



Motivation

- **Sample efficiency** (# of human annotations) is crucial in online RLHF.
- Previous works focus on strategic exploration, while we study from a underexplored perspective—**transfer learning**.
- Rich scenarios for transfer learning:
 - Reward models (RMs) from relevant tasks
 - Easy-to-access metric other than human feedback
 - Guidance from advanced LLMs

Key Question: How to improve sample efficiency in online RLHF by leveraging those imperfect source RMs?

Setting and Assumptions

Standard Contextual Bandit Framework

- \mathcal{S} : prompt space; \mathcal{A} : response space,
- $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$: LLM as a policy. W.L.O.G., $\pi(\cdot|\cdot) > 0$ everywhere.
- r^* : unknown true (human intrinsic) RM,
- Learning objective:

$$\pi_{r^*}^* \leftarrow \arg \max_{\pi} J_{\beta}(\pi) := \mathbb{E}_{s \sim \rho, a \sim \pi} [r^*(s, a)] - \beta \text{KL}(\pi \| \pi_{\text{ref}}), \quad (1)$$

with ρ as prompt distribution, π_{ref} as the reference policy, and yields:

$$\pi_{r^*}^*(a|s) \propto \pi_{\text{ref}}(a|s) \exp\left(\frac{r^*(s, a)}{\beta}\right), \quad (2)$$

- Bradley-Terry preference model

$$\mathbb{P}_{r^*}(\mathbb{I}[a \succ \tilde{a}] | s, a, \tilde{a}) = \text{Sigmoid}(r(s, a) - r(s, \tilde{a})).$$

Standard Assumptions

- **Bounded rewards:** $r^* \in [0, R]$,
- **Function approximation:** A policy class Π is available.
 - (i) $\pi_{r^*}^* \in \Pi$.
 - (ii) $\forall \pi \in \Pi, \|\log \frac{\pi}{\pi_{\text{ref}}}\|_{\infty} \leq \frac{R}{\beta}$.

Online RLHF with Reward Transfer

Besides human feedback, we access to W source RMs $r^1, \dots, r^W \in [0, R]$:

- no prior knowledge on their quality
- due to Eq. (2), any LLM policy can be converted as a RM.

**Transfer RL has been studied for decades.
But the KL-regularization in Eq. (1) makes something different!**

Blessing of Regularization: A Policy Coverage Perspective

Policy Coverage: The coverage coefficient of policy $\tilde{\pi}$ by another policy π :

$$\text{Cov}^{\tilde{\pi}|\pi} := \mathbb{E}_{s \sim \rho, a \sim \tilde{\pi}} \left[\frac{\tilde{\pi}(a|s)}{\pi(a|s)} \right].$$

Why Policy Coverage Perspective?

- Fundamental complexity measure in online [1] and offline [2] RLHF.
- Optimization/exploration on policy (LLM) space is more efficient in RLHF.

Key Lemma: special structure due to KL regularization ($\beta > 0$)

Lemma 3.1: For any $\pi \in \text{conv}(\Pi) \cup \{\pi_{r^w}^*\}_{w=1}^W$,

$$\text{Cov}^{\pi_{r^*}^*|\pi} = 1 + O(1) \cdot \frac{J_{\beta}(\pi_{r^*}^*) - J_{\beta}(\pi)}{\beta}.$$

Interpretation

- $\text{Cov}^{\pi_{r^*}^*|\pi}$ can be identified by π 's value gap
 - vastly distinguished from pure reward maximization
- KL-reg “reconciles” exploration and exploitation
 - exploiting policies with high policy value coincides with exploration!

Transfer Policy Optimization (TPO): Provably Efficient Online Transfer Learning

Main Idea in TPO Algorithm Design: transfer from π with the lowest $\text{Cov}^{\pi_{r^*}^*|\pi}$

- **Principle 1:** Transfer from the source policy with the highest policy value.
- **Principle 2:** “Self-Transfer Learning”: treat the offline policy distilled from collected data as a transfer candidate (see paper for details).

Theorem 4.4: In sharp contrast to $\tilde{O}(\sqrt{\text{Complex}(\Pi) \cdot T})$ regret in previous works, our TPO achieves:

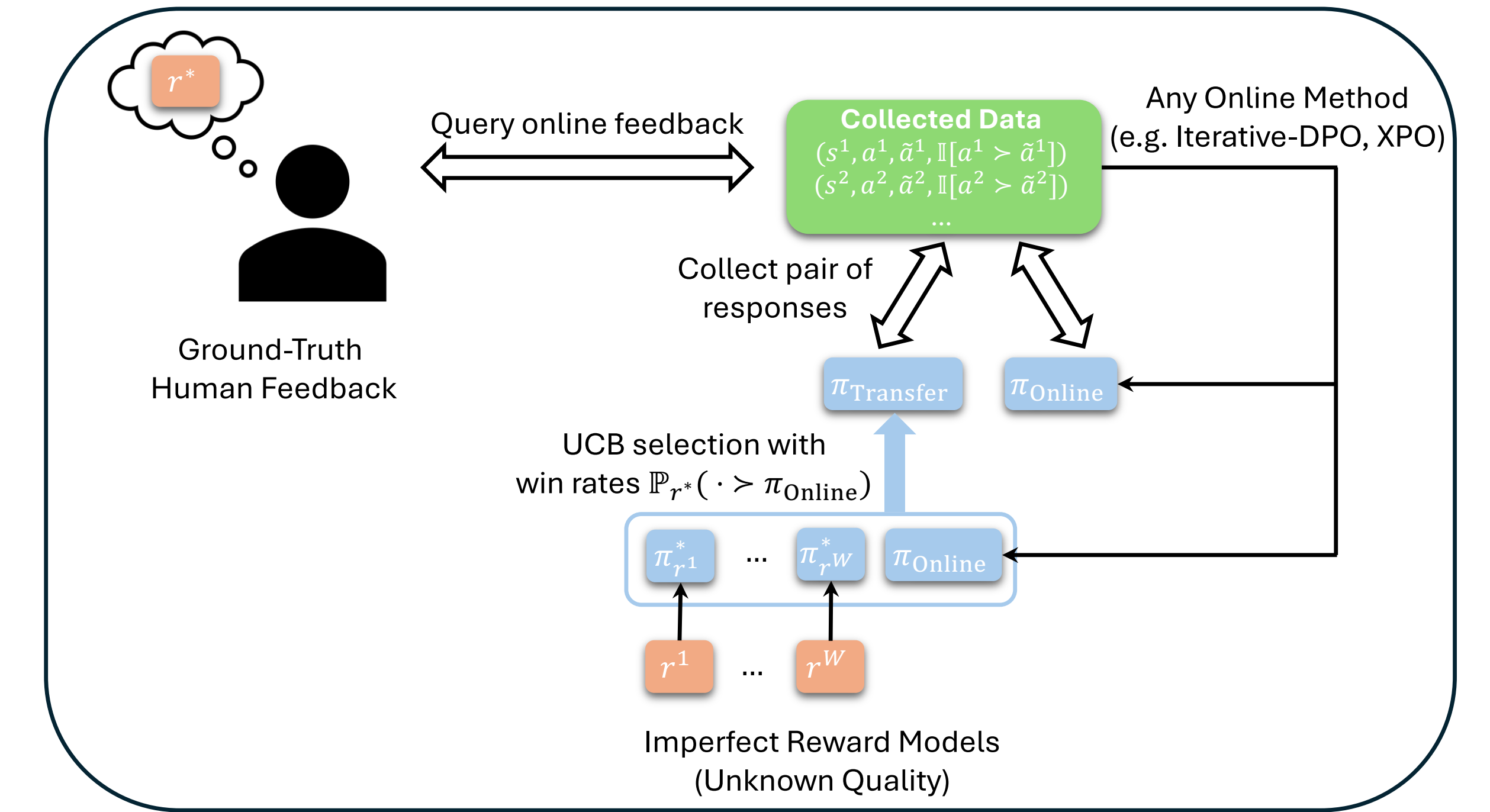
- $\tilde{O}(\sqrt{WT})$, when T is small
 - Reduce dependence on complexity of Π to W
- $\tilde{O}(\sqrt{T})$, when T is large enough
 - No dependence on complexity of Π

Empirical TPO: From Theory to Practice

- Estimating policy value can be computationally expensive.
- Is there a more accessible indicator for $\text{Cov}^{\pi_{r^*}^*|\pi}$? **Yes, the win rates!**
- **Lemma 5.1 [Informal]** A lower bound for $\text{Cov}^{\pi_{r^*}^*|\pi}$:

$$\text{Cov}^{\pi_{r^*}^*|\pi} \geq \max_{\gamma > 0} \frac{1}{\sqrt{\gamma + 2 \cdot \mathbb{P}_{r^*}(\pi \succ \pi_{r^*}^*) \log \frac{1+\gamma}{\gamma}}} \quad (3)$$

- Empirical algorithm design



Main Idea: “Transfer from the expert until beat it”

- $\pi_{r^*}^*$ is unknown, but the online learning policy π_{Online} converges to it.
- Transfer from $\arg \max_{\pi} \mathbb{P}_{r^*}(\pi \succ \pi_{\text{Online}})$.
- Formulate as a multi-armed bandit problem since win rates are unknown.
- Scalable in practice!

Experiments in Summarization Tasks with T5

- Fine-tuning T5-small (60M) on XSum dataset.
- 4 source RMs:
 - 2 metrics of similarity with human summary: (i) ROUGE (ii) BERTScore
 - 2 advanced LLMs: (iii) T5-Base (250M) (iv) T5-Large (770M)
- Llama3-8B to simulate human feedback.

	Without Transfer	Purely Exploit ROUGE	Purely Exploit T5-Large
Iter 1	52.1 ± 1.2	53.1 ± 1.1	49.5 ± 0.9
Iter 2	53.3 ± 1.6	54.5 ± 1.3	49.1 ± 0.4
Iter 3	54.0 ± 1.2	53.3 ± 1.5	50.6 ± 0.3

Table 1. Win rates (%) of the policies trained by empirical TPO competed with 3 baselines.

References

- [1] Xie et al., Exploratory preference optimization: Harnessing implicit q*-approximation for sample-efficient rlhf.
 [2] Liu et al., Provably mitigating overoptimization in rlhf: Your sft loss is implicitly an adversarial regularizer.