# The Information Geometry of RL from Human Feedback

## Gokul Swamy

# Outline for Today

1. What is the fine-tuning problem?

2. End-to-end, what is the two-stage RLHF process doing?

3. What are direct alignment algorithms?
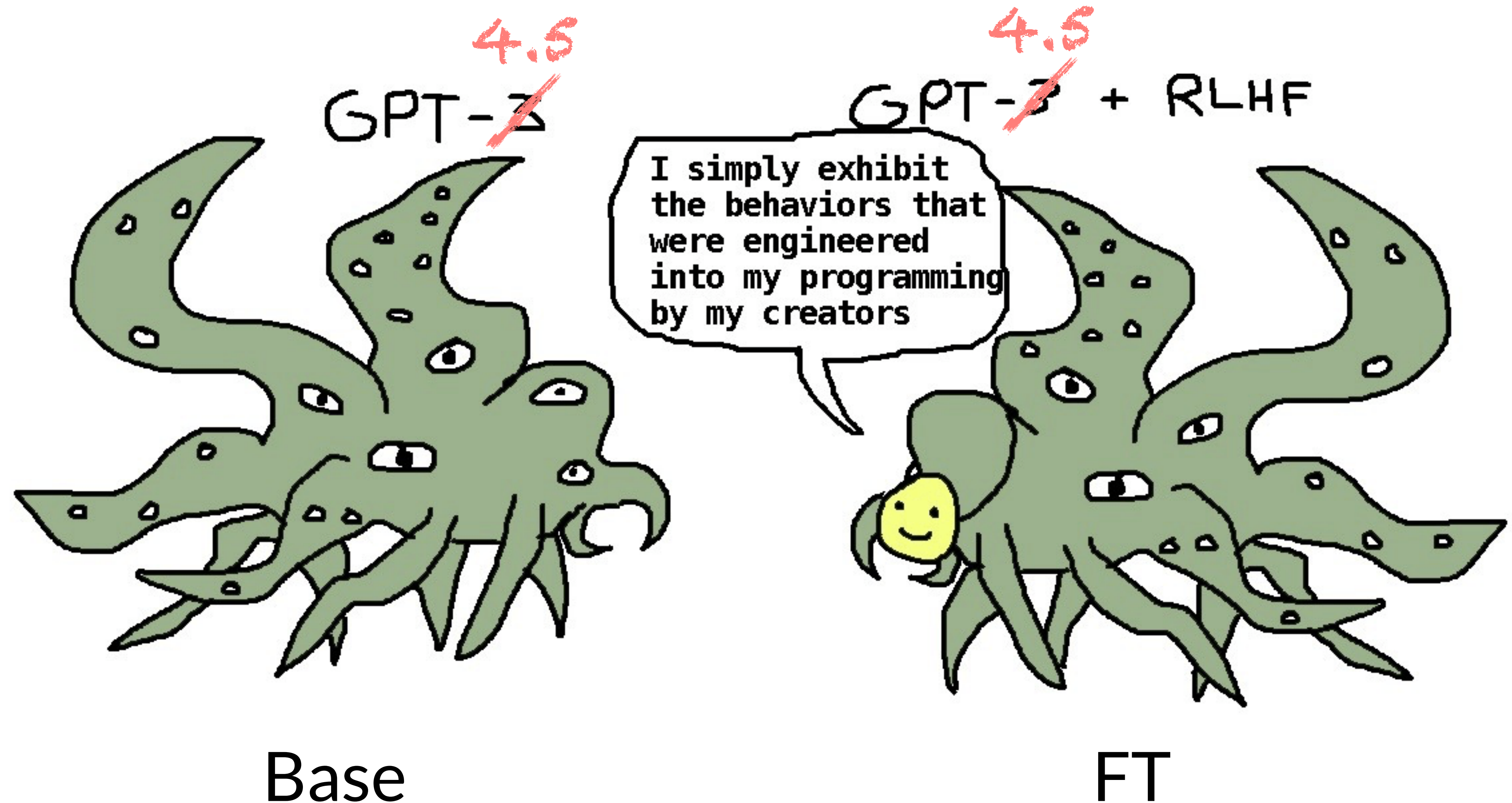
# Outline for Today

1. What is the fine-tuning problem?
   *A: Regularized maximum likelihood estimation.*

2. End-to-end, what is the two-stage RLHF process doing?

3. What are direct alignment algorithms?

# We Live in the Era of Fine-Tuning

# We Live in the Era of Fine-Tuning



Base

FT

[Oertell et al.]

# We Live in the Era of Fine-Tuning

**Prompt:**
What is the purpose of the list C in the code below?

```
def binomial_coefficient(n, r):
    C = [0 for i in range(r + 1)];
    C[0] = 1;
    for i in range(1, n + 1):
        j = min(i, r);
        while j > 0:
            C[j] += C[j - 1];
            j -= 1;
    return C[r]
```

**GPT-3 175B completion:**
A. to store the value of C[0]
B. to store the value of C[1]
C. to store the value of C[i]
D. to store the value of C[i - 1]

**InstructGPT 175B completion:**
The list C in this code is used to store the values of the binomial coefficient as the function iterates through the values of n and r. It is used to calculate the value of the binomial coefficient for a given value of n and r, and stores the result in the final return value of the function.

Base                                        FT

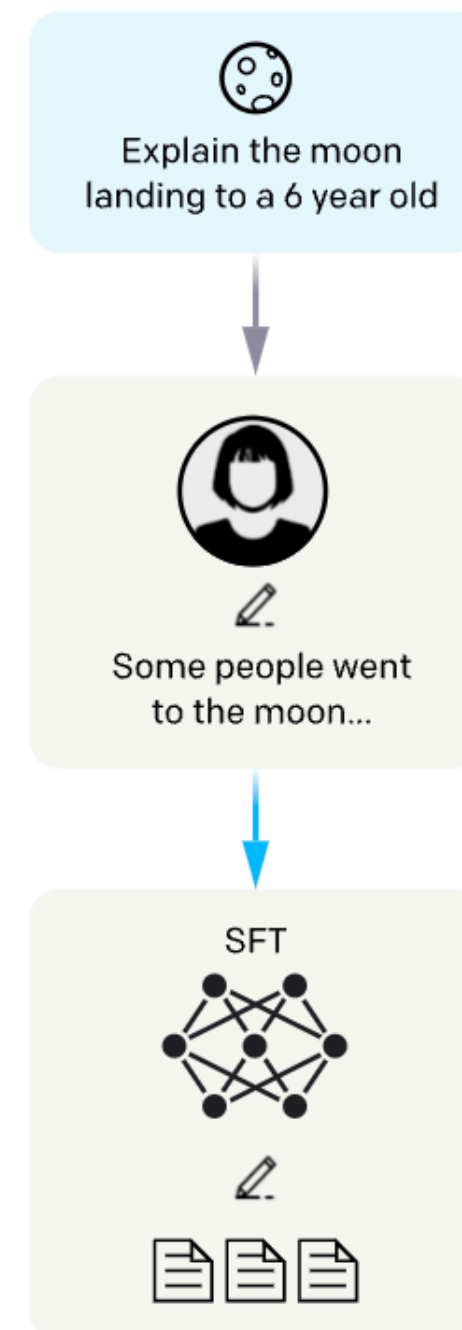[Ouyang et al.]

# We Live in the Era of Fine-Tuning



*Handwritten annotations on the diagram:*

SFT / IL

RM / Classifier

RL

PFT

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...

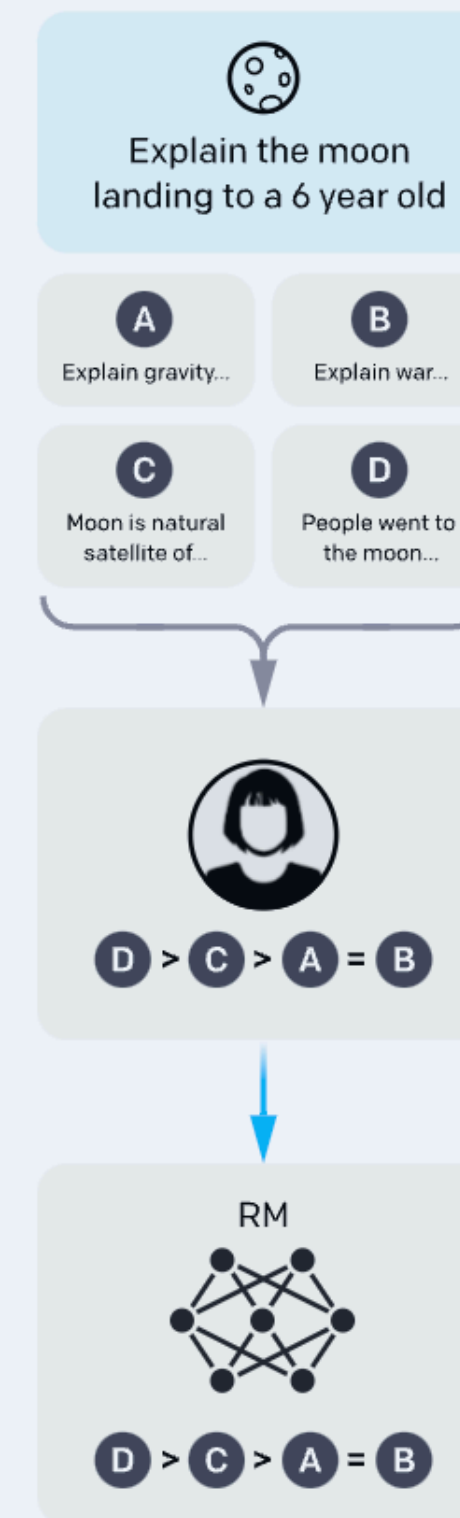This data is used to fine-tune GPT-3 with supervised learning.

SFT

BC

$\pi_{ref}$

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.**
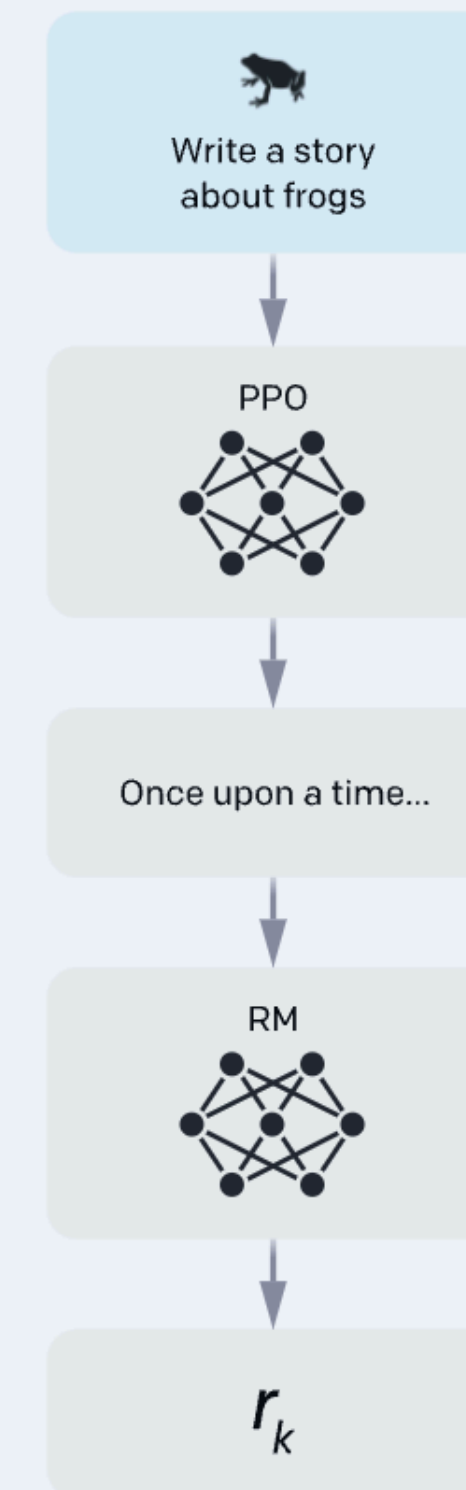
A new prompt is sampled from the dataset.

Write a story about frogs

The policy generates an output.

PPO

→ REBEL / REINFORCE / 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.

$r_k$

# Language Modeling as an MDP

$$s_1 = [s_0, a_1], \; s_2 = [s_0, \; a_1, \; a_2] \quad \cdots \quad s_H = [s_0, a_1 \cdots a_H]$$

$\hookrightarrow$ reset = generating from a prefix



**Prompt**

$s_0 \sim \rho_0$

**Completion** $(s \sim \pi \mid s_0)$

$T(s' \mid s, a) = 1$ if $s' = s \cdot a$

$0$ o/w

$\hookrightarrow$ deterministic, tree-structured

**Reward**

$r(s_H)$

| Token 1 | Token 2 | ... | Token $H$ |
|---|---|---|---|
| $a_1 \in A$ | $a_2 \in A$ | | $a_H \in A$ |

# What makes the Language MDP Special

1. Dynamics are deterministic, known, and tree-structured.

2. Resets are just generating from a prefix — easy to do.

3. The reward function is non-Markovian and doesn't decompose into token-wise rewards.

# Preference Fine-Tuning

$D = \{s_0, \xi^+, \xi^-\}$

Prompt
$(s_0 \sim \rho_0)$

Completion 1 $(\xi_1 \sim \pi_{ref} \mid s_0)$

| Token 1 $(a_1 \in \mathscr{A})$ | Token 2 $(a_2 \in \mathscr{A})$ | Token 3 $(a_3 \in \mathscr{A})$ |
|---|---|---|

Completion 2 $(\xi_2 \sim \pi_{ref} \mid s_0)$

| Token 1 $(a_1' \in \mathscr{A})$ | Token 2 $(a_2' \in \mathscr{A})$ | Token 3 $(a_3' \in \mathscr{A})$ |
|---|---|---|

(+) easier data collection

(−) one bit of information

Preference
$\xi_1 \succ \xi_2$

# *Preference* Fine-Tuning

***Goal****: Maximize the relative likelihood of preferred to dis-preferred completions.*

FKL

RKL

$$\pi^\star = \arg\min_{\pi \in \Pi} \mathbb{D}_{KL}\left(\mathscr{D} \,||\, \pi\right) \;+\; \mathbb{D}_{KL}\left(\pi \,||\, \pi_{\text{ref}}\right)$$

*(Data Likelihood)*

l initul
lonrusl

*(Prior Reg.)*

# Outline for Today

1. What is the fine-tuning problem?
   *A: Regularized maximum likelihood estimation.*

2. End-to-end, what is the two-stage RLHF process doing?
   *A: MLE over reward models followed by MaxEnt over policies.*

3. What are direct alignment algorithms?

# Notation

For simplicity, we're going to assume the "Bradley-Terry" model of preferences:

*(handwritten annotation: assuming transitivity ✓ / all raters mostly agree)*

$$\mathbb{P}_r(\xi_1 > \xi_2 \,|\, s_0) = \sigma(r(\xi_1) - r(\xi_2))$$

Also, let's denote the empirical preference distribution as:

$$\mathbb{P}_{\mathscr{D}}(\xi_1 > \xi_2 \,|\, s_0)$$

i.e. how often raters preferred $\xi_1$ to $\xi_2$ given prompt $s_0$.

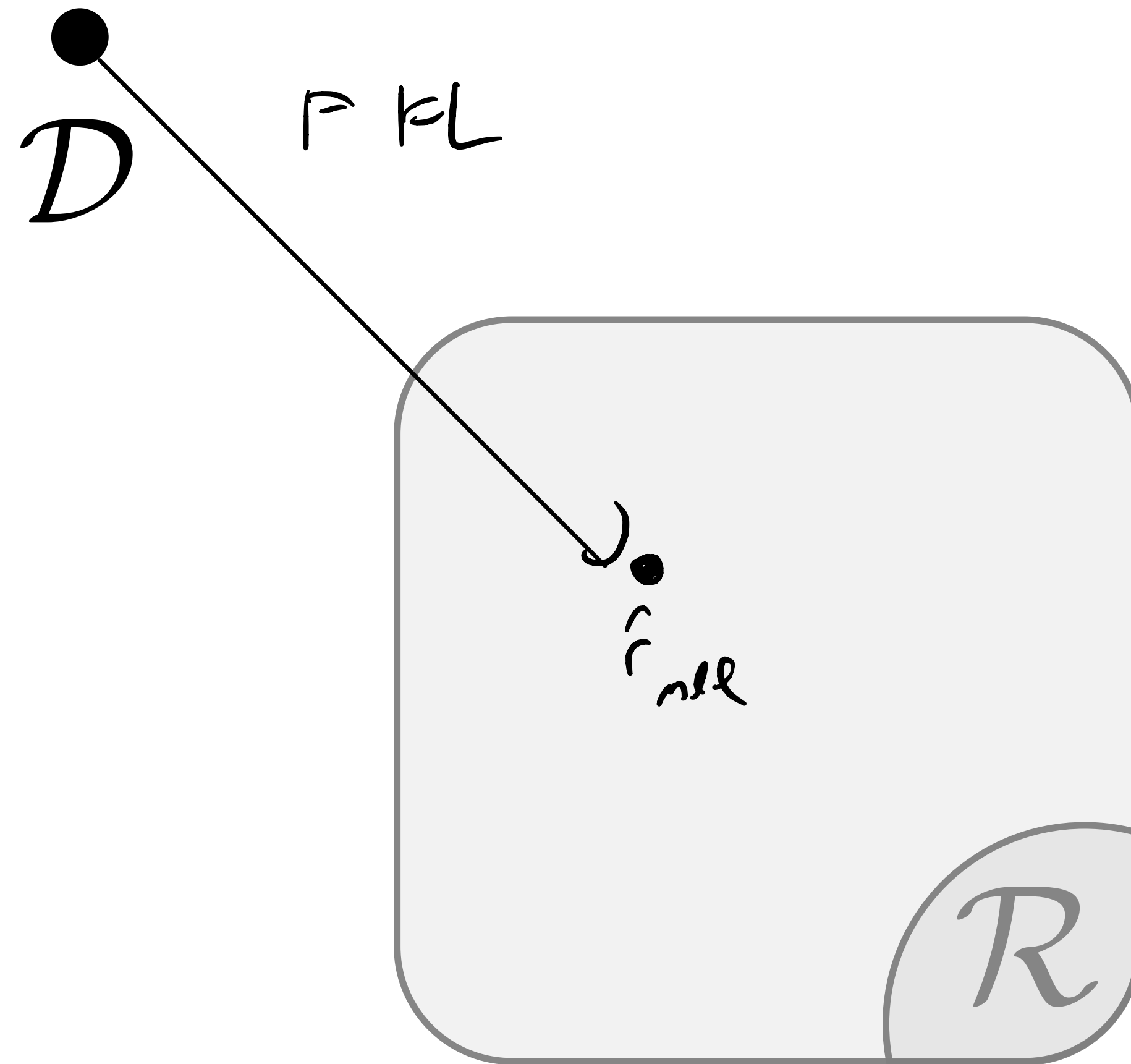# Reward Modeling is MLE

Then,

$$\hat{r}_{\text{mle}} = \arg\min_{r \in \mathcal{R}} \mathbb{E}_{s_0 \sim \mathcal{D}}[\mathbb{D}_{KL}(\mathbb{P}_{\mathcal{D}} \,||\, \mathbb{P}_r)]$$

*(handwritten annotations, blue: "$\leftarrow$ Ftormul PL")*

*(handwritten, blue): $D_{KL}(P||Q) = \mathbb{E}_P[\log P] - \mathbb{E}_P[\log Q]$*

$$= \arg\max_{r \in \mathcal{R}} \mathbb{E}_{(s_0, \xi^+, \xi^-) \sim \mathcal{D}}[\log \mathbb{P}_r(\xi^+ \succ \xi^- \,|\, s_0)]$$

*(handwritten, blue): $BT$*

$$= \arg\max_{r \in \mathcal{R}} \mathbb{E}_{(s_0, \xi^+, \xi^-) \sim \mathcal{D}}[\log \sigma(r(\xi^+) - r(\xi^-))]$$

*This is just logistic regression / classification!*

# Reward Modeling is a FKL Projection onto $\mathcal{R}$
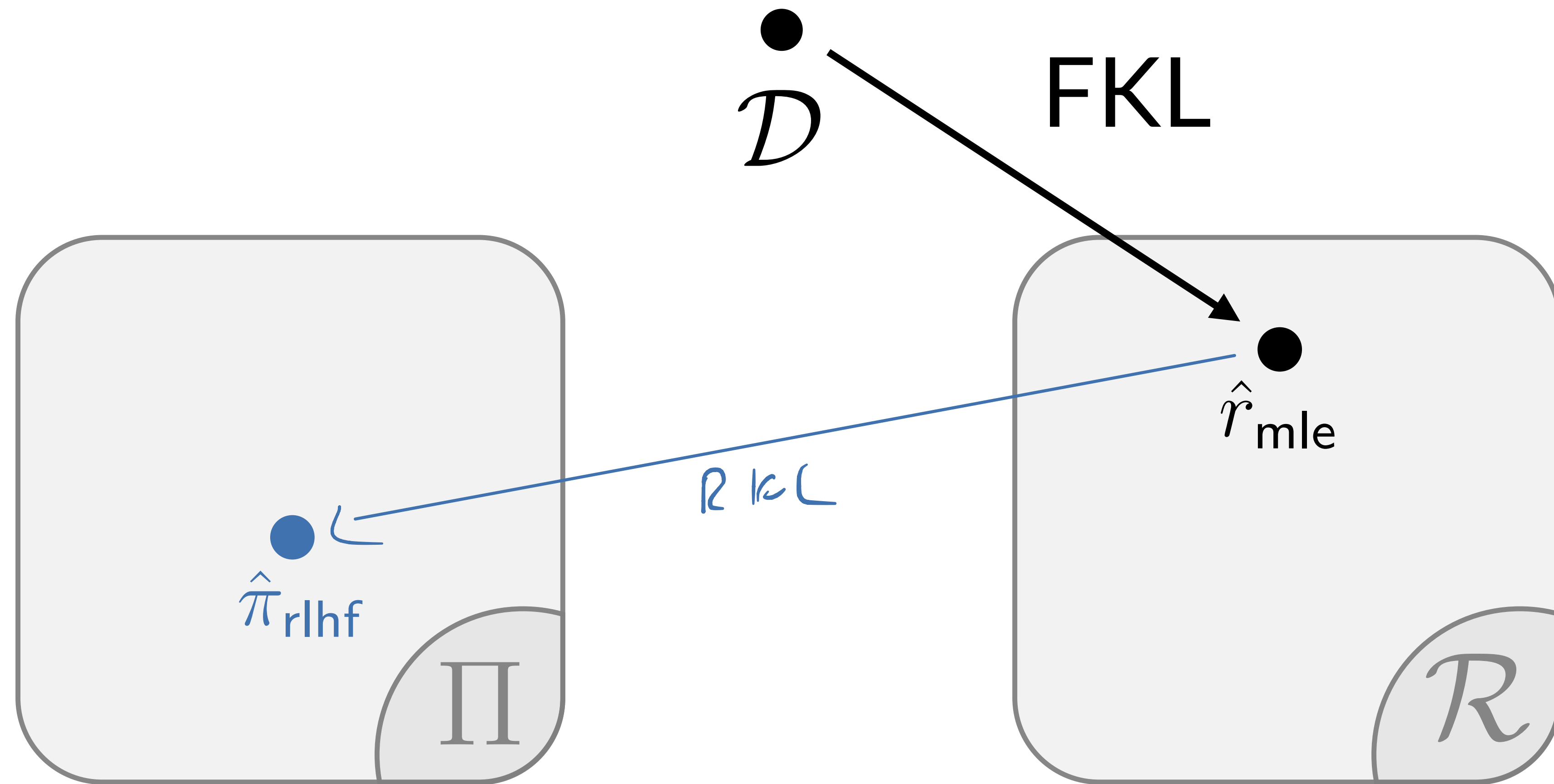
# *Recap*: "Soft" / Entropy Regularized RL

$$\hat{\pi}_{\mathsf{rlhf}} = \arg\max_{\pi \in \Pi} \mathbb{E}_{\xi \sim \pi}[\hat{r}_{\mathsf{mle}}(\xi)] + \underbrace{\mathbb{D}_{KL}(\pi \,||\, \pi_{\mathsf{ref}})}$$

$$\mathbb{E}_{\xi \sim \pi}\left[ \sum_{h}^{H} \log\left( \frac{\pi(a_h \,|\, s_h)}{\pi_{\mathsf{ref}}(a_h \,|\, s_h)} \right) \right]$$

deterministic dynamics

$$\prod_{h}^{H} \pi_r^{\star}(a_h \,|\, s_h) = \mathbb{P}_{\hat{r}}^{\star}(\xi \,|\, s_0) = \frac{\mathbb{P}_{\mathsf{ref}}(\xi) \cdot \exp(\hat{r}(\xi))}{\sum_{\xi' \in \Xi | s_0} \mathbb{P}_{\mathsf{ref}}(\xi') \cdot \exp(\hat{r}(\xi'))}$$

# Soft RL is a *Reverse* KL Projection onto $\Pi$



**E2E**, *(1) RLHF is FKL to $\mathscr{R}$ and (2) RKL to $\Pi$*

# *If* ⏱️: Soft RL is a *Reverse* KL Projection onto $\Pi$

$$\hat{\pi}_{\mathsf{rlhf}} = \arg\min_{\pi \in \Pi} \mathbb{D}_{KL}(\mathbb{P}_{\pi} || \mathbb{P}_{\hat{r}}^{\star})$$

$$= \arg\min_{\pi \in \Pi} \mathbb{E}_{\xi \sim \mathbb{P}_{\pi}} \left[ \log \left( \frac{\mathbb{P}_{\pi}(\xi)}{\mathbb{P}_{\hat{r}}^{\star}(\xi)} \right) \right]$$

$$= \arg\min_{\pi \in \Pi} \sum_{\xi \in \Xi} \mathbb{P}_{\pi}(\xi)(\log \mathbb{P}_{\pi}(\xi) - \log \mathbb{P}_{\hat{r}}^{\star}(\xi))$$

$$= \arg\min_{\pi \in \Pi} \sum_{\xi \in \Xi} \mathbb{P}_\pi(\xi)(\log \mathbb{P}_\pi(\xi) - \hat{r}(\xi) + \log Z_{\hat{r}}^\star)$$

$$= \arg\min_{\pi \in \Pi} \sum_{\xi \in \Xi} \mathbb{P}_\pi(\xi)(\log \mathbb{P}_\pi(\xi) - \hat{r}(\xi))$$

$$= \arg\max_{\pi \in \Pi} \sum_{\xi \in \Xi} \mathbb{P}_\pi(\xi)(-\log \mathbb{P}_\pi(\xi) + \hat{r}(\xi))$$

$$= \arg\max_{\pi \in \Pi} \mathbb{E}_{\xi \sim \pi}[\hat{r}(\xi)] + \mathbb{H}(\pi)$$

# Outline for Today

1. What is the fine-tuning problem?
   *A: Regularized maximum likelihood estimation.*

2. End-to-end, what is the two-stage RLHF process doing?
   *A: MLE over reward models followed by MaxEnt over policies.*

3. What are direct alignment algorithms?
   *A: Algorithms like DPO directly maximize likelihood over $\Pi$ without passing through $\mathscr{R}$.*
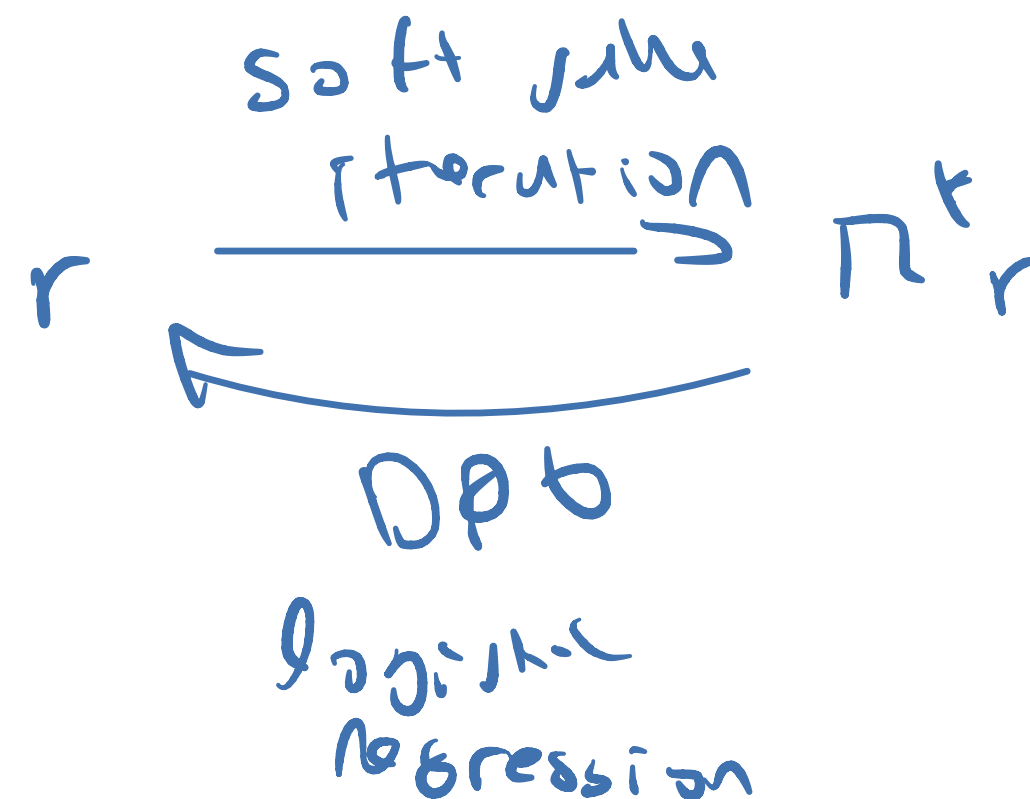
# The DPO "Reparameterization Trick"

$$\prod_h^H \pi_r^\star(a_h \mid s_h) = \frac{\prod_h^H \pi_{\mathsf{ref}}(a_h \mid s_h) \cdot \exp(r(\xi))}{Z(s_0)}$$

taking log of both sides

$$\sum_h^H \log \pi_r^\star(a_h \mid s_h) = \sum_h^H \log \pi_{\mathsf{ref}}(a_h \mid s_h) + r(\xi) - \log Z(s_0)$$

$$r(\xi) = \sum_{h}^{H} \log \pi_r^{\star}(a_h \,|\, s_h) - \log \pi_{\mathsf{ref}}(a_h \,|\, s_h) + \log Z(s_0)$$

$$\triangleq r_{\pi}(\xi)$$



We can express the reward model that makes a policy (soft) optimal in terms of said policy by "inverting" the MaxEnt RL equations!

More explicitly, consider the soft-optimal policy for $r_\pi$:

$$\mathbb{P}^{\star}_{r_\pi}(\xi) \propto \exp(r_\pi(\xi))$$

$$\propto \exp\left(\sum_h^H \log \pi(a_h \mid s_h) + \log Z(s_0)\right)$$

$$\propto \exp\left(\sum_h^H \log \pi(a_h \mid s_h)\right)$$

$$\propto \prod_h^H \pi(a_h \mid s_h)$$

The soft optimal policy for $r_\pi$ is $\pi$, which means we can optimize over $r_\pi$ and get the soft optimal policy "for free"!

Now, we proceed by MLE *directly* over policies:

$$\hat{\pi}_{\mathsf{dpo}} = \arg\max_{\pi \in \Pi} \mathbb{E}_{(s_0, \xi^+, \xi^-) \sim \mathscr{D}} [\log \sigma(r_\pi(\xi^+) - r_\pi(\xi^-))]$$

$$= \arg\max_{\pi \in \Pi} \mathbb{E}_{(s_0, \xi^+, \xi^-) \sim \mathscr{D}} \left[ \log \sigma \left( \sum_h^H \log \frac{\pi(a_h^+ \mid s_h^+)}{\pi_{\mathsf{ref}}(a_h^+ \mid s_h+)} - \log \frac{\pi(a_h^- \mid s_h^-)}{\pi_{\mathsf{ref}}(a_h^- \mid s_h^-)} \right) \right]$$

So, we end up with a single-step MLE procedure!

# DPO is a FKL Projection onto $\Pi$