# Align-Pro: A Principled Approach to Prompt Optimization for LLM Alignment

Prashant Trivedi[1], Souradip Chakraborty[2], Avinash Reddy[1], Vaneet Aggarwal[3], Amrit Singh Bedi[1], and George K. Atia[1]

[1]University of Central Florida
[2]University of Maryland, College Park
[3]Purdue University

## Abstract

The alignment of large language models (LLMs) with human values is critical as these models become increasingly integrated into various societal and decision-making processes. Traditional methods, such as Reinforcement Learning from Human Feedback (RLHF), achieve alignment by fine-tuning model parameters, but these approaches are often computationally expensive and impractical when models are frozen or inaccessible for parameter modification. In contrast, in this work, we consider prompt optimization as a viable alternative to RLHF for LLM alignment, leveraging the fact that the output of an LLM is a direct function of its input prompt. While prompt optimization has shown empirical promise, its theoretical underpinning remains under-explored. To address this gap, we propose Align Pro, a unified theoretical framework consists of a prompter, designed to rigorously study LLM alignment via prompt optimization. To analyze the prompt optimization process, we derive a closed-form expression for the optimal prompter distribution, and establish theoretical bounds that compare the performance of prompt optimization against RLHF-based fine-tuning. Finally, we provide empirical validation through experiments on the various datasets, demonstrating that prompt optimization can effectively align LLMs, even when parameter fine-tuning is not feasible. Align-Pro achieves approximately a 50% increase in mean rewards and a 40% improvement in win rate over the baseline without parameter fine-tuning, demonstrating its effectiveness across three datasets and various model configurations.

## 1 Introduction

The quest to align large language models (LLMs) with human values is not just an academic pursuit but a practical necessity [35, 16]. As these AI models (e.g., ChatGPT, Llamma2, etc.) increasingly become an essential part of various aspects of daily life and decision-making processes, ensuring their outputs reflect ethical considerations and societal norms becomes
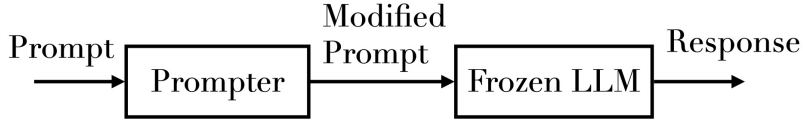
Figure 1: The figure illustrates our setup for aligning a frozen large language model (LLM) through prompt optimization. In this setup, input prompts are refined by a prompter to produce aligned responses from the frozen LLM.

crucial [19, 9]. The standard approach to aligning LLMs has been through fine-tuning parameters via Reinforcement Learning from Human Feedback (RLHF) [42, 3, 43], which involves three main steps: Supervised Fine-Tuning (SFT), reward learning, and RL fine-tuning. However, this process can be resource-intensive, as it necessitates significant updates to the model parameters [5, 26]. A further complication to alignment arises when models are either 'frozen' or operate as 'black boxes,' where direct access to tweak parameters is restricted [10, 31]. These scenarios pose a critical question: *How can we ensure LLM alignment when parameter modification is not allowed or possible?*

One promising solution lies in the concept of *prompt optimization* [22, 21, 18]. This technique leverages the idea that the output of an LLM is fundamentally a function of the input prompt—thereby turning the prompt into a powerful tool to elicit desired responses to align with specific rewards (cf. Figure 1). Various empirical studies in the literature have shown the significant benefits of prompt optimization techniques for LLM alignment [31, 17, 37]. But they often lack a systematic approach grounded in theory. This raises an important question about the efficacy of prompt optimization compared to traditional fine-tuning: *Can prompt optimization achieve comparable or better performance as compared to parameter fine-tuning for LLM alignment?*

Addressing the above question requires a theoretical framework to understand the mechanics of prompt optimization. True to the best of our knowledge, there is a notable absence of literature focusing on a theoretical formulation of prompt optimization specifically for LLM alignment. This paper aims to fill this gap by developing a unified theoretical optimization framework (called Align Pro) to analyze prompt optimization for LLM alignment rigorously. We explore its theoretical performance, particularly in terms of suboptimality bounds, which measure how closely the responses generated via prompt optimization approach to the outcomes achieved through fine-tuning methods. Furthermore, we substantiate our theoretical insights with empirical evidence, demonstrating the practical efficacy of our proposed framework. By bridging the gap between theory and practice, this work not only advances our understanding of prompt optimization but also provides a viable pathway for aligning LLMs in scenarios where direct modification of model parameters is infeasible. We summarize our main contributions as follows.

- **A Theoretical Framework to prompt optimization for LLM alignment.** We propose Align Pro: a novel prompt optimization framework wherein a prompter language model takes an initial prompt as input and generates an optimized prompt. This

optimized prompt is then fed into the target frozen language model to produce the desired response. This setup allows us to theoretically study the prompt optimization for LLM alignment in general settings.

- **Prompt Optimization Objective.** We derive an optimization objective that a prompter should optimize and also derive a closed-form expression for the optimal prompt distribution. The closed form expression provides a direct mathematical formula to identify the optimal prompt, simplifying the process of prompt optimization and providing theoretical insights into the nature of optimal prompts.

- **Theoretical Bounds for performance with respect to fine tuning via RLHF.** We establish theoretical bounds on the difference between the expected rewards obtained from the RLHF policy, which represents the benchmark for model performance, and the optimal policy derived from our prompt optimization approach.

- **Experimental Validity:** We conduct a series of experiments on the three datasets to support the insights we obtain from the theoretical analysis. Align-Pro demonstrates approximately a 50% increase in mean rewards and a 40% improvement in win rate over the baseline without fine-tuning, showcasing its effectiveness across three datasets and diverse model configurations.

## 2  Related Works

**RLHF and LLM fine-tuning:** RLHF has become the most widely used method for aligning LLM responses with human values [11, 26, 44]. For a more comprehensive discussion on RLHF, refer to some recent surveys [5, 7]. Recently, some methods have been developed to bypass the need for RL, directly utilizing a preference dataset for alignment, including Direct Preference Optimization (DPO) [29], SLiC [39], and other extensions [1, 2, 13, 23, 24, 34, 36]. The recent work of [12] has demonstrated the potential of efficient exploration methods to improve LLM responses based on human preference feedback. Moreover, methods such as ORPO [14] align the model without using a reference model. Furthermore, Intuitive Fine Tuning (IFT) conducts alignment solely relying on positive samples and a single policy, starting from a pre-trained base model [15]. But all of these approaches are focused on alignment via parameter fine tuning.

**Prompt optimization for Alignment:** The field of prompt optimization has grown substantially in recent years. Early work focused on optimizing prompts for white-box LLMs, such as AutoPrompt [31] and FluentPrompt [30]. These also include methods based on soft prompts, like those in [18, 20, 40]. However, these model majorly use the gradient-based methods to update the prompts from labeled data. Recent attention has shifted toward optimizing prompts for black-box LLMs. Methods such as Clip-Tuning [6], BBT [33], and BBTv2 [32], which rely on access to input embeddings and output logits of the black-box LLM, were among the most popular. We are also interested modifying the prompt without tinkering the frozen LLM model, but with the theoretical understanding in terms of how well we can do

with respect to fine tuning methods. Some prompt optimization techniques use reinforcement learning methods, such as BDPL [10] and PRewrite [17], while others, like PromptAgent [37] are purely planning-based. Unlike our prompt optimization framework, these method use iterative update of the prompt by first initializing it using an LLM, and then training the LLM using RL. Recently, the APOHF approach [22] was introduced, utilizing dueling bandits theory to select pairs of prompts for preference feedback in each iteration. To the best of our knowledge, there is little to no existing research that provides a theoretical foundation for prompt optimization frameworks in alignment and compares them with fine-tuning methods, which is the focus of our work. We want to emphasize that our goal is not to design the best prompter but rather to develop a unified framework for prompt optimization.

# 3 Preliminaries and Background

This section provides the essential background and foundational concepts relevant to alignment. We start by defining the notation, followed by a quick overview of the RLHF framework, which involves three key steps: (i) supervised fine-tuning (SFT), (ii) reward learning, and (iii) fine-tuning with RL.

**Language Models.** We start by defining the language model mathematically. Let us denote the vocabulary set by $\mathcal{V}$, and we denote the language model by $\pi(y|x)$, which takes in the sequence of tokens $x := \{x_1, x_2, \cdots, x_N\}$ (with each $x_i \in \mathcal{V}$) as an input, and generate response $y := \{y_1, y_2, \cdots, y_M\}$ (with each $y_i \in \mathcal{V}$) as the output. At instant $t$, each output token $y_t \sim \pi(\cdot|x_t)$.

**Supervised Fine-Tuning (SFT).** SFT is the initial step in the RLHF process. It involves fine-tuning a pre-trained LLM on a vast dataset of human-generated text in a supervised manner.

**Reward Learning.** This stage involves learning the reward model by gathering preferences from experts/human feedback or an oracle based on outputs generated by the SFT model denoted by $\pi_{\text{sft}}$. The optimization is generally performed under the Bradley-Terry model for pairwise comparison [4], which seeks to minimize the loss formulated as:

$$\mathcal{L}(r, D_r) = -\mathbb{E}_{(x,y_u,y_v)\sim D_r} \left[\log\left(\sigma(r(x, y_u) - r(x, y_v))\right)\right] \tag{1}$$

where $D_r$ denotes the dataset of response pairs $(y_u, y_v)$, with $y_u$ and $y_v$ representing the winning and the losing responses, respectively, which are generated by the policy $\pi_{\text{sft}}$ optimized under the reward $r(x, y)$, and evaluated by human experts or an oracle function $p^*(\cdot|y_u, y_v, x)$, and $\sigma(\cdot)$ is the sigmoid function.

**Fine-tuning with RL.** In this step, we obtain the aligned model which maximizes the reward model $r(x, y)$ (trained in the previous step) by solving a KL-regularized optimization problem:

$$\pi_{\text{RLHF}} := \arg\max_{\pi} \mathbb{E}_{x\sim P, y\sim\pi(\cdot|x)} \left[r(x, y) - \beta\mathbb{D}_{KL}(\pi(\cdot|x)\|\pi_{\text{sft}}(\cdot|x))\right], \tag{2}$$
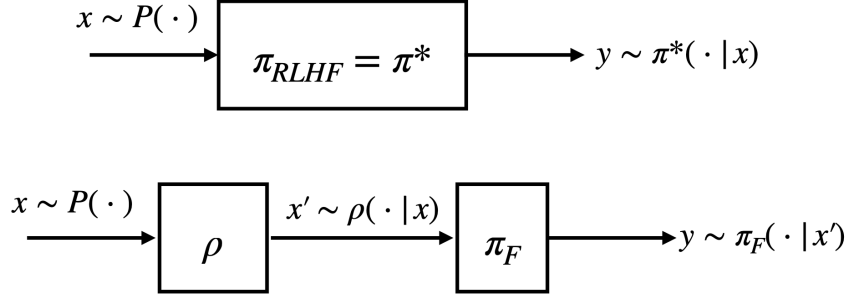
Figure 2: The top figure depicts the optimal RLHF framework, in which responses are generated by the RLHF policy, $\pi_{\text{RLHF}}$. In contrast, the bottom figure illustrates our prompt optimization framework: an input prompt $x \sim P(\cdot)$ is modified by the prompter $\rho$ which updates the prompt. This refined prompt is then passed to the frozen LLM, $\pi_F$, to generate aligned responses.

where, $\beta > 0$ is a parameter controlling the deviation from the baseline policy $\pi_{\text{sft}}$. This iterative process alternates between updating the policy and reward models until convergence, as detailed in prior works [16, 42].

# 4  Prompt Optimization Framework for LLM Alignment

In this section, we provide a mathematical formulation for the framework of prompt optimization for LLM alignment. In traditional LLM alignment, as described in (2), the model parameters are fine-tuned to adjust the response distributions in a way that maximizes the reward function. However, in our setting, we operate under a different regime, starting with a pre-trained language model, denoted by $\pi_F$, whose parameters remain frozen. In this case, direct modification of the model to align with a reward function is not possible. Therefore, an alternative and widely adopted approach in the literature is to optimize the input prompt itself to yield better-aligned responses [17, 31, 41]. Typically, this process involves iterative refinement of prompts, where the outputs of the model are evaluated, compared to human preferences, and the prompts are adjusted accordingly. However, such iterative fine-tuning can be computationally expensive and time-intensive.

To address this challenge, we propose a general prompt optimization framework that introduces a secondary language model—referred to as the "prompter," denoted by $\rho$. This prompter model is typically smaller than the frozen LLM $\pi_F$ and operates on the input prompt. Instead of directly feeding the original prompt $x$ to the frozen LLM $\pi_F$, we first pass $x$ through the prompter model, which generates a modified prompt $x' \sim \rho(\cdot|x)$. This modified prompt $x'$ is then fed into $\pi_F$, which produces the final response $y \sim \pi_F(\cdot|x')$. This process is illustrated in Figure 2, where we visually depict the sequence of steps described above. Interestingly, although we cannot fine-tune the frozen model $\pi_F$, we have the flexibility to fine-tune the prompter model $\rho$ in any desired manner. However, a fundamental challenge arises: what

should be the objective for optimizing the prompter? This remains an open question. While there is substantial empirical evidence in the literature showing that prompt optimization can significantly enhance response generation and thus improve alignment [31, 17, 41], there is no established mathematical framework to guide this process. Addressing this gap is the focus of this our work.

**Optimization Objective for Prompter Design.** To explore this, we revisit the fundamentals of LLM alignment. For a given prompt $x$, the probability of generating a response $y$ from the frozen model is represented by $\pi_F(y|x)$. After introducing the prompter model $\rho$, the probability of generating response $y$ given input $x$ (denoted by $\widetilde{\pi}_\rho$) is now expressed as:

$$\widetilde{\pi}_\rho(y|x) = \sum_{x'} \pi_F(y|x')\rho(x'|x), \tag{3}$$

which captures the average probability of generating the response $y$ for a given $x$ under the influence of the prompter $\rho$. Let us consider the ideal scenario: if we were able to fine-tune the language model $\pi_F$, we would solve the optimization problem in (2) and obtain the RLHF optimal solution $\pi^*$, which is given by [27, 28]

$$\pi^*(y|x) = \frac{1}{Z^*(x)} \pi_F(y|x) \exp\left(\frac{r^*(x,y)}{\beta}\right) , \tag{4}$$

where $Z^*(x) = \sum_y \pi_F(y|x) \exp(r^*(x,y)/\beta)$ is the normalizing constant, and $\beta$ is the alignment tuning parameter, and reward $r^*$ is obtained from solving (1). We emphasize that if we have a prompter $\rho$ that performs as well as the RLHF-optimal policy $\pi^*$, it should be a sufficient indicator of a good prompter. With this understanding, we consider the following prompter suboptimality gap given by

$$\triangle(\rho) := |J(\pi^*) - J(\widetilde{\pi}_\rho)|, \tag{5}$$

which captures how well our prompter is doing with respect to fine-tuned optimal policy $\pi^*$. Mathematically, it holds that

$$J(\pi^*) - J(\widetilde{\pi}_\rho) = \mathbb{E}_{x \sim P, y \sim \pi^*(\cdot|x)}[r^*(x,y)] - \mathbb{E}_{x \sim P, y \sim \widetilde{\pi}_\rho(\cdot|x)}[r^*(x,y)]$$

$$= \mathbb{E}_{x \sim P} \left[ \mathbb{E}_{y \sim \pi^*(\cdot|x)}[r^*(x,y)] - \mathbb{E}_{\substack{x' \sim \rho(\cdot|x) \\ y \sim \pi_F(\cdot|x')}}[r^*(x,y)] \right]. \tag{6}$$

The Equation (6) evaluates the difference in expected return between the optimal RLHF policy $\pi^*$ and our prompt optimization policy $\widetilde{\pi}_\rho$, indicating how much better (or worse) $\pi^*$ performs compared to $\widetilde{\pi}_\rho$. We highlight that this performance gap is clearly influenced by the choice of the prompt distribution $\rho$; a non-optimal $\rho$ can result in a significant gap. This leads us to the following questions:

- **Q1**: Can we design an optimal prompter $\rho^*$ that closes the suboptimality gap between the fine-tuned policy $\pi^*$, and the prompt optimization policy $\widetilde{\pi}_{\rho^*}$ as mentioned in Equation (6)?

- **Q2**: If such a $\rho^*$ exists, then can $\widetilde{\pi}_{\rho^*}$ outperform the fine-tuned optimal policy $\pi^*$?

To the best of our knowledge, the above questions remain unanswered in the existing literature. We address these questions in the next section.

# 5 Proposed Approach: Align-Pro

In this section, we address both questions Q1 and Q2. Let us start by addressing Q1 and develop a general prompt optimization framework to design an optimal prompter $\rho^*$. But then the first question arises: in what sense is $\rho^*$ optimal? In order to see that, let us reconsider $J(\pi^*) - J(\widetilde{\pi}_\rho)$ and after adding-subtracting $\mathbb{E}_{y \sim \pi_F(\cdot|x)}[r^*(x,y)]$ in the right hand side of Equation (6), we get

$$J(\pi^*) - J(\widetilde{\pi}_\rho) = \mathbb{E}_{x \sim P}[\Delta_1 + \Delta_2], \tag{7}$$

where $\Delta_1$ and $\Delta_2$ are defined as

$$\Delta_1 := \mathbb{E}_{y \sim \pi^*(\cdot|x)}[r^*(x,y)] - \mathbb{E}_{y \sim \pi_F(\cdot|x)}[r^*(x,y)]$$

$$\Delta_2 := \mathbb{E}_{y \sim \pi_F(\cdot|x)}[r^*(x,y)] - \mathbb{E}_{y \sim \widetilde{\pi}_\rho(\cdot|x)}[r^*(x,y)]$$

$$= \mathbb{E}_{y \sim \pi_F(\cdot|x)}[r^*(x,y)] - \mathbb{E}_{\substack{x' \sim \rho(\cdot|x) \\ y \sim \pi_F(\cdot|x')}}[r^*(x,y)].$$

We remark that in (7), $\Delta_1$ is the suboptimality gap between the optimal fine-tuned policy, and the frozen model $\pi_F$. Thus, it captures the effectiveness of the optimal RLHF policy with respect to the frozen model. In other words, it quantifies how good or bad our frozen model is with respect to the optimally aligned model. We note that $\Delta_1$ is constant for a given $\pi_F$ and does not depend upon prompter $\rho$, hence we cannot improve this part with the prompter. Another insight is that since $\pi^*$ is the optimal RLHF policy, $\Delta_1 \geqslant 0$, i.e., is always positive. On the other hand, the second term, $\Delta_2$, depends upon our prompter $\rho$ and can be controlled by designing a prompter. This observation leads to the formulation of an optimization problem for the prompter as follows.

## 5.1 Optimization Problem for Prompter

We recall from the definition of $\Delta_2$ that we would need to learn a $\rho$ such that $\Delta_2$ is minimized. To achieve that, we recognize that the only term involving the prompter $\rho$ in $\Delta_2$ is $\mathbb{E}_{x' \sim \rho(\cdot|x), y \sim \pi_F(\cdot|x')}[r^*(x,y)]$, and minimizing $\Delta_2$, we should solve the following optimization problem

$$\max_\rho \mathbb{E}_{x' \sim \rho(\cdot|x), y \sim \pi_F(\cdot|x')}[r^*(x,y)]. \tag{8}$$

However, at the same time, since our prompter is also another language model, we will already have access to a baseline supervised fine-tuned prompter $\rho_{\mathrm{sft}}$, and we want to ensure that our prompter $\rho^*$ does not deviate significantly from $\rho_{\mathrm{sft}}$, which motivates us to include a

known and supervised fine-tuned prompter, denoted by $\rho_{\text{sft}}$. Thus, we solve the following optimization problem:

$$\max_{\rho} \mathbb{E}_{\substack{x' \sim \rho(\cdot|x) \\ y \sim \pi_F(\cdot|x')}} [r^*(x, y)] - \lambda \mathbb{D}_{KL}(\rho(\cdot|x) \| \rho_{\text{sft}}(\cdot|x)). \tag{9}$$

We have introduced a KL-divergence-based regularizer above between the prompter $\rho$ and a reference supervised fine-tuned prompter $\rho_{\text{sft}}$. This helps with the development of a proper optimization problem with a closed-form expression, and enables control over proximity to the initial prompter $\rho_{\text{sft}}$ through the tuning parameter $\lambda$.

**Interpretation of $\lambda$.** Another interesting interpretation of $\lambda$ is that it controls the extent of prompt optimization we want to introduce into the pipeline, hence we also refer to it as the prompt tuning parameter. For instance, $\lambda \to \infty$ means no prompt optimization, while $\lambda \to 0$, drives the optimization toward maximizing the prompter reward, albeit at the cost of deviating from $\rho_{\text{sft}}$ which might be important in certain cases. Therefore, $\lambda$ provides a meaningful trade-off, and its effects will be further elucidated in the theoretical section that follows.

The following Lemma 5.1 provides the optimal solution to the optimization problem (9).

**Lemma 5.1.** *Let $R(x, x') := \mathbb{E}_{y \sim \pi_F(\cdot|x')}[r^*(x, y)]$, and $\lambda > 0$ be the prompter tuning parameter. The optimal prompt distribution $\rho^*$ that maximizes the objective function of the optimization problem (9) is given by:*

$$\rho^*(x'|x) = \frac{1}{Z(x)} \rho_{\text{sft}}(x'|x) \exp\left(\frac{1}{\lambda} R(x, x')\right), \tag{10}$$

*where $Z(x)$ is the log partition function given by*

$$Z(x) = \sum_{x'} \rho_{\text{sft}}(x'|x) \exp\left(\frac{1}{\lambda} R(x, x')\right).$$

The proof of Lemma 5.1 is deferred to the technical Appendix A. Now, we have answered Q1. Next, we move to answer Q2 in which we utilize the optimal prompter $\rho^*(x'|x)$ to obtain a bound on the suboptimality gap. Notably, the integration of this optimal prompter with the frozen model will lead to the refined performance expressed in terms of the modified optimal policy $\widetilde{\pi}_\rho^*(y|x) = \sum_{x'} \rho^*(x'|x)\pi_F(y|x')$. This refined bound captures the effectiveness of the prompt optimization process and offers a clearer understanding of how closely the modified policy $\widetilde{\pi}_{\rho^*}$ approximates the true optimal policy $\pi^*$.

# 6  Theoretical Results w.r.t Fine Tuning

We begin by establishing a bound on the suboptimality gap for the optimal prompter. The following theorem bounds the suboptimality gap $\triangle(\rho^*) = |J(\pi^*) - J(\widetilde{\pi}_{\rho^*})|$ when the optimal prompter $\rho^*$ as obtained in Lemma 5.1 is used. We present our result in Theorem 6.1 as follows. The proof of this theorem is deferred to the technical Appendix B.

**Theorem 6.1.** *Let the optimal prompter $\rho^*(x'|x)$ be given as in (10). Then, the suboptimality gap is given by*

$$\triangle(\rho^*) \leqslant r_{\max} \mathbb{E}_{x \sim P} \sqrt{2\mathbb{E}_{\pi^\star}\left[\log\left(\frac{\pi^*(y|x)}{\pi_F(y|x)}\right)\right]} + \lambda \; \mathbb{E}_{x \sim P}[\mathbb{D}_{KL}(\rho^*(\cdot|x)\|\rho_{\mathrm{sft}}(\cdot|x))], \qquad (11)$$

*where $P$ denotes the prompt distribution, $\lambda$ is the prompter tuning parameter, and $\triangle(\rho) = |J(\pi^*) - J(\widetilde{\pi}_\rho)|$.*

*Remark* 6.2. If $\pi_F = \pi_{\mathrm{sft}}$, that is, if the frozen model is just the base or supervised fined-tune model, then we have a tighter suboptimality gap. Moreover, the reward $R(x, x')$ will be defined in terms of the SFT policy $\pi_{\mathrm{sft}}$, i.e., $R(x, x') = \mathbb{E}_{y \sim \pi_{\mathrm{sft}}(\cdot|x')}[r^*(x, y)]$.

Theorem 6.1 provides an upper bound on the suboptimality gap between an optimal RLHF policy $\pi^*$ and the optimal policy obtained by our prompt optimization approach Align-Pro that uses the optimal prompter $\rho^*$ and generates the response from the policy $\widetilde{\pi}_{\rho^*}$. This upper bound is controlled by a term involving the KL divergence, $\lambda$, reward model, and the parameter $\beta$. We now provide the interpretations to each term of the suboptimality gap given in Theorem 6.1.

- **Significance of first term in RHS of** (11)**:** The first term in Equation (11) is always non-negative. Thus, it captures the intrinsic difficulty of obtaining the optimal RLHF policy via a prompt optimization setup when the frozen model is not fully aligned. Therefore, the suboptimality gap is non-negative, irrespective of the prompter used. We note that when $\pi_F = \pi^*$, the first term in Theorem 6.1 becomes zero. However, this scenario is not relevant to our prompt optimization framework, as it necessitates fine-tuning the frozen LLM.

- **Significance of second term in RHS of** (11)**:** The second term captures the KL divergence between the optimal prompter $\rho^*$ and the supervised fine-tuned prompter $\rho_{\mathrm{sft}}$. This ensures that the optimal $\rho^*$ does not deviate significantly from the SFT prompter, and allowing control over their proximity.

# 7 Experimental Evaluations

In this section, we present the experimental setup and results to validate the effectiveness of our proposed prompt optimization framework, Align-Pro. We begin by outlining our experimental setup, including the dataset, model architecture, and evaluation metrics. Following this, we present our results and provide a detailed analysis of our findings.

## 7.1 Experimental Setup

We evaluate the performance of our Align-Pro framework using two distinct prompter models, denoted as P1 (Phi-3.5-Instruct) and P2 (Qwen-2.5-1.5B-Instruct), which modifies

and updates the original prompt. Additionally, we use two frozen models, denoted as F1 (Llama-3.1-8B-Instruct) and F2 (Llama-3.1-8B-Instruct) to generate the final responses. This setup results in four unique model architectures, each representing a combination of the prompter and frozen models. For each architecture, we assess performance for the following three different configurations.

- **No Fine-Tuning**: In this configuration, the prompter is not used, and only the frozen model is used to generate responses without any fine-tuning or prompt modifications.
- **Align-Pro**: In this setup, a fine-tuned prompter is placed before a frozen. The prompter refines the input prompt, and the frozen model generates the response based on the optimized prompt.
- **RLHF**: In this configuration, the frozen model undergoes fine-tuning through Reinforcement Learning from Human Feedback (RLHF), and the response is generated directly from this fine-tuned model.

**Datasets:** To capture the diversity in our experimental evaluations, we evaluate the performance over different datasets:

- **UltraFeedback** [8] : A large-scale, high-quality, and diversified AI feedback dataset which contains feedback from user-assistant conversations from various aspects. This dataset evaluates the coherence of the prompt-response pairs.

- **HelpSteer** [38]: A multi-attribute helpfulness dataset annotated for correctness, coherence, complexity, and verbosity in addition to overall helpfulness of responses.

- **Orca** [25]: This dataset features responses with detailed explanations for each prompt, promoting thinking and effective instruction-following capabilities in the models.

**Evaluation Criteria.** The primary objective of our experiments is to optimize the input prompt to guide the frozen LLM that produces the desired response effectively. We fine-tune the prompter using Proximal Policy Optimization (PPO) within the RLHF framework to achieve this. The reward signal for this fine-tuning process is derived from the quality of the enhanced prompt and the output generated by the frozen LLM. We assess the performance of our Align-Pro framework based on three key metrics: mean reward, variance, and win-rate comparison against the no-fine-tuning baseline.

**Computational Resources.** Since we do not alter the parameters of the frozen model, our experiments require relatively fewer computational resources. Consequently, we were able to conduct all our experiments using a machine equipped with an INTEL(R) XEON(R) GOLD 6526Y processor with a Nvidia H100 GPU. We used Python 3.11 to execute the experiments. we used the *PPOTrainer* variant from Hugging Face TRL library to run the RLHF and Prompt Optimization pipeline experiments.

**Hyper-parameters used.** All of our experiments use the open-access TRL library, which is publicly available. The library can be accessed using the following link. For our experiments, we do not perform any extra hyper-parameter tuning; rather, we use the parameters *learning rate* $= 1.41e - 5$ given in the above-mentioned link. Moreover, we use the

following generation configurations to generate the response for evaluation in all experiments: temperature $= 1.5$, top $P = 0.6$ and top $K = 20$.

## 7.2  Results

**Mean reward and variance comparison:** We calculate mean rewards and variances to assess the quality of preferred response generation and the diversity of the language model for all configurations and different model architectures. To associate the reward to each response, we use the available reward model[1], which scores the response. This reward model is trained to assign higher scores to the responses that comply with the off-target attributes.

We also compared Align-Pro with an oracle model, where the LLM is fine-tuned using reinforcement learning from human feedback (RLHF). Figure 3 presents the mean rewards across all three datasets for each model configuration, while Figure 4 shows the corresponding reward variances. Interestingly, Align-Pro consistently outperforms the baseline (no fine-tuning) in terms of mean reward, demonstrating its ability to generate more preferred and stable responses. Moreover, the variance in reward for Align-Pro is the lowest, indicating that it produces more reliable and stable outputs. In each figure, we employ two prompters, denoted as P1 (Phi-3.5-Instruct) and P2 (Qwen-2.5-1.5B-Instruct), along with two frozen LLMs, denoted as F1 (Llama-3.1-8B-Instruct) and F2 (Llama-3.1-8B-Instruct).
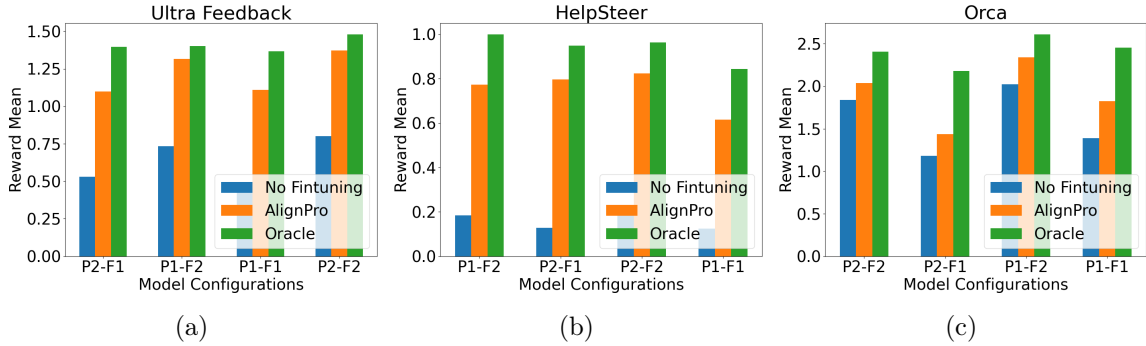


Figure 3: **Reward mean comparisons.** Figure shows the reward mean across the chosen datasets. Align-Pro shows a significant improvement over the no-fine-tuning approach. We employ two prompters P1 (Phi-3.5-Instruct) and P2 (Qwen-2.5-1.5B-Instruct), along with two frozen LLMs, denoted as F1 (Llama-3.1-8B-Instruct) and F2 (Llama-3.1-8B-Instruct). The oracle is fine-tuned LLM via RLHF.

**Win rate comparison:** We evaluate the performance of our Align-Pro method by comparing it to the no fine-tuning configuration using win rate as the primary performance metric. We rely on GPT-4 as an external, impartial judge to ensure unbiased evaluation. The evaluation criteria focus on critical aspects of the response: helpfulness, harmlessness, relevance, accuracy, depth, creativity, and level of detail. To update the prompt, we use a standardized system prompt template as detailed in Appendix C.2. Table 1 presents the win rates for Align-Pro

---

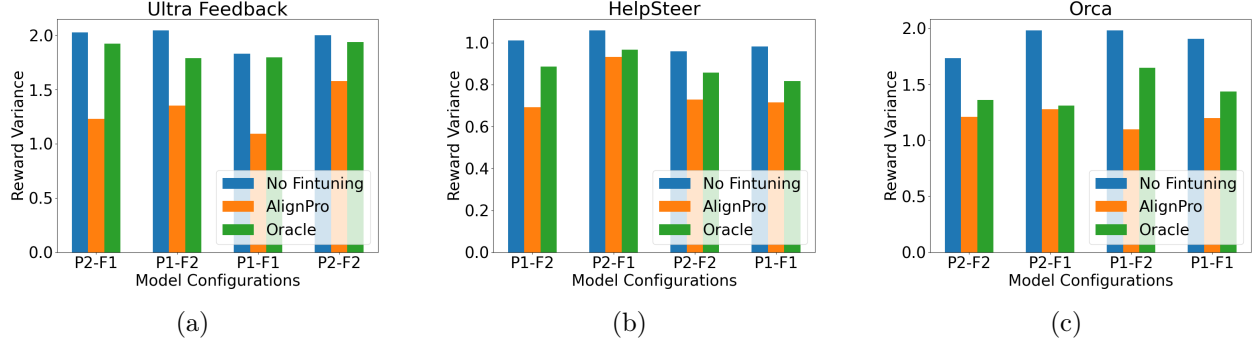[1] https://huggingface.co/weqweasdas/RM-Gemma-2B

11

Figure 4: **Reward variance comparisons.** Align-Pro has the least variance compared to Oracle and no finetuning approach. Due to the prompter's precise guidance, frozen LLM generates almost similar responses in terms of helpfulness and coherence, which results in less diverse responses. We use the following terminologies for the prompters and the frozen models: P1 (Phi-3.5-Instruct), P2 (Qwen-2.5-1.5B-Instruct), F1 (Llama-3.1-8B-Instruct), and F2 (Llama-3.1-8B-Instruct), respectively.

(denoted by A) against the no fine-tuning baseline (denoted by B). The results clearly show that, on average, Align-Pro significantly outperforms the no fine-tuning approach across all model architectures and datasets. These findings demonstrate the effectiveness of our Align-Pro framework, which enhances performance by optimizing the input prompt while keeping the large language model frozen.

| Model Architectures | UltraFeedback | | | HelpSterer | | | Orca | | |
|---|---|---|---|---|---|---|---|---|---|
| Prompter, Frozen LLM | A win | B win | Tie | A win | B win | Tie | A win | B win | Tie |
| Phi-3.5-Instruct, Llama-3.1-8B-Instruct | **60** | 24 | 16 | **46** | 37 | 17 | **63** | 26 | 11 |
| Qwen-2.5-1.5B-Instruct, Llama-3.1-8B-Instruct | **65** | 23 | 12 | **67** | 23 | 10 | **63** | 30 | 4 |
| Phi-3.5-Instruct, Qwen-2.5-7B-Instruct | **59** | 27 | 14 | **58** | 27 | 15 | **46** | **46** | 8 |
| Qwen-2.5-1.5B-Instruct, Qwen-2.5-7B-Instruct | **56** | 30 | 14 | **59** | 25 | 16 | **59** | 27 | 14 |

Table 1: The table presents the win rates (for 100 samples) of our Align-Pro method, denoted by **A**, compared to the baseline no fine-tuning method, denoted by **B**. A higher win rate indicates superior performance. Bolded numbers highlight the higher win rates. Across all model architectures and datasets, Align-Pro consistently outperforms the no fine-tuning baseline, demonstrating its effectiveness in improving response quality.

**Summary:** Our experiments confirm that using a prompter alongside a frozen LLM significantly enhances alignment. Moreover, the expected reward and the win-rate differences are affected by the degree to which the prompter and frozen model align with human preferences. These experimental results, therefore, strongly support our theoretical findings.

*Remark* 7.1. Our aim is not to present the best prompt optimizer, but to develop a unified theoretical framework and a principled approach to prompt optimization. We seek to

understand its theoretical performance relative to RLHF and fine-tuning methods, hence we did not compare our approach with other prompt optimization methods in our experimental section.

# 8    Conclusion, Limitations and Future Work

This work introduces a novel theoretical framework for prompt optimization by utilizing a smaller, trainable model to generate optimized prompts for a frozen large language model (LLM). This approach significantly reduces computational costs while preserving the LLM's pre-trained capabilities. We provide a closed-form expression for the optimal prompter and use it to establish theoretical bounds on suboptimality gap that compares the optimized prompt policy with the standard RLHF policy. We demonstrate the effectiveness of our method across three datasets and various model configurations. In each scenario, we observe an approximate 50% improvement in mean rewards and a 40% higher win rate compared to the baseline with no fine-tuning.

**Limitations:** Our prompt optimization framework is inherently limited by the capabilities of the frozen language model. Another limitation includes the sensitivity of the prompt to the final response; a slight change in the prompt can lead to profound changes in the final responses. We will consider some of these issues as part of our future work.

**Future Work:** Some potential future directions of our work include analyzing the robustness of the optimal prompter in the presence of noise in the frozen model and exploring the use of multiple prompters in sequence before inputting them into the frozen model.

# References

[1] Afra Amini, Tim Vieira, and Ryan Cotterell. Direct preference optimization with an offset. *arXiv preprint arXiv:2402.10571*, 2024.

[2] Mohammad Gheshlaghi Azar, Zhaohan Daniel Guo, Bilal Piot, Remi Munos, Mark Rowland, Michal Valko, and Daniele Calandriello. A general theoretical paradigm to understand learning from human preferences. In *International Conference on Artificial Intelligence and Statistics*, pages 4447–4455. PMLR, 2024.

[3] Mohammad Gheshlaghi Azar, Mark Rowland, Bilal Piot, Daniel Guo, Daniele Calandriello, Michal Valko, and Rémi Munos. A general theoretical paradigm to understand learning from human preferences. *arXiv preprint arXiv:2310.12036*, 2023.

[4] Ralph Allan Bradley and Milton E Terry. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952.

[5] Stephen Casper, Xander Davies, Claudia Shi, Thomas Krendl Gilbert, Jérémy Scheurer, Javier Rando, Rachel Freedman, Tomasz Korbak, David Lindner, Pedro Freire, et al.

Open problems and fundamental limitations of reinforcement learning from human feedback. *arXiv preprint arXiv:2307.15217*, 2023.

[6] Yekun Chai, Shuohuan Wang, Yu Sun, Hao Tian, Hua Wu, and Haifeng Wang. Clip-tuning: Towards derivative-free prompt learning with a mixture of rewards. In *Proc. EMNLP (Findings)*, pages 108–117, 2022.

[7] Shreyas Chaudhari, Pranjal Aggarwal, Vishvak Murahari, Tanmay Rajpurohit, Ashwin Kalyan, Karthik Narasimhan, Ameet Deshpande, and Bruno Castro da Silva. RLHF Deciphered: A Critical Analysis of Reinforcement Learning from Human Feedback for LLMs. *arXiv preprint arXiv:2404.08555*, 2024.

[8] Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Bingxiang He, Wei Zhu, Yuan Ni, Guotong Xie, Ruobing Xie, Yankai Lin, et al. Ultrafeedback: Boosting language models with scaled ai feedback. In *Forty-first International Conference on Machine Learning*, 2024.

[9] Josef Dai, Xuehai Pan, Ruiyang Sun, Jiaming Ji, Xinbo Xu, Mickel Liu, Yizhou Wang, and Yaodong Yang. Safe rlhf: Safe reinforcement learning from human feedback. *arXiv preprint arXiv:2310.12773*, 2023.

[10] Shizhe Diao, Zhichao Huang, Ruijia Xu, Xuechun Li, LIN Yong, Xiao Zhou, and Tong Zhang. Black-box prompt learning for pre-trained language models. *Transactions on Machine Learning Research*, 2023.

[11] Yann Dubois, Chen Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin, Percy S Liang, and Tatsunori B Hashimoto. Alpacafarm: A simulation framework for methods that learn from human feedback. *Advances in Neural Information Processing Systems*, 36, 2024.

[12] Vikranth Dwaracherla, Seyed Mohammad Asghari, Botao Hao, and Benjamin Van Roy. Efficient exploration for llms. *arXiv:2402.00396*, 2024.

[13] Qi Gou and Cam-Tu Nguyen. Mixed preference optimization: Reinforcement learning with data selection and better reference model. *arXiv preprint arXiv:2403.19443*, 2024.

[14] Jiwoo Hong, Noah Lee, and James Thorne. Orpo: Monolithic preference optimization without reference model, 2024.

[15] Ermo Hua, Biqing Qi, Kaiyan Zhang, Yue Yu, Ning Ding, Xingtai Lv, Kai Tian, and Bowen Zhou. Intuitive fine-tuning: Towards simplifying alignment into a single process, 2024.

[16] Timo Kaufmann, Paul Weng, Viktor Bengs, and Eyke Hüllermeier. A survey of reinforcement learning from human feedback. *arXiv preprint arXiv:2312.14925*, 2023.

[17] Weize Kong, Spurthi Amba Hombaiah, Mingyang Zhang, Qiaozhu Mei, and Michael Bendersky. PRewrite: Prompt Rewriting with Reinforcement Learning. *arXiv preprint arXiv:2401.08189*, 2024.

[18] Brian Lester, Rami Al-Rfou, and Noah Constant. The Power of Scale for Parameter-Efficient Prompt Tuning. In *Proc. EMNLP*, pages 3045–3059, 2021.

[19] Aaron J Li, Satyapriya Krishna, and Himabindu Lakkaraju. More rlhf, more trust? on the impact of human preference alignment on language model trustworthiness. *arXiv preprint arXiv:2404.18870*, 2024.

[20] Xiang Lisa Li and Percy Liang. Prefix-Tuning: Optimizing Continuous Prompts for Generation. In *Proc. ACL*, pages 4582–4597, 2021.

[21] Yuanzhi Li and Yingyu Liang. Learning overparameterized neural networks via stochastic gradient descent on structured data. *Advances in neural information processing systems*, 31, 2018.

[22] Xiaoqiang Lin, Zhongxiang Dai, Arun Verma, See-Kiong Ng, Patrick Jaillet, and Bryan Kian Hsiang Low. Prompt optimization with human feedback. *arXiv preprint arXiv:2405.17346*, 2024.

[23] Tianqi Liu, Zhen Qin, Junru Wu, Jiaming Shen, Misha Khalman, Rishabh Joshi, Yao Zhao, Mohammad Saleh, Simon Baumgartner, Jialu Liu, et al. LiPO: Listwise Preference Optimization through Learning-to-Rank. *arXiv preprint arXiv:2402.01878*, 2024.

[24] Tetsuro Morimura, Mitsuki Sakamoto, Yuu Jinnai, Kenshi Abe, and Kaito Air. Filtered direct preference optimization. *arXiv preprint arXiv:2404.13846*, 2024.

[25] Subhabrata Mukherjee, Arindam Mitra, Ganesh Jawahar, Sahaj Agarwal, Hamid Palangi, and Ahmed Awadallah. Orca: Progressive learning from complex explanation traces of gpt-4. *arXiv preprint arXiv:2306.02707*, 2023.

[26] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.

[27] Xue Bin Peng, Aviral Kumar, Grace Zhang, and Sergey Levine. Advantage-weighted regression: Simple and scalable off-policy reinforcement learning. *arXiv preprint arXiv:1910.00177*, 2019.

[28] Jan Peters and Stefan Schaal. Reinforcement learning by reward-weighted regression for operational space control. In *Proceedings of the 24th international conference on Machine learning*, pages 745–750, 2007.

[29] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36, 2024.

[30] Weijia Shi, Xiaochuang Han, Hila Gonen, Ari Holtzman, Yulia Tsvetkov, and Luke Zettlemoyer. Toward human readable prompt tuning: Kubrick's the shining is a good movie, and a good prompt too? In *Proc. EMNLP*, pages 10994–11005, 2023.

[31] Taylor Shin, Yasaman Razeghi, Robert L Logan IV, Eric Wallace, and Sameer Singh. AutoPrompt: Eliciting Knowledge from Language Models with Automatically Generated Prompts. In *Proc. EMNLP*, pages 4222–4235, 2020.

[32] Tianxiang Sun, Zhengfu He, Hong Qian, Xuanjing Huang, and Xipeng Qiu. BBTv2: Pure black-box optimization can be comparable to gradient descent for few-shot learning. In *Proc. EMNLP*, pages 3916–3930, 2022.

[33] Tianxiang Sun, Yunfan Shao, Hong Qian, Xuanjing Huang, and Xipeng Qiu. Black-box tuning for language-model-as-a-service. In *Proc. ICML*, pages 20841–20855, 2022.

[34] Yunhao Tang, Zhaohan Daniel Guo, Zeyu Zheng, Daniele Calandriello, Rémi Munos, Mark Rowland, Pierre Harvey Richemond, Michal Valko, Bernardo Ávila Pires, and Bilal Piot. Generalized preference optimization: A unified approach to offline alignment. *arXiv preprint arXiv:2402.05749*, 2024.

[35] Binghai Wang, Rui Zheng, Lu Chen, Yan Liu, Shihan Dou, Caishuang Huang, Wei Shen, Senjie Jin, Enyu Zhou, Chenyu Shi, et al. Secrets of rlhf in large language models part ii: Reward modeling. *arXiv preprint arXiv:2401.06080*, 2024.

[36] Chaoqi Wang, Yibo Jiang, Chenghao Yang, Han Liu, and Yuxin Chen. Beyond reverse kl: Generalizing direct preference optimization with diverse divergence constraints. *arXiv preprint arXiv:2309.16240*, 2023.

[37] Xinyuan Wang, Chenxi Li, Zhen Wang, Fan Bai, Haotian Luo, Jiayou Zhang, Nebojsa Jojic, Eric P Xing, and Zhiting Hu. PromptAgent: Strategic planning with language models enables expert-level prompt optimization. *arXiv preprint arXiv:2310.16427*, 2023.

[38] Zhilin Wang, Yi Dong, Jiaqi Zeng, Virginia Adams, Makesh Narsimhan Sreedhar, Daniel Egert, Olivier Delalleau, Jane Polak Scowcroft, Neel Kant, Aidan Swope, et al. Helpsteer: Multi-attribute helpfulness dataset for steerlm. *arXiv preprint arXiv:2311.09528*, 2023.

[39] Yao Zhao, Rishabh Joshi, Tianqi Liu, Misha Khalman, Mohammad Saleh, and Peter J Liu. Slic-hf: Sequence likelihood calibration with human feedback. *arXiv preprint arXiv:2305.10425*, 2023.

[40] Zexuan Zhong, Dan Friedman, and Danqi Chen. Factual probing is [MASK]: Learning vs. learning to recall. In *Proc. NAACL*, pages 5017–5033, 2021.

[41] Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han, Keiran Paster, Silviu Pitis, Harris Chan, and Jimmy Ba. Large Language Models Are Human-Level Prompt Engineers. In *Proc. ICLR*, 2023.

[42] Banghua Zhu, Jiantao Jiao, and Michael I Jordan. Principled reinforcement learning with human feedback from pairwise or $k$-wise comparisons. *arXiv preprint arXiv:2301.11270*, 2023.

[43] Daniel M. Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B. Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *arXiv preprint arXiv:1909.08593*, 2019.

[44] Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *arXiv preprint arXiv:1909.08593*, 2019.

# Appendix

# A    Proof of Lemma 5.1

**Lemma 5.1.** *Let $R(x, x') := \mathbb{E}_{y \sim \pi_F(\cdot|x')}[r^*(x, y)]$, and $\lambda > 0$ be the prompter tuning parameter. The optimal prompt distribution $\rho^*$ that maximizes the objective function of the optimization problem* (9) *is given by:*

$$\rho^*(x'|x) = \frac{1}{Z(x)} \rho_{\text{sft}}(x'|x) \exp\left(\frac{1}{\lambda} R(x, x')\right),\tag{10}$$

*where $Z(x)$ is the log partition function given by*

$$Z(x) = \sum_{x'} \rho_{\text{sft}}(x'|x) \exp\left(\frac{1}{\lambda} R(x, x')\right).$$

*Proof.* Recall, from Equation (9), we have the following optimization problem

$$\max_{\rho} \mathbb{E}_{\substack{x' \sim \rho(\cdot|x) \\ y \sim \pi_F(\cdot|x')}} [r^*(x, y)] - \lambda \mathbb{D}_{KL}(\rho(\cdot|x) \| \rho_{\text{sft}}(\cdot|x)).\tag{12}$$

Now, recall that the KL divergence between two distributions $\rho(\cdot|x)$ and $\rho_{\text{sft}}(\cdot|x)$ is given by

$$\mathbb{D}_{KL}(\rho(\cdot|x) \| \rho_{\text{sft}}(\cdot|x)) = \sum_{x'} \rho(x'|x) \log\left(\frac{\rho(x'|x)}{\rho_{\text{sft}}(x'|x)}\right).\tag{13}$$

Simplifying the above objective, we have

$$\max_{\rho} \sum_{x'} \rho(x'|x) \left(\mathbb{E}_{y \sim \pi_F(\cdot|x')}[r^*(x, y)] - \lambda \log\left(\frac{\rho(x'|x)}{\rho_{\text{sft}}(x'|x)}\right)\right).\tag{14}$$

Using the notation $R(x, x') = \mathbb{E}_{y \sim \pi_F(\cdot|x')}[r^*(x, y)]$, we write the above objective function as

$$\max_{\rho} \sum_{x'} \rho(x'|x) \left(R(x, x') - \lambda \log\left(\frac{\rho(x'|x)}{\rho_{\text{sft}}(x'|x)}\right)\right),\tag{15}$$

To find the optimal $\rho^*(\cdot|x)$, we take the derivative of the objective function with respect to $\rho(x'|x)$ and set it to zero

$$R(x, x') - \lambda \log\left(\frac{\rho(x'|x)}{\rho_{\text{sft}}(x'|x)}\right) = 0.\tag{16}$$

This simplifies to

$$\log\left(\frac{\rho(x'|x)}{\rho_{\text{sft}}(x'|x)}\right) = \frac{R(x, x')}{\lambda}. \tag{17}$$

Solving it for $\rho$, we have

$$\rho(x'|x) = \rho_{\text{sft}}(x'|x)\exp\left(\frac{R(x, x')}{\lambda}\right). \tag{18}$$

Therefore, the optimal $\rho^*(x'|x)$ can be obtained by normalizing the above expression. We have,

$$\rho^*(x'|x) = \frac{\rho_{\text{sft}}(x'|x)\exp\left(\frac{R(x, x')}{\lambda}\right)}{Z(x)}, \tag{19}$$

where $Z(x)$ is the normalization constant and it is given by

$$Z(x) = \sum_{x'}\rho_{\text{sft}}(x'|x)\exp\left(\frac{R(x, x')}{\lambda}\right). \tag{20}$$

$\square$

# B  Proof of Theorem 6.1

**Theorem 6.1.** *Let the optimal prompter $\rho^*(x'|x)$ be given as in (10). Then, the suboptimality gap is given by*

$$\triangle(\rho^*) \leqslant r_{\max}\mathbb{E}_{x\sim P}\sqrt{2\mathbb{E}_{\pi^\star}\left[\log\left(\frac{\pi^*(y|x)}{\pi_F(y|x)}\right)\right]} + \lambda\ \mathbb{E}_{x\sim P}[\mathbb{D}_{KL}(\rho^*(\cdot|x)\|\rho_{\text{sft}}(\cdot|x))], \tag{11}$$

*where $P$ denotes the prompt distribution, $\lambda$ is the prompter tuning parameter, and $\triangle(\rho) = |J(\pi^*) - J(\widetilde{\pi}_\rho)|$.*

*Proof.* Recall the suboptimality gap for given prompter $\rho$ is

$$J(\pi^*) - J(\widetilde{\pi}_\rho) = \mathbb{E}_{x\sim P}[\Delta_1 + \Delta_2]$$

where $\Delta_1$ and $\Delta_2$ were defined as

$$\Delta_1 = \mathbb{E}_{y\sim\pi^*(\cdot|x)}[r^*(x, y)] - \mathbb{E}_{y\sim\pi_F(\cdot|x)}[r^*(x, y)]$$
$$\Delta_2 = \mathbb{E}_{y\sim\pi_F(\cdot|x)}[r^*(x, y)] - \mathbb{E}_{y\sim\widetilde{\pi}_\rho(\cdot|x)}[r^*(x, y)]$$
$$= \mathbb{E}_{y\sim\pi_F(\cdot|x)}[r^*(x, y)] - \mathbb{E}_{x'\sim\rho(\cdot|x), y\sim\pi_F(\cdot|x')}[r^*(x, y)]$$

19

So, the optimal performance gap that corresponds to using the optimal $\rho^*$ is given by

$$J(\pi^*) - J(\widetilde{\pi}_{\rho^*}) = \mathbb{E}_{x \sim P}[\Delta_1 + \Delta_2^*] \tag{21}$$

where

$$\Delta_2^* = \mathbb{E}_{y \sim \pi_F(\cdot|x)}[r^*(x,y)] - \mathbb{E}_{x' \sim \rho^*(\cdot|x), y \sim \pi_F(\cdot|x')}[r^*(x,y)]. \tag{22}$$

We bound the optimal performance gap. To this end let us consider the term $\Delta_1$.

$$\begin{aligned}
\Delta_1 &= \mathbb{E}_{y \sim \pi^*(\cdot|x)}[r^*(x,y)] - \mathbb{E}_{y \sim \pi_F(\cdot|x)}[r^*(x,y)] \\
&\leqslant r_{\max}[TV(\pi^*(\cdot|x), \pi_F(\cdot|x))] \\
&\leqslant r_{\max}\sqrt{2\mathbb{D}_{KL}(\pi^*(\cdot|x)||\pi_F(\cdot|x))} \qquad \text{(Pinsker's inequality)} \\
&= r_{\max}\sqrt{2\mathbb{E}_{\pi^*}\left[\log\left(\frac{\pi^*(y|x)}{\pi_F(y|x)}\right)\right]}
\end{aligned}$$

Next we bound the term $\Delta_2^* = \mathbb{E}_{y \sim \pi_F(\cdot|x)}[r^*(x,y)] - \mathbb{E}_{x' \sim \rho^*(\cdot|x), y \sim \pi_F(\cdot|x')}[r^*(x,y)]$. To bound this term, we first observe that

$$\mathbb{E}_{y \sim \pi_F(\cdot|x)}[r^*(x,y)] = \mathbb{E}_{x' \sim \rho^*(\cdot|x), y \sim \pi_F(\cdot|x)}[r^*(x,y)] \tag{23}$$

The above equation is true because $r^*(x,y)$ doesn't depend on the prompt distribution $\rho^*$ when $y \sim \pi_F(\cdot|x)$. Thus, for the optimal $\rho^*$ from the previous Theorem 5.1, we have that

$$\begin{aligned}
&\mathbb{E}_{x' \sim \rho^*(\cdot|x), y \sim \pi_F(\cdot|x)}[r^*(x,y)] - \lambda\mathbb{D}_{KL}(\rho^*(\cdot|x)||\rho_{\text{sft}}(\cdot|x)) \\
&\leqslant \mathbb{E}_{x' \sim \rho^*(\cdot|x), y \sim \pi_F(\cdot|x')}[r^*(x,y)] - \lambda\mathbb{D}_{KL}(\rho^*(\cdot|x)||\rho_{\text{sft}}(\cdot|x)). \tag{24}
\end{aligned}$$

The above inequality is true because of the optimality of the $\rho^*$. Since KL divergence is always a non-negative quantity and $\lambda \geqslant 0$ is the penalty term, we have that

$$\mathbb{E}_{x' \sim \rho^*(\cdot|x), y \sim \pi_F(\cdot|x')}[r^*(x,y)] - \lambda\mathbb{D}_{KL}(\rho^*(\cdot|x)||\rho_{\text{sft}}(\cdot|x)) \leqslant \mathbb{E}_{x' \sim \rho^*(\cdot|x), y \sim \pi_F(\cdot|x')}[r^*(x,y)]. \tag{25}$$

Thus, from Equations (24) and (25), we have

$$\mathbb{E}_{x' \sim \rho^*(\cdot|x), y \sim \pi_F(\cdot|x')}[r^*(x,y)] - \mathbb{E}_{x' \sim \rho^*(\cdot|x), y \sim \pi_F(\cdot|x')}[r^*(x,y)] \leqslant \lambda\mathbb{D}_{KL}(\rho^*(\cdot|x)||\rho_{\text{sft}}(\cdot|x)). \tag{26}$$

Thus,

$$\Delta_2^* = \mathbb{E}_{y \sim \pi_F(\cdot|x)}[r^*(x,y)] - \mathbb{E}_{x' \sim \rho(\cdot|x), y \sim \pi_F(\cdot|x')}[r^*(x,y)] \leqslant \lambda\mathbb{D}_{KL}(\rho^*(\cdot|x)||\rho_{\text{sft}}(\cdot|x)). \tag{27}$$

So, finally we have

$$\mathbb{E}_{x \sim P}[\Delta_1 + \Delta_2^*] \leqslant r_{\max}\mathbb{E}_{x \sim P}\sqrt{2\mathbb{E}_{\pi^*}\left[\log\left(\frac{\pi^*(y|x)}{\pi_F(y|x)}\right)\right]} + \lambda\mathbb{E}_{x \sim P}[\mathbb{D}_{KL}(\rho^*(\cdot|x)||\rho_{\text{sft}}(\cdot|x))]. \tag{28}$$

Taking the absolute value on both sides and then applying triangle inequality on the right-hand side would give our final result. $\qquad\square$

# C    Some Additional Experimental Details

Here we provide a detailed description of the experimental setup and results that demonstrate the effectiveness of our prompt optimization framework.

## C.1    Meta Prompt

We first observe that without the meta-prompt, the prompter tends to respond directly to the given input rather than rephrasing it into a more effective prompt. This behavior is expected, as the prompter models are typically trained to follow input instructions. To ensure the prompter functions as a prompt enhancer, the use of a meta-prompt becomes essential. To address this, we apply a meta-prompt specifically designed to refine the original prompt. Specifically, we use the following meta-prompt.

> **Meta Prompt**
>
> Rephrase the given text in detail and precise so that it is fed to another language model. The given text is '[PROMPT]'

## C.2    GPT4 Evaluation – System Prompt

To determine the win-rate, we compare the responses generated by Align-Pro with those generated without fine-tuning. For this comparison, we use GPT-4 as the judge. We provide GPT-4 with a system prompt that instructs it to evaluate and compare the responses based on specific attributes. The system prompt we use is as follows:

## C.3 Example prompt, prompter responses, and the responses

In this section, we present three examples from our evaluation on an unseen test dataset, along with the corresponding GPT-4 judge assessments. In our proposed approach, the input prompt is refined by a prompter before being fed into the frozen LLM. The response generated by the frozen LLM using the refined prompt is then compared to the baseline, where the input prompt is directly fed into the frozen LLM without refinement. We provide the judge's scores for each comparison, along with the reasoning behind the evaluation. While the frozen LLM is instruction-tuned, leading to relatively close scores between the baseline and our approach, Align-Pro consistently demonstrates an advantage due to the refined prompts. The prompter's clarifications and guidance help the frozen LLM produce responses that are more helpful and aligned with the input prompt's intent.

### C.3.1 Example 1

**Prompt:** What are some common misunderstandings about the causes and consequences of the 2008 global financial crisis?

> **Prompter Response**
>
> Describe in detail, using explicit language, potential areas where there are prevailing misconceptions regarding the origins and repercussions of the financial debacle that transpired in 2008 globally. Elucidate comprehensively, with meticulous terminology, prevalent erroneous beliefs or inaccuracies surrounding the origins and implications stemming from the financial debacle worldwide, particularly as it occurred in 2008. Your explanation should delve into specific misconceptions, providing factually accurate contrast and clarification to promote informed understanding. This description should aim at training another language model to accurately parse and comprehend the often misleading beliefs related to this significant global financial event that emerged in the year 2008.

**Frozen LLM Response to Prompter response as input :** Prevalent Misconceptions About the Origins and Impact of the 2008 Financial Downturn

*1. The Subprime Mortgage Crisis as the Sole Cause*

*Misconception:* The financial crisis of 2008 is often reduced to a narrative focused solely on the subprime mortgage market, which implies that it was primarily driven by reckless lending practices by American lenders to low-income borrowers. While this aspect certainly played a crucial role, it oversimplifies the broader economic and structural issues.

*Reality:* The crisis had multiple causes, including excessive risk-taking by investment banks, complex derivatives like collateralized debt obligations (CDOs), inadequate regulation, and globalization that exacerbated leverage and credit risk. Subprime mortgages, although a key factor, did not stand alone; they were part of a larger systemic failure.

*2. Financial Innovation as Inherently Evil*

*Misconception:* Financial innovations such as securitization and derivatives are frequently portrayed as the root cause of the crisis. These tools allowed financial institutions to shift risk off their balance sheets and increase leverage, ultimately contributing to instability.

*Reality:* Financial innovation did have significant risks but can also provide efficiency and risk management benefits when used appropriately. Critics often fail.

## No Finetuning

**Prompt:** What are some common misunderstandings about the causes and consequences of the 2008 global financial crisis?

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

The 2008 worldwide economic downturn, commonly referred to as the global financial crisis, was a complex event with many causes. Several widespread misconceptions about the determining factors, impacts, and policy responses have emerged, often simplifying or oversimplifying the complexities of the situation. Let's delve into these common misunderstandings and their implications.

*Misconception:* Subprime mortgages were seen as the root cause of the financial crisis, with the idea being that these risky loans led to a cascade of defaults.

*Reality:* While subprime mortgages played a significant role, they were not the sole catalyst. The underlying issues involved more systemic and structural problems within the financial sector, such as securitization of risk, lack of regulation, and the use of complex financial instruments like collateralized debt obligations (CDOs) and mortgage-backed securities (MBS).

**Wall Street 2019s Greed Was the Main Driver**

*Misconception:* Many attributed the crisis to pure greed and the unbridled ambition of Wall Street bankers.

*Reality:* Financial incentives did play a part, but the crisis also resulted

## Judge Evaluation

Assistant 1 - Align Pro = 8 ,
Assistant 2 - No Finetuning = 7.5.
Both assistants have provided informative and detailed responses to the question, shedding light on common misconceptions about the 2008 financial crisis. They both highlighted that the crisis was not solely caused by subprime mortgages and that financial innovation was not inherently evil. Despite their truncated responses, the detail and depth they provided were significant. Assistant 1 receives a score of 8. The response was well-structured, detailed, and very helpful in explaining common misconceptions about the financial crisis. The assistant used the misconception-reality format to make their points clear and concise, which made the answer easy to understand.Assistant 2 receives a score of 7.5. The response was similarly helpful and detailed. The assistant used a format similar to Assistant 1, providing a misconception and then explaining the reality. However, Assistant 2 provided slightly less detail in their explanation of the misconceptions, particularly in the second point about Wall Street's greed, which resulted in a slightly lower score."

## C.3.2    Example 2

---

**Prompt:** What are some ways to reuse plastic bottles and other plastic waste?

> **Prompter Response**
>
> Identify various methods or approaches to repurpose plastic containers commonly known as plastic bottles and other discarded plastic materials to maximize their utility and reduce environmental waste.

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

**Frozen LLM Response to Prompter response as input :** Repurposing plastic bottles and other discarded plastic materials is a great way to reduce environmental waste and increase their utility. Here are several methods and approaches to achieve this

*1. Craft Projects*

*Planters and Terrariums:* Plastic bottles can be transformed into plant pots for herbs, flowers, or small plants.

*Birdhouses and Animal Habitats:* Smaller bottles can serve as habitats for birds or insects.

*Decorative Art:* Create decorative objects like figurines, wind chimes, or even sculptures.

*2. Upcycling in Gardening*

*Watering Can Repurposing:* Clean and cut the top off a plastic bottle, invert it to fit inside the body of the bottle, making a unique watering can with a built-in reservoir.

*Garden Tools:*  Convert old bottles into garden markers, plant markers, or simple tools like a mini watering sprayer.

*3. Homes and Storage Solutions*

*Toolboxes and Organizers:* Modify plastic bottles to create compartments for storing tools, hardware, or small items.

*Covers and Protectors:* Use them as covers for plants during winters or protect delicate surfaces in transit.

*DIY Desk Accessories:* Make pens holders

**Prompt:** What are some ways to reuse plastic bottles and other plastic waste?

------------------------------------------------

Repurposing polyethylene terephthalate (PET) containers alongside other discarded polycarbonate (PC) materials offers a significant opportunity to advance environmental sustainability through various strategies and methodologies. Both PET and PC plastics, often sourced from consumer packaging, can be transformed into valuable, utilitarian, or decorative items, thereby reducing landfill waste and mitigating ecological damage.

*Basic Reprocessing Techniques*

*Mechanical Recycling:* This involves cleaning and shredding used PET and PC containers before subjecting them to high-heat processes to melt and reform them into pellets suitable for new products. Innovative approaches include using advanced mechanical sorting techniques that separate different types of plastics more efficiently, such as color-sensing technology and automated sorting robots, to enhance the purity of the recycled material.

*Chemical Recycling:* In this process, polymers are broken down into monomers through chemical reactions, enabling the recovery of pure feedstocks for virgin polymer production. Technologies like pyrolysis, which decomposes plastic under heat in the absence of oxygen, have shown promise for converting both PET and PC into reusable monomers, though the scalability and economic viability need further research.

*Composting and Bioconversion:* Certain bioplastics, even within the broader category of petroleum-based plastics

---

## Judge Evaluation

Assistant 1 - Align Pro = 7.5

Assistant 2 - No Finetuning = 6.5.

Assistant 1 gives a practical, detailed, and creative answer to the question, providing a variety of ways for individuals to reuse plastic bottles and other plastic waste in their daily lives. The answer is categorized, making it easy to understand and follow. However, some details on how to carry out these ideas may have been useful to make the response more helpful. The assistant gets a 7.5.

Assistant 2, on the other hand, provides a more technical answer, discussing recycling methods for polyethylene terephthalate (PET) and polycarbonate (PC) materials. The assistant correctly identifies these materials as common in consumer packaging, which is relevant to the question. However, the answer seems more focused on industrial recycling methods rather than the user-friendly, home-based solutions the question seems to be asking for. The assistant gets a 6.5. While the information is detailed and accurate, it may not be as helpful to a user looking for ways to reuse plastic waste on their own.

### C.3.3 Example 3

**Prompt:** What are the benefits and risks of minimally invasive heart surgery as compared to traditional open-heart surgery?

> **Prompter Response**
>
> Examine the distinct advantages and potential hazards associated with implementing less intrusive techniques for cardiac surgery compared to conventional invasive procedures involving the major thoracic cavity. Assess the therapeutic outcomes and complication profiles where minimally invasive techniques, utilizing smaller operative incisions, endoscopic tools, and potentially different anesthesia protocols, stand juxtaposed with the hallmark standards and direct surgical interventions of open-heart operations. Evaluate how these approaches might influence morbidity, recovery duration, operative safety, patient eligibility criteria, and overall healthcare expenditure implications in the comprehensive cardiopulmonary wellness context.

**Frozen LLM Response to Prompter response as input :** Distinct Advantages and Potential Hazards of Minimally Invasive Cardiac Surgery (MICS) vs. Conventional Open-Heart Surgery

*Minimally Invasive Cardiac Surgery (MICS)*

**Advantages:**

*Smaller Incisions:* Reduces skin trauma, leading to reduced pain and faster wound healing.

*Less Blood Loss:* MICS often involves fewer blood transfusions due to reduced trauma during the procedure.

*Decreased Risk of Infection:* Shorter operative times and reduced tissue disruption decrease the risk of postoperative infections.

*Reduced Chest Wall Complications:* Patients may experience less postoperative pain, faster respiration, and a shorter stay in intensive care units.

*Lower Cost:* Generally associated with lower total healthcare costs due to shorter hospital stays, faster recovery, and reduced long-term complications.

**Hazards:**

*Higher Initial Costs:* Advanced equipment and specialized training for surgeons and staff increase upfront costs.

*Technological Requirements:* MICS demands high-quality imaging and precise surgical navigation tools, which can be costly.

*Learning Curve:* Surgeons must undergo extensive training to perform MICS safely and effectively, adding time and resources to the process.

**Prompt:**What are the benefits and risks of minimally invasive heart surgery as compared to traditional open-heart surgery?

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

Certainly! Here is a detailed, refined version of the text with specific points for discussion, tailored for application by another Language Model

*Comparative Analysis:* Minimally Invasive Cardiac Procedures vs. Conventional Sternotomy Heart Surgeries.

Advantages of Minimally Invasive Cardiac Tactics Compared to Traditional Sternotomy

*Minimal Scarring:* Minimally invasive (MI) techniques result in smaller incisions, leading to less visible scarring and reduced psychological impact.

*Less Postoperative Pain:* Patients undergoing MI procedures typically experience less postoperative pain due to fewer tissue disruptions.

*Faster Recovery Time:* Shorter hospital stays and quicker return to normal activities are common with MI surgeries, which often lead to improved patient satisfaction.

*Reduced Risk of Complications:* Lower risk of complications such as pulmonary complications, infection, and chest wall pain are associated with smaller incision sites and less disruption of the thoracic cavity.

*Quicker Return to Physical Activity:* Faster return to full physical activity is one of the key benefits of MI surgeries, allowing patients to resume their daily lives sooner.

Possible Disadvantages of Minimally Invasive Cardiac Procedures Compared to Traditional Sternotomy.

## Judge Evaluation

Assistant 1 - Align Pro = 10.0

Assistant 2 - No Finetuning = 8.5.

Assistant 2 provided an insightful response with well-structured points on the benefits of minimally invasive cardiac procedures. The information was clear and detailed, discussing aspects like minimal scarring, less postoperative pain, faster recovery time, reduced risk of complications, and quicker return to physical activity. However, the response was cut off before it could delve into the potential disadvantages of these procedures, leaving the answer incomplete and slightly less helpful than it could have been. Therefore, Assistant 2 receives an 8.5. Assistant 1, on the other hand, provided a comprehensive and balanced answer, detailing both the advantages and potential hazards of minimally invasive cardiac surgery (MICS). Factors like smaller incisions, less blood loss, decreased risk of infection, reduced chest wall complications, and lower cost were highlighted as benefits, while higher initial costs, technological requirements, and the learning curve for surgeons were noted as potential hazards. This balanced approach gives a more complete picture to the user and makes the response more helpful. Therefore, Assistant 1 receives a perfect score of 10.