

Doubly Robust Alignment for Large Language Models

Erhan Xu^{†1}, Kai Ye^{†1}, Hongyi Zhou^{†2}, Luhan Zhu³, Francesco Quinzan^{‡4}, Chengchun Shi^{‡1}

¹London School of Economics, ²Tsinghua University, ³University of the Arts London, ⁴University of Oxford
[†]equal contribution, [‡]joint senior contributors

Problem: Model Misspecification in RLHF

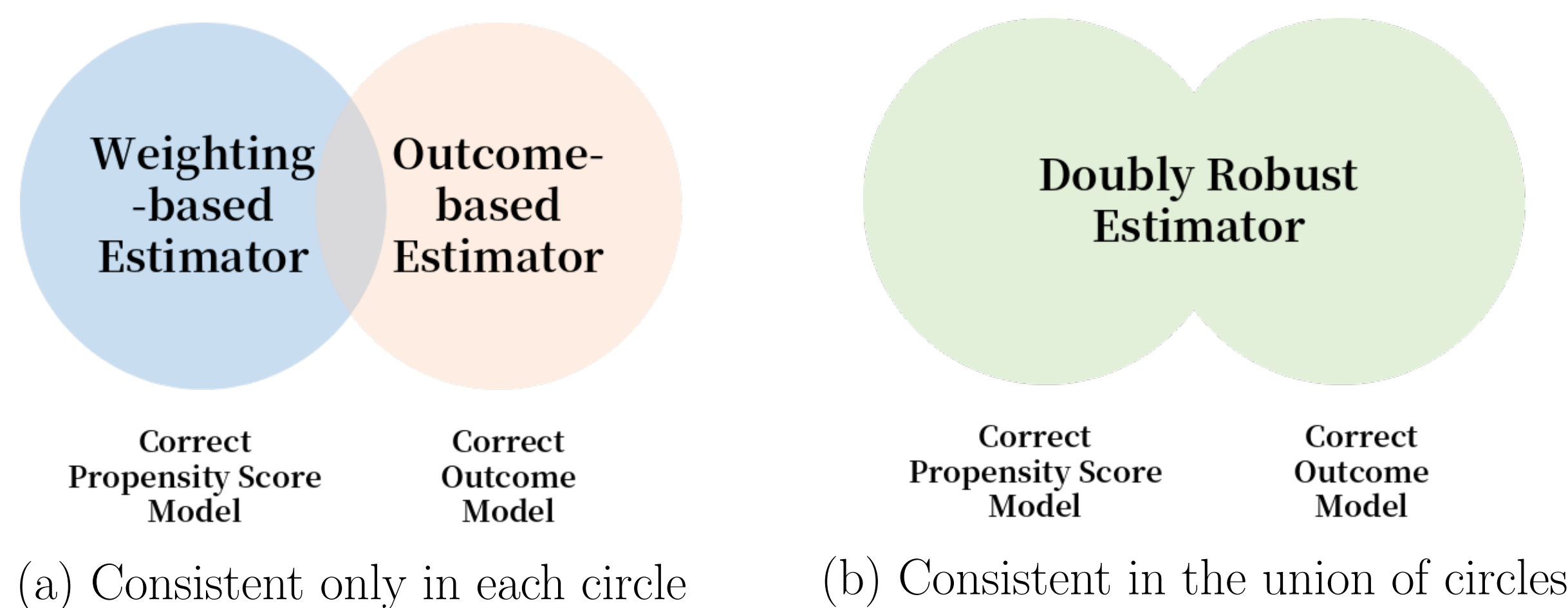
- **PPO-based algorithms:** Sensitive to reward model misspecification
 - Can lead to reward hacking and misguided policy learning
- **DPO-based algorithms:** Sensitive to reference policy misspecification
 - Performance degrades when reference policy is inaccurate
- **Preference-based algorithms:** Rely on correct preference model specification
 - Bradley-Terry (BT) model often violated due to intransitivity in human preferences

Robust to misspecified: preference model reward model reference policy				
RLHF	Reward-based	PPO-type	✗	✗
		DPO-type	✗	✗
		IPO	✓	✗
	Preference-based	GPM	✗	✓
		DRPO	✓	✓

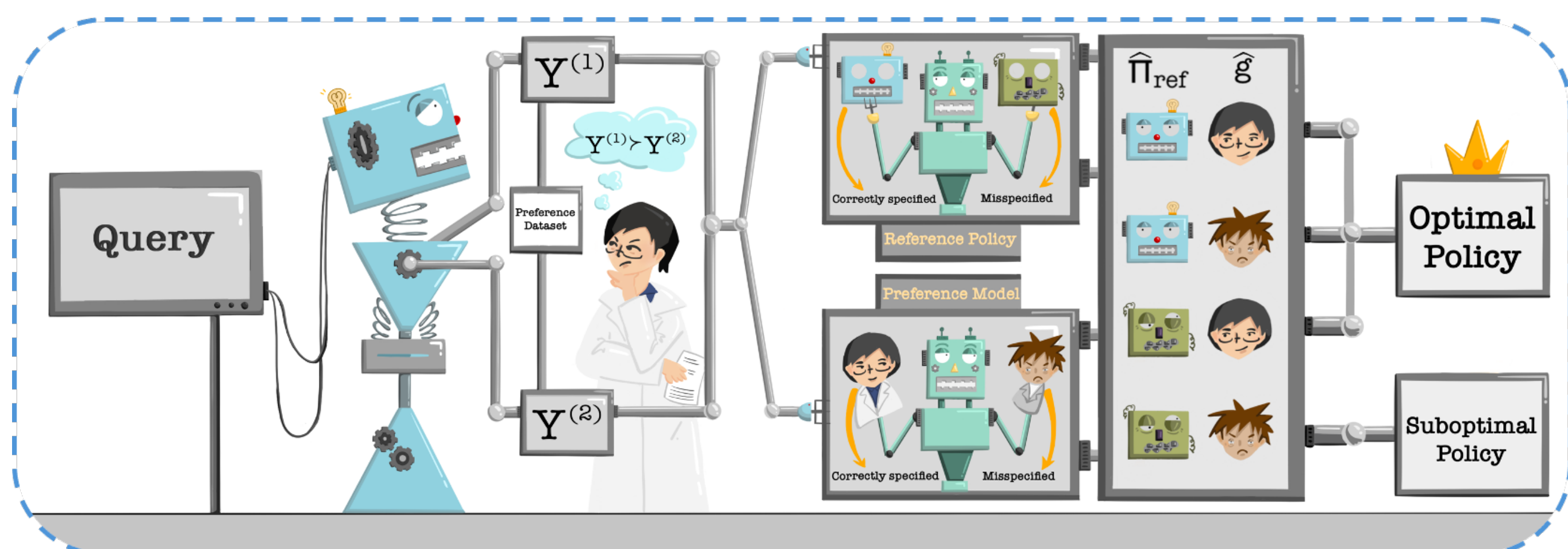
Our Solution: Use Double Robust Method

Doubly robust methods originate from the **missing data** and **causal inference** literature. Consider the estimation of **average treatment effect** (ATE) in causal inference. These methods estimate two models:

- A **propensity score** model for treatment assignment mechanism
- An **outcome regression** model for subject's outcome given treatment
- Similar to **reference policy** in LLMs
- Similar to **reward model** in LLMs



When DR methods meet LLMs:



Key Insight:

- Estimate two models: preference model \hat{g} and reference policy $\hat{\pi}_{\text{ref}}$
- Construct estimator that remains consistent when *either* model is correct

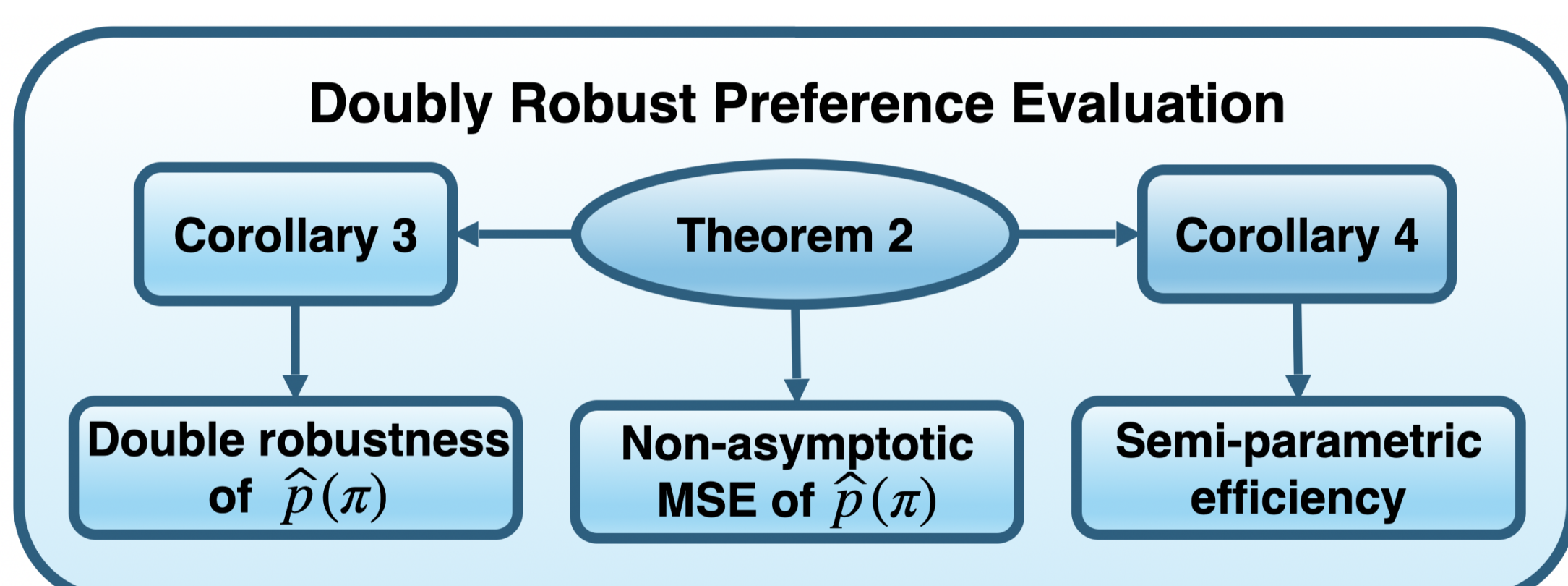
Doubly Robust Preference Evaluation

Goal: Estimate total preference $p^*(\pi) = \mathbb{E}_{y \sim \pi, y' \sim \pi_{\text{ref}}} g^*(X, y, y')$

Our DR Estimator:

$$\hat{p}(\pi) = \frac{1}{2} \mathbb{E}_{(X, Y^{(1)}, Y^{(2)}, Z) \sim \mathcal{D}} \left\{ \sum_{a=1}^2 \mathbb{E}_{y \sim \pi(\cdot|X)} [\hat{g}(X, y, Y^{(a)})] + \sum_{a=1}^2 (-1)^{a-1} \frac{\pi(Y^{(a)}|X)}{\hat{\pi}_{\text{ref}}(Y^{(a)}|X)} [Z - \hat{g}(X, Y^{(1)}, Y^{(2)})] \right\}$$

Main Theoretical Results:



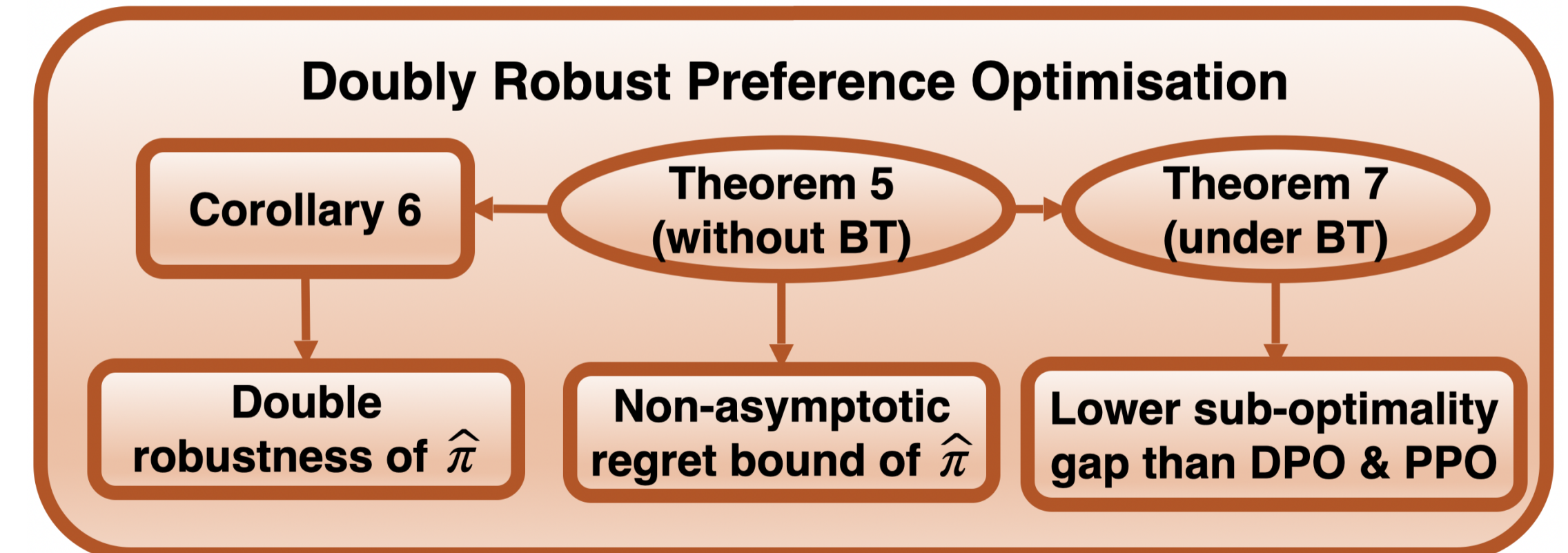
- **Double Robustness:** $\text{MSE} \rightarrow 0$ when either $\hat{\pi}_{\text{ref}}$ or \hat{g} is correct
- **Semiparametric Efficiency:** Achieves smallest-possible MSE when both correct

Doubly Robust Preference Optimization

For Preference Optimization:

$$\hat{\pi} = \arg \max_{\pi} \hat{p}(\pi) - \beta \text{KL}(\pi, \hat{\pi}_{\text{ref}})$$

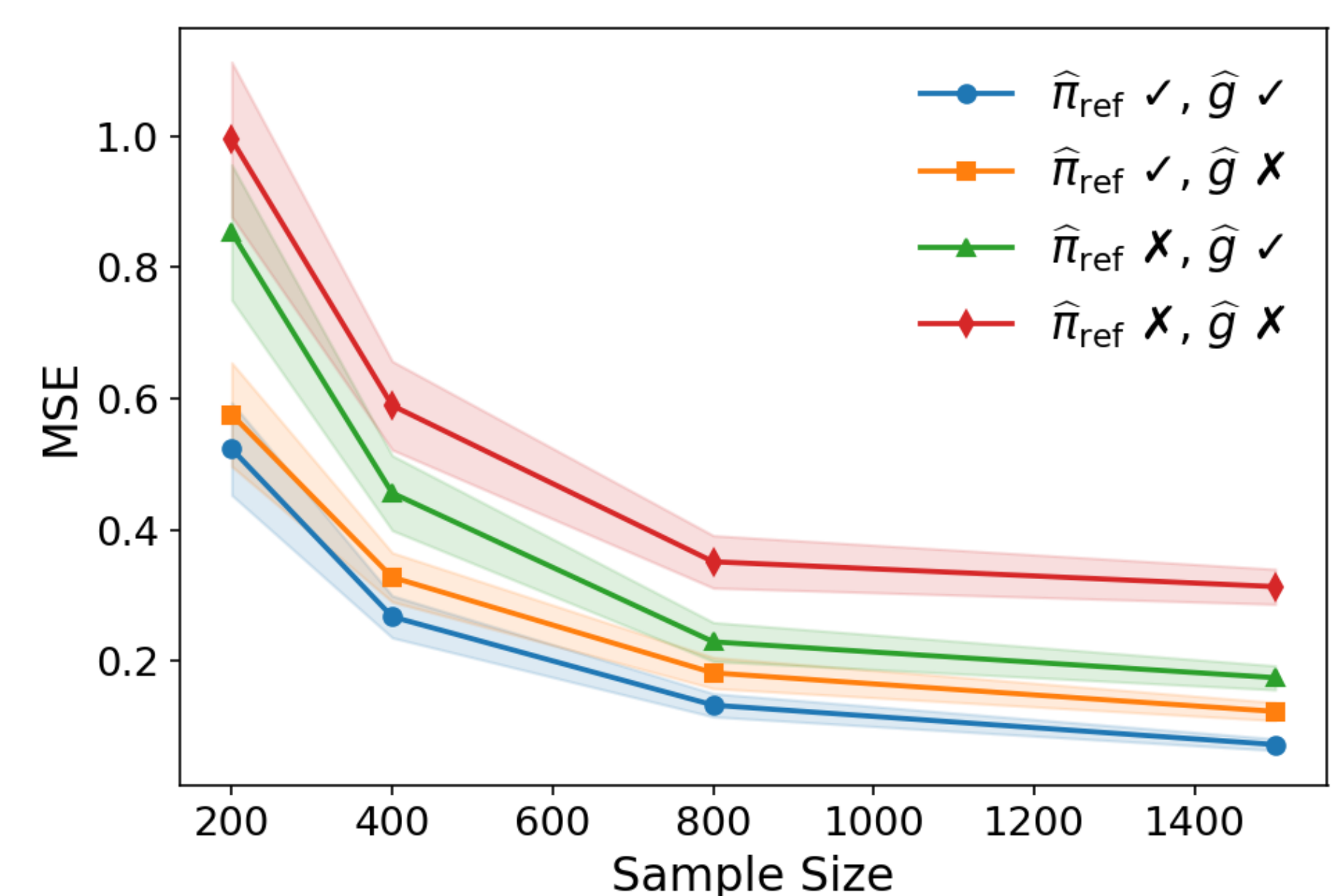
Main Theoretical Results:



- **Double robustness:** Regret of $\hat{\pi}$ decays to zero when *either* reference policy *or* preference model (not necessarily both) is correct
- **Sub-optimality gaps:**
 - PPO: $O(n^{-1/2} + \|\hat{r} - r\|)$
 - DPO: $O(n^{-1/2} + \|\hat{\pi}_{\text{ref}} - \pi_{\text{ref}}\|)$
 - DRPO: $O(n^{-1/2} + \|\hat{r} - r\| \|\hat{\pi}_{\text{ref}} - \pi_{\text{ref}}\|)$

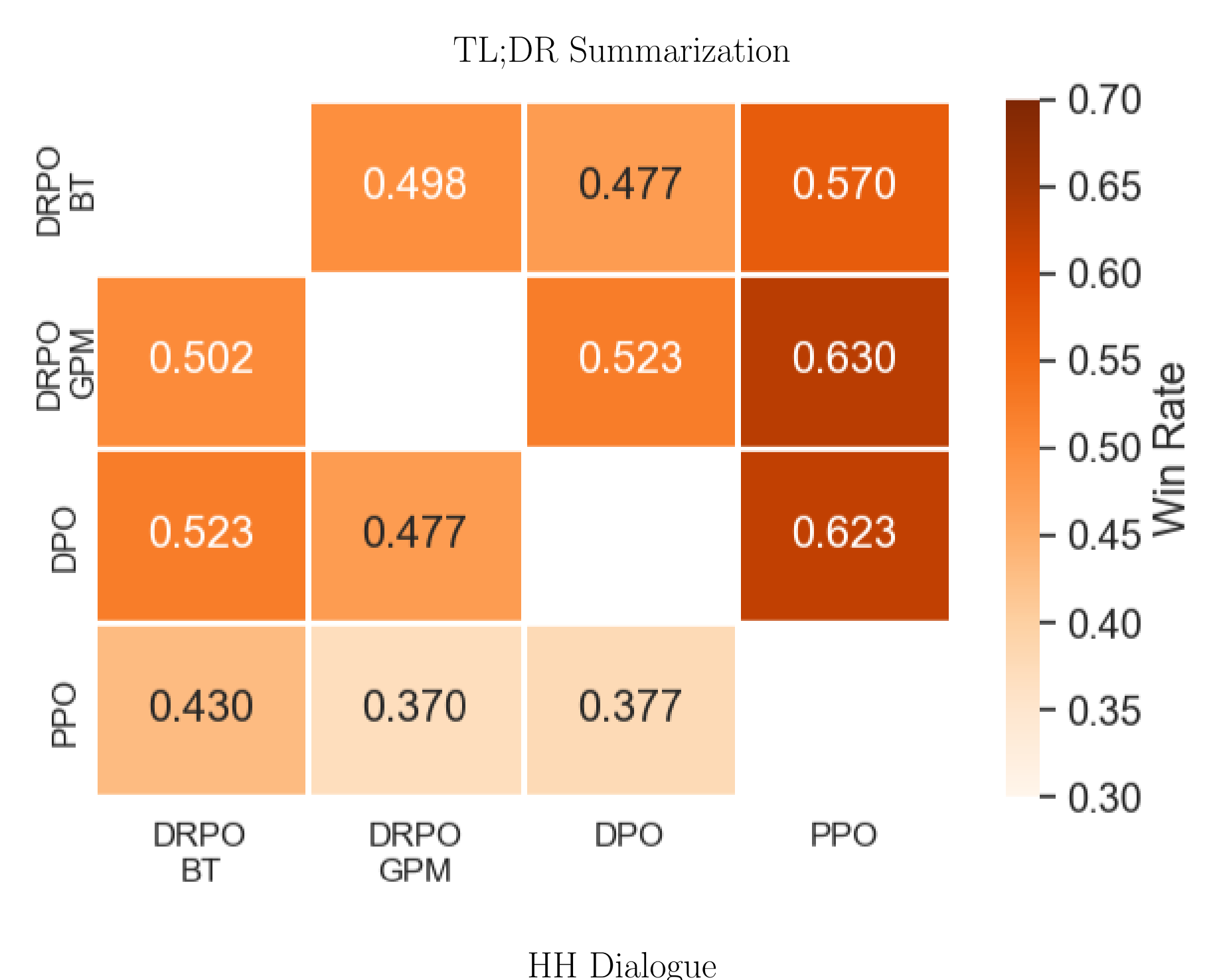
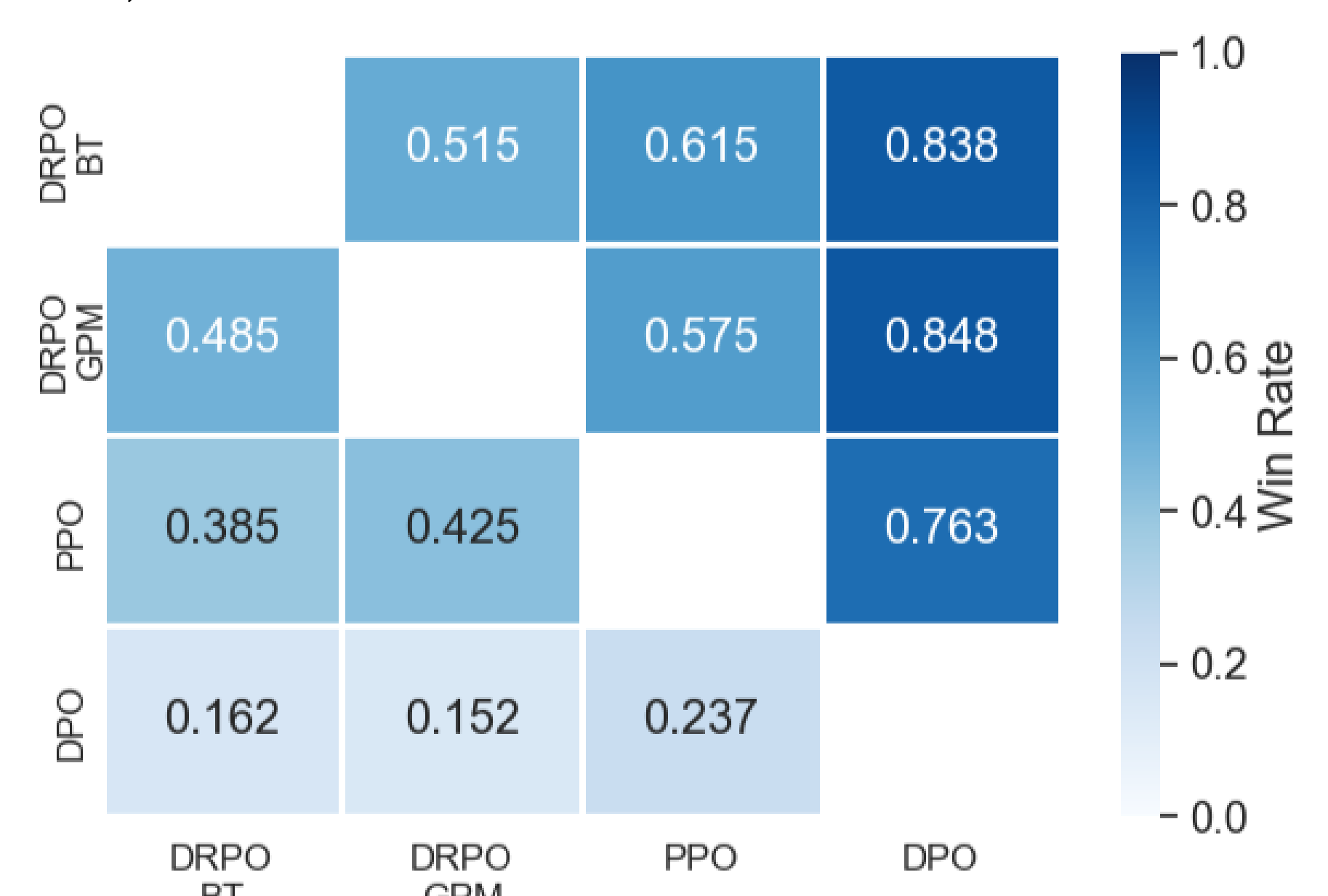
Empirical Results

Applications to IMDb Dataset:



- Evaluate total preference of DPO-trained policy over SFT reference
- Ground truth: 0.681 (computed via Monte Carlo)
- MSE converges to zero when either model is correct

Applications to TL;DR and HH Datasets:



- DRPO (both BT and General Preference Model variants) achieves more robust and often superior performance to PPO and DPO under GPT-4o-mini evaluation
- Robust performance without extensive hyperparameter tuning