

Robust Reinforcement Learning from Human Feedback for Large Language Models Fine-Tuning

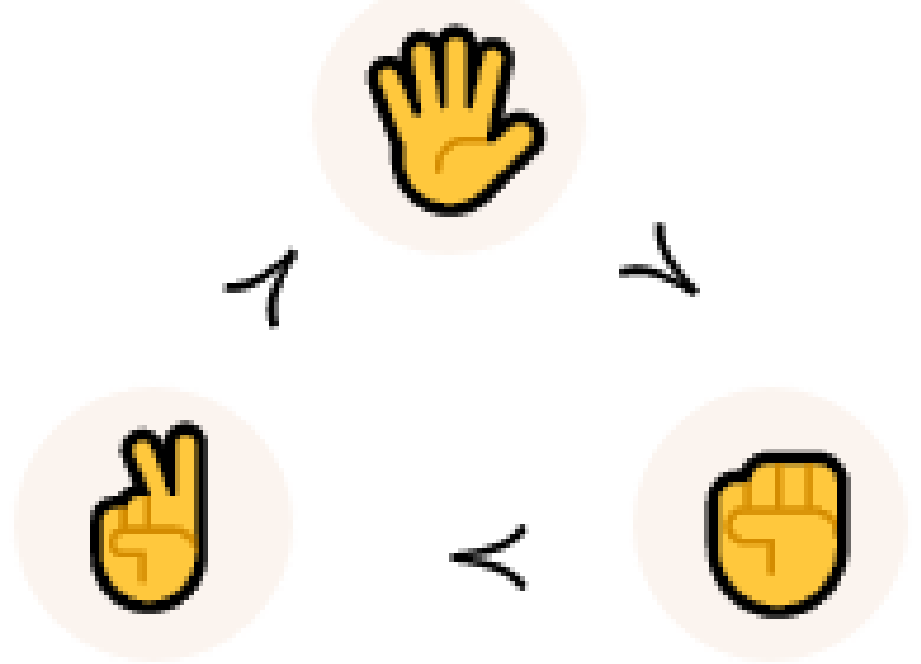
Kai Ye¹ Hongyi Zhou² Jin Zhu¹ Francesco Quinzan³ Chengchun Shi¹
¹LSE ²Tsinghua University ³University of Oxford

Abstract

Reinforcement learning from human feedback (RLHF) has emerged as a key technique for aligning the output of large language models (LLMs) with human preferences. To learn the reward function, most existing RLHF algorithms use the Bradley-Terry model, which relies on assumptions about human preferences that may not reflect the complexity and variability of real-world judgments. In this paper, we propose a robust algorithm to enhance the performance of existing approaches under such reward model misspecifications. Theoretically, our algorithm reduces the variance of reward and policy estimators, leading to improved regret bounds. Empirical evaluations on LLM benchmark datasets demonstrate that the proposed algorithm consistently outperforms existing methods, with 77-81% of responses being favored over baselines on the Anthropic Helpful and Harmless dataset.

Problem: preference model misspecification

RLHF algorithms for fine-tuning large language models typically require a human preference model, with the Bradley-Terry (BT) model being the most widely adopted. However, it relies on a reward-based preference assumption that imposes several idealized conditions on human judgments:



- **Transitivity:** human preferences are logically ordered ($A \succ B \succ C$ implies $A \succ C$).
- **Context-independence:** preferences between two responses are based solely on the prompts and responses themselves.
- **Perfect rationality:** users provide consistent and deterministic feedback.

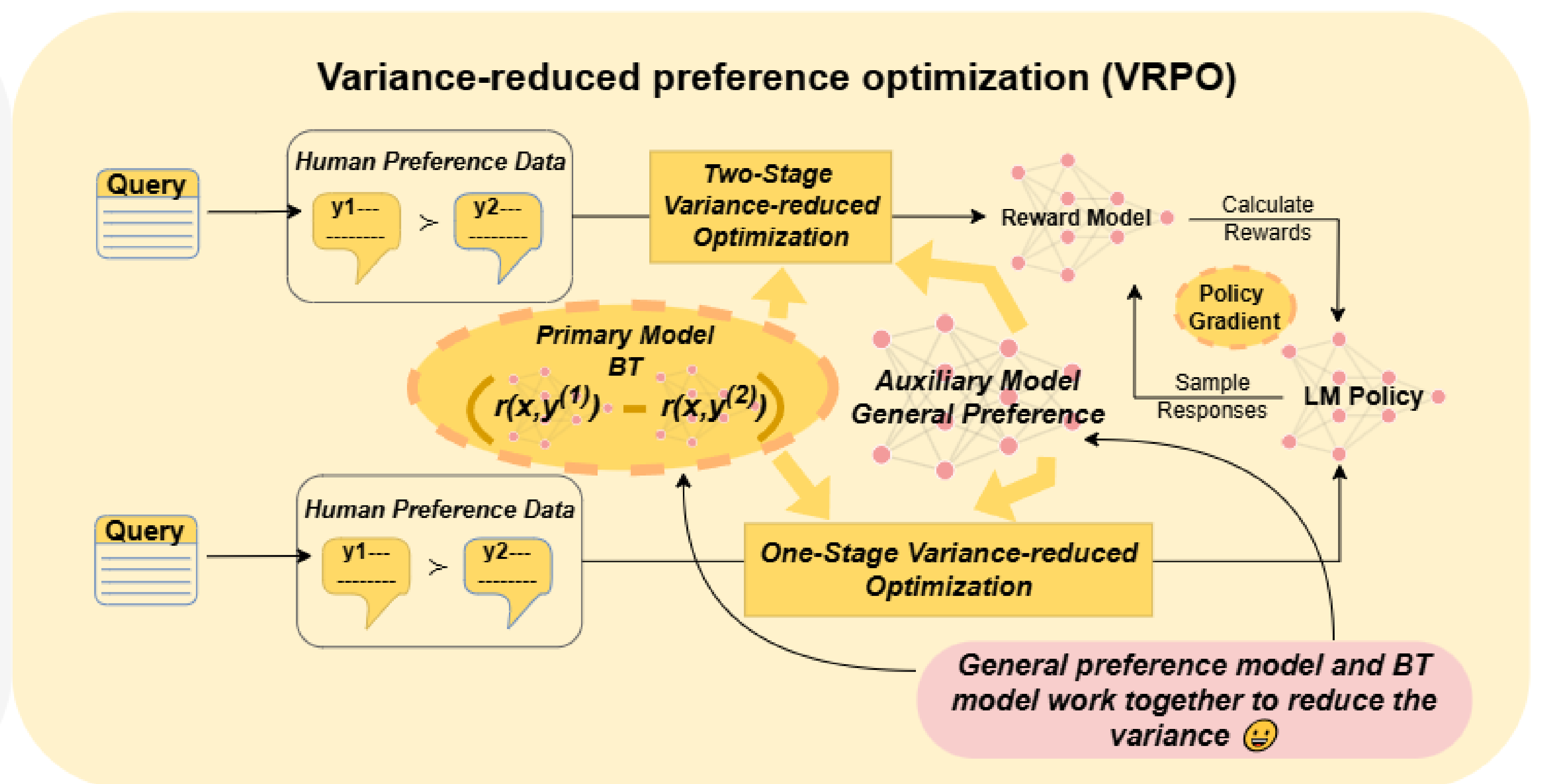
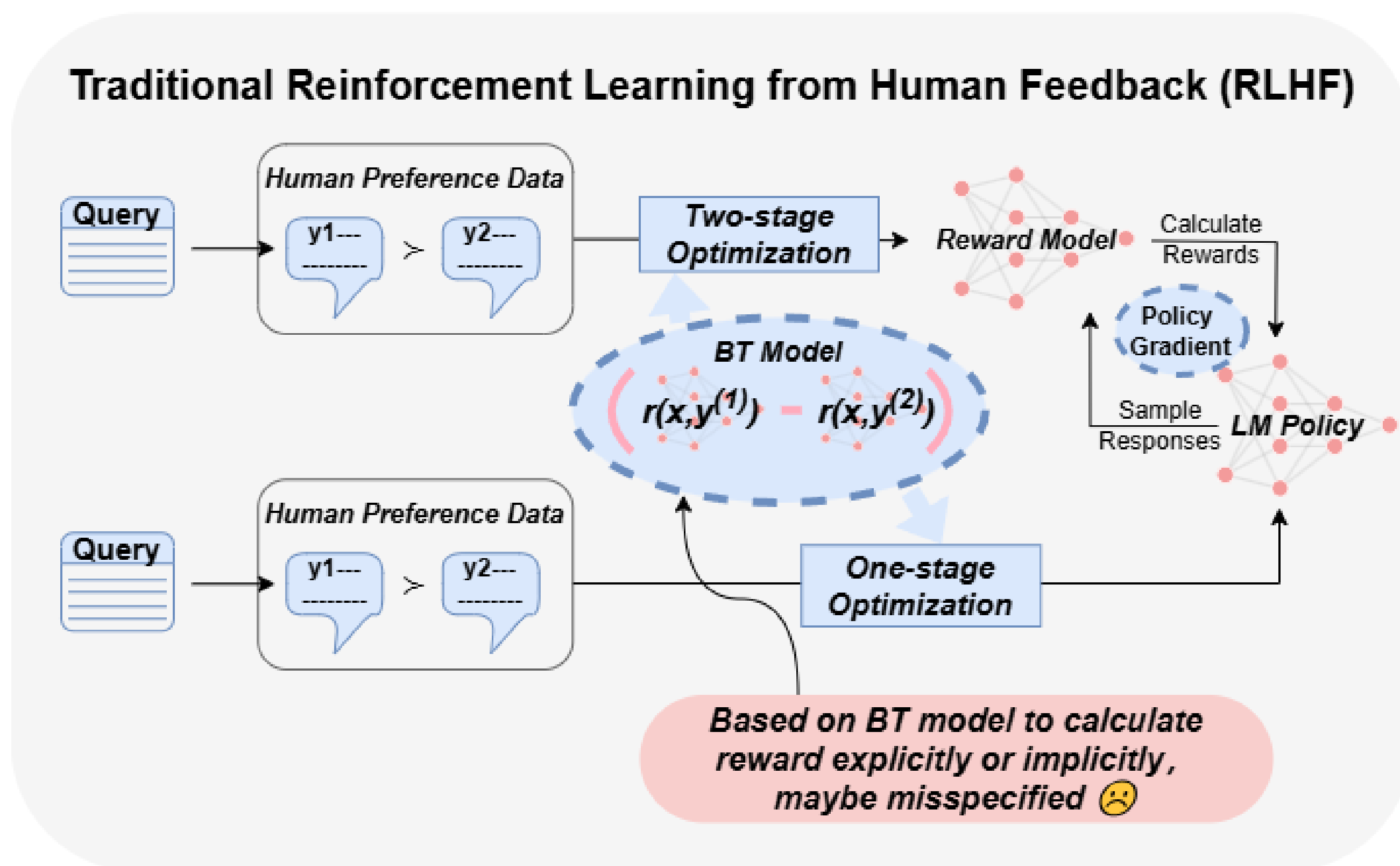
Those assumptions may be misspecified or unrealistic in practice.

Our solution: VRPO

- A flexible pipeline applicable to a variety of existing RLHF algorithms to enhance their sample efficiency under human preference model misspecification.
- Estimate two models for the preference function p^* : (i) A primary, simpler reward-based model p_θ , as in existing algorithms; and (ii) an auxiliary, more complex preference model p_η , designed to enhance the accuracy of the primary model.
- Loss function:

$$\tilde{\mathcal{L}}(\theta) = \mathbb{E}_n \left[\ell(X, Y^{(1)}, Y^{(2)}, Z; \theta) - \sum_{u=0}^1 \ell(X, Y^{(1)}, Y^{(2)}, u; \theta) p_\eta(X, Y^{(1)}, Y^{(2)}, u) \right. \\ \left. + \sum_{u=0}^1 \mathbb{E}_{y^{(1)*}, y^{(2)*} \sim \pi_{\text{ref}}(\cdot | X)} \ell(X, y^{(1)*}, y^{(2)*}, u; \theta) p_\eta(X, y^{(1)*}, y^{(2)*}, u) \right].$$

Variance-reduced preference optimization (VRPO) pipeline



VRPO incorporates an auxiliary preference model to reduce the variance of the estimated primary model. **Left:** The classic one-stage and two-stage optimization schemes in RLHF. Both approaches require fitting a reward model, either explicitly or implicitly, which may lead to model misspecification. **Right:** In contrast, VRPO employs an auxiliary reward-free preference model to better capture human preferences. It works jointly with the primary model for variance reduction and policy improvement.

Theoretical Guarantees

MODEL SETTING	VARIANCE OF ESTIMATOR	MSE OF ESTIMATOR	SUBOPTIMALITY GAP
MISSPECIFIED	↓	↓	↓
CORRECTLY SPECIFIED	↓	↓	↓

Double Robustness

In the correctly specified setting, the target parameter $\bar{\theta} = \arg \min_{\theta} \mathbb{E}[\tilde{\mathcal{L}}(\theta)]$, when either the reference policy π_{ref} or the auxiliary preference model p_η is correctly specified.

Variance and MSE reductions

regardless of whether the model is correctly specified or misspecified, we have

$$\|\mathbb{E}(\hat{\theta}) - \bar{\theta}\|_2 = O\left(\frac{d}{n\lambda_{\min}}\right), \quad \|\mathbb{E}(\tilde{\theta}) - \bar{\theta}\|_2 = O\left(\frac{d}{n\lambda_{\min}}\right),$$

and

$$\text{Var}(\hat{\theta}) - \text{Var}(\tilde{\theta}) = \underbrace{\frac{1}{n} A^{-1}(\bar{\theta}) H A^{-1}(\bar{\theta})}_{\text{variance reduction}} + O\left(\frac{d^{3/2}}{n^{3/2} \lambda_{\min}^2}\right) + O\left(\frac{\|p_\eta - p^*\|_\infty^2}{n \lambda_{\min}^2}\right),$$

where H denotes certain positive semi-definite matrix, $A(\bar{\theta}) := -\mathbb{E}\left\{\frac{\partial^2}{\partial \theta^2} \mathcal{L}(y, A, x; \bar{\theta})\right\}$, λ_{\min} denote the minimum eigenvalue of $A(\bar{\theta})$ and $\|p_\eta - p^*\|_\infty$ denotes the difference between p_η and p^* in supremum norm.

Variance and MSE reductions

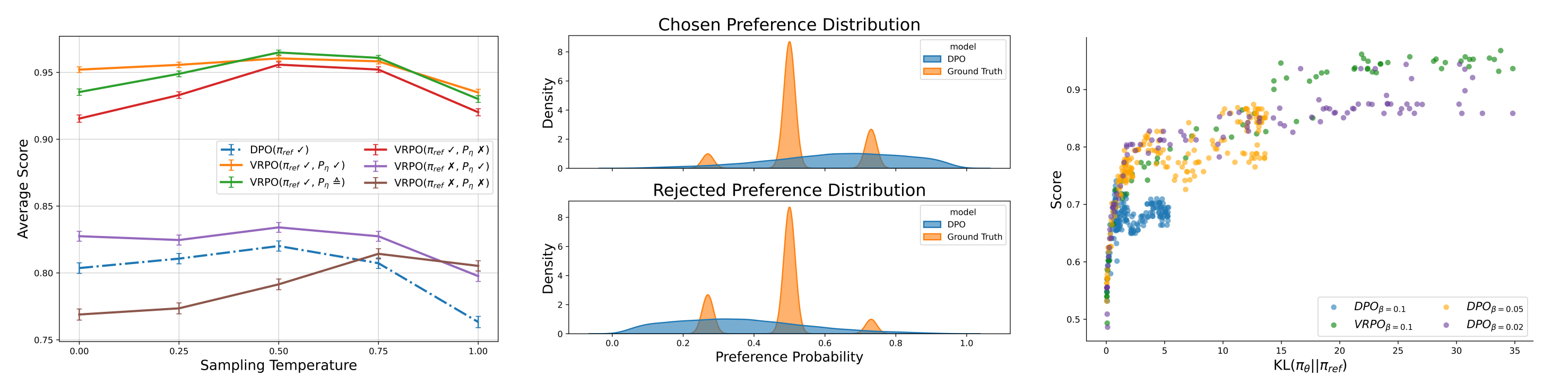
regardless of whether the model is correctly specified or misspecified, we have

$$\mathbb{E}\{R(\hat{\theta})\} = R(\bar{\theta}) + \text{trace}\left(\text{Var}(\hat{\theta})(-\nabla_{\bar{\theta}}^2 J(\bar{\theta}))\right) + O\left(\frac{d^{3/2}}{n^{3/2} \lambda_{\min}^3}\right),$$

$$\mathbb{E}\{R(\tilde{\theta})\} = R(\bar{\theta}) + \text{trace}\left(\text{Var}(\tilde{\theta})(-\nabla_{\bar{\theta}}^2 J(\bar{\theta}))\right) + O\left(\frac{d^{3/2}}{n^{3/2} \lambda_{\min}^3}\right),$$

where $\text{Var}(\hat{\theta})$ and $\text{Var}(\tilde{\theta})$ denote the covariance matrices of $\hat{\theta}$ and $\tilde{\theta}$, respectively.

Experiments performance



Comparisons in IMDb dataset. **Left** panel represents the expected reward in different VRPO setting compared to DPO, for example $(\pi_{\text{ref}} \checkmark, P_\eta \times)$ means the reference model is correctly specified and the preference model is misspecified, and $P_\eta \hat{=}$ means the preference model is estimated, demonstrating the robustness of our method. **Middle** plane illustrates the difference in preference probability distributions between the ground truth and the DPO estimation for both the Chosen and Rejected responses. **Right** The panel reports the expected reward versus KL-divergence for VRPO and DPO, demonstrating the quality of the optimization.

VRPO	0.500	0.565	0.572
DPO	0.435	0.500	0.538
SFT	0.428	0.462	0.500
VRPO	0.500	0.793	0.948
DPO	0.207	0.500	0.891
SFT	0.052	0.109	0.500

Head-to-head comparisons between VRPO, DPO, SFT. Win rates are evaluated by GPT-4o-mini. **Left** panel displays the win rate in summarization task. In both tasks, VRPO outperforms DPO by defeating it directly and demonstrating a higher win rate against SFT. **Right** panel displays the win rate in Anthropic HH one-step dialogue task.