

# RL from Human Feedback as Game-Solving

Gokul Swamy

# Outline for Today

1. When is the Bradley-Terry assumption inaccurate and what happens to online / offline PFT as a result?
2. What is a more robust criterion for preference aggregation and how can we efficiently optimize it?

# Outline for Today

1. When is the Bradley-Terry assumption inaccurate and what happens to online / offline PFT as a result?

*A: BT is violated when when a reward function can't explain (aggregate) preferences, leading to mode collapse in RLHF.*

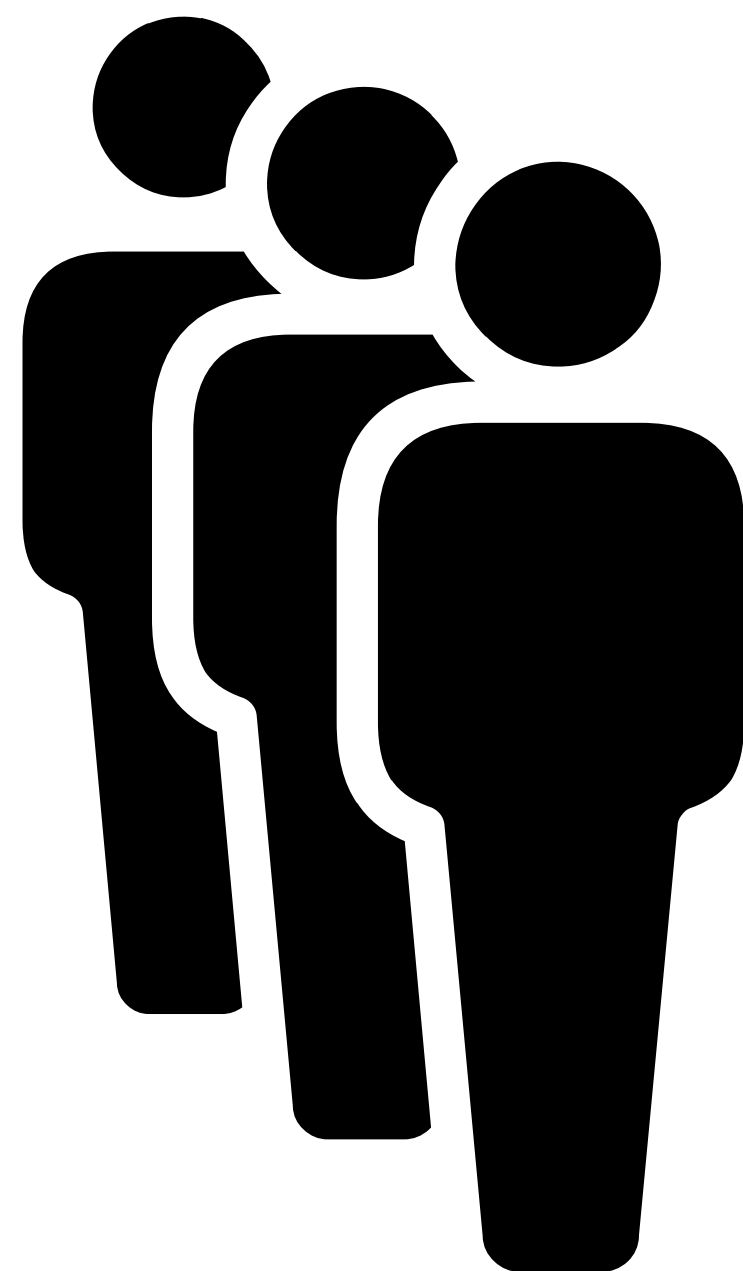
2. What is a more robust criterion for preference aggregation and how can we efficiently optimize it?

$$r \in [1, 0, 0]$$

# Preference Matrix

$$r(\xi_1) > r(\xi_2)$$

$$\xi_1 > \xi_2$$



$$r(\xi_1) > r(\xi_3)$$

$$\xi_1 \succ \xi_3$$

$$r(\xi_2) = r(\xi_3)$$

$$\xi_2 \sim \xi_3$$

$\mathcal{D}$  w.l.o.g.

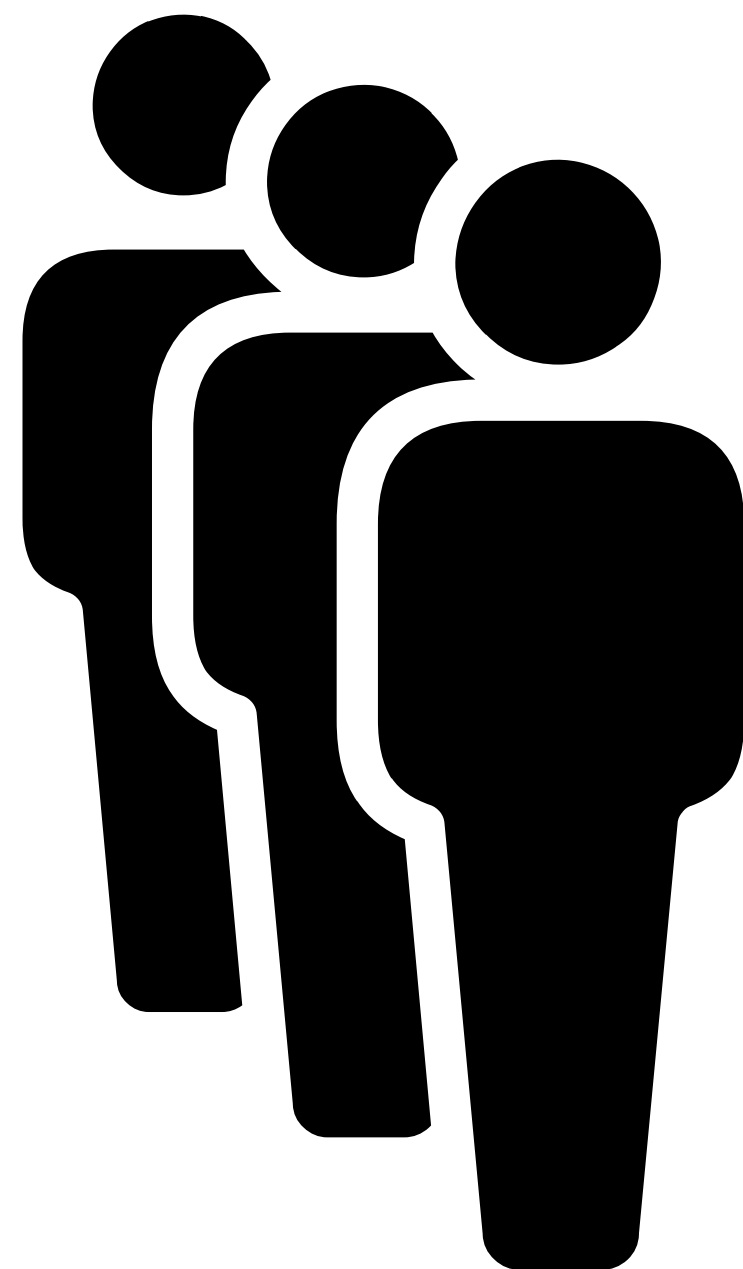
$$\mathcal{P}(\xi_1 > \xi_2) \triangleq \mathbb{P}_{\mathcal{D}}(\xi_1 > \xi_2)$$

$\mathcal{P}$	$\xi_1$	$\xi_2$	$\xi_3$
$\xi_1$	0.5	1	1
$\xi_2$	0	0.5	0.5
$\xi_3$	0	0.5	0.5

0.5

# (Partial) Preference Matrix

$\xi_1 > \xi_2$

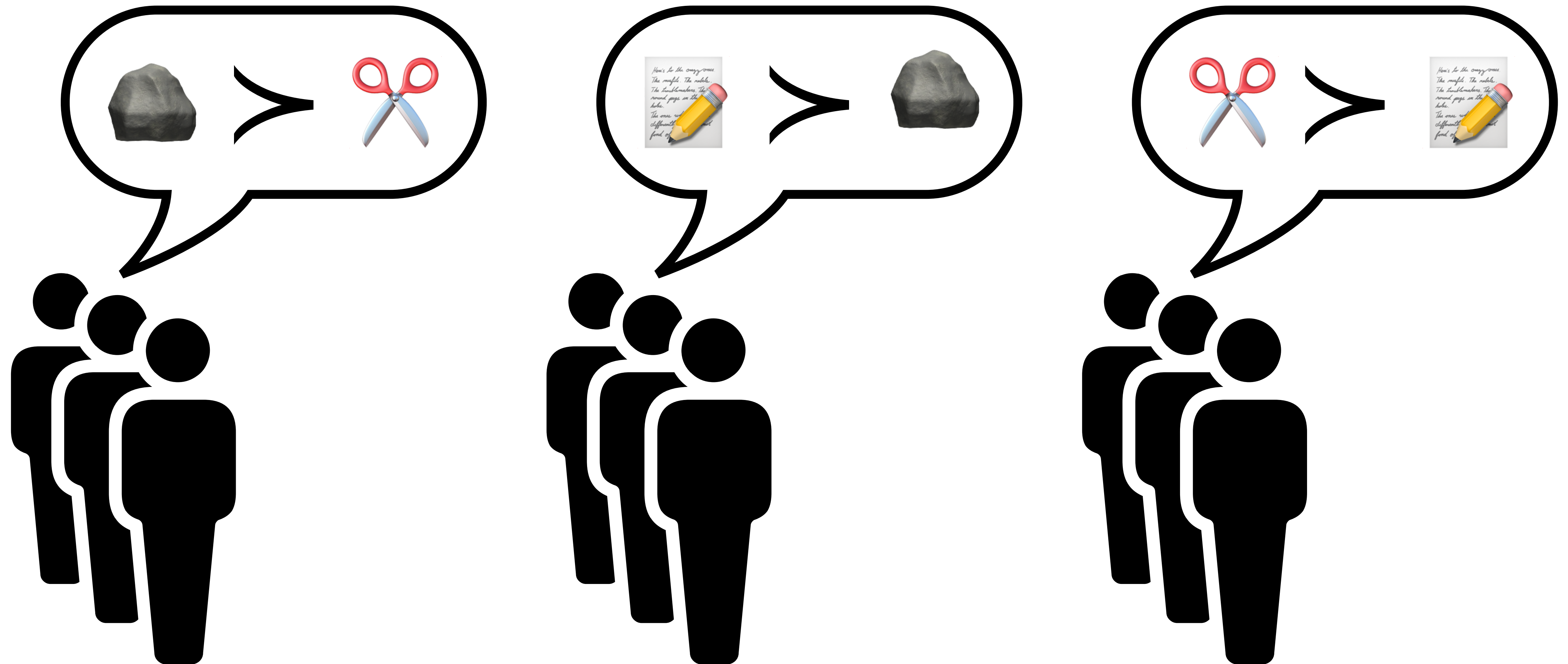


$\mathcal{D}$

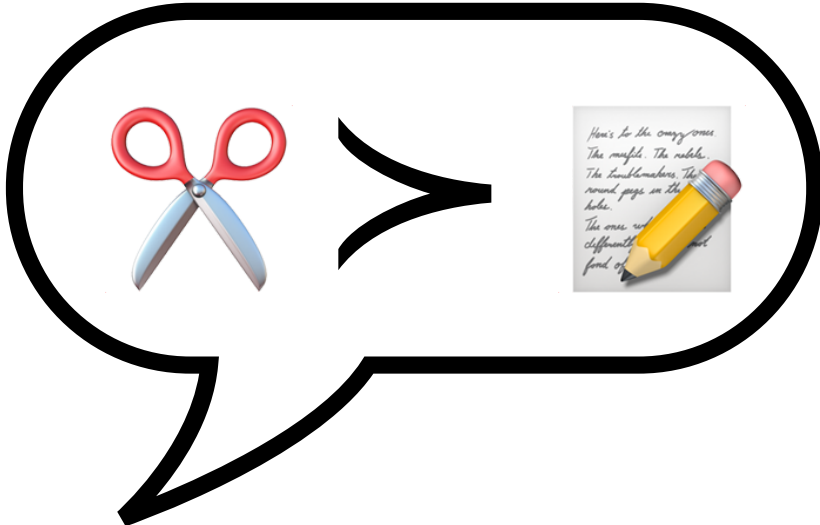
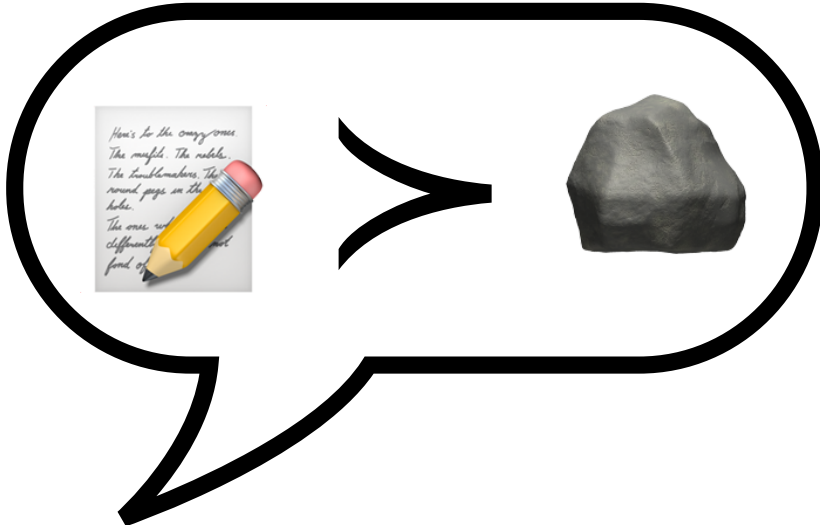
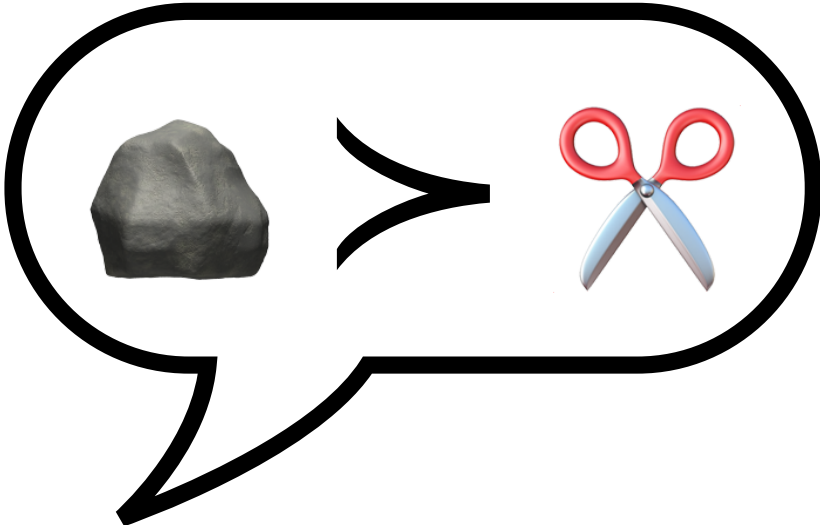
$$\mathcal{P}(\xi_1 > \xi_2) \triangleq \mathbb{P}_{\mathcal{D}}(\xi_1 > \xi_2)$$







$\mathcal{P}$	$\xi_1$	$\xi_2$	$\xi_3$
$\xi_1$	0.5	1	?
$\xi_2$	0	0.5	?
$\xi_3$	?	?	0.5







# Intransitivity from Preference Aggregation









# Intransitivity from Preference Aggregation



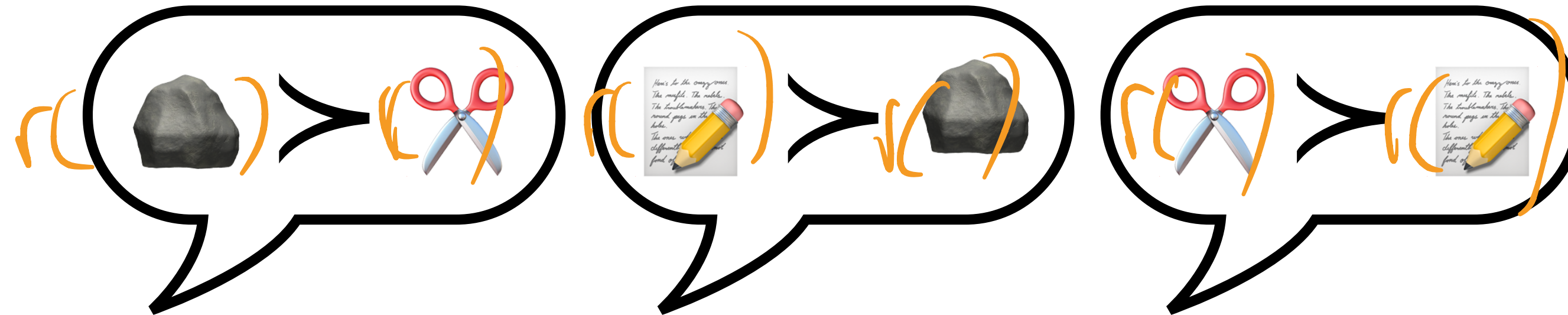
			
	0.5	?	1
	?	0.5	?
	0	?	0.5


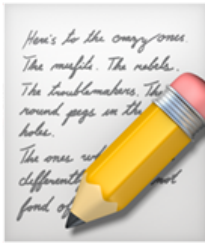


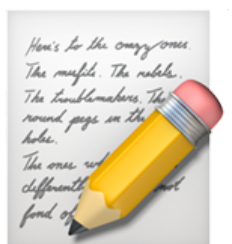

			
	0.5	0	?
	1	0.5	?
	?	?	0.5

			
	0.5	?	?
	?	0.5	0
	?	1	0.5



# Intransitivity from Preference Aggregation


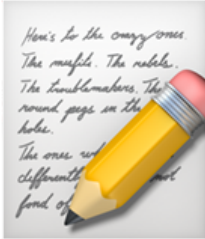


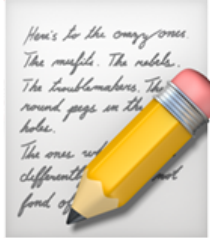



				
	0.5	0	1	$\Rightarrow 1.5$
	1	0.5	0	$\Rightarrow 1.5$
	0	1	0.5	$\Rightarrow 1.5$

There is no  $r^\star$  that explains these preferences!



# No $r^\star$ Leads to Mode Collapse in RLHF

				
	0.5	0.7	0	$\approx 1.2$
	0.3	0.5	1	$\approx 1.8$
	1	0	0.5	$\approx 1.5$


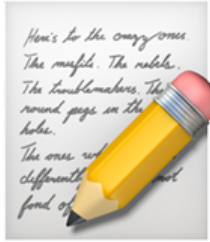


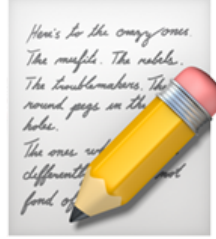

$$\pi_{\text{ref}} = \left[ \frac{1}{3}, \frac{1}{3}, \frac{1}{3} \right]$$

$$\hat{r}_{\text{mle}} = [1.2, 1.8, 1.5]$$

$$\hat{\pi}_{\text{rlhf}} \neq \pi_{\text{ref}} \cdot \exp\left(\frac{1}{\beta} \hat{r}_{\text{mle}}\right)$$

Loss

# No $r^\star$ Leads to Mode Collapse in RLHF

				...
	0.5	0.7	0	...
	0.3	0.5	1	...
	1	0	0.5	...
...	...	...	...	...

This problem gets worse with a larger  $\Xi$ !

# Outline for Today




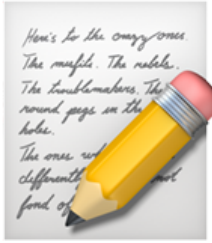

1. When is the Bradley-Terry assumption inaccurate and what happens to online / offline PFT as a result?

*A: BT is violated when when a reward function can't explain (aggregate) preferences, leading to mode collapse in RLHF.*


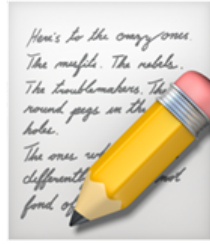


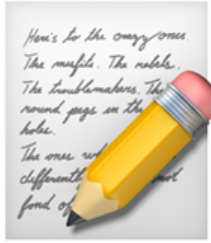

2. What is a more robust criterion for preference aggregation and how can we efficiently optimize it?

*A: The minimax winner doesn't assume transitivity of preferences. We can use a self-play algorithm to compute it.*

# Beyond Bradley-Terry in RLHF

			
	0.5	0.7	0
	0.3	0.5	1
	1	0	0.5

$$2 \cdot \mathcal{P} - 1$$

			
	0	0.4	-1
	-0.4	0	1
	1	-1	0

anti-symmetric

$2 \cdot \mathcal{P} - 1$  defines a symmetric 2p0s game!

# Von Neumann / Minimax Winners

$$\boxed{\pi_1^\star}, \pi_2^\star = \arg \max_{\pi_1 \in \Pi} \arg \min_{\pi_2 \in \Pi} \mathbb{E}_{\xi_1 \sim \pi_1, \xi_2 \sim \pi_2} [2\mathcal{P}(\xi_1 > \xi_2) - 1]$$

$\hookrightarrow$  robust against the worst case comparator policy





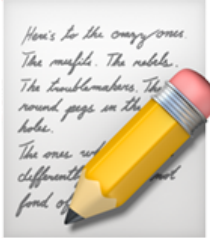

$$\max_{\pi_1} (\pi_1^\star)^\top (2P - 1) \pi_2^\star \leq \pi_1^\top (2P - 1) \pi_2 = 0 \quad \frac{0+1}{2} = 0.5$$

1. Pick a policy that is robust against worst case comparator.
2. Preferred to any other policy w.p. 1/2.
3. No assumptions on underlying shared  $r^\star$  required!



# Von Neumann / Minimax Winners

$$\pi_1^\star, \pi_2^\star = \arg \max_{\pi_1 \in \Pi} \arg \min_{\pi_2 \in \Pi} \mathbb{E}_{\xi_1 \sim \pi_1, \xi_2 \sim \pi_2} [2\mathcal{P}(\xi_1 > \xi_2) - 1]$$

				
	0.5	0.7	0	1.2
	0.3	0.5	1	
	1	0	0.5	1.5

$$\pi^\star = \left[ \frac{5}{12}, \frac{5}{12}, \frac{1}{6} \right]$$

# If : Computing Minimax Winners

$$\pi_1^\star, \pi_2^\star = \arg \max_{\pi_1 \in \Pi} \arg \min_{\pi_2 \in \Pi} \mathbb{E}_{\xi_1 \sim \pi_1, \xi_2 \sim \pi_2} [2\mathcal{P}(\xi_1 > \xi_2) - 1]$$

We can define a sequence of losses for each NR player:

$$\ell_t^1(\pi) = \mathbb{E}_{\xi \sim \pi, \xi' \sim \pi_2^t} [2\mathcal{P}(\xi > \xi') - 1] \quad \ell_t^2(\pi) = \mathbb{E}_{\xi \sim \pi_t^1, \xi' \sim \pi} [-(2\mathcal{P}(\xi > \xi') - 1)]$$

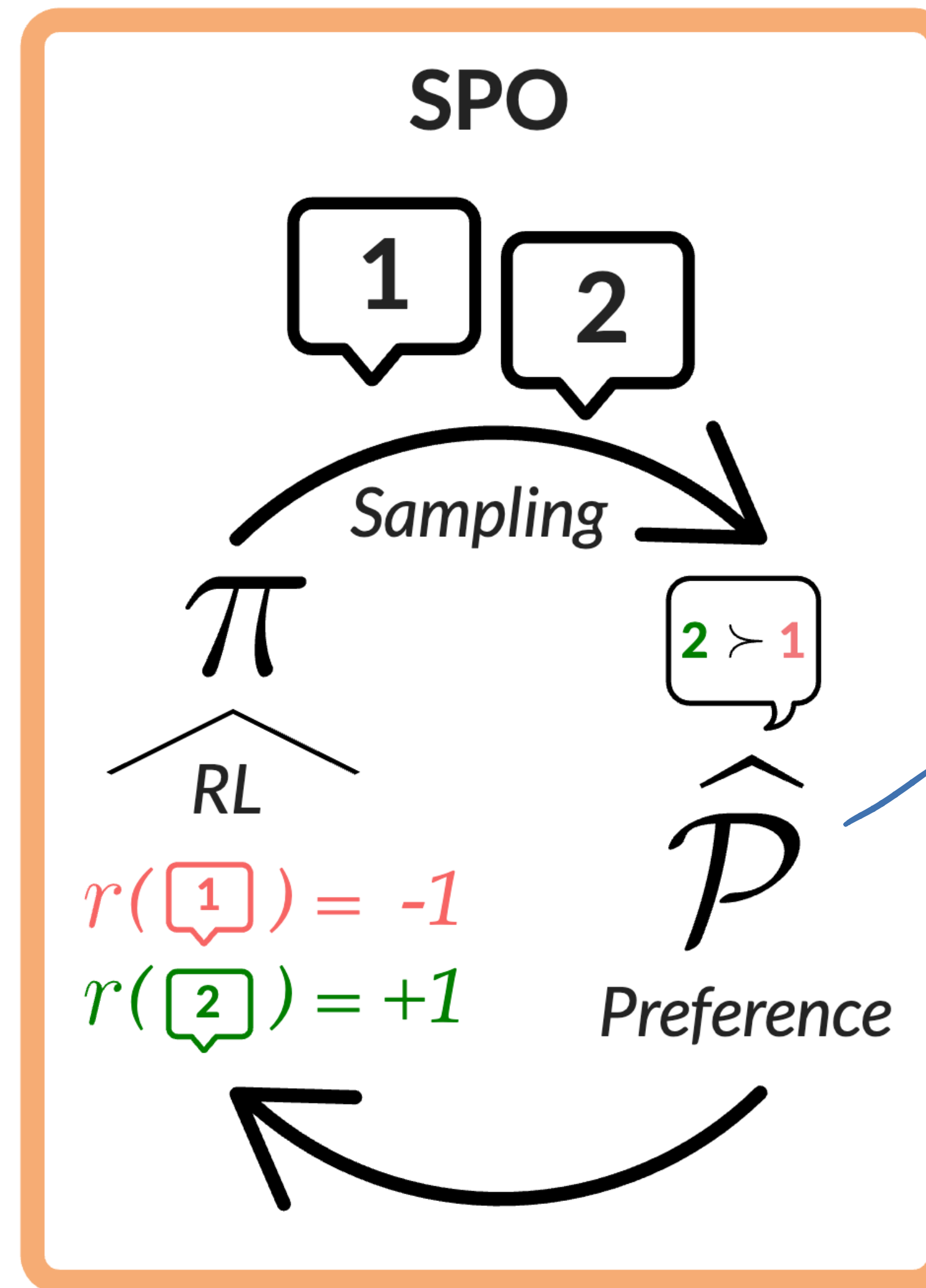
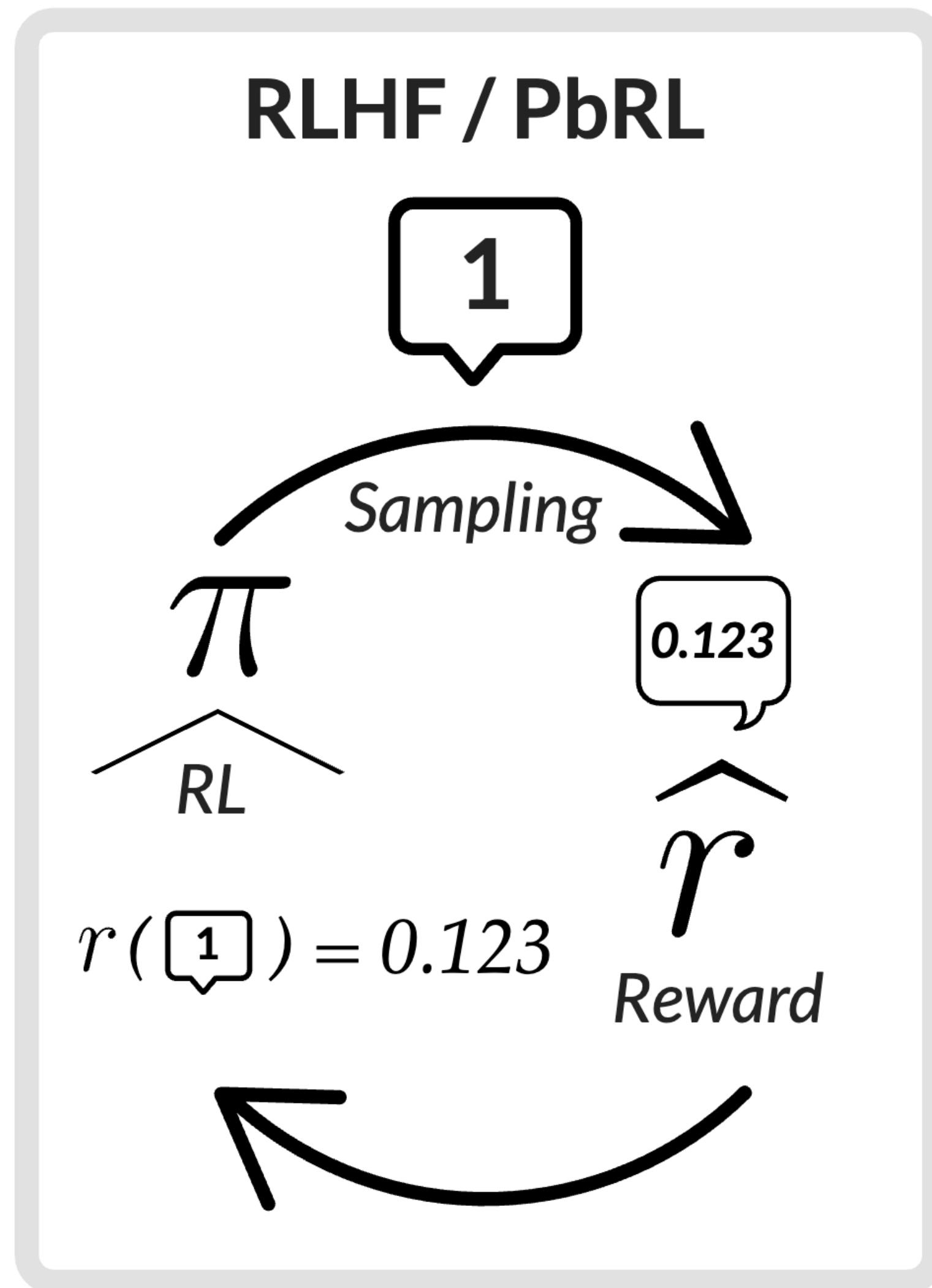
Without loss of generality, assume that  $\pi_1^0 = \pi_2^0$ . Then,

$$\begin{aligned} \ell_0^1(\pi) &= \mathbb{E}_{\xi \sim \pi, \xi' \sim \pi_2^0} [2\mathcal{P}(\xi > \xi') - 1] \\ &= \mathbb{E}_{\xi \sim \pi, \xi' \sim \pi_1^0} [2\mathcal{P}(\xi > \xi') - 1] \\ &= \mathbb{E}_{\xi \sim \pi_1^0, \xi' \sim \pi} [-(2\mathcal{P}(\xi > \xi') - 1)] = \ell_0^2(\pi) \\ &\Rightarrow \forall t \in [T], \pi_1^t = \pi_2^t \end{aligned}$$

Can just do self-play, no adversarial training required!



# SPO: Self-Play Preference Optimization



$$r_t = P \cdot \Pi_t$$

$$\sim N(\xi, \sigma^2)$$

win rate against  $\pi$  as the DM

# Instantiations of SPO on Large Language Models

## Offline Dataset

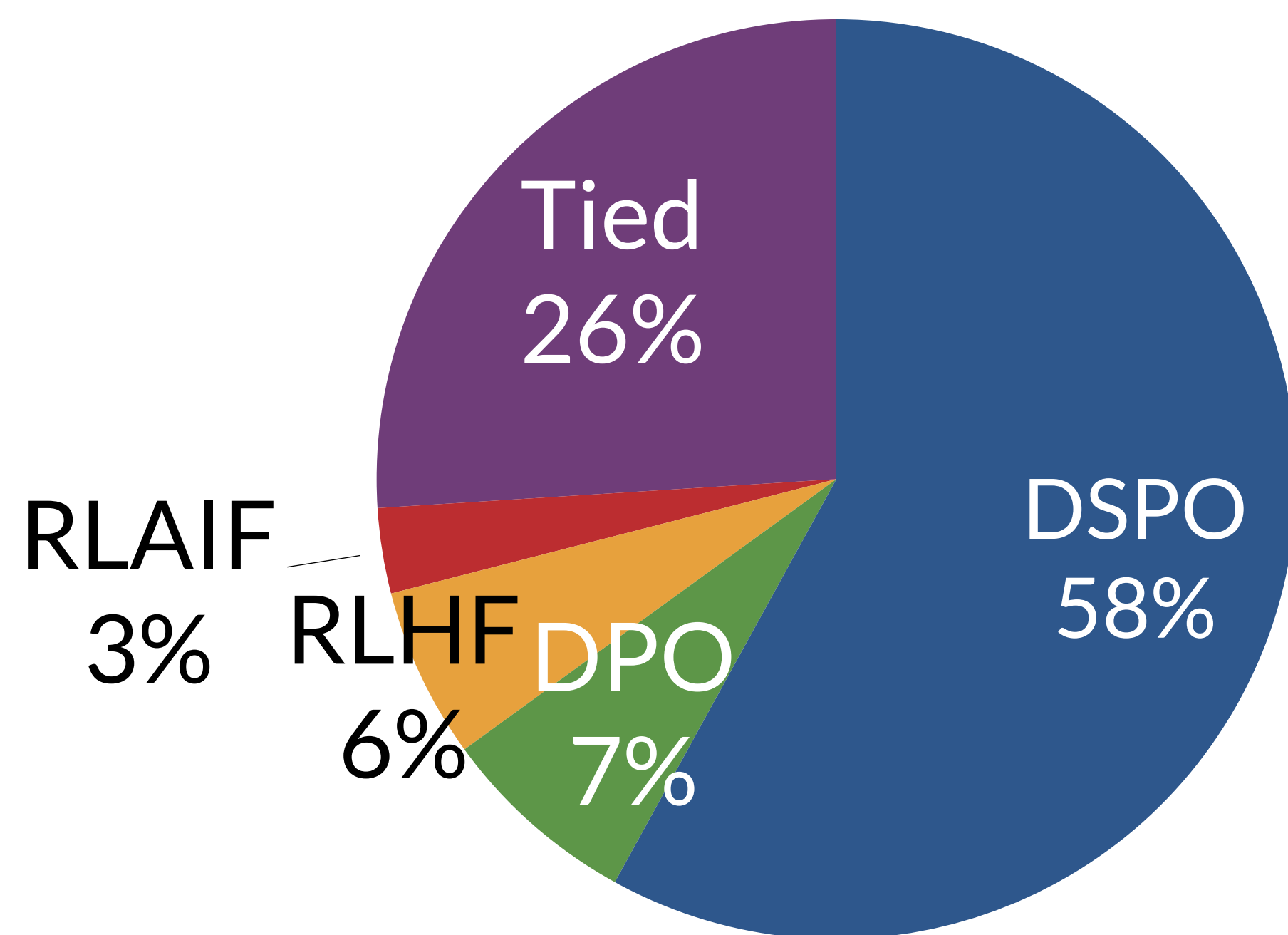
$P(r > c)$	SFT	RLHF	SPO	$P(r > c)$	RLHF	ISPO	DSPO
SFT	0.5	0.02	0.02	RLHF	0.5	0.21	0.24
RLHF	0.98	0.5	0.25	ISPO	0.79	0.5	0.61
SPO	0.98	0.75	0.5	DSPO	0.76	0.39	0.5

[Munos+'23]

[Calandriello+'24]

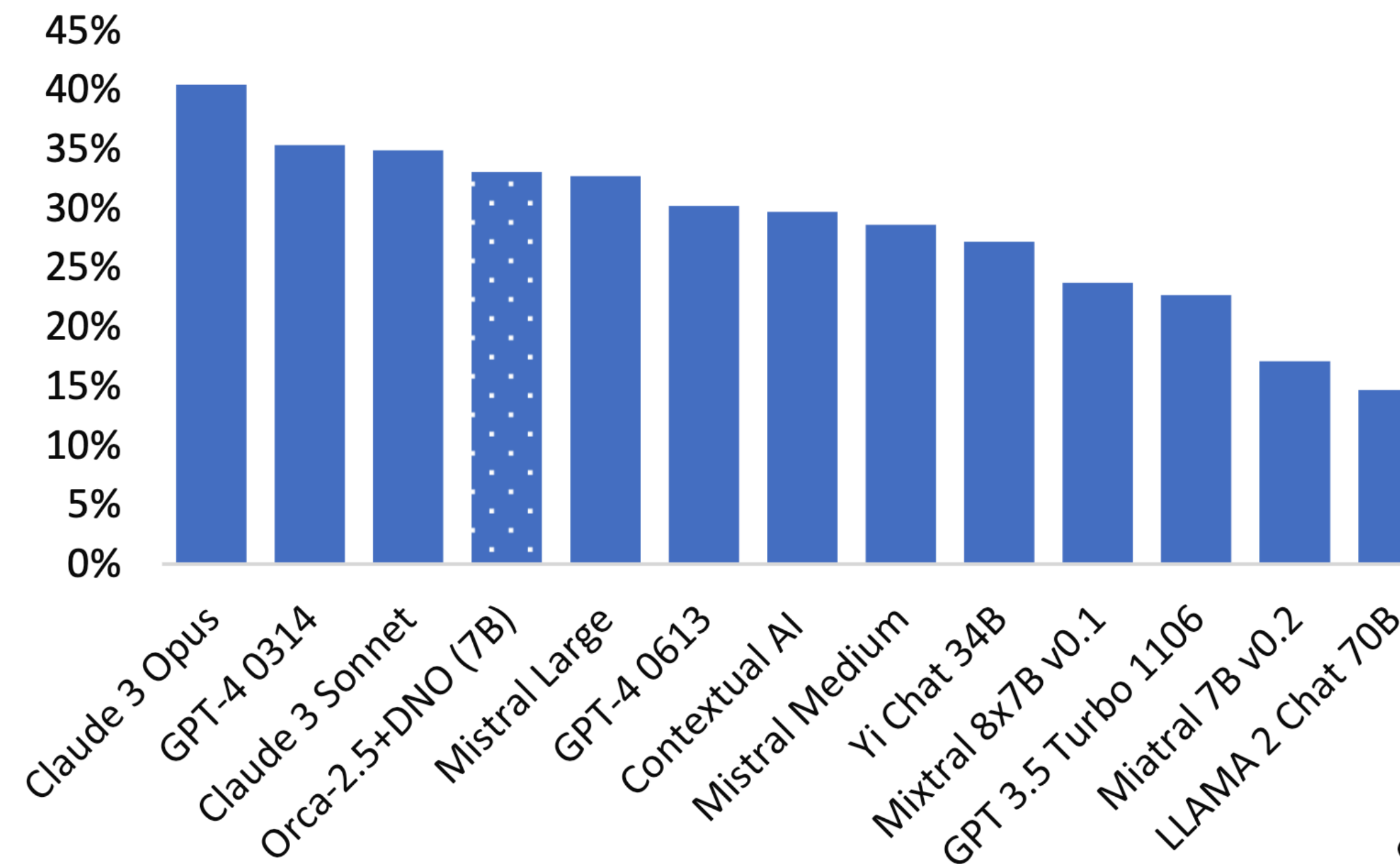
# Instantiations of SPO on Large Language Models

## Online Oracle



[Guo+'24]

LC Win Rate



[Rosset+'24]

# Outline for Today

1. When is the Bradley-Terry assumption inaccurate and what happens to online / offline PFT as a result?

*A: BT is violated when when a reward function can't explain (aggregate) preferences, leading to mode collapse in RLHF.*

2. What is a more robust criterion for preference aggregation and how can we efficiently optimize it?

*A: The minimax winner doesn't assume transitivity of preferences. We can use a self-play algorithm to compute it.*