# Modulated Intervention Preference Optimization (MIPO): Keep the Easy, Refine the Difficult

## **Cheolhun Jang**

NC Research 9405hun@ncsoft.com

#### **Abstract**

Preference optimization methods typically begin training with a well-trained SFT model as a reference model. In RLHF and DPO, a regularization term is used during the preference optimization process to prevent the policy model from deviating too far from the reference model's distribution, thereby avoiding the generation of anomalous responses. When the reference model is already well-aligned with the given data or only requires slight adjustments, this approach can produce a well-aligned model. However, if the reference model is not aligned with the given data and requires significant deviation from its current state, a regularization term may actually hinder the model alignment. In this study, we propose **Modulated Intervention Preference Optimization (MIPO)** to address this issue. MIPO modulates the degree of intervention from the reference model based on how well the given data is aligned with it. If the data is well-aligned, the intervention is increased to prevent the policy model from diverging significantly from reference model. Conversely, if the alignment is poor, the interference is reduced to facilitate more extensive training. We compare the performance of MIPO and DPO using Mistral-7B and Llama3-8B in Alpaca Eval 2.0 and MT-Bench. The experimental results demonstrate that MIPO consistently outperforms DPO across various evaluation scenarios.

#### 1 Introduction

As the performance of Large Language Models (LLMs) trained with a large amount of data has been attracting attention, methods (Chowdhery et al. 2023; Touvron et al. 2023; Brown et al. 2020) for training them have been actively studied. The commonly used LLM training pipeline is to pretrain LLM using a large amount of data, and then use the instruction-tuning method (Wei et al. 2021) to allow LLM to follow the human-provided instruction.

However, it is difficult to train LLM to produce the desired output (helpful, harmless) or to prevent LLM from producing the output that LLM should not produce (Bai et al. 2022a). Therefore, LLM alignment methods employing human feedback have started to gain significant attention.

Among these methods, Reinforcement Learning from Human Feedback (RLHF) (Christiano et al. 2017; Askell et al. 2021) received significant attention. Models trained with RLHF are well-aligned with human feedback, demonstrating reliable performance as a result (Korbak et al. 2023;

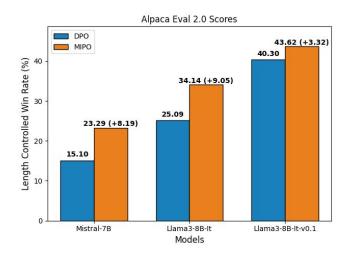


Figure 1: Alpacaeval 2.0 scores for MIPO and DPO implementations on Mistral-7B-Base and Llama-8B-Instruct. v0.1 is a model trained on different dataset.

Havrilla et al. 2024; Dai et al. 2023). However, the RLHF approach involves a complex training process, including the training of a reward model, which has posed significant challenges in the implementation and training (Casper et al. 2023; Peng et al. 2023).

Direct Preference Optimization (DPO) (Rafailov et al. 2024) is a method designed to overcome these limitations. In DPO, the optimization problem of RLHF is modified to eliminate the reward model and train only the policy model. This makes it easier to train DPO compared to RLHF, and DPO also effectively learned human preferences, demonstrating strong performance.

In DPO and RLHF, the policy model is trained to align with the instance while ensuring its distribution does not move significantly away from the reference model's distribution to prevent it from generating anomalous responses (ex. hallucinations). Therefore, if the reference model is moderately aligned with the given preference pair, it could be possible to train a well-aligned policy model for the given data without significantly diverging from the reference model's distribution. However, if the reference model is not

aligned with the given preference pair, it will be difficult for the policy model to align with the data through minor adjustments, without moving far from the reference model's distribution. Therefore, it is crucial to adjust the training objective based on how well the reference model is aligned.

In this paper, we propose a preference optimization algorithm, **Modulated Intervention Preference Optimization** (MIPO), to address this issue. As seen in Figure 2, MIPO utilizes the average log likelihood to measure how well the reference model is aligned with the given preference pair. Through this value, the MIPO objective is configured to modulate the intervention of the reference model, allowing more extensive training on pairs that are judged to be poorly aligned with the reference model. We use Alpaca Eval 2.0 and MT-Bench to compare the performance of MIPO with DPO and other preference optimization methods. Across diverse experimental settings, MIPO consistently achieves outstanding performance. To summarize, **MIPO has the following properties:** 

- Modulate the intervention of the Reference Model: MIPO is a novel approach that modulates the intervention of the reference model for each instance. It determines the extent of the reference model's intervention based on the degree of alignment. MIPO maintains performance on pairs where the reference model already well-aligned, while simultaneously achieving substantial performance gains on pairs where the reference model previously underperformed (Section \$4).
- Outstanding Benchmark Performance: We conduct experiments using Llama3-8B-Instruct (AI@Meta 2024) and Mistral-7B-Base (Jiang et al. 2023) to verify the effectiveness of MIPO in various models. On Alpaca Eval 2.0, our proposed method consistently outperforms DPO. As we can see in Figure 1, in Llama3-8B-Instruct, it outperforms DPO by approximately 9 points (+36.07%), and in Mistral-7B-Base, it outperforms about 8 points (+54.24%). In most cases, MIPO achieves the best performance not only compared to DPO but also when compared to other methods. On MT-Bench, MIPO also exhibits the best performance among the compared approaches (Section \$6.1).
- Simple and Effective Training: The high-performance model can be found in MIPO by tuning only the hyper-parameter  $\beta$ . Moreover, consistently outstanding performance is achieved within a specific range of  $\beta$ , independent of model architecture or dataset. Thus, unlike other methods that require extensive tuning, this approach allows for easy acquisition of a high-performance model with minimal tuning effort (Section \$6.2).

## 2 Related Works

After being pretrained on a large amount of data (Chowdhery et al. 2023) and fine-tuned (Chung et al. 2024; Ramamurthy et al. 2022), LLMs have achieved notable performance across many tasks (Touvron et al. 2023; Brown et al. 2020; Thoppilan et al. 2022). However, LLMs that could generate responses that were even more helpful and

harmless were needed, leading to the development of preference optimization methods (Christiano et al. 2017; Bai et al. 2022a,b) that fine-tune LLMs more closely to human feedback.

RLHF (Askell et al. 2021; Ouyang et al. 2022) is one such preference optimization method for LLM alignment. In RLHF, preference data is used to train a reward model, which is then utilized to optimize the policy model by Proximal Policy Optimization (Schulman et al. 2017). RLHF effectively aligns models with human feedback, resulting in good performance (Korbak et al. 2023; Havrilla et al. 2024). However, there are challenges, such as the difficulty of obtaining scored data, ensuring stable training, and the necessity of training a reward model (Casper et al. 2023; Peng et al. 2023; Wang et al. 2024).

DPO is a preference optimization method that solves optimization problem of RLHF in a more easier and efficient manner. (Rafailov et al. 2024) proposed DPO to eliminate the reward model in RLHF and train only the policy model with preference data. It is simple compared to RLHF, and the training phase is more stable. So it has become one of the widely used method for aligning language models. However, DPO also has its drawbacks like dependency on the reference model and issues with length exploitation (Liu, Liu, and Cohan 2024; Gorbatovski et al. 2024; Xu et al. 2024). Therefore, new model alignment methods such as KTO (Ethayarajh et al. 2024), IPO (Azar et al. 2024) and ORPO (Hong, Lee, and Thorne 2024) continue to emerge.

However, most methods including DPO does not take into account the differences in the degree of alignment of the reference model between preference pairs. As mentioned earlier, if the reference model is already well-aligned, only minimal training will be needed to achieve alignment. Conversely, if the reference model is completely misaligned, extensive training will be required. However, DPO does not account for these differences (Section \$3.3).

To address this issue, we propose **MIPO**, which varies the learning weights among instances by modulates the degree of intervention from the reference model (Section \$4).

# 3 Background

In this section, we will review the DPO in Section \$3.2, and analyze the ineffective aspects of DPO in Section \$3.3.

## 3.1 Terminology

 $D = \left\{x^i, y_w^i, y_l^i\right\}_{i=1}^N$  is for pairwise-preference dataset, where  $x^i$  is prompt and  $y_w^i$  is chosen (preferred) response and  $y_l^i$  is rejected (dis-preferred) response for that prompt.  $\pi_{ref}$  is reference model, initial LLM that we start training from.  $\pi_{\theta}$  is policy model, which is a model we train.

# 3.2 **DPO**

DPO employs the Bradley-Terry (BT) model (Bradley and Terry 1952) to represent the distribution of human preference. BT model represents human preference distribution for  $y_w$ ,  $y_l$  by the reward function as follows:

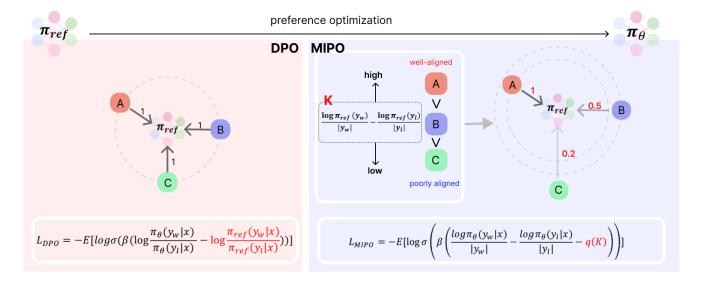


Figure 2: **Optimization process of MIPO**. In DPO, the objective utilizes a consistent regularization term (red part in DPO objective) for the reference model across all instances (A, B, C in Figure), regardless of the degree of alignment of each instance. However, in MIPO, the alignment of each instance with the reference model is first assessed by using the difference in average log likelihood. Based on this value, K, the extent to which the reference model will intervene in the learning process is determined and subsequently reflected in the MIPO objective.

$$p(y_w > y_l | x) = \frac{exp(r(x, y_w))}{exp(r(x, y_w)) + exp(r(x, y_l))}$$
(1)

DPO's reward function is reparameterized from the RLHF's objective as following equation.

$$r(x,y) = \beta \log \frac{\pi_{\theta}(y|x)}{\pi_{ref}(y|x)} + \beta \log Z(x)$$
 (2)

From equations (1) and (2), we can formulate preference distribution by using  $\pi_{ref}$  and  $\pi_{\theta}$ . Subsequently, the DPO objective is derived as expressed in (3)

$$L_{DPO}(\pi_{\theta}; \pi_{ref})$$

$$= E_{(x, y_w, y_l)} \sim_D \left[ -log\sigma(\beta log \frac{\pi_{\theta}(y_w|x)}{\pi_{ref}(y_w|x)}) - \beta log \frac{\pi_{\theta}(y_l|x)}{\pi_{ref}(y_l|x)} \right]$$
(3)

# 3.3 Ineffective Aspects of DPO

**DPO** does not consider how well the preference pairs are aligned. Looking at the reward of DPO in Eq (2) without Z(x). It can be seen that the reward is the difference between the log likelihood of the policy model and the log likelihood of the reference model. This implies that DPO allows for high rewards to be obtained solely by increasing the log likelihood of a response, without considering the degree to which the reference model already performs well on that response. Consequently, the training process proceeds

without taking into account the extent to which the reference model is aligned with the give preference data.

For example, consider  $pair_1$ , preference data where the reference model already well-aligned, and  $pair_2$ , where it does not. Ideally, model will require to train slightly on  $pair_1$  to maintain its current performance, while it will require substantial training for  $pair_2$  compared to  $pair_1$ .

Let's assume that the policy model has been trained so that the log likelihood of the chosen response increases by  $\alpha$  compared to the reference model, while the log likelihood of the rejected response remains unchanged in both pairs  $(log\pi_{\theta}(y_w|x) - log\pi_{ref}(y_w|x) = \alpha, log\pi_{\theta}(y_l|x) - log\pi_{ref}(y_l|x) = 0)$ . In DPO, both pairs would yield the same loss by Eq (3). This implies that the improvement in log likelihood for  $pair_1$  and  $pair_2$  holds equal significance in DPO.

Consequently, DPO trains the model without discriminating between instances of strong and weak alignment with the reference model. This uniform approach can result in insufficient training for pairs where the reference model needs improvement and excessive training for pairs where preferences are already adequately captured. Therefore, this issue can negatively impact the performance of the trained model.

#### 4 Methodology

In this section, we explain why we use average log likelihood to determine how well reference model is aligned to data in Section \$4.1. Then we introduce **Modulated Intervention Preference Optimization** (MIPO), an algorithm that adjusts the degree of intervention from the reference model based on the level of alignment in Section \$4.2.

## 4.1 Measuring the Alignment Degree

To solve the problem of DPO mentioned above Section \$3.3, we first need to measure which pairs are well-aligned to reference model and which pairs are poorly aligned.

In the context of preference learning, being "well-aligned" can be interpreted as the model being more likely to generate a chosen response  $y_w$  than a rejected response  $y_l$  for a given input x. However, using the difference in log likelihoods between chosen and rejected responses to measure alignment is not feasible, as log likelihood is highly sensitivity to response length. If the lengths of the chosen and rejected responses differ significantly, the longer response's log likelihood will be disproportionately lower, regardless of individual token probabilities.

Therefore, we decide to use of average log likelihood. It allows for a more fairer comparison of generation probabilities between chosen and rejected responses, mitigating the impact of length discrepancies. We have decided to use the difference in average log likelihood, K, as a metric to assess the alignment of the reference model with a given pair.

$$K = \frac{log\pi_{ref}(y_w|x)}{|y_w|} - \frac{log\pi_{ref}(y_l|x)}{|y_l|}$$
(4)

We interpret a high K value as indicative of strong alignment in the reference model, whereas a low K value suggest insufficient alignment. Based on this assumption, we propose our objective as follows:

# 4.2 Deriving the MIPO Objective

$$L_{MIPO}(\pi_{\theta}; \pi_{ref}) = E_{(x, y_w, y_l) \sim D} - \log \sigma(\beta \underbrace{\left(\frac{\log \pi_{\theta}(y_w|x)}{|y_w|} - \frac{\log \pi_{\theta}(y_l|x)}{|y_l|}\right)}_{f(\theta)} - \beta \underbrace{\log(1 + e^K)}_{q(K)})$$
(5)

For the reasons mentioned above, the MIPO objective is designed to enhance the alignment of the policy model by using average log likelihood,  $f(\theta)$ . Additionally, it is adjusted based on the degree of alignment through q(K), which acts as a modulator for the degree of intervention from the reference model.

Let's examine the MIPO objective in two cases:

When reference model is well aligned for a given pair It means K is large enough. Then, q(K) converges to K and the objective of MIPO can be expressed as follows.

$$L_{MIPO} = -log\sigma(\beta(\frac{log\pi_{\theta}(y_w|x)}{|y_w|} - \frac{log\pi_{\theta}(y_l|x)}{|y_l|}) - \beta(\frac{log\pi_{ref}(y_w|x)}{|y_w|} - \frac{log\pi_{ref}(y_l|x)}{|y_l|})).$$
(6)

The objective is calculated based on the difference between the policy model's average log likelihood difference,  $f(\theta)$ , and this values of reference model, K. Therefore, as

 $f(\theta)$  diverges further from K, the loss decreases, preventing the policy model from significantly diverging from the reference model.

When reference model is poorly aligned for a given pair It means K is low. In this case, q(K) approaches to 0 and objective can be expressed as follows.

$$L_{MIPO} = -log\sigma(\beta(\frac{log\pi_{\theta}(y_w|x)}{|y_w|} - \frac{log\pi_{\theta}(y_l|x)}{|y_l|})) \quad (7)$$

Since the MIPO objective does not include a term for the reference model, it only considers the  $f(\theta)$  for alignment, focusing solely on increasing this value. When compared to the case where q(K) = K, it is clear that the MIPO loss significantly greater because  $f(\theta)$  is less than  $f(\theta) - K(\because K < 0)$ . Consequently, the policy model can be trained while diverging further from the distribution of the reference model.

In summary, the MIPO assesses how well the reference model is aligned with the given instance through the metric K. This metric is then used to calculate q(K), which determines the extent to which the reference model's influence on the policy model's learning. When K is high, it indicates strong alignment with the given data. In this case, q(K) takes on the value of K, thereby increasing the intervention of the reference model. Consequently, the policy model train without diverging significantly from the reference model. Conversely, if K is low, q(K) becomes zero, allowing the policy model to train without intervention from the reference model.

More detailed explanations about objective are provided in the Section \$6.4 and gradient analysis can be found in Appendix A.

## 5 Experimental Settings

## 5.1 Datasets

**Binarized UltraFeedback** We train models with Binarized UltraFeedback Dataset (Cui et al. 2023). It consists of 64K preference pairs from diverse resources.

**Llama3 UltraFeedback**<sup>1</sup> Because there is a possibility that Binarized Ultrafeedback data was used in the training phase of Llama3-8B-instruct, (Meng, Xia, and Chen 2024) proposed new dataset. The data is created base on responses generated by Llama3-8B-Instruct by using the Binarized Ultrafeedback prompts. Among these responses, the highest scoring response and the lowest scoring response, which are scored by reward model (Jiang, Ren, and Lin 2023), are used to form preference pairs. In this study, models trained using this dataset is labeled with the **v0.1** tag.

# 5.2 Evaluation

The trained models are evaluated on AlpacaEval2.0 and MT-Bench.

https://huggingface.co/datasets/princeton-nlp/ llama3-ultrafeedback.

Method	Mistral-7B-Base			Llama3-8B-Instruct			Llama3-8B-Instruct-v0.1		
	Alpaca Eval 2.0		MT-Bench	Alpaca Eval 2.0		MT-Bench	Alpaca Eval 2.0		MT-Bench
	LC(%)	WR(%)	Avg. Score	LC(%)	WR(%)	Avg. Score	LC(%)	WR(%)	Avg. Score
ORPO	14.7*	12.2*	-	-	-	_	28.5*	27.4*	-
KTO	13.1*	9.1*	-	-	-	-	33.1*	31.8*	_
SimPO	21.4*	20.8*	7.05	-	-	-	44.7*	40.5*	7.72
DPO	15.1*	12.5*	7.01	25.09	21.18	7.95	40.3*	37.9*	7.79
MIPO	22.02	17.50	7.12	34.14	30.00	7.97	43.62	40.74	7.92

Table 1: AlpacalEval 2.0 and MT-Bench scores for preference optimization methods in Mistral-7B, Llama3-8B. The v0.1 tag refers to a model trained using **Llama3 Ultrafeedback** data, and the others are all trained with **Binarized UltraFeedback**. Results denoted with (\*) are sourced from (Meng, Xia, and Chen 2024).

Alpaca Eval 2.0 Alpaca Eval 2.0 (Li et al. 2023; Dubois et al. 2024) consists of 805 prompts. The responses generated using these prompts are compared against those produced by GPT-4-Turbo. Through this comparison, Alpaca Eval 2.0 quantify the model's performance by calculating the percentage of instances where its response surpass those of GPT-4-Turbo, expressed as a win rate (WR). AlpacaEval 2.0 also provides length controlled win rate (LC) that considers bias due to length.

MT-Bench MT-Bench (Zheng et al. 2023) is a multi-turn benchmark consisting of 80 distinct instructions to evaluate model performance. Model generated responses from these prompts are scored by using GPT-4. The benchmark's strength lies in its diverse category coverage, enabling comprehensive model assessment across multiple dimensions.

## 5.3 Models and Baselines

To compare across different model families, we use Mistral-7B-Base (Jiang et al. 2023) and Llama3-8B-Instruct (AI@Meta 2024) as base model for preference optimization. We compare MIPO with DPO and also with SimPO, which utilizes average log likelihood. Additionally, results are compared with offline preference optimization methods, such as ORPO and KTO.

We implement MIPO, DPO and SimPO by using TRL (von Werra et al. 2020) and the alignment book (Tunstall et al.). When the Alpaca Eval 2.0 scores for models trained with DPO and SimPO are lower than those reported in the reference<sup>2</sup>, we adapts the reference values for a fair comparison. For MT-Bench evaluations, we utilize the checkpoints in reference to generate responses and evaluate. Additionally, we reference results from it for ORPO and KTO.

# 6 Result and Analysis

#### 6.1 Benchmark Results

As shown in Table 1, MIPO consistently achieves higher scores compared to DPO and demonstrates outstanding performance relative to other methods in the most cases.

Comparative analysis using Alpaca Eval 2.0 reveals that MIPO consistently and significantly outperforms DPO

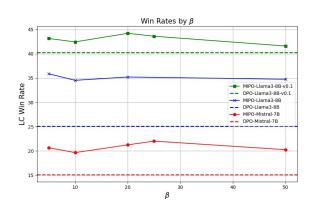


Figure 3: Alpaca Eval 2.0 scores in Mistral and Llama3 based on  $\beta$ . The dotted line represents the performance of DPO.

across all experiments. Moreover, MIPO achieves performance levels comparable to SimPO, which had previously demonstrated the highest performance.

In MT-Bench, MIPO consistently exhibits enhanced performance relative to DPO and SimPO across all experiments.

## **6.2** Performance Based on $\beta$

One of the advantages of MIPO is the ease of hyperparameter tuning. MIPO objective contains only a single hyperparameter,  $\beta$ , allowing for optimal model training by adjusting just this one. Figure 3 illustrates how the model's performance varies with different  $\beta$  in Mistral-7B and Llama-8B. As depicted in Figure 3, MIPO maintains exceptionally high performance across a similar beta range ([5,50]), demonstrating robustness across various models and datasets. The optimal model configuration is consistently identified within this range.

In conclusion, MIPO demonstrates a significant advantage: it consistently produces models that substantially outperform DPO and approach optimal performance levels, achieved through the tuning of a single hyperparameter,  $\beta$ , within a moderate range. This capability persists

<sup>&</sup>lt;sup>2</sup> https://github.com/princeton-nlp/SimPO

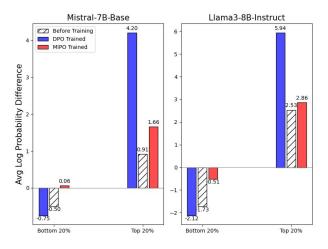


Figure 4: The difference in average log likelihood changes after training for both MIPO and DPO, as applied to Mistral-7B-Base and Llama3-8B-Instruct.

across diverse model architectures and datasets, underscoring MIPO's robustness and effectiveness.

## 6.3 Analysis about Average Log Likelihood

Figure 4, represents the average log likelihood difference between chosen and rejected responses for the model on the evaluation dataset, showing how this difference changes after training. It specifically highlights how the values for instances in the top 20% and bottom 20% of average log likelihood differences in reference model have evolved.

At this point, the top 20% are instances with a large average log likelihood difference in reference model, indicating they are already well-aligned data, while the bottom 20% are poorly aligned and require more training.

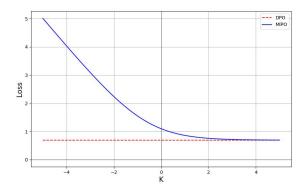
In the bottom 20%, the average log likelihood difference for DPO actually decrease, whereas for MIPO, the average log likelihood clearly increase. Conversely, in the top 20%, the average log likelihood for DPO increase significantly, while for MIPO, it only increase slightly. This pattern is observed in both the Llama3-8B and Mistral-7B.

This indicates that in DPO, the data that is already well-aligned continued to be better aligned through further training, while the data that is not well-aligned do not see significant improvement. However, in MIPO, the training is operated to maintain performance on well-aligned data while significantly improving the alignment of poorly aligned data. MIPO achieves the intended outcome described in Section \$4.2, thereby effectively enhancing model alignment.

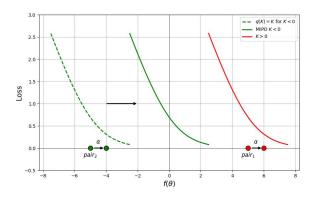
#### 6.4 Analysis about MIPO objective function

As seen in Eq (5), the MIPO objective can be expressed as the difference between the average log likelihood of the chosen response and rejected response in policy model and minus q(K) consists of values calculated from the reference model.

Let's examine how the MIPO objective behaves during the training process in two scenarios.



(a) MIPO loss in early stages of training



(b) MIPO loss in high K and low K

Early Stage in Training In the early stage of training, there is minimal difference between the reference model and the policy model. Therefore, the average log likelihood difference of the policy model does not significantly diverge from that of the reference model ( $\pi_{ref} \approx \pi_{\theta}$ ). Consequently, the MIPO loss can be written as  $\log\left(1+e^{-\left(K-\ln\left(1+e^K\right)\right)}\right)$ . However, DPO loss for all instance initially 0. This can be observed in Figure 5a.

**Loss Reflection During Training** Next, let's examine how the loss for two pairs,  $pair_1$  which has high K value and  $pair_2$  which has low K value, behave during training. Suppose that the average log likelihood difference of the policy model,  $f(\theta)$ , increases by  $\alpha>0$  compared to the reference model for both pairs  $(f(\theta)=K+\alpha)$ .

In Figure 5b, the red section represents  $pair_1$ . Since  $pair_1$  has a high K, the MIPO objective is expressed as  $-logsigmoid(f(\theta)-K)$  (the red line in the figure). Therefore, the MIPO loss is  $-logsigmoid(\alpha)$ , as we can be seen in the graph. Next,  $pair_2$  is represented by the green section. Since K is low, the MIPO objective is expressed as  $-logsigmoid(f(\theta))$  (the green line in the figure). Therefore, the MIPO loss is  $-logsigmoid(K+\alpha)$ , which is larger than the loss for  $pair_1$ . Thus, even with the same amount of increase,  $pair_2$  has a larger loss, indicating that training is

accelerated for pairs with lower K.

Additionally, the figure's dotted line facilitates a comparative analysis between the scenarios where the q(K) is simply K and 0. In dotted line, even if K is low, the loss is calculated based on the K. Thus, when the same increase occurs, the loss is calculated equally for both high and low K pairs, causing the model to train with the same weight for both pairs.

As a result, the MIPO objective results in a relatively large loss when K is low, indicating data that is poorly aligned. Thus, more extensive training can occur on poorly aligned data. Conversely, in the case of well-aligned data, the intervention from the reference model is substantial, causing the objective to be calculated based on the values of the reference model. This prevents the policy model from diverging significantly from the reference model.

#### 7 Conclusion

In DPO, rewards are calculated based on the reference model for all pair data without considering how well the reference model is aligned with the given pair data. Therefore, DPO does not distinguish between instances that require more training ant those that only need minimal training. In this paper, we proposed **Modulated Intervention Preference Optimization (MIPO)** as a method to address and improve upon this issue. MIPO adjusts the objective based on the degree of alignment of the reference model on the given instances. For pairs that require more learning, MIPO reduces the intervention of the reference model, allowing the policy model to diverge from it and find better weights. Conversely, for pairs that are better aligned, the intervention of the reference model is maintained, ensuring that the policy model does not significantly diverge from the reference model.

Through experiments, we found that models trained using MIPO demonstrated significantly improved performance compared to those trained using DPO. Moreover, we observed a notable increase in the average log likelihood difference for instances with initially small differences from the reference model, aligning with our expectations compared to DPO.

**Limitations & Future Work** While MIPO achieved significantly better performance than DPO, there are still several areas for improvement and further investigation

Average log likelihood is not an absolute measure of the degree of alignment The degree of preference between the chosen and rejected responses can vary for each preference pair. In some cases, the chosen and rejected responses might be decided by a very subtle difference, while in others, the difference could be significant. If a given preference pair has only a slight difference, the model may be well-aligned, but the average log probability difference (K) is unlikely to be large. Therefore, it is difficult to accurately assert that a large K indicates superiority on a particular preference pair. The K alone does not provide an absolute measure of performance across different preference pairs.

Although MIPO does not account for the difficulty differences between preference pairs, it is likely that pairs where the model was poorly aligned improved more, as higher av-

erage log likelihoods typically indicate better performance for each pair.

**Experiments on Various Adaptation Terms** While we used  $ln(1+e^K)$  to construct the loss for MIPO, there are various functions that could be used as adaptation term. To achieve the same effect as MIPO, the function q should use an alternative function that converges to K when K is large and to 0 when K is small. However, experiments exploring all possible functions for this have not yet been conducted.

#### References

AI@Meta. 2024. Llama 3 Model Card.

Askell, A.; Bai, Y.; Chen, A.; Drain, D.; Ganguli, D.; Henighan, T.; Jones, A.; Joseph, N.; Mann, B.; DasSarma, N.; et al. 2021. A general language assistant as a laboratory for alignment. *arXiv preprint arXiv:2112.00861*.

Azar, M. G.; Guo, Z. D.; Piot, B.; Munos, R.; Rowland, M.; Valko, M.; and Calandriello, D. 2024. A general theoretical paradigm to understand learning from human preferences. In *International Conference on Artificial Intelligence and Statistics*, 4447–4455. PMLR.

Bai, Y.; Jones, A.; Ndousse, K.; Askell, A.; Chen, A.; Das-Sarma, N.; Drain, D.; Fort, S.; Ganguli, D.; Henighan, T.; et al. 2022a. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv* preprint arXiv:2204.05862.

Bai, Y.; Kadavath, S.; Kundu, S.; Askell, A.; Kernion, J.; Jones, A.; Chen, A.; Goldie, A.; Mirhoseini, A.; McKinnon, C.; et al. 2022b. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*.

Bradley, R. A.; and Terry, M. E. 1952. Rank analysis of incomplete block designs: I. The method of paired comparisons. *Biometrika*, 39(3/4): 324–345.

Brown, T.; Mann, B.; Ryder, N.; Subbiah, M.; Kaplan, J. D.; Dhariwal, P.; Neelakantan, A.; Shyam, P.; Sastry, G.; Askell, A.; et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33: 1877–1901.

Casper, S.; Davies, X.; Shi, C.; Gilbert, T. K.; Scheurer, J.; Rando, J.; Freedman, R.; Korbak, T.; Lindner, D.; Freire, P.; et al. 2023. Open problems and fundamental limitations of reinforcement learning from human feedback. *arXiv* preprint arXiv:2307.15217.

Chowdhery, A.; Narang, S.; Devlin, J.; Bosma, M.; Mishra, G.; Roberts, A.; Barham, P.; Chung, H. W.; Sutton, C.; Gehrmann, S.; et al. 2023. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240): 1–113.

Christiano, P. F.; Leike, J.; Brown, T.; Martic, M.; Legg, S.; and Amodei, D. 2017. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30.

Chung, H. W.; Hou, L.; Longpre, S.; Zoph, B.; Tay, Y.; Fedus, W.; Li, Y.; Wang, X.; Dehghani, M.; Brahma, S.; et al. 2024. Scaling instruction-finetuned language models. *Journal of Machine Learning Research*, 25(70): 1–53.

- Cui, G.; Yuan, L.; Ding, N.; Yao, G.; Zhu, W.; Ni, Y.; Xie, G.; Liu, Z.; and Sun, M. 2023. UltraFeedback: Boosting Language Models with High-quality Feedback. arXiv:2310.01377.
- Dai, J.; Pan, X.; Sun, R.; Ji, J.; Xu, X.; Liu, M.; Wang, Y.; and Yang, Y. 2023. Safe rlhf: Safe reinforcement learning from human feedback. *arXiv preprint arXiv:2310.12773*.
- Dubois, Y.; Galambosi, B.; Liang, P.; and Hashimoto, T. B. 2024. Length-Controlled AlpacaEval: A Simple Way to Debias Automatic Evaluators. *arXiv preprint arXiv:2404.04475*.
- Ethayarajh, K.; Xu, W.; Muennighoff, N.; Jurafsky, D.; and Kiela, D. 2024. Kto: Model alignment as prospect theoretic optimization. *arXiv preprint arXiv:2402.01306*.
- Gorbatovski, A.; Shaposhnikov, B.; Malakhov, A.; Surnachev, N.; Aksenov, Y.; Maksimov, I.; Balagansky, N.; and Gavrilov, D. 2024. Learn your reference model for real good alignment. *arXiv preprint arXiv:2404.09656*.
- Havrilla, A.; Du, Y.; Raparthy, S. C.; Nalmpantis, C.; Dwivedi-Yu, J.; Zhuravinskyi, M.; Hambro, E.; Sukhbaatar, S.; and Raileanu, R. 2024. Teaching large language models to reason with reinforcement learning. *arXiv preprint arXiv:2403.04642*.
- Hong, J.; Lee, N.; and Thorne, J. 2024. Orpo: Monolithic preference optimization without reference model. *arXiv* preprint arXiv:2403.07691, 2(4): 5.
- Jiang, A. Q.; Sablayrolles, A.; Mensch, A.; Bamford, C.; Chaplot, D. S.; Casas, D. d. l.; Bressand, F.; Lengyel, G.; Lample, G.; Saulnier, L.; et al. 2023. Mistral 7B. *arXiv* preprint arXiv:2310.06825.
- Jiang, D.; Ren, X.; and Lin, B. Y. 2023. LLM-Blender: Ensembling Large Language Models with Pairwise Comparison and Generative Fusion. In *Proceedings of the 61th Annual Meeting of the Association for Computational Linguistics (ACL 2023)*.
- Korbak, T.; Shi, K.; Chen, A.; Bhalerao, R. V.; Buckley, C.; Phang, J.; Bowman, S. R.; and Perez, E. 2023. Pretraining language models with human preferences. In *International Conference on Machine Learning*, 17506–17533. PMLR.
- Li, X.; Zhang, T.; Dubois, Y.; Taori, R.; Gulrajani, I.; Guestrin, C.; Liang, P.; and Hashimoto, T. B. 2023. AlpacaEval: An Automatic Evaluator of Instruction-following Models. https://github.com/tatsu-lab/alpaca\_eval.
- Liu, Y.; Liu, P.; and Cohan, A. 2024. Understanding Reference Policies in Direct Preference Optimization. *arXiv* preprint arXiv:2407.13709.
- Meng, Y.; Xia, M.; and Chen, D. 2024. Simpo: Simple preference optimization with a reference-free reward. *arXiv* preprint arXiv:2405.14734.
- Ouyang, L.; Wu, J.; Jiang, X.; Almeida, D.; Wainwright, C.; Mishkin, P.; Zhang, C.; Agarwal, S.; Slama, K.; Ray, A.; et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35: 27730–27744.
- Peng, B.; Song, L.; Tian, Y.; Jin, L.; Mi, H.; and Yu, D. 2023. Stabilizing RLHF through advantage model and selective rehearsal. *arXiv preprint arXiv:2309.10202*.

- Rafailov, R.; Sharma, A.; Mitchell, E.; Manning, C. D.; Ermon, S.; and Finn, C. 2024. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36.
- Ramamurthy, R.; Ammanabrolu, P.; Brantley, K.; Hessel, J.; Sifa, R.; Bauckhage, C.; Hajishirzi, H.; and Choi, Y. 2022. Is reinforcement learning (not) for natural language processing: Benchmarks, baselines, and building blocks for natural language policy optimization. *arXiv preprint arXiv:2210.01241*.
- Schulman, J.; Wolski, F.; Dhariwal, P.; Radford, A.; and Klimov, O. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.
- Thoppilan, R.; De Freitas, D.; Hall, J.; Shazeer, N.; Kulshreshtha, A.; Cheng, H.-T.; Jin, A.; Bos, T.; Baker, L.; Du, Y.; et al. 2022. Lamda: Language models for dialog applications. *arXiv preprint arXiv:2201.08239*.
- Touvron, H.; Lavril, T.; Izacard, G.; Martinet, X.; Lachaux, M.-A.; Lacroix, T.; Rozière, B.; Goyal, N.; Hambro, E.; Azhar, F.; et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Tunstall, L.; Beeching, E.; Lambert, N.; Rajani, N.; Huang, S.; Rasul, K.; Bartolome, A.; M. Rush, A.; and Wolf, T. ???? The Alignment Handbook.
- von Werra, L.; Belkada, Y.; Tunstall, L.; Beeching, E.; Thrush, T.; Lambert, N.; and Huang, S. 2020. TRL: Transformer Reinforcement Learning. https://github.com/huggingface/trl.
- Wang, B.; Zheng, R.; Chen, L.; Liu, Y.; Dou, S.; Huang, C.; Shen, W.; Jin, S.; Zhou, E.; Shi, C.; et al. 2024. Secrets of rlhf in large language models part ii: Reward modeling. *arXiv preprint arXiv:2401.06080*.
- Wei, J.; Bosma, M.; Zhao, V. Y.; Guu, K.; Yu, A. W.; Lester, B.; Du, N.; Dai, A. M.; and Le, Q. V. 2021. Finetuned language models are zero-shot learners. *arXiv preprint arXiv:2109.01652*.
- Xu, H.; Sharaf, A.; Chen, Y.; Tan, W.; Shen, L.; Van Durme, B.; Murray, K.; and Kim, Y. J. 2024. Contrastive preference optimization: Pushing the boundaries of llm performance in machine translation. *arXiv* preprint arXiv:2401.08417.
- Zheng, L.; Chiang, W.-L.; Sheng, Y.; Zhuang, S.; Wu, Z.; Zhuang, Y.; Lin, Z.; Li, Z.; Li, D.; Xing, E. P.; Zhang, H.; Gonzalez, J. E.; and Stoica, I. 2023. Judging LLM-as-a-judge with MT-Bench and Chatbot Arena. arXiv:2306.05685.