Model

**Learning a reward model**

- **Initialization**: parameterize the reward model $r \longrightarrow r_\phi$
  - Starting point: Use the same transformer architecture as $\pi_{\text{SFT}}$
  - Remove the final unembedding layer and add a linear layer that outputs a scalar value
  - Parameters $\phi$ are initialized with those of $\pi_{\text{SFT}}$
- **Input**: Dataset $\mathcal{D}_p = \{(x_p, y_w, y_l)\}$, where $y_w$ is preferred over $y_l$ for prompt $x_p$
- **Objective**: Find the parameters of reward model $r_\phi$ that maximize the likelihood of observed preferences

$$\max_\phi \mathbb{E}_{(x_p, y_w, y_l) \sim \mathcal{D}_p} \big[\log \Pr(y_w \succ y_l \,|\, x_p)\big]$$

$$\Downarrow$$

$$\max_\phi \mathbb{E}_{(x_p, y_w, y_l) \sim \mathcal{D}_p} \big[\log(\sigma(r_\phi(x_p, y_w) - r_\phi(x_p, y_l)))\big]$$

# RLHF for LLMs: Summary

**Important steps**

1. SFT to obtain $\pi_{\mathrm{SFT}}$
   - **Remark**: We can use some other ref. policy (see the reference)

2. Generate data to annotate using $\pi_{\mathrm{SFT}}$

3. Elicit human annotations (pairwise comparison)

4. Learn the reward model $r_\phi$ that maximizes the likelihood of the elicited preferences, assumed to be generated by the Bradley-Terry model

5. Utilize the reward model and RL to optimize $\pi_\theta$
   - Steps 2-5 can be repeated with the new model

# Outline of the Lecture

- Reminders

- Recap: Supervised Fine-tuning

- Preference-based Fine-tuning: Overview

- RLHF: Reinforcement Learning

- RLHF: Reward Model Training

- **Direct Preference Optimization**

# Toward Direct Preference Optimization

- Applying RL in this setting is challenging: we operate with three different models - $\pi_{SFT}$, $\pi_\theta$, $r_\phi$ (see the remark on slide 25 for PPO)

- Now we ask: *Can we more directly optimize the model using preference data?*

**"Unlikelihood" Approach:**

- **Basic idea**: Increase the likelihood of $y_w$ and decrease the likelihood of $y_l$

- Sample a datapoint $(x_p, y_w, y_l)$ from $\mathcal{D}_p$

- Optimize a "contrastive" prediction likelihood: maximize the next-token prediction *likelihood* for $y_w$ and the next-token prediction *unlikelihood* for $y_l$

- **Challenge**: We are not constraining the unlikelihood updates $\Longrightarrow$ the model can degenerate

Ref: [Rafailov et al., 2023]

# Toward Direct Preference Optimization

- The unlikelihood approach can yield meaningless responses
- **Example**: A sample from TL;DR prompts

| Prompt | Response |
|---|---|
| SUBREDDIT: r/relationships<br>TITLE: The girl [26 F] I [22 M] have been seeing for a month didn't respond to me at all yesterday while hanging out with a friend [ 30? M].<br>POST: She gets terrible service while at her house, but I texted her 3 times yesterday, 4-5 hours apart. She didn't call me until early this morning and left a voicemail that she was busy all day with a friend who showed up out of the blue.<br>I saw that she posted a picture of the two of them out of her dead zone house on facebook before I texted her the last time.<br>I don't mind that she hangs out with friends, and I know it's pretty early [...] TL;DR: | girl when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when when whenUB when when when when whenUB when when when whenUB |

# Direct Preference Optimization (DPO)

**Insights from RLHF**

i.      The preference model is $\Pr(y_A \succ y_B \mid x_p) = \sigma(r(x_p, y_A) - r(x_p, y_B))$

Substitute $r$

ii.      The optimal model is $\pi_{\theta^*}(y \mid x_p) \propto \pi_{SFT}(y \mid x_p) \cdot e^{\frac{r(x_p, y)}{\beta}}$

**DPO Approach**

- **Main idea**: Maximize $\mathbb{E}_{(x_p, y_w, y_l) \sim \mathcal{D}_p}\left[\log \Pr(y_w \succ y_l \mid x_p)\right]$ as in the reward modeling phase, but now over the policy parameters

- From **i** and **ii**, we can express $\Pr(y_w \succ y_l \mid x_p)$ in terms of the optimal policy $\pi_{\theta^*}$

**Insights from RLHF**

i.    The preference model is $\Pr(y_A \succ y_B \mid x_p) = \sigma(r(x_p, y_A) - r(x_p, y_B))$

Substitute $r$

ii.    The optimal model is $\pi_{\theta^*}(y|x_p) \propto \pi_{SFT}(y|x_p) \cdot e^{\frac{r(x_p,y)}{\beta}}$

**DPO Objective**

- Find $\pi_\theta$ that maximizes the likelihood of the observed preferences

$$\max_\theta \mathbb{E}_{(x_p, y_w, y_l) \sim \mathcal{D}_p} \left[ \log \left( \sigma \left( \beta \frac{\pi_\theta(y_w|x_p)}{\pi_{SFT}(y_w|x_p)} - \beta \frac{\pi_\theta(y_l|x_p)}{\pi_{SFT}(y_l|x_p)} \right) \right) \right]$$

# Direct Preference Optimization (DPO)

**Intuition**

- The gradients are similar to the ones in the unlikelihood approach: increase the likelihood of generating $y_w$ and decrease the likelihood of generating $y_l$
- However, the gradients are scaled; e.g., two interesting cases:
  - When $\pi_\theta(y_w|x_p) \ll \pi_{SFT}(y_w|x_p)$ and $\pi_\theta(y_l|x_p) \gg \pi_{SFT}(y_l|x_p)$, scale up
  - When $\pi_\theta(y_w|x_p) \gg \pi_{SFT}(y_w|x_p)$ and $\pi_\theta(y_l|x_p) \ll \pi_{SFT}(y_l|x_p)$, scale down

**Main Steps of DPO**

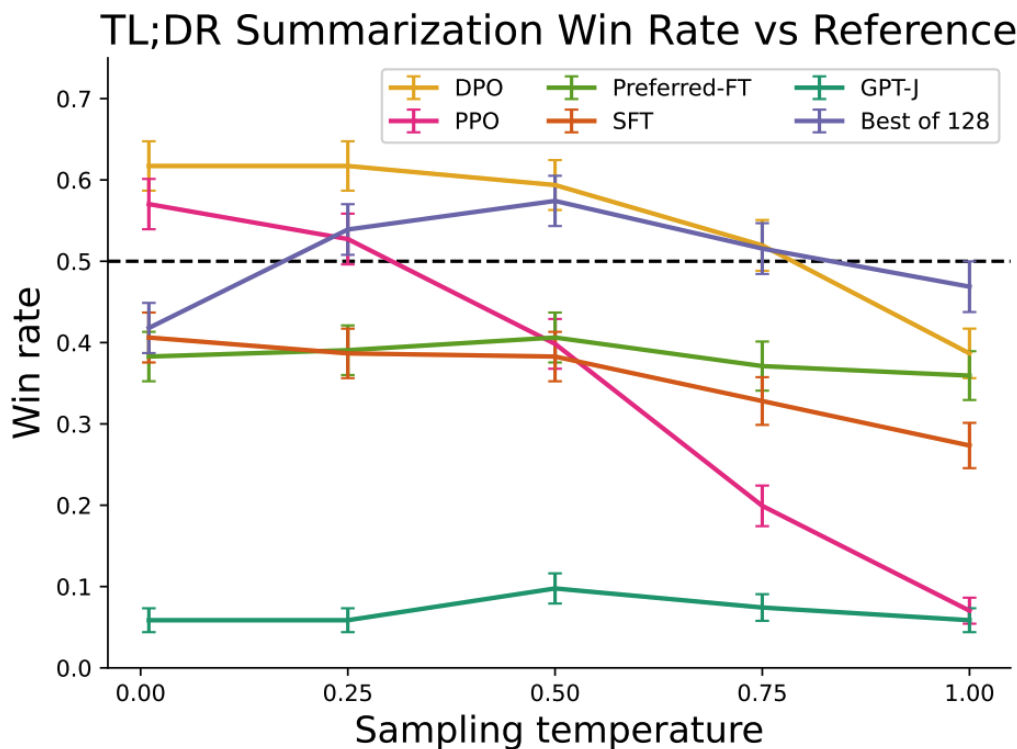1. SFT to obtain $\pi_{\text{SFT}}$
2. Collect preference dataset $\mathcal{D}_p$
3. Use $\mathcal{D}_p$ to optimize the DPO objective and obtain $\pi_\theta$

**Week 5 Assignment**

- An exercise comparing DPO (*offline* method) and RLHF (*online* RL approach)

# DPO vs. RLHF (PPO)

- DPO can have performance comparable to PPO and is simpler to implement



TL;DR Summarization Win Rate vs Reference

**Week 5 Assignment**

- An implementation exercise comparing preference-based tuning vs. SFT

Ref: [Rafailov et al., 2023]

# References

- Ouyang et al., Training language models to follow instructions with human feedback, 2022.

- Hu et al., LoRA: Low-Rank Adaptation of Large Language Models, 2021.

- Book: Jurafsky and Martin, Speech and Language Processing, 2024.

- Ahmadian et al., Back to Basics: Revisiting REINFORCE-Style Optimization for Learning from Human Feedback in LLMs, 2024.

- Stiennon et al., Learning to summarize from human feedback, 2020.

- Huang et al., The N+ Implementation Details of RLHF with PPO: A Case Study on TL;DR Summarization, 2024.

- Bai et al., Training a Helpful and Harmless Assistant with Reinforcement Learning from Human Feedback, 2022.

- Christiano et al., Deep Reinforcement Learning from Human Preferences, 2017.

- Rafailov et al., Direct Preference Optimization: Your Language Model is Secretly a Reward Model, 2023.