# The Value of RL in Fine-Tuning

## Gokul Swamy

# Outline for Today

1. What assumption on the preference dataset did we make in the DPO derivation and what happens when it breaks?

2. When are two-stage RLHF and DPO equivalent?

3. Why does two-stage RLHF work *much* better in practice?

# *Recap*: Forward and Reverse KL

= mode covering "

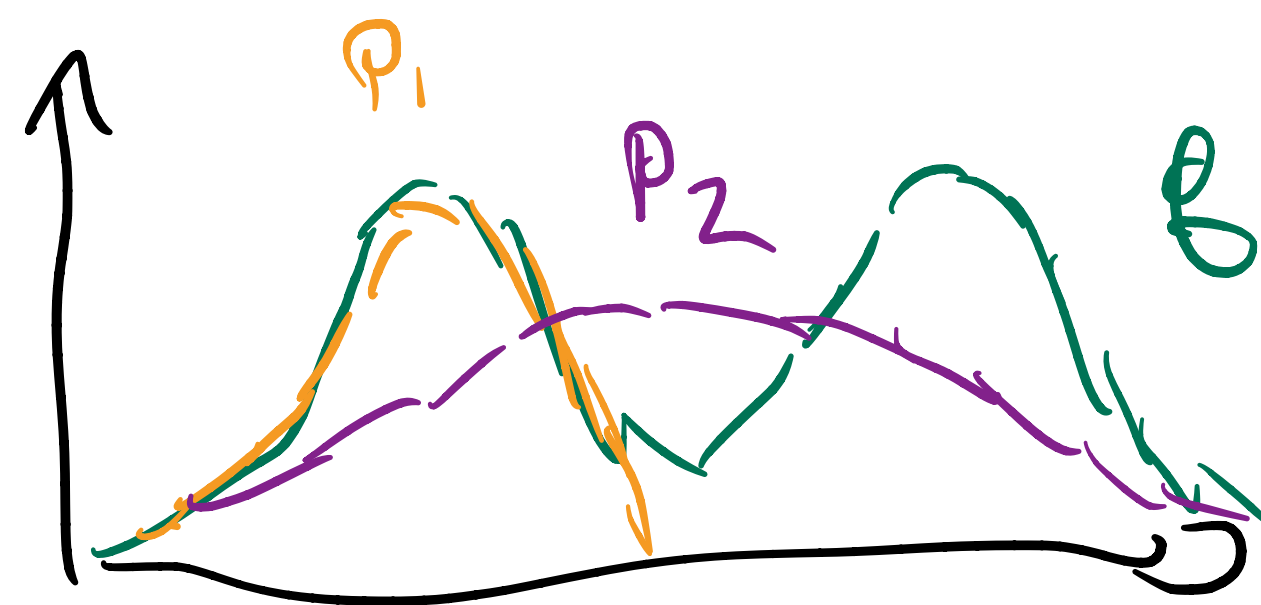**FKL**: $\min_p \mathbb{D}_{\mathsf{KL}}(q||p) = \min_p \sum_x q(x)\log\left(\dfrac{q(x)}{p(x)}\right)$   MLE

**RKL**: $\min_p \mathbb{D}_{\mathsf{KL}}(p||q) = \min_p \sum_x p(x)\log\left(\dfrac{p(x)}{q(x)}\right)$   Soft RC
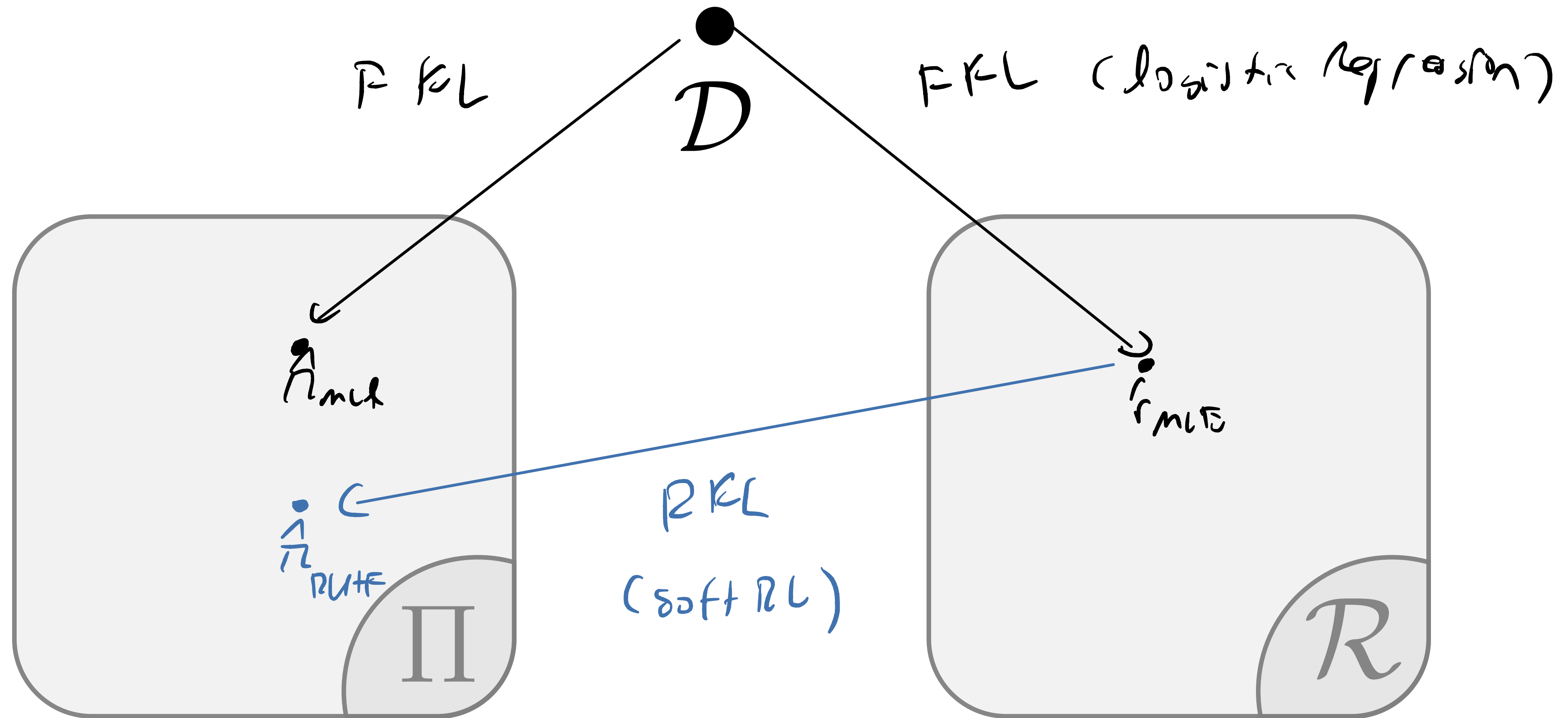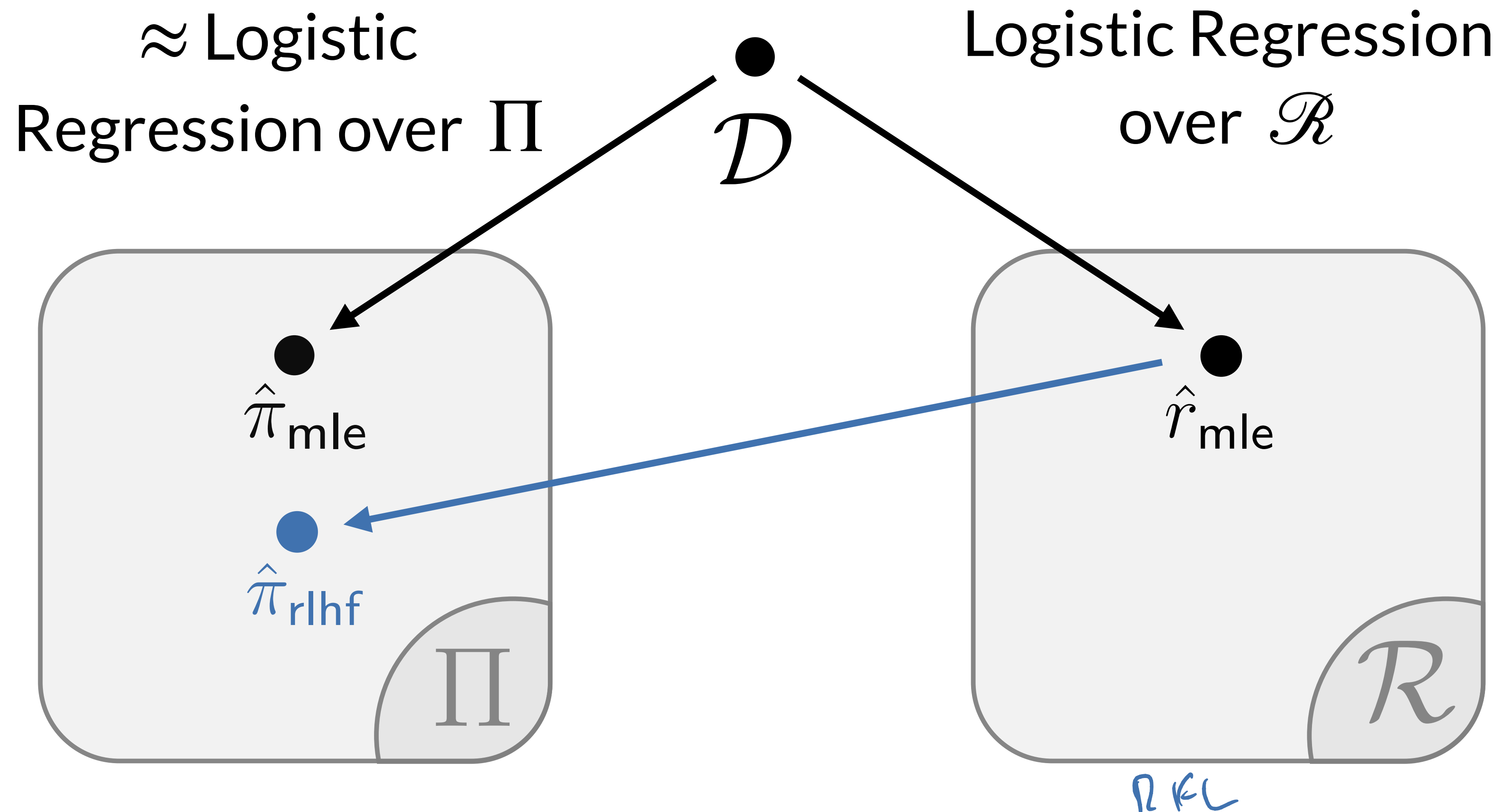
= mode seeking

↓

= mode selection



*FKL is "mode-covering" while RKL is "mode-seeking"*

# *Recap*: Information Geometry of RLHF

# *Recap*: Information Geometry of RLHF

$\approx$ Logistic
Regression over $\Pi$

Logistic Regression
over $\mathscr{R}$

$\mathcal{D}$

$\hat{\pi}_{\mathsf{mle}}$

$\hat{r}_{\mathsf{mle}}$

$\hat{\pi}_{\mathsf{rlhf}}$

$\Pi$

$\mathcal{R}$

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

# Outline for Today

1. What assumption on the preference dataset did we make in the DPO derivation and what happens when it breaks?
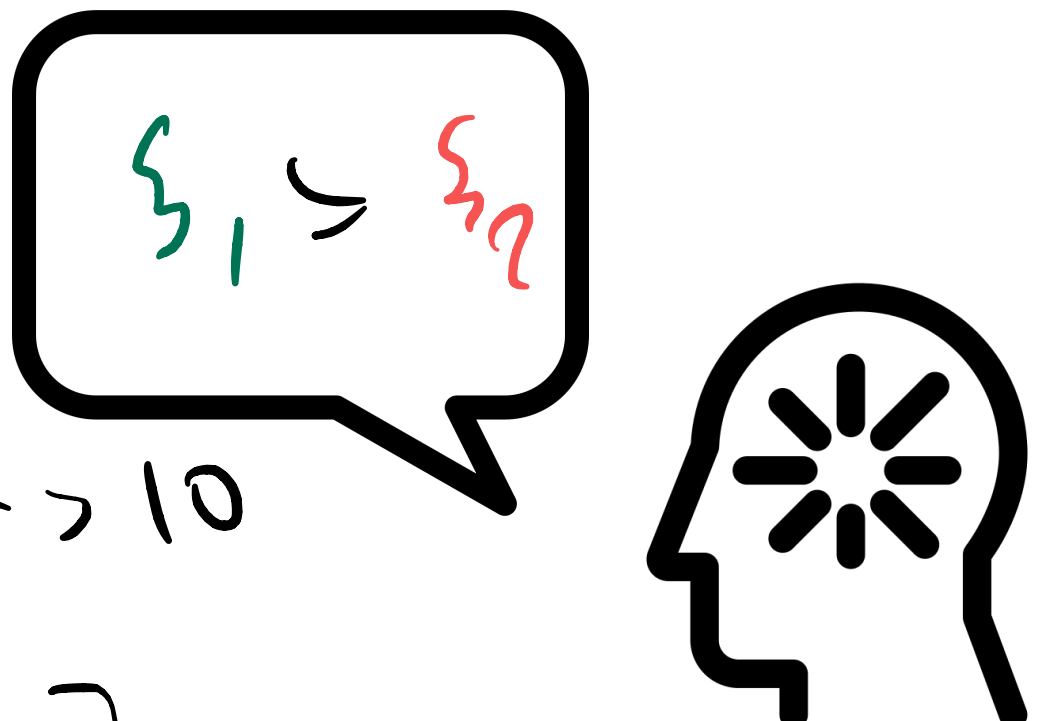
   *A: Full coverage of $\mathcal{D}$. Without it, we can't control the RKL.*

2. When are two-stage RLHF and DPO equivalent?

3. Why does two-stage RLHF work *much* better in practice?

# Why does DPO break with Partial Coverage?

| | $\xi_1$ | $\xi_2$ | $\xi_3$ |
|---|---|---|---|
| $\pi_{\text{ref}}$ | 0.5 | 0.5 | 0 |
| $\hat{\pi}_{\text{rlhf}}$ | | | |
| $\hat{\pi}_{\text{dpo}}$ | | | |

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

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

*DPO Doesn't Regularize to $\pi_{ref}$ and can produce OOD responses.*

# Outline for Today

1. What assumption on the preference dataset did we make in the DPO derivation and what happens when it breaks?
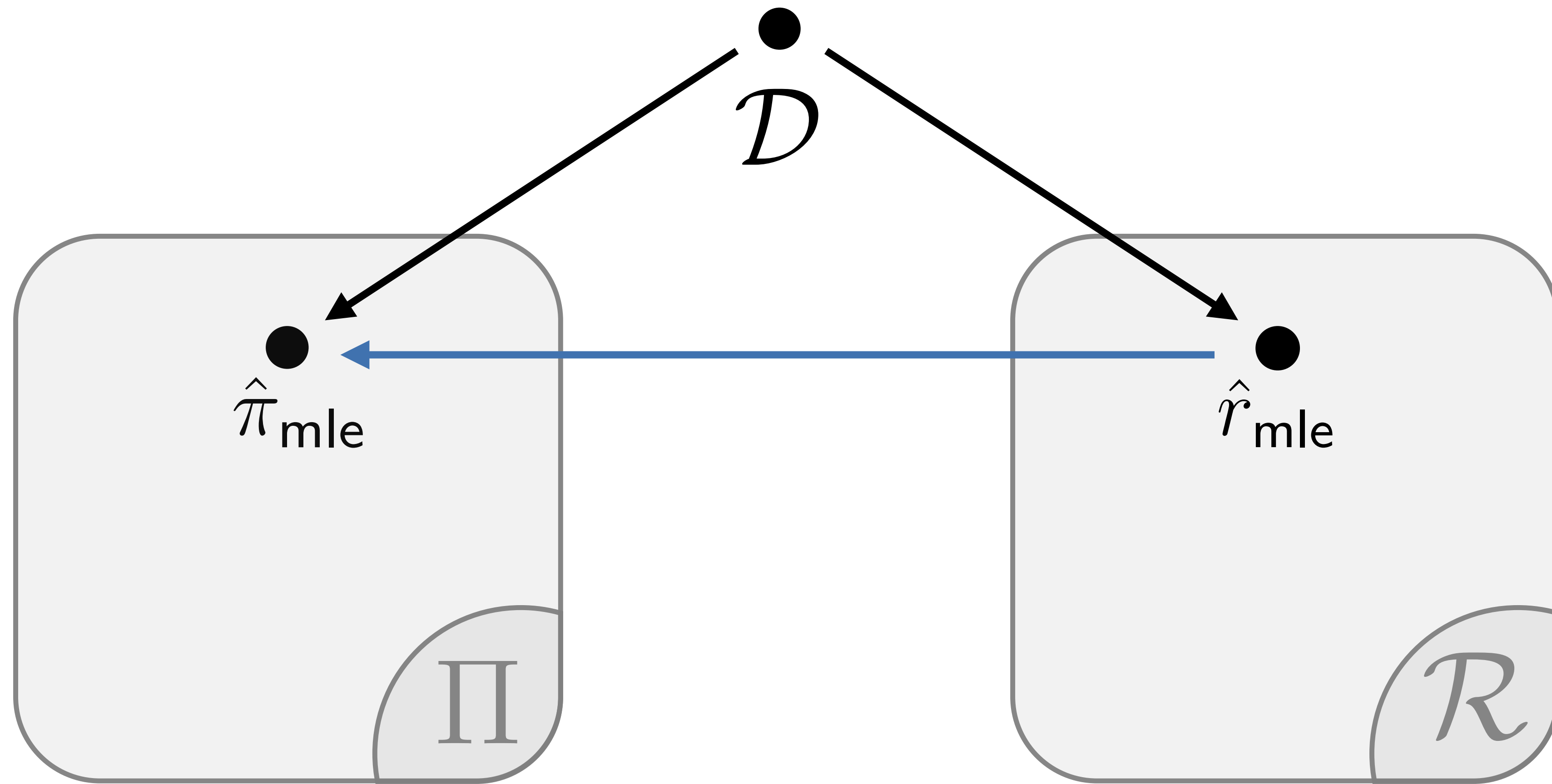
*A: Full coverage of $\mathcal{D}$. Without it, we can't control the RKL.*

2. When are two-stage RLHF and DPO equivalent?

*A: When $\Pi$ and $\mathcal{R}$ are isomorphic and all projections are exact.*

3. Why does two-stage RLHF work *much* better in practice?

# When are two-stage RLHF and DPO equivalent?



**A:** *If (1) $\Pi \Leftrightarrow \mathcal{R}$ and (2) all projections exact.*

# If ⏱: RLHF = DPO when $\Pi \Leftrightarrow \mathscr{R}$

$$\mathbb{E}_{\mathscr{D}}\left[\log\sigma\left(\sum_h^H \log\hat{\pi}_{\mathsf{mle}}(a_h^+\,|\,s_h^+) - \log\hat{\pi}_{\mathsf{mle}}(a_h^-\,|\,s_h^-)\right)\right] = \mathbb{E}_{\mathscr{D}}\left[\log\sigma\left(\hat{r}_{\mathsf{mle}}(\xi^+) - \hat{r}_{\mathsf{mle}}(\xi^-)\right)\right]$$

($\hat{r}_{\mathsf{mle}}$

(minimizing same fn.
over the same class)

(uniqueness
of minimizer)

$$\Rightarrow \forall \xi \in \Xi, \ \sum_h^H \log\hat{\pi}_{\mathsf{mle}}(a_h\,|\,s_h) = \hat{r}_{\mathsf{mle}}(\xi)$$

(soft RL = RL)

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

the same
function

$$= \arg\min_{\pi\in\Pi} \mathbb{D}_{KL}\left(\mathbb{P}_\pi\,||\,\mathbb{P}_{r_{\hat{\pi}}}\right)$$

$$= \hat{\pi}_{\mathsf{mle}}$$

$\hat{\pi}_{\mathsf{mle}} \in \Pi$

NPG
covariant

*MLE is invariant to reparameterization.*

# *Policies vs. Reward Models*

**Policies**: $\pi : \mathcal{S} \rightarrow \Delta(\mathscr{A}) \in \Pi$

*(Prefixes)*     *(Tokens)*

**Rewards**: $r : \Xi \rightarrow \mathbb{R} \in \mathscr{R}$

*(Completions)*

**Both of these are fine-tuned from the same SFT checkpoint!**

# Outline for Today

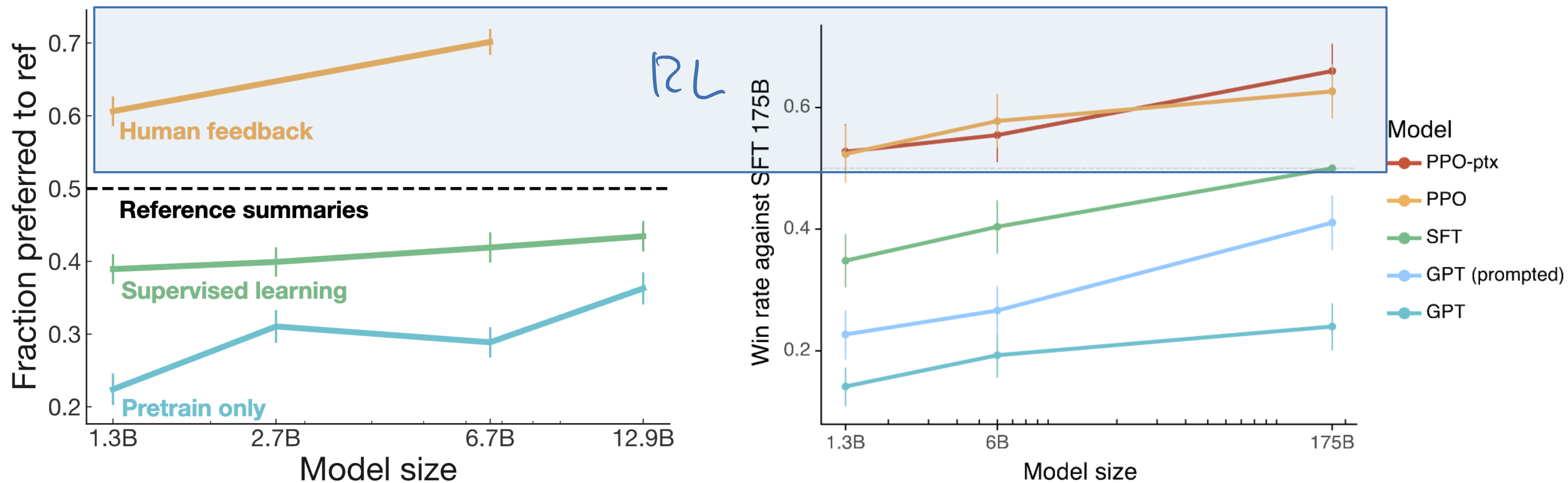1. What assumption on the preference dataset did we make in the DPO derivation and what happens when it breaks?

   **A:** *Full coverage of $\mathcal{D}$. Without it, we can't control the RKL.*

2. When are two-stage RLHF and DPO equivalent?

   **A:** *When $\Pi$ and $\mathcal{R}$ are isomorphic and all projections are exact.*
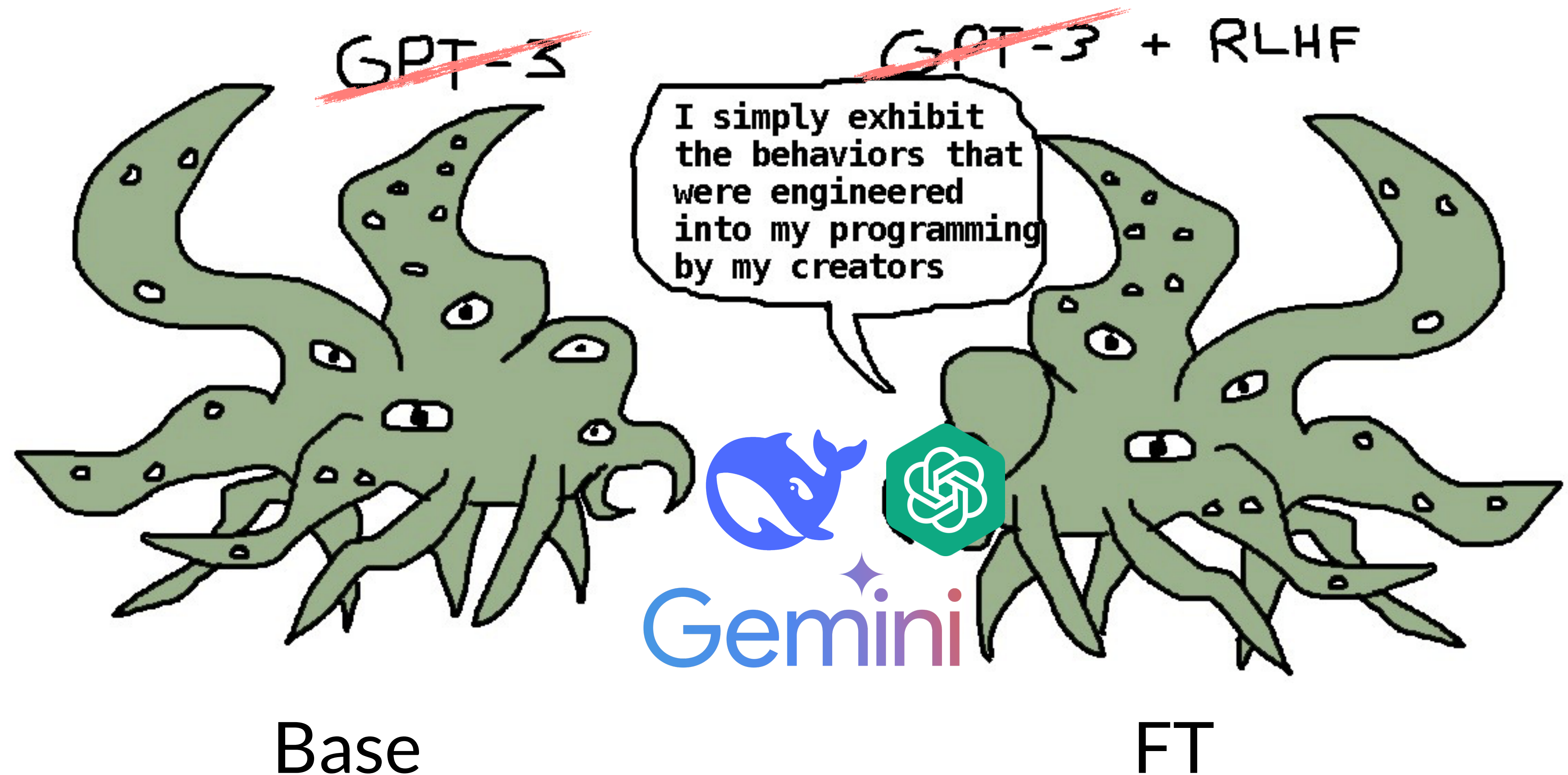
3. Why does two-stage RLHF work *much* better in practice?

# Two Stage RLHF > DPO



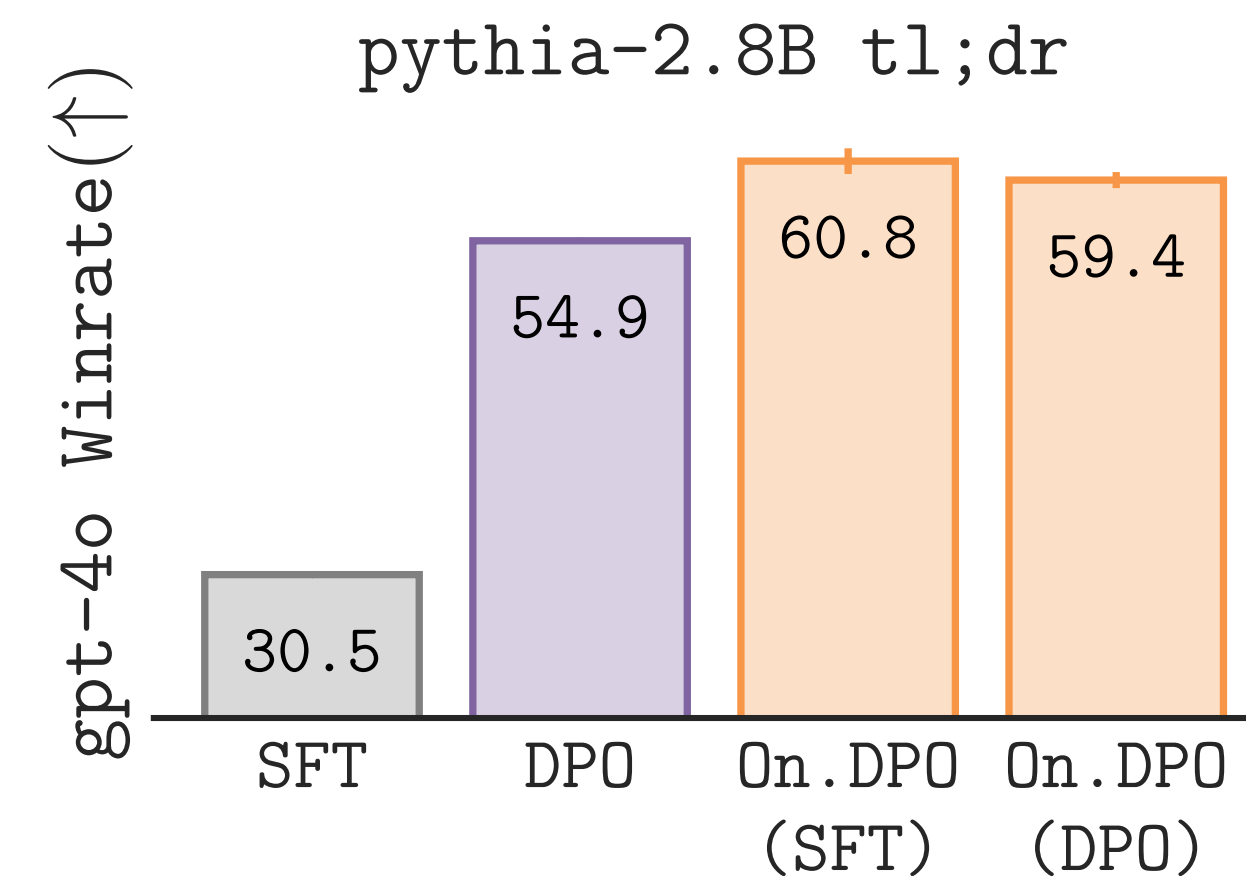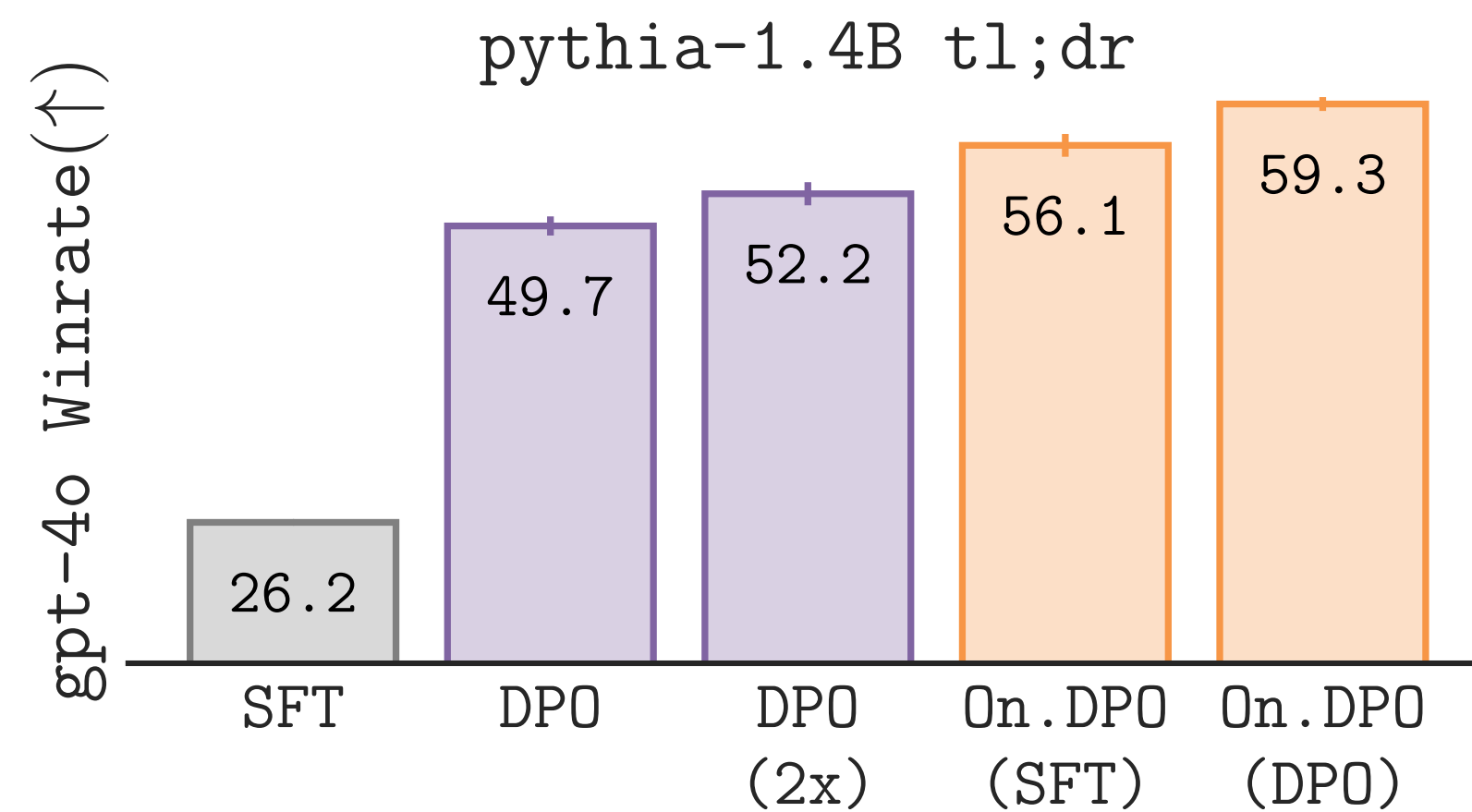[Stiennon et al., Ouyang et al.]

# Two Stage RLHF > DPO

# Grounding the Gap in Summarization

1. We will focus on the task of summarization of Reddit posts, using models from the Pythia family pre-trained on the Pile.

2. We will use the *same* dataset to train both policies and reward models.

3. We will start from the *same* SFT checkpoint to train both.

4. We will use the *same* optimizer (DPO) for both online and offline PFT with the *same* hyperparameters.

# Gap Appears in "🍎s to 🍎s" Comparison



pythia-1.4B tl;dr

gpt-4o Winrate(↑)

| SFT | DPO | DPO (2x) | On.DPO (SFT) | On.DPO (DPO) |
|-----|-----|----------|--------------|--------------|
| 26.2 | 49.7 | 52.2 | 56.1 | 59.3 |

pythia-2.8B tl;dr

gpt-4o Winrate(↑)

| SFT | DPO | On.DPO (SFT) | On.DPO (DPO) |
|-----|-----|--------------|--------------|
| 30.5 | 54.9 | 60.8 | 59.4 |

# 6 Hypotheses for the Online-Offline Gap

$H_1$

$H_2$

$H_3$

$H_4$

$H_5$

$H_6$

# $\mathbb{H}_1$: Intrinsic Value of On-Policy Feedback

*... but the on-policy labels are just imputed*

*... from an RM trained on the same data as the policy*

*... and we can't create any new info via sampling.*

**Sasha Rush** ✓
@srush_nlp

Lot of pitches this week for "perpetual data machines". Either laundering self-generated data or attributing prescience to reward models. Just want to caution that is a common trap smart people fall for.

9:58 AM · Dec 15, 2023 · **139.4K** Views

# $\mathbb{H}_2$: Failure of Offline Regularization to $\pi_{\mathsf{ref}}$

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

*Reverse KL has an on-policy expectation.*

*... but we used the same regularizer for all experiments.*

[Song et al.]

# $\mathbb{H}_5$ : RMs Generalize Better OOD



## ... but they also generalize better ID

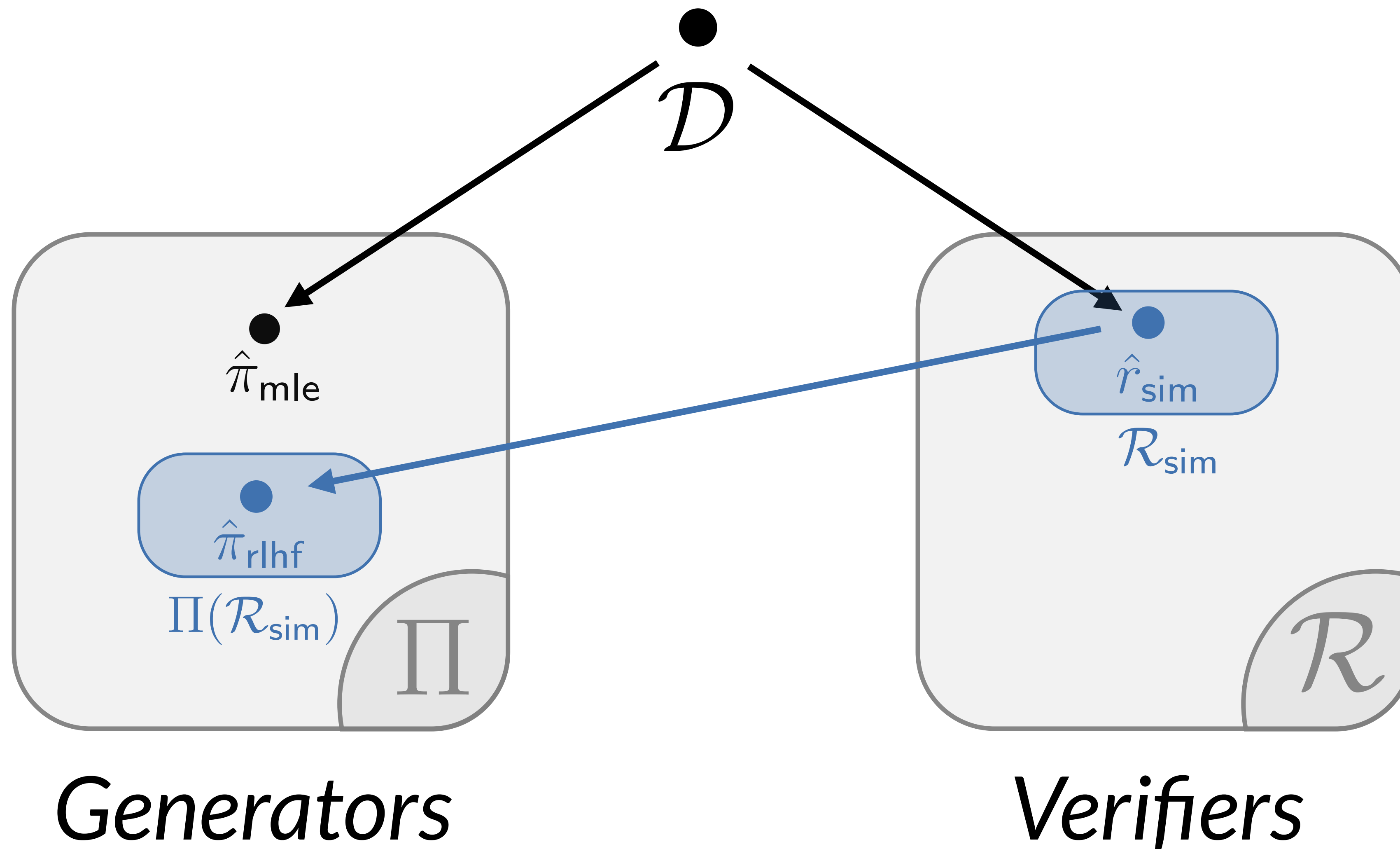# Generation-Verification Gaps

$\hookrightarrow$ P != NP



*GV Gap = easier to check than to solve!*

# $\mathbb{H}_6$: *Proper* Learning w/ a *Generation-Verification Gap*



*Only need to search over $\Pi\left(\mathscr{R}_{sim}\right) \subset \Pi$!*

$\mathbb{H}_6$: *Proper* Learning w/ a *Generation-Verification Gap*

$\mathcal{D}$

$\hat{\pi}_{\mathsf{mle}}$

$\hat{r}_{\mathsf{sim}}$

$\mathcal{R}_{\mathsf{sim}}$

$\hat{\pi}_{\mathsf{rlhf}}$

$\Pi(\mathcal{R}_{\mathsf{sim}})$
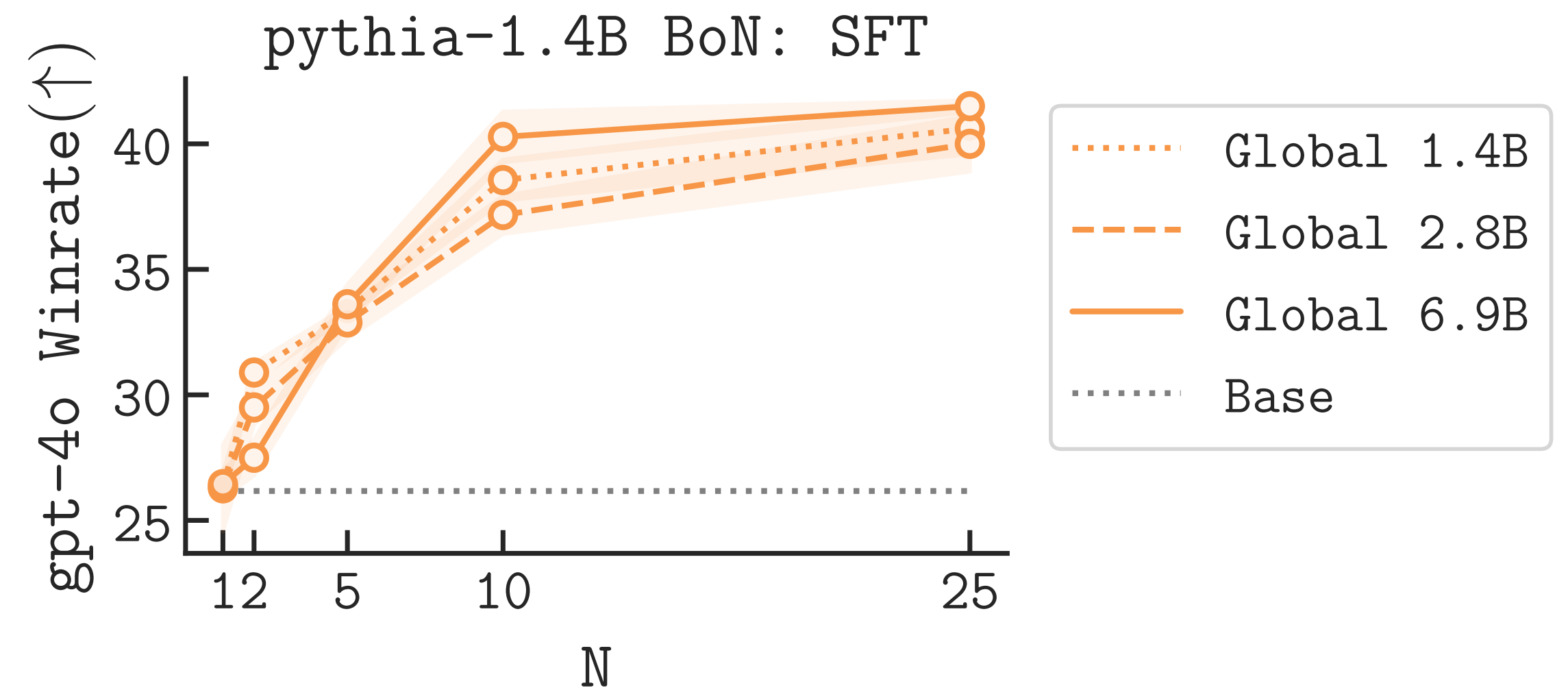
$\Pi$

$\mathcal{R}$

*Generators*

*Verifiers*

*All roads lead to likelihood, but RL takes a shortcut!*

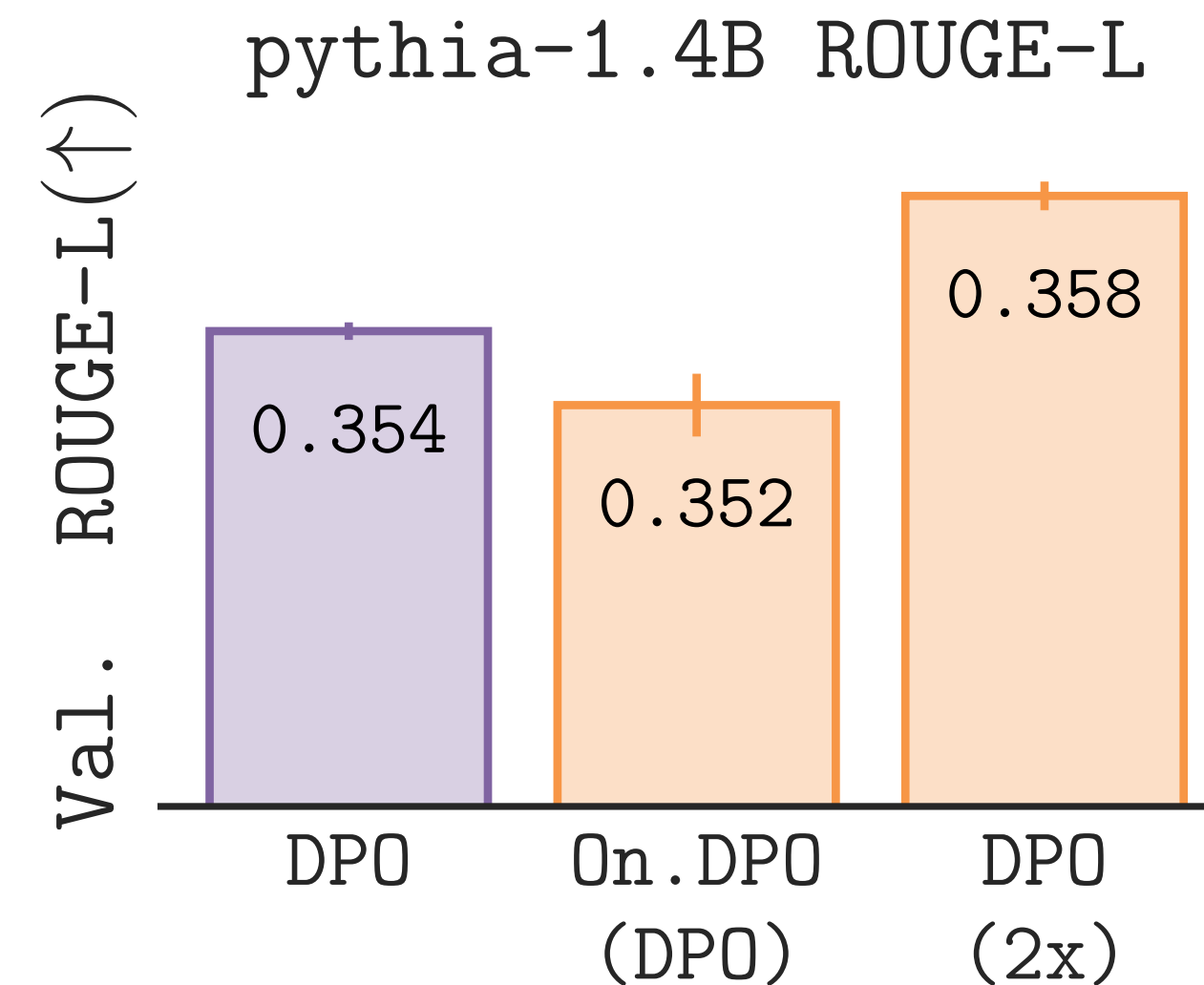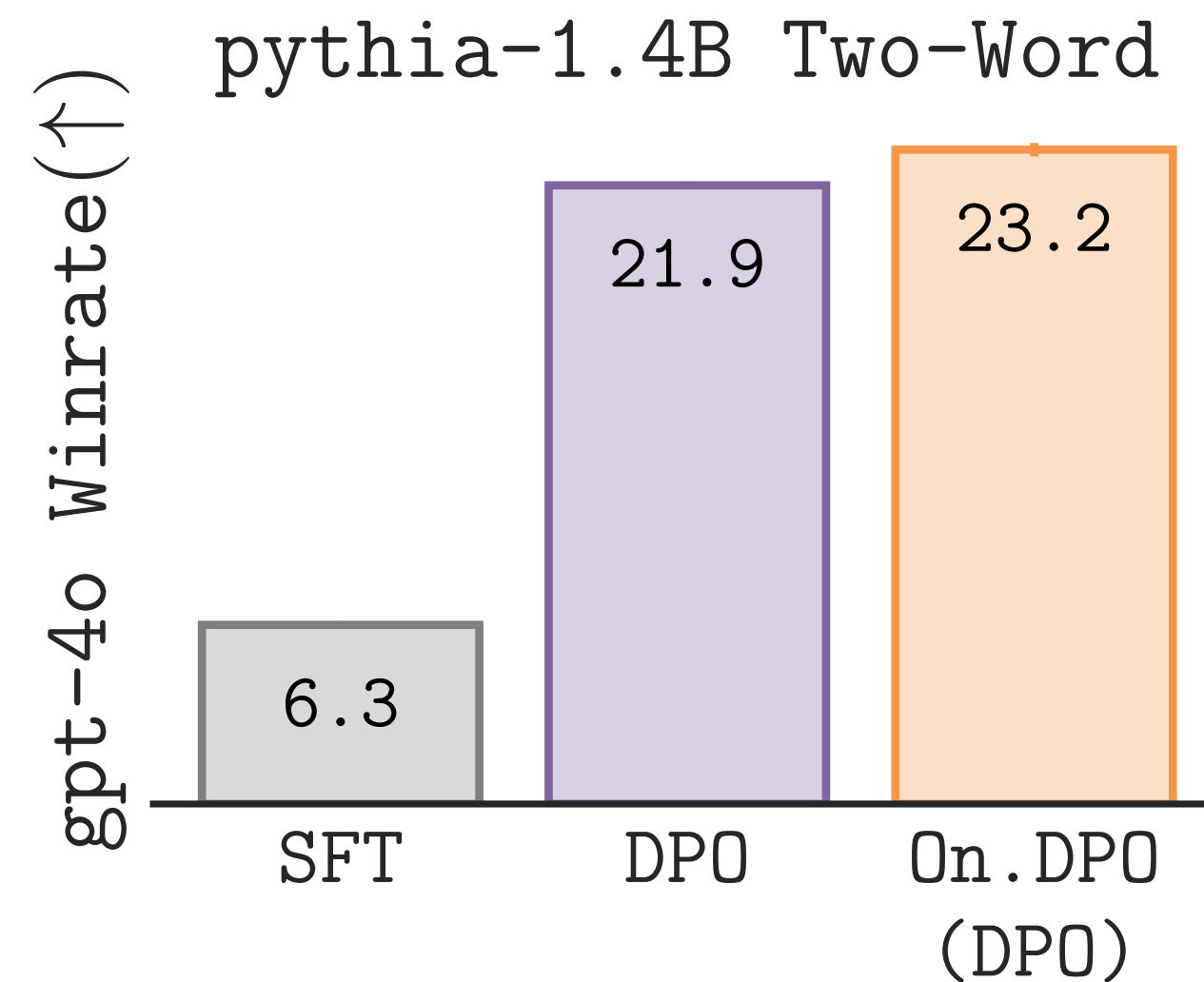# Evidence for Generation-Verification Gap



*Using a much smaller RM than policy doesn't hurt.*

*Using a much larger RM than policy doesn't help.*

# Closing the Generation-Verification Gap

*Simplify Policy:*   *Complicate Reward:*



*Online PFT ≈ Offline PFT with no generation-verification gap!*

# Outline for Today

1. What assumption on the preference dataset did we make in the DPO derivation and what happens when it breaks?

*A: Full coverage of $\mathcal{D}$. Without it, we can't control the RKL.*

2. When are two-stage RLHF and DPO equivalent?

*A: When $\Pi$ and $\mathcal{R}$ are isomorphic and all projections are exact.*

3. Why does two-stage RLHF work *much* better in practice?

*A: RLHF only has to search over policies (generators) that are optimal for simple rewards (verifiers) rather than over all of $\Pi$.*