

---

# A Meta-Algorithm for Aligning LLMs with General Preferences

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Many alignment methods, including reinforcement learning from human feedback  
2 (RLHF), rely on the Bradley-Terry reward assumption, which is insufficient to cap-  
3 ture the full range of general human preferences. To achieve robust alignment with  
4 general preferences, we model the alignment problem as a two-player zero-sum  
5 game, where the Nash equilibrium policy guarantees a 50% win rate against any  
6 competing policy. However, previous algorithms for finding the Nash policy either  
7 diverge or converge to a Nash policy in a modified game, even in a simple synthetic  
8 setting, thereby failing to maintain the 50% win rate guarantee against all other  
9 policies. We propose a meta-algorithm for language model alignment with general  
10 preferences, inspired by convergent algorithms in game theory. Theoretically, we  
11 prove that our meta-algorithm converges to an exact Nash policy. Additionally, our  
12 meta-algorithm is simple and can be integrated with many existing methods de-  
13 signed for RLHF and preference optimization with minimal changes. Experimental  
14 results demonstrate the effectiveness of the proposed framework when combined  
15 with existing preference policy optimization methods.

## 16 1 Introduction

17 Large Language Models (LLMs) [Brown et al., 2020, OpenAI, 2023, Dubey et al., 2024] have  
18 fundamentally transformed the fields of natural language processing and artificial intelligence. They  
19 excel in tasks ranging from text generation and translation to complex question answering and  
20 interactive dialogue systems. As these models become more integrated into daily life, a key challenge  
21 is ensuring they achieve high levels of alignment with human values and preferences.

22 One of the most widely adopted approaches to addressing this challenge is Reinforcement Learning  
23 from Human Feedback (RLHF) [Christiano et al., 2017, Ouyang et al., 2022]. This framework  
24 consists of two steps: first, learning a reward model from a dataset containing human preferences,  
25 and second, optimizing the LLM using the proximal policy optimization (PPO) algorithm [Schulman  
26 et al., 2017]. Recently, Rafailov et al. [2024] observed that the first step can be bypassed, proposing  
27 the direct preference optimization (DPO) algorithm, which directly optimizes the LLM from the  
28 dataset.

29 However, the aforementioned approaches crucially rely on the assumption that human preferences  
30 can be expressed using the Bradley-Terry (BT) model [Bradley and Terry, 1952]. Unfortunately, the  
31 BT model is too restrictive to capture the richness and complexity of human preferences. Specifically,  
32 the BT model can only induce *transitive* preferences—i.e., if more people favor A over B, and B  
33 over C, then more people must favor A over C. Such transitivity may not hold in the presence of  
34 diverse populations and is also incompatible with evidence from human decision-making [May, 1954,  
35 Tversky, 1969].

36 To overcome this limitation, recent research has begun to explore alignment under general preferences.  
 37 [Munos et al. \[2024\]](#) formulate this alignment problem as a symmetric two-player zero-sum game,  
 38 where both players’ strategies are LLMs, and their payoffs are determined by the win rate against the  
 39 opponent’s LLM according to the preference model. The objective is to identify a Nash equilibrium  
 40 policy that guarantees at least a 50% win rate against any other policy [[Azar et al., 2024](#), [Munos](#)  
 41 [et al., 2024](#), [Calandriello et al., 2024](#)]. However, the trajectory of all the proposed algorithms either  
 42 diverge or converge to the Nash policy of a modified game, thereby failing to maintain the 50% win  
 43 rate guarantee against all other policies.

44 **Our Contribution.** We introduce a novel meta-algorithm, **Last-Iterate Nash Equilibrium Policy**  
 45 **Optimization (LINE-PO)**, inspired by the proximal point method, a convergent algorithm for  
 46 solving two-player zero-sum games [[Nemirovski, 2004](#)]. Our first observation is that many existing  
 47 algorithms, including PPO [[Schulman et al., 2017](#)], DPO [[Rafailov et al., 2024](#)], IPO [[Azar et al.,](#)  
 48 [2024](#)], SPPO [[Wu et al., 2024](#)], INPO [[Zhang et al., 2024](#)], etc., can be interpreted as implementations  
 49 of the Prox operator [[Nemirovski, 2004](#)]. LINE-PO employs the Prox operator as its fundamental  
 50 building block and provably *converges* to the Nash equilibrium policy in the *last iterate*, assuming  
 51 the Prox operator can be computed exactly. This approach allows us to leverage many existing  
 52 algorithms in a black-box manner. While several algorithms in the literature demonstrate average-  
 53 iterate convergence to the Nash equilibrium policy, they all diverge in the last iterate. Unfortunately,  
 54 iterate averaging can be cumbersome, particularly when deep-learning components are involved, as  
 55 it may not be feasible to average the outputs of LLMs.<sup>1</sup> Compared to these algorithms, LINE-PO  
 56 achieves the more desirable last-iterate convergence.

57 Additionally, we validate the effectiveness of LINE-PO in both synthetic and LLM settings.

58 **Synthetic Setting.** We construct a  $3 \times 3$  two-player zero-sum preference game, and compare  
 59 LINE-PO with a wide range of algorithms proposed in the literature. The result clearly shows that  
 60 LINE-PO is the only algorithm that converges to the Nash equilibrium of the game in the last iterate.

61 **LLM Setting.** Furthermore, we evaluate the performance of LINE-PO against existing preference  
 62 optimization algorithms under a real-world setting, where a pre-trained LLM, Qwen2-1.5B [[Yang](#)  
 63 [et al., 2024](#)], is fine-tuned using different algorithms on the UltraFeedback [[Cui et al., 2023](#)] dataset,  
 64 which is commonly used for alignment fine-tuning of LLMs. Our experimental results demonstrate  
 65 the advantages of LINE-PO: it achieves at least 55% win rate compared against baseline algorithms  
 66 including iterative algorithms such as iterative IPO [[Azar et al., 2024](#)] and INPO [[Zhang et al., 2024](#)].

## 67 2 Backgrounds

68 We use  $\Delta(\mathcal{Z})$  to denote a distribution over a set  $\mathcal{Z}$ . We denote  $x \in \mathcal{X}$  as an instruction where  $\mathcal{X}$  is  
 69 the instruction set. We assume a fixed distribution  $\rho \in \Delta(\mathcal{X})$  over the instruction set. We denote  $\mathcal{Y}$  as  
 70 the response set and  $y \in \mathcal{Y}$  as one response. Given any instruction  $x \in \mathcal{X}$ , an LLM policy  $\pi$  specifies  
 71 the output distribution  $\pi(\cdot | x) \in \Delta(\mathcal{Y})$ . For distributions  $p, q \in \Delta(\mathcal{Z})$ , the Kullback-Leibler (KL)  
 72 divergence is defined as  $\text{KL}(p||q) := \sum_{z \in \mathcal{Z}} p(z) \log \frac{p(z)}{q(z)}$ . The sigmoid function is  $\sigma(x) := \frac{e^x}{e^x + 1}$ .  
 73 We use  $\text{supp}(p)$  to denote the support of a distribution  $p$ .

74 **Preference Models** In this paper, we focus on general preference models.

75 **Definition 1** (General Preference Model). A general preference model  $\mathbb{P} : \mathcal{X} \times \mathcal{Y} \times \mathcal{Y} \rightarrow [0, 1]$   
 76 satisfies  $\mathbb{P}(y_1 \succ y_2 | x) = 1 - \mathbb{P}(y_2 \succ y_1 | x)$ . When we query  $\mathbb{P}$  with  $(x, y_1, y_2)$ , it outputs 1 with  
 77 probability  $\mathbb{P}(y_1 \succ y_2 | x)$  meaning  $y_1$  is preferred over  $y_2$ , and it outputs 0 otherwise.

78 We define  $\mathbb{P}(\pi_1 \succ \pi_2) := \mathbb{E}_{x \sim \rho} [\mathbb{E}_{y_1 \sim \pi_1, y_2 \sim \pi_2} [\mathbb{P}(y_1 \succ y_2 | x)]]$  as the win rate of  $\pi_1$  over  $\pi_2$  under  
 79 preference model  $\mathbb{P}$ . A special case of the general preference model is the Bradley-Terry (BT) model,  
 80 which assumes a reward function parameterizes the preference. We review alignment under the BT  
 81 model in [Appendix A](#).

<sup>1</sup>Storing all LLMs produced during training could solve this, but it is highly space-inefficient and, to our knowledge, has not been implemented.

Table 1: Property comparison of different preference optimization algorithms. (\*) Means convergence in the original game  $J(\pi_1, \pi_2)$

Algorithm	General Preference	Regularized Game Solver	Last-Iterate Convergence*
DPO [Rafailov et al., 2024]	✗	✗	✗
IPO [Azar et al., 2024]	✓	✗	✗
SPPO [Wu et al., 2024]	✓	✗	✗
INPO [Zhang et al., 2024]	✓	✓	✗
LINE-PO	✓	✓	✓

82 **Definition 2** (Bradley-Terry Model). *A preference model  $\mathbb{P}$  satisfies the Bradley-Terry (BT) assumption*  
 83 *if there exists a reward function  $r^* : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$  such that*

$$\mathbb{P}(y_1 \succ y_2 \mid x) = \frac{\exp(r^*(x, y_1))}{\exp(r^*(x, y_1)) + \exp(r^*(x, y_2))} = \sigma(r^*(x, y_1) - r^*(x, y_2)).$$

## 84 2.1 Alignment with General Preference Models

85 The BT model assumption is insufficient to capture the full range of general human preferences  
 86 [Munos et al., 2024, Swamy et al., 2024]. To achieve robust alignment with general preferences, we  
 87 model the policy optimization problem as a two-player zero-sum game with the objective function as  
 88 follows:<sup>2</sup>

$$J(\pi_1, \pi_2) := \mathbb{P}(\pi_1 \succ \pi_2) - \frac{1}{2} = \mathbb{E}_{x \sim \rho} [\mathbb{E}_{y_1 \sim \pi_1, y_2 \sim \pi_2} [\mathbb{P}(y_1 \succ y_2 \mid x)]] - \frac{1}{2}. \quad (1)$$

89 In this game, the max-player controls  $\pi_1$  and tries to maximize  $J(\pi_1, \pi_2)$  while the min-player  
 90 controls  $\pi_2$  and tries to minimize  $J(\pi_1, \pi_2)$ . We focus only on policies with  $\Pi := \{\pi : \text{supp}(\pi) \subseteq$   
 91  $\text{supp}(\pi_{\text{sft}})\}$  in the support of the initial SFT policy. A Nash equilibrium policy  $(\pi_1^*, \pi_2^*)$  satisfies

$$\pi_1^*, \pi_2^* \in \underset{\pi_1 \in \Pi}{\text{argmax}} \underset{\pi_2 \in \Pi}{\text{argmin}} J(\pi_1, \pi_2), \quad J(\pi_1, \pi_2^*) \leq J(\pi_1^*, \pi_2^*) \leq J(\pi_1^*, \pi_2), \forall \pi_1, \pi_2 \in \Pi.$$

92 Since  $J(\pi_1, \pi_2)$  is symmetric, the game has a symmetric Nash equilibrium  $(\pi^*, \pi^*)$ . Moreover,  
 93 the Nash equilibrium policy  $\pi^*$  guarantees that for any other policy  $\pi$ , its win rate is at least  
 94  $\mathbb{P}(\pi^* \succ \pi) \geq \mathbb{P}(\pi^* \succ \pi^*) = 50\%$ . We call this property *robust alignment*. Our goal is to find a  
 95 policy with robust alignment.

96 Existing online iterative preference optimization methods designed for or applicable to the original  
 97 game including iterative IPO [Azar et al., 2024] and SPPO [Wu et al., 2024], are based on Multi-  
 98 plicative Weights Update, and thus *diverge* as we show in Section 4. There is also a line of works  
 99 including Nash-MD [Munos et al., 2024, Ye et al., 2024], Online IPO [Calandriello et al., 2024],  
 100 INPO [Zhang et al., 2024] aim to find the Nash equilibrium of a modified KL-regularized game:

$$J_\tau(\pi_1, \pi_2, \pi_{\text{ref}}) := J(\pi_1, \pi_2) - \tau \mathbb{E}_{x \sim \rho} [\text{KL}(\pi_1(\cdot \mid x) \parallel \pi_{\text{ref}}(\cdot \mid x))] + \tau \mathbb{E}_{x \sim \rho} [\text{KL}(\pi_2(\cdot \mid x) \parallel \pi_{\text{ref}}(\cdot \mid x))].$$

101 The additional KL regularization terms in the objective are introduced for training stability. However,  
 102 the Nash equilibrium of the modified game no longer achieves robust alignment, i.e., has a win rate  
 103 of at least 50% against any competing policy.

104 Moreover, most existing theoretical convergence guarantees only hold for the average iterate, i.e., the  
 105 uniform mixture of training iterates, which is not used in practice. We focus on designing algorithms  
 106 with provable last-iterate convergence to Nash equilibrium, which aligns with practice and is more  
 107 space-efficient [Munos et al., 2024].

108 In the next section, we propose a meta-algorithm that uses algorithms designed for the regular-  
 109 ized game  $J_\tau(\pi_1, \pi_2, \pi_{\text{ref}})$  or other preference optimization methods as black-boxes to find Nash  
 110 equilibrium of  $J(\pi_1, \pi_2)$  (1), thereby achieving robust alignment.

## 111 3 Last-Iterate Nash Equilibrium Policy Optimization

112 We propose an extremely simple meta-algorithm, Last-Iterate Nash Equilibrium Policy Optimization  
 113 (LINE-PO, Algorithm 1), for robustly aligning LLMs with general preferences. LINE-PO is an

<sup>2</sup>We introduce the constant  $\frac{1}{2}$  only to ensure the game is zero-sum and it has no effect on its Nash equilibria.

---

**Algorithm 1:** Last-Iterate Convergent Nash Equilibrium Policy Optimization (LINE-PO)

---

**Input:** Initial policy  $\pi_{\text{sft}}$ , preference oracle  $\mathbb{P}$ , regularization  $\tau > 0$

```
1 Initialize  $\pi^1, \pi_{\text{ref}} \leftarrow \pi_{\text{sft}}$ 
2 for  $t = 1, 2, \dots, T - 1$  do
3    $\pi^{t+1} \leftarrow \operatorname{argmax}_{\pi_1} \min_{\pi_2} J_\tau(\pi_1, \pi_2, \pi_{\text{ref}})$  using Algorithm 2
4    $\pi_{\text{ref}} \leftarrow \pi^{t+1}$ 
5 return  $\pi^T$ 
```

---

114 online iterative algorithm inspired by the classic Conceptual Prox method [Nemirovski, 2004] first  
115 introduced in the optimization theory community. This method has recently been applied to finding  
116 a Nash equilibrium in zero-sum games [Perolat et al., 2021, Abe et al., 2024] and has had notable  
117 success in training advanced game AI models [Perolat et al., 2022].

### 118 3.1 LINE-PO

119 In each iteration  $t$ , LINE-PO updates the next-iteration policy  $\pi^{t+1}$  as the Nash equilibrium policy of  
120 a regularized game  $J_\tau(\pi_1, \pi_2, \pi_{\text{ref}})$  using the current policy as reference  $\pi_{\text{ref}} = \pi^t$ . The rationale  
121 behind LINE-PO is simple: update the reference policy when there is no improvement in the  
122 regularized game. Denote  $\pi^*$  the Nash equilibrium of the original game. We show that KL divergence  
123 to  $\pi^*$  is monotonically decreasing:  $\text{KL}(\pi^* || \pi^{t+1}) \leq \text{KL}(\pi^* || \pi^t)$ . Since  $\pi^{t+1}$  is closer to the Nash  
124 equilibrium than  $\pi^t$ , LINE-PO updates the reference policy from  $\pi^t$  to  $\pi^{t+1}$  for further optimization.  
125 We also remark that in LINE-PO, the regularization amount  $\tau > 0$  does not need to decrease and  
126 could be kept constant.

127 Each iteration of LINE-PO requires solving a zero-sum game with additional KL regularization  
128  $J_\tau(\pi_1, \pi_2, \pi_{\text{ref}})$ . We will show momentarily that many existing policy optimization methods for  
129 alignment can be applied to the KL regularized game and have exponentially fast convergence. We  
130 prove the meta-algorithm LINE-PO achieves last-iterate convergence to a Nash equilibrium with  
131 robust alignment property, which appears to be the first in the context of LLM alignment.

132 **Theorem 1.** *We assume that there exists a Nash equilibrium  $\pi^*$  of  $J(\pi_1, \pi_2)$  (defined in (1)) such  
133 that  $\text{supp}(\pi^*) = \text{supp}(\pi_{\text{sft}})$ . In every iteration  $t \geq 1$ , it holds that  $\text{KL}(\pi^* || \pi^{t+1}) \leq \text{KL}(\pi^* || \pi^t)$ .  
134 Moreover, LINE-PO has last-iterate convergence, i.e.,  $\lim_{t \rightarrow \infty} \pi^t$  exists and is a Nash equilibrium.*

### 135 3.2 Solving a Regularized Game

136 We show how to solve the Nash equilibrium of the regularized game  $J_\tau(\pi_1, \pi_2, \pi_{\text{ref}})$  using the Mirror  
137 Descent (MD) algorithm and how to implement MD using existing policy optimization algorithms.  
138 For simplicity, we consider policy  $\pi \in \Delta(\mathcal{Y})$  and omit the dependence on the instruction  $x$ . All  
139 discussions can be extended to the contextual setting in a straightforward way.

140 **Mirror Descent and Multiplicative Weights Update** Mirror Descent (MD) is a classical family of  
141 optimization algorithms. An important member of this family is the Multiplicative Weights Update  
142 (MWU) algorithm, which is MD with negative entropy regularization. For a maximization problem  
143  $\max_{\pi} f(\pi)$ , given an existing policy  $\pi^t$ , MWU computes the update  $\pi^{t+1}$  as follows:

$$\pi^{t+1} := \operatorname{argmax}_{\pi} \langle \nabla f(\pi^t), \pi \rangle - \eta^{-1} \cdot \text{KL}(\pi || \pi^t). \quad (2)$$

144 Note that RLHF in (4) is equivalent to one step of MWU if we interpret the reward  $r$  as the gradient  
145  $\nabla f(\pi_{\text{ref}})$ .

146 **Prox operator.** The update rule of MWU can be compactly written using the *prox operator* as  
147 shown in Algorithm 2.<sup>3</sup> Fix a 1-strongly convex function  $\varphi : \mathcal{Z} \rightarrow \mathbb{R}$  over a closed convex set  
148  $\mathcal{Z} \subset \mathbb{R}^n$ . The *Bregman divergence* induced by  $\varphi$  is

$$D_\varphi(\cdot || \cdot) : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}_{\geq 0},$$
$$D_\varphi(z || z') := \varphi(z) - \varphi(z') - \langle \nabla \varphi(z'), z - z' \rangle.$$

---

<sup>3</sup>The prox operator is also called the prox-mapping [Nemirovski, 2004].

---

**Algorithm 2:** Regularized game solver

---

**Input:** Reference policy  $\pi_{\text{ref}}$ , preference oracle  $\mathbb{P}$ , regularization  $\tau > 0$ , step size  $\eta > 0$ , number of iterations  $K \geq 1$

**Output:** A Nash policy  $\text{argmax}_{\pi_1} \min_{\pi_2} J_\tau(\pi_1, \pi_2, \pi_{\text{ref}})$

```
1 Initialize  $\mu^1 \leftarrow \pi_{\text{ref}}$ 
2 for  $k = 1, 2, \dots, K - 1$  do
3    $g_\tau^k \leftarrow \nabla_\mu(\mathbb{P}(\mu \succ \mu_k) - \tau \text{KL}(\mu || \pi_{\text{ref}})) = \mathbb{P}(\cdot \succ \mu_k) - \tau(\log \frac{\mu_k(\cdot)}{\pi_{\text{ref}}(\cdot)} + 1)$ 
4    $\mu^{k+1} \leftarrow \text{Prox}(\mu_k, \eta g_\tau^k)$ 
5 return  $\mu_K$ 
```

---

149 Given a reference point  $z \in \mathcal{Z}$  and a vector  $g \in \mathbb{R}^n$ , the prox operator  $\text{Prox}(z, g)$  generalizes the  
150 notion of a gradient ascent step from  $z$  in the direction of  $g$ .

151 **Definition 3** (Prox Operator). For a strongly convex regularizer  $\varphi$ , the prox operator is defined as

$$\text{Prox}(z, g) := \underset{z'}{\text{argmax}} \langle g, z' \rangle - D_\varphi(z' || z) = \underset{z'}{\text{argmax}} \langle g + \nabla \varphi(z), z' \rangle - \varphi(z'). \quad (3)$$

152 When  $\varphi(z) = \frac{1}{2} \|z\|_2^2$  is the  $\ell_2$  regularizer, the prox operator  $\text{Prox}(z, g) = \Pi_{\mathcal{Z}}[z + g]$  is the  
153 exactly the projected gradient ascent step. In this paper, without additional notes, we choose  
154  $\varphi = \sum_{i=1}^n z[i] \ln z[i]$  as the negative entropy regularizer and the corresponding Bregman divergence  
155  $D_\varphi$  is the KL divergence.

156 The flexibility of the prox operator lies in the choice of  $g$  for different objectives. In RLHF,  $g$  is  
157 the reward model  $r$  and we compute the optimal policy  $\pi^* = \text{Prox}(\pi_{\text{ref}}, \eta r)$ . For vanilla MWU,  $g$   
158 is the gradient  $\nabla f(\pi^t)$  and we update  $\pi^{t+1} = \text{Prox}(\pi^t, \eta \nabla f(\pi^t))$ . When a preference model  $\mathbb{P}$  is  
159 available, we can choose  $g$  as the preference function  $\mathbb{P}(\cdot \succ \pi)$  over the current policy  $\pi$ . For our  
160 theoretical results, we assume the prox operator  $\text{Prox}$  can be evaluated exactly or approximately.  
161 Practically, we can use many existing preference optimization methods to compute the prox operator  
162 as shown in the next section.

163 **Exponentially Fast Convergence** Denote  $\pi_\tau^*$  the Nash equilibrium of the KL regularized game  
164  $J_\tau(\pi_1, \pi_2, \pi_{\text{ref}})$ , which is  $\tau$ -strongly monotone. We can apply classical results to show that MWU  
165 (Algorithm 2) achieves linear last-iterate convergence rate: the distance to the Nash equilibrium  $\pi_\tau^*$   
166 decreases exponentially fast.

167 **Theorem 2.** For appropriate step size  $\eta > 0$ , Algorithm 2 guarantees for every  $k \geq 1$ ,  
168  $\text{KL}(\pi_\tau^* || \mu^{k+1}) \leq (1 - \frac{\eta\tau}{2})^k \text{KL}(\pi_\tau^* || \pi_{\text{ref}})$ .

### 169 3.3 Computing the prox operator

170 We show how to compute the prox operator in practical large-scale applications like LLM alignment.  
171 Specifically, we show that many existing algorithms designed for RLHF and preference optimization  
172 with neural network parameters can be adapted to solve the prox operator  $\text{Prox}(\pi, \eta g)$  ( $\eta > 0$  is  
173 the step size). These algorithms include RL algorithms like PPO, and loss-minimization algorithms  
174 like DPO, IPO, SPPO, DRO, each of which may be preferred in certain settings. Our contribution  
175 here is not proposing new algorithms but unifying existing diverse preference methods through  
176 the perspective of computing the prox operator. Due to space limit, we defer the discussion to  
177 Appendix F. This perspective opens the possibility of applying other algorithms from online learning  
178 and optimization to robust LLM alignment and we include implementation for two other algorithms  
179 in Appendix H.

## 180 4 Synthetic Experiments

181 We conduct experiments on a simple bandit problem with  $\mathcal{Y} = \{y_a, y_b, y_c\}$  and non-BT preference  
182 model over  $\mathcal{Y}$ . Specifically, we set  $\mathbb{P}[y_b \succ y_a] = \mathbb{P}[y_c \succ y_b] = 0.9$  and  $\mathbb{P}[y_a \succ y_c] = 0.8$ . We can  
183 observe that the preference is intransitive and exhibits a preference cycle  $y_c \succ y_b \succ y_a \succ y_c$ .

184 **Experiments using noiseless gradient** We  
 185 present numerical results of mirror-descent  
 186 (MD) algorithms (equivalent to MWU) and  
 187 LINE-PO (Algorithm 1) in Figure 1. We can see  
 188 that the MD algorithm diverges from the unique  
 189 Nash equilibrium and suffers a large equilibrium  
 190 gap, while LINE-PO achieves fast last-iterate  
 191 convergence to the Nash equilibrium, aligned  
 192 with our theoretical results (Theorem 1).

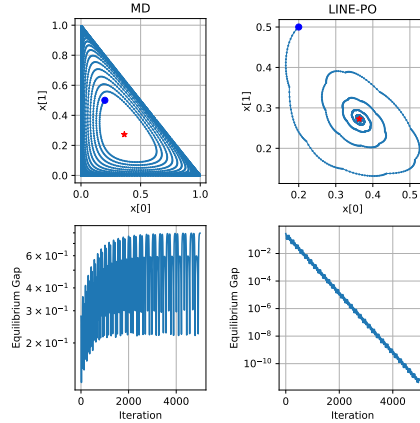


Figure 1: Dynamics on a simple 3-dimensional preference game. The unique Nash equilibrium is  $[4/11, 3/11, 3/11]$  represented as red star. We initialize all algorithms at the blue dot point  $[0.2, 0.5, 0.3]$ .

193 **Experiments using preference samples**  
 194 Since the popular iterative DPO algorithm does  
 195 not contain a gradient step, we also conduct ex-  
 196 periments with only Oracle query access to the  
 197 preference model. We compare the performance  
 198 of various algorithms, including iterative DPO,  
 199 iterative IPO, SPPO, and INPO and present re-  
 200 sults in Figure 2. We remark that iterative DPO  
 201 and iterative IPO both diverge in the last iterate;  
 202 INPO converges to a point that is not Nash equi-  
 203 librium and does not guarantee robust alignment;  
 204 LINE-PO is the only algorithm that achieves  
 205 last-iterate convergence to the Nash equilibrium.  
 206 We defer a more detailed discussion to Appendix I.

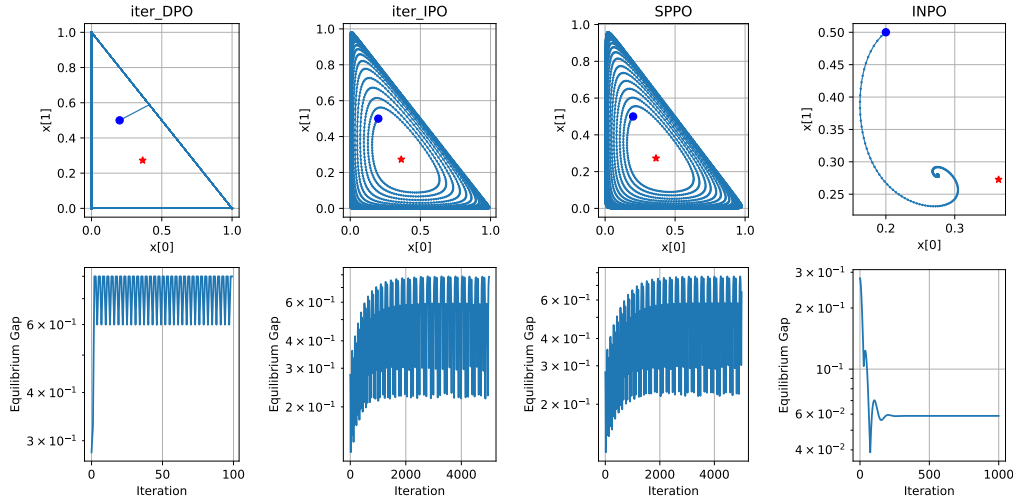


Figure 2: Dynamics on a simple 3-dimensional preference game. The unique Nash equilibrium is  $[4/11, 3/11, 3/11]$  represented as red star. We initialize all algorithms at the blue dot point  $[0.2, 0.5, 0.3]$ .

## 207 5 Real-World Experiments

208 Apart from the controlled synthetic experiments, we conduct experiments with a pre-trained LLM,  
 209 Qwen2-1.5B [Yang et al., 2024], on a commonly used dataset UltraFeedback [Cui et al., 2023] to  
 210 show the effectiveness of LINE-PO under the real-world preference optimization setting.

211 **5.1 Experimental Settings**

212 **Datasets** We use the UltraFeedback dataset, specifically its binarized version for preference fine-  
213 tuning.<sup>4</sup> It contains 64K data examples consisting of a user instruction and a positive-negative output  
214 pair annotated by GPT-4. The instructions contained in this dataset cover a wide range of instruction  
215 types, making it suitable to study preference optimization in a real-world setting. Since we focus on  
216 online and iterative preference optimization, only the instructions are used because the output pairs  
217 will be generated and annotated online. In addition, to reduce the computation cost, the instructions  
218 are randomly split into 6 equal-size subsets. Each subset therefore contains around 10K instructions  
219 and is used in one training iteration.

220 **Preference Oracle** The preference oracle we chose is Llama-3-OffsetBias-8B [Park et al., 2024a],  
221 which is a pairwise preference model that predicts which output is better given an instruction and  
222 an output pair. Fine-tuned from Meta-Llama-3-8B-Instruct [Dubey et al., 2024], it achieves strong  
223 performance on various human preference alignment benchmarks on RewardBench [Lambert et al.,  
224 2024]. We chose it as the preference oracle since it strikes a balance between computation efficiency  
225 and alignment with human preferences, making it suitable for iterative preference optimization.

226 **Online Preference Data Generation** To construct the preference data, i.e., output pairs with a  
227 preference annotation specifying which one is better, we adopt the setting of Zhang et al. [2024]  
228 by sampling 5 candidate outputs for each instruction with a temperature of 0.8 and applying the  
229 preference oracle to compare all the output pairs constructed. The best and the worst candidate  
230 outputs, derived from the pairwise comparison results, are then selected to form a data point.

231 **Baselines** We include the following baselines for comparisons with LINE-PO: (1) SFT, which  
232 fine-tunes the pre-trained Qwen2-1.5B on the UltraChat dataset, and the resulted checkpoint serves  
233 as the start point and/or the reference policy for the other training algorithms; (2) vanilla online  
234 DPO [Rafailov et al., 2024] and (3) vanilla online IPO [Azar et al., 2024], where one training iteration  
235 is performed over the entire instruction set of UltraFeedback; (4) INPO [Zhang et al., 2024], where  
236 at each iteration the training is performed on one data split; (5) iterative IPO, which has a similar  
237 training setting to INPO but without the KL constraint from the reference policy.

238 **Evaluations** We use the instructions in a widely used benchmark, AlpacaEval [Li et al., 2023], to  
239 construct the test set, since these instructions are diverse and cover various task scenarios. However,  
240 we choose not to use the default evaluator of the AlpacaEval benchmark, GPT-4, to perform the  
241 evaluation, but instead use the same preference oracle used in data generation, Llama-3-OffsetBias-  
242 8B, as the evaluator. This is because we aim for a controlled experimental setting – the preference  
243 oracle that the model learns to fit should also be the one used to evaluate the model performance.

244 **Training Details** We follow the training recipe proposed in Tunstall et al. [2023] for the experiments.  
245 Specifically, at each training iteration, the models are fine-tuned for 3 epochs with the batch size  
246 setting to 32 and with a linear learning rate scheduler. The checkpoints are selected based on their  
247 validation loss on the UltraFeedback dataset. As for the hyper-parameters, we perform a grid search  
248 for the strength of the KL regularization,  $\eta^{-1}$ , in vanilla DPO and IPO. Specifically, we found that  
249 DPO achieves the best performance when it is set to 0.01, while IPO achieves the best performance  
250 when  $\eta^{-1}$  is set within the range of 0.01 - 0.002. We then choose the value of  $\eta$  to be 0.002 to  
251 encourage larger learning steps. This value of  $\eta$  is also used for iterative IPO and INPO. INPO has  
252 another hyper-parameter  $\tau$  which controls the strength of the KL regularization from the reference  
253 policy. We determine its value following the setting of Zhang et al. [2024], where  $\eta\tau$  is set to a fixed  
254 ratio,  $1/3$ . Regarding LINE-PO, the second training round starts when the first training round based  
255 on INPO begins to converge/overfit, and  $\eta^{-1}$  is set to 0.01 for the second round for training stability.

256 **5.2 Result Analysis**

257 Figure 3 presents the training dynamics of three iterative preference optimization algorithms we com-  
258 pared: iterative IPO (Iter-IPO), INPO, and LINE-PO, which are demonstrated by their checkpoints’  
259 win rates against the SFT checkpoint and the average length of their outputs. For INPO and LINE-PO,  
260 the model is trained for up to 18 iterations, which are equivalent to 3 training rounds over the entire  
261 instruction set since it has been split into 6 subsets. We note that:

---

<sup>4</sup>[https://huggingface.co/datasets/HuggingFaceH4/ultrafeedback\\_binarized](https://huggingface.co/datasets/HuggingFaceH4/ultrafeedback_binarized).

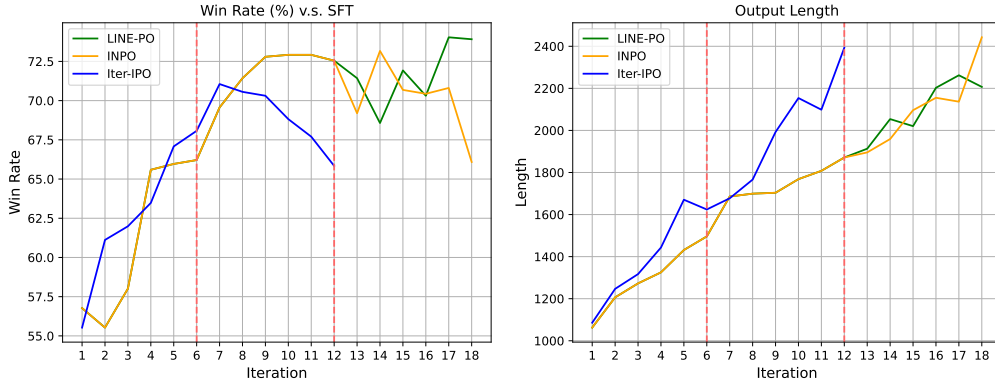


Figure 3: Comparisons of Iterative IPO (Iter-IPO), INPO, and LINE-PO. The win rate of the trained checkpoints against the SFT checkpoint, and the average length of the outputs are compared. The red vertical lines mark the end of one training round – a complete iteration over the 6 data splits.

Table 2: Performance comparison of different training algorithms. The row v.s. column win rate (%) is reported. For INPO, we report its performance with 2-round (R2) and 3-round (R3) training.

Row/Col	SFT	DPO	IPO	Iter-IPO	INPO-R2	INPO-R3	LINE-PO	Avg
Iter-IPO	65.84	56.40	54.04	50.00	47.83	46.21	39.01	51.33
INPO-R2	72.55	60.25	58.39	52.17	50.00	49.32	41.37	54.86
INPO-R3	66.09	60.25	58.51	53.79	50.68	50.00	44.97	54.90
LINE-PO	<b>73.91</b>	<b>66.71</b>	<b>66.21</b>	<b>60.99</b>	<b>58.63</b>	<b>55.03</b>	<b>50.00</b>	<b>61.64</b>

262 (1) Iter-IPO shows a quicker improvement rate at the beginning of the training, but its performance  
 263 against the SFT checkpoint starts to degrade in the second training round, which indicates the inherent  
 264 instability of this training algorithm.

265 (2) INPO archives a relatively stable win rate against SFT at the end of the second training round.  
 266 However, its win rate starts to slightly degrade in the third training round. We suspect this suggests  
 267 that INPO has started to converge and/or overfit. Therefore, for LINE-PO, which shares the same  
 268 training trajectory as INPO for the first two training rounds, we update the reference policy at the  
 269 beginning of the third training round, following the optimization process described in Algorithm 1.

270 (3) LINE-PO is able to further improve the model performance with the updated reference policy.  
 271 Notably, it also results in the shortest outputs compared to Iter-IPO and INPO, suggesting that it is  
 272 more robust to the length bias of the preference models which preference optimization algorithms  
 273 tend to exploit [Park et al., 2024b].

274 Table 2 provides pairwise comparisons between the final checkpoints of the iterative preference  
 275 optimization algorithms and a few baselines. It demonstrates the clear advantage of LINE-PO, which  
 276 is able to achieve an above 50% win rate against all the other checkpoints. In contrast, Iter-IPO can  
 277 only outperform the vanilla DPO and IPO settings. Regarding INPO, we found that the average win  
 278 rate of its checkpoint after the third training round (INPO-R3) is only slightly higher than that of  
 279 its intermediate checkpoint at the end of the second training round (INPO-R2) (54.90 vs. 54.86),  
 280 suggesting that its performance plateaued by the end of the second training round.

## 281 6 Conclusion

282 We have proposed LINE-PO, a meta-algorithm for preference optimization that provably converges  
 283 to the Nash equilibrium policy in the last iterate. We have provided a theoretical analysis of the  
 284 properties of LINE-PO and have empirically demonstrated its effectiveness under both synthetic  
 285 and real-world experimental settings. We believe LINE-PO has significant potential to enhance the  
 286 performance of LLMs in the alignment fine-tuning setting, due to its theoretical guarantees and  
 287 flexibility, as it can be integrated with existing learning algorithms while overcoming their limitations.



## 288 References

- 289 Kenshi Abe, Kaito Ariu, Mitsuki Sakamoto, and Atsushi Iwasaki. Adaptively perturbed mirror  
290 descent for learning in games. In *Proceedings of the 41st International Conference on Machine*  
291 *Learning*, 2024.
- 292 Mohammad Gheshlaghi Azar, Zhaohan Daniel Guo, Bilal Piot, Remi Munos, Mark Rowland, Michal  
293 Valko, and Daniele Calandriello. A general theoretical paradigm to understand learning from  
294 human preferences. In *International Conference on Artificial Intelligence and Statistics*, pages  
295 4447–4455. PMLR, 2024.
- 296 Heinz H. Bauschke, Walaa M. Moursi, and Xianfu Wang. Generalized monotone operators and their  
297 averaged resolvents. *Math. Program.*, 189(1):55–74, 2021. doi: 10.1007/S10107-020-01500-6.  
298 URL <https://doi.org/10.1007/s10107-020-01500-6>.
- 299 Ralph Allan Bradley and Milton E. Terry. Rank analysis of incomplete block designs: I. the method  
300 of paired comparisons. *Biometrika*, 39:324, 1952. URL <https://api.semanticscholar.org/CorpusID:125209808>.
- 302 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhari-  
303 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agar-  
304 wal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh,  
305 Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Ma-  
306 teusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCand-  
307 lish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot  
308 learners. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Ad-  
309 vances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Asso-  
310 ciates, Inc., 2020. URL [https://proceedings.neurips.cc/paper\\_files/paper/  
311 2020/file/1457c0d6bfcb4967418bfb8ac142f64a-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfcb4967418bfb8ac142f64a-Paper.pdf).
- 312 Yang Cai, Argyris Oikonomou, and Weiqiang Zheng. Finite-time last-iterate convergence for learning  
313 in multi-player games. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- 314 Daniele Calandriello, Zhaohan Daniel Guo, Remi Munos, Mark Rowland, Yunhao Tang,  
315 Bernardo Avila Pires, Pierre Harvey Richemond, Charline Le Lan, Michal Valko, Tianqi Liu,  
316 Rishabh Joshi, Zeyu Zheng, and Bilal Piot. Human alignment of large language models through  
317 online preference optimisation. In *Forty-first International Conference on Machine Learning*, 2024.  
318 URL <https://openreview.net/forum?id=2RQqg2Y7Y6>.
- 319 Paul F. Christiano, Jan Leike, Tom B. Brown, Miljan Martic, Shane Legg, and Dario Amodei.  
320 Deep reinforcement learning from human preferences. In Isabelle Guyon, Ulrike von Luxburg,  
321 Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett,  
322 editors, *Advances in Neural Information Processing Systems 30: Annual Conference on Neu-  
323 ral Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages  
324 4299–4307, 2017. URL [https://proceedings.neurips.cc/paper/2017/hash/  
325 d5e2c0adad503c91f91df240d0cd4e49-Abstract.html](https://proceedings.neurips.cc/paper/2017/hash/d5e2c0adad503c91f91df240d0cd4e49-Abstract.html).
- 326 Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Wei Zhu, Yuan Ni, Guotong Xie, Zhiyuan Liu,  
327 and Maosong Sun. Ultrafeedback: Boosting language models with high-quality feedback. *arXiv*  
328 *preprint arXiv:2310.01377*, 2023.
- 329 Constantinos Daskalakis and Ioannis Panageas. The limit points of (optimistic) gradient descent in  
330 min-max optimization. In *the 32nd Annual Conference on Neural Information Processing Systems*  
331 *(NeurIPS)*, 2018.
- 332 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha  
333 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn,  
334 Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston  
335 Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron,  
336 Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris  
337 McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton  
338 Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David

339 Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes,  
340 Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip  
341 Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme  
342 Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu,  
343 Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov,  
344 Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah,  
345 Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu  
346 Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph  
347 Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani,  
348 Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz  
349 Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yearly, Laurens van der Maaten, Lawrence  
350 Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas  
351 Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri,  
352 Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis,  
353 Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov,  
354 Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan  
355 Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan,  
356 Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy,  
357 Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohit Girdhar, Rohit  
358 Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou,  
359 Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia  
360 Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Rapparthi, Sheng Shen, Shengye Wan,  
361 Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla,  
362 Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek  
363 Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao,  
364 Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent  
365 Gonguet, Virginie Do, Vish Vogeti, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu,  
366 Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia,  
367 Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen  
368 Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe  
369 Papakipos, Aaditya Singh, Aaron Grattafiori, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya  
370 Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex  
371 Vaughan, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Anam Yunus, Andrei  
372 Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew  
373 Ryan, Ankit Ramchandani, Annie Franco, Aparajita Saraf, Arkabandhu Chowdhury, Ashley  
374 Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin  
375 Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu,  
376 Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt  
377 Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Changan Wang, Changkyu Kim, Chao  
378 Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Damon  
379 Civin, Dana Beaty, Daniel Kreymer, Daniel Li, Danny Wyatt, David Adkins, David Xu, Davide  
380 Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dinggang Wang, Duc Le,  
381 Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily  
382 Hahn, Emily Wood, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix  
383 Kreuk, Feng Tian, Firat Ozgenel, Francesco Caggioni, Francisco Guzmán, Frank Kanayet, Frank  
384 Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern,  
385 Govind Thattai, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid  
386 Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen  
387 Suk, Henry Aspegren, Hunter Goldman, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-  
388 Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Japhet Asher, Jean-Baptiste  
389 Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul,  
390 Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie,  
391 Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik  
392 Prasad, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraghavan, Kelly  
393 Michelena, Keqian Li, Kun Huang, Kunal Chawla, Kushal Lakhota, Kyle Huang, Lailin Chen,  
394 Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu,  
395 Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Maria  
396 Tsimpoukelli, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev,  
397 Maxim Naumov, Maya Lathi, Meghan Keneally, Michael L. Seltzer, Michal Valko, Michelle

- 398 Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang,  
399 Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam,  
400 Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier,  
401 Nikolay Pavlovich Laptev, Ning Dong, Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia  
402 Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro  
403 Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani,  
404 Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy,  
405 Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan  
406 Maheswari, Russ Howes, Ruty Rinott, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara  
407 Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh  
408 Ramaswamy, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha,  
409 Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe,  
410 Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan  
411 Govindaprasad, Sumit Gupta, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury,  
412 Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe  
413 Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi,  
414 Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vitor Albiero, Vlad Ionescu,  
415 Vlad Poenaru, Vlad Tiberiu Mihalescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang,  
416 Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang,  
417 Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang,  
418 Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait,  
419 Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. The llama 3 herd  
420 of models, 2024. URL <https://arxiv.org/abs/2407.21783>.
- 421 Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. KTO: model  
422 alignment as prospect theoretic optimization. *CoRR*, abs/2402.01306, 2024. doi: 10.48550/ARXIV.  
423 2402.01306. URL <https://doi.org/10.48550/arXiv.2402.01306>.
- 424 Francisco Facchinei and Jong-Shi Pang. *Finite-dimensional variational inequalities and complemen-*  
425 *tarity problems*. Springer, 2003.
- 426 Zhaolin Gao, Jonathan D. Chang, Wenhao Zhan, Owen Oertell, Gokul Swamy, Kianté Brantley,  
427 Thorsten Joachims, J. Andrew Bagnell, Jason D. Lee, and Wen Sun. REBEL: reinforcement  
428 learning via regressing relative rewards. *CoRR*, abs/2404.16767, 2024. doi: 10.48550/ARXIV.  
429 2404.16767. URL <https://doi.org/10.48550/arXiv.2404.16767>.
- 430 Noah Golowich, Sarath Pattathil, and Constantinos Daskalakis. Tight last-iterate convergence rates  
431 for no-regret learning in multi-player games. *Advances in neural information processing systems*  
432 (*NeurIPS*), 2020a.
- 433 Noah Golowich, Sarath Pattathil, Constantinos Daskalakis, and Asuman Ozdaglar. Last iterate is  
434 slower than averaged iterate in smooth convex-concave saddle point problems. In *Conference on*  
435 *Learning Theory (COLT)*, 2020b.
- 436 Eduard Gorbunov, Adrien B. Taylor, and Gauthier Gidel. Last-iterate convergence of opti-  
437 mistic gradient method for monotone variational inequalities. In Sanmi Koyejo, S. Mo-  
438 hamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh, editors, *Advances in Neu-*  
439 *ral Information Processing Systems 35: Annual Conference on Neural Information Pro-*  
440 *cessing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December*  
441 *9, 2022*, 2022. URL [http://papers.nips.cc/paper\\_files/paper/2022/hash/893cd874ba98afa54ae9e385a24a83ac-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/893cd874ba98afa54ae9e385a24a83ac-Abstract-Conference.html).
- 443 Yu-Guan Hsieh, Kimon Antonakopoulos, and Panayotis Mertikopoulos. Adaptive learning in  
444 continuous games: Optimal regret bounds and convergence to nash equilibrium. In *Conference on*  
445 *Learning Theory*, pages 2388–2422. PMLR, 2021.
- 446 A. N. Iusem, T. Pennanen., and B. F. Svaiter. Inexact variants of the proximal point algorithm  
447 without monotonicity. *SIAM Journal on Optimization*, 13(4):1080–1097, 2003. doi: 10.1137/  
448 S1052623401399587. URL <https://doi.org/10.1137/S1052623401399587>.
- 449 G. M. Korpelevich. The extragradient method for finding saddle points and other problems. *Matecon*,  
450 12:747–756, 1976. URL <https://ci.nii.ac.jp/naid/10017556617/>.

- 451 Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu,  
452 Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, et al. Rewardbench: Evaluating reward models  
453 for language modeling. *arXiv preprint arXiv:2403.13787*, 2024.
- 454 Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy  
455 Liang, and Tatsunori B. Hashimoto. AlpacaEval: An automatic evaluator of instruction-following  
456 models. [https://github.com/tatsu-lab/alpaca\\_eval](https://github.com/tatsu-lab/alpaca_eval), 5 2023.
- 457 Kenneth O May. Intransitivity, utility, and the aggregation of preference patterns. *Econometrica:  
458 Journal of the Econometric Society*, pages 1–13, 1954.
- 459 Panayotis Mertikopoulos, Christos H. Papadimitriou, and Georgios Piliouras. Cycles in adversarial  
460 regularized learning. In Artur Czumaj, editor, *Proceedings of the Twenty-Ninth Annual ACM-  
461 SIAM Symposium on Discrete Algorithms, SODA 2018, New Orleans, LA, USA, January 7-10,  
462 2018*, pages 2703–2717. SIAM, 2018. doi: 10.1137/1.9781611975031.172. URL <https://doi.org/10.1137/1.9781611975031.172>.
- 464 Aryan Mokhtari, Asuman Ozdaglar, and Sarath Pattathil. A unified analysis of extra-gradient and  
465 optimistic gradient methods for saddle point problems: Proximal point approach. In *International  
466 Conference on Artificial Intelligence and Statistics (AISTATS)*, 2020a.
- 467 Aryan Mokhtari, Asuman E Ozdaglar, and Sarath Pattathil. Convergence rate of  $\mathcal{O}(1/k)$  for optimistic  
468 gradient and extragradient methods in smooth convex-concave saddle point problems. *SIAM  
469 Journal on Optimization*, 30(4):3230–3251, 2020b.
- 470 Remi Munos, Michal Valko, Daniele Calandriello, Mohammad Gheshlaghi Azar, Mark Rowland,  
471 Zhaohan Daniel Guo, Yunhao Tang, Matthieu Geist, Thomas Mesnard, Côme Fiegel, et al. Nash  
472 learning from human feedback. In *Forty-first International Conference on Machine Learning*,  
473 2024.
- 474 Arkadi Nemirovski. Prox-method with rate of convergence  $\mathcal{O}(1/t)$  for variational inequalities with  
475 Lipschitz continuous monotone operators and smooth convex-concave saddle point problems. *SIAM  
476 Journal on Optimization*, 15(1):229–251, 2004.
- 477 OpenAI. GPT-4 technical report. *CoRR*, abs/2303.08774, 2023. doi: 10.48550/ARXIV.2303.08774.  
478 URL <https://doi.org/10.48550/arXiv.2303.08774>.
- 479 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin,  
480 Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser  
481 Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan  
482 Leike, and Ryan Lowe. Training language models to follow instructions with human feedback.  
483 In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh, editors,  
484 *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information  
485 Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December  
486 9, 2022*, 2022. URL [http://papers.nips.cc/paper\\_files/paper/2022/hash/blefde53be364a73914f58805a001731-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/blefde53be364a73914f58805a001731-Abstract-Conference.html).
- 488 Junsoo Park, Seungyeon Jwa, Meiying Ren, Daeyoung Kim, and Sanghyuk Choi. Offsetbias:  
489 Leveraging debiased data for tuning evaluators. *arXiv preprint arXiv:2407.06551*, 2024a.
- 490 Ryan Park, Rafael Rafailov, Stefano Ermon, and Chelsea Finn. Disentangling length from quality in  
491 direct preference optimization. *arXiv preprint arXiv:2403.19159*, 2024b.
- 492 Julien Perolat, Remi Munos, Jean-Baptiste Lespiau, Shayegan Omidshafiei, Mark Rowland, Pedro  
493 Ortega, Neil Burch, Thomas Anthony, David Balduzzi, Bart De Vylder, et al. From Poincaré  
494 recurrence to convergence in imperfect information games: Finding equilibrium via regularization.  
495 In *International Conference on Machine Learning*, pages 8525–8535. PMLR, 2021.
- 496 Julien Perolat, Bart De Vylder, Daniel Hennes, Eugene Tarassov, Florian Strub, Vincent de Boer,  
497 Paul Muller, Jerome T Connor, Neil Burch, Thomas Anthony, et al. Mastering the game of Stratego  
498 with model-free multiagent reinforcement learning. *Science*, 378(6623):990–996, 2022.

- 499 Leonid Denisovich Popov. A modification of the Arrow-Hurwicz method for search of saddle points.  
500 *Mathematical notes of the Academy of Sciences of the USSR*, 28(5):845–848, 1980. Publisher:  
501 Springer.
- 502 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea  
503 Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances*  
504 *in Neural Information Processing Systems*, 36, 2024.
- 505 Sasha Rakhlin and Karthik Sridharan. Optimization, learning, and games with predictable sequences.  
506 *Advances in Neural Information Processing Systems*, 2013.
- 507 Pierre Harvey Richemond, Yunhao Tang, Daniel Guo, Daniele Calandriello, Mohammad Gheshlaghi  
508 Azar, Rafael Rafailov, Bernardo Avila Pires, Eugene Tarassov, Lucas Spangher, Will Ellsworth,  
509 et al. Offline regularised reinforcement learning for large language models alignment. *arXiv*  
510 *preprint arXiv:2405.19107*, 2024.
- 511 R. Tyrrell Rockafellar. Monotone operators and the proximal point algorithm. *SIAM Journal*  
512 *on Control and Optimization*, 14(5):877–898, 1976. doi: 10.1137/0314056. URL <https://doi.org/10.1137/0314056>.
- 514 Corby Rosset, Ching-An Cheng, Arindam Mitra, Michael Santacrose, Ahmed Awadallah, and  
515 Tengyang Xie. Direct nash optimization: Teaching language models to self-improve with general  
516 preferences. *CoRR*, abs/2404.03715, 2024. doi: 10.48550/ARXIV.2404.03715. URL <https://doi.org/10.48550/arXiv.2404.03715>.
- 518 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy  
519 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 520 Gokul Swamy, Christoph Dann, Rahul Kidambi, Zhiwei Steven Wu, and Alekh Agarwal. A mini-  
521 maximalist approach to reinforcement learning from human feedback, 2024.
- 522 Vasilis Syrgkanis, Alekh Agarwal, Haipeng Luo, and Robert E Schapire. Fast convergence of  
523 regularized learning in games. *Advances in Neural Information Processing Systems (NeurIPS)*,  
524 2015.
- 525 Yunhao Tang, Zhaohan Daniel Guo, Zeyu Zheng, Daniele Calandriello, Remi Munos, Mark Rowland,  
526 Pierre Harvey Richemond, Michal Valko, Bernardo Avila Pires, and Bilal Piot. Generalized  
527 preference optimization: A unified approach to offline alignment. In *Forty-first International*  
528 *Conference on Machine Learning*, 2024.
- 529 Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada,  
530 Shengyi Huang, Leandro von Werra, Clémentine Fourrier, Nathan Habib, et al. Zephyr: Direct  
531 distillation of lm alignment. *arXiv preprint arXiv:2310.16944*, 2023.
- 532 Amos Tversky. Intransitivity of preferences. *Psychological review*, 76(1):31, 1969.
- 533 Chen-Yu Wei, Chung-Wei Lee, Mengxiao Zhang, and Haipeng Luo. Linear last-iterate convergence  
534 in constrained saddle-point optimization. In *International Conference on Learning Representations*  
535 *(ICLR)*, 2021.
- 536 Yue Wu, Zhiqing Sun, Huizhuo Yuan, Kaixuan Ji, Yiming Yang, and Quanquan Gu. Self-play  
537 preference optimization for language model alignment. *arXiv preprint arXiv:2405.00675*, 2024.
- 538 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li,  
539 Chengyuan Li, Dayiheng Liu, Fei Huang, et al. Qwen2 technical report. *arXiv preprint*  
540 *arXiv:2407.10671*, 2024.
- 541 Chenlu Ye, Wei Xiong, Yuheng Zhang, Nan Jiang, and Tong Zhang. A theoretical analysis of  
542 nash learning from human feedback under general kl-regularized preference. *arXiv preprint*  
543 *arXiv:2402.07314*, 2024.
- 544 Yuheng Zhang, Dian Yu, Baolin Peng, Linfeng Song, Ye Tian, Mingyue Huo, Nan Jiang, Haitao  
545 Mi, and Dong Yu. Iterative nash policy optimization: Aligning llms with general preferences via  
546 no-regret learning. *arXiv preprint arXiv:2407.00617*, 2024.

547	<b>Contents</b>	
548	<b>1 Introduction</b>	<b>1</b>
549	<b>2 Backgrounds</b>	<b>2</b>
550	2.1 Alignment with General Preference Models . . . . .	3
551	<b>3 Last-Iterate Nash Equilibrium Policy Optimization</b>	<b>3</b>
552	3.1 LINE-PO . . . . .	4
553	3.2 Solving a Regularized Game . . . . .	4
554	3.3 Computing the prox operator . . . . .	5
555	<b>4 Synthetic Experiments</b>	<b>5</b>
556	<b>5 Real-World Experiments</b>	<b>6</b>
557	5.1 Experimental Settings . . . . .	7
558	5.2 Result Analysis . . . . .	7
559	<b>6 Conclusion</b>	<b>8</b>
560	<b>A Alignment under the Bradley-Terry Model Assumption</b>	<b>14</b>
561	<b>B Related Work</b>	<b>15</b>
562	<b>C Properties of the Prox Operator</b>	<b>16</b>
563	<b>D Proof of Theorem 1</b>	<b>16</b>
564	<b>E Proof of Theorem 2</b>	<b>18</b>
565	<b>F Computing the Prox Operator using Preference Learning Methods</b>	<b>18</b>
566	<b>G Practical Implementation of Algorithms</b>	<b>20</b>
567	<b>H Implementation of Mirror-Prox and Optimistic Multiplicative Weights Update</b>	<b>20</b>
568	<b>I Detailed Discussion on Synthetic Experiments</b>	<b>21</b>

569 **A Alignment under the Bradley-Terry Model Assumption**

570 **RLHF** The canonical formulation of Reinforcement Learning from Human Feedback (RLHF) is to  
571 first learn a reward function  $r$  under the BT model and then find the optimal KL regularized policy  
572  $\pi^*$  with respect to the learned reward function  $r$ :

$$\pi^* := \arg \max_{\pi} \mathbb{E}_{x \sim \rho, y \sim \pi(\cdot|x)} [r(x, y) - \eta^{-1} \text{KL}(\pi(\cdot|x) || \pi_{\text{ref}}(\cdot|x))], \quad (4)$$

573 where  $\eta^{-1} > 0$  controls the regularization, and  $\pi_{\text{ref}}$  is the initial reference model, usually the policy  
574  $\pi_{\text{sft}}$  obtained from pre-training and supervised fine-tuning.

575 **DPO** Rafailov et al. [2024] observe that the regularized optimization problem (4) has a closed-form  
 576 solution : for any  $x$  and  $y$ ,

$$\pi^*(y | x) = \frac{\pi_{\text{ref}}(y | x) \exp(\eta r(x, y))}{Z_x}, \quad (5)$$

577 where  $Z_x = \mathbb{E}_{y \sim \pi_{\text{ref}}(\cdot | x)}[\exp(\frac{1}{\eta} r(y, x))]$  is the normalization constant known as the partition function.  
 578 In (5), we see that  $\pi^*$  implicitly parameterizes the reward function  $r$ . Rafailov et al. [2024] propose  
 579 direct preference optimization (DPO) to learn the optimal policy using the maximum likelihood  
 580 objective directly:

$$\ell_{\text{DPO}}(\pi; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[ \log \sigma \left( \eta^{-1} \log \frac{\pi(y_w | x)}{\pi_{\text{ref}}(y_w | x)} - \eta^{-1} \log \frac{\pi(y_l | x)}{\pi_{\text{ref}}(y_l | x)} \right) \right],$$

581 where  $\mathcal{D}$  is a data set containing win-loss pair of responses  $\{y_w, y_l\}$  given prompt  $x$ .

## 582 B Related Work

583 **Alignment under Preference models** Most existing approaches adopt the Bradley-Terry (BT)  
 584 preference model [Bradley and Terry, 1952, Christiano et al., 2017], which involves first learning a  
 585 preference model and then optimizing the objective function with a KL divergence penalty relative to  
 586 the original language model. For example, RLHF [Ouyang et al., 2022] aims to ensure that LLMs  
 587 follow instructions by initially learning a BT model and subsequently fine-tuning the model based on  
 588 the learned reward while regularizing it with the original LLM.

589 Building on this framework, Rafailov et al. [2024] introduces Direct Preference Optimization (DPO)  
 590 that maintains the assumption of the BT model for preferences but eliminates the preference learning  
 591 step by reformulating the objective and optimizing it directly. Additionally, Ethayarajh et al. [2024]  
 592 diverges from the traditional BT-based methods by deriving algorithms that bypass the preference  
 593 modeling step altogether. Instead, they model user preferences based on Kahneman and Tversky’s  
 594 utility theory.

595 **Alignment Solution Concepts under General Preferences** Azar et al. [2024] are the first to  
 596 consider general preferences and propose a family of optimization objectives that optimize a function  
 597 of the preferences probabilities regularized by the KL divergence with respect to the original model.  
 598 They propose the IPO algorithm, an offline algorithm that directly optimizes the win rate of the  
 599 model penalized by the KL divergence with respect to the original model. Munos et al. [2024]  
 600 also consider general preferences and aim to find the *von Neumann winner*, which corresponds to  
 601 the Nash equilibrium of a game played between the two LLMs over the win rate. They propose a  
 602 variant of the Mirror Descent (MD) algorithm called Nash-MD and show last-iterate convergence  
 603 in the KL-regularized game. Concurrently, Swamy et al. [2024] study the same solution concept  
 604 focusing more on sequential games. Calandriello et al. [2024] proved that the objective of the the  
 605 IPO algorithm coincides with the Nash policy under a proper choice of the parameter that controls  
 606 the regularization.

607 **Iterative Self-Play Algorithms** Apart from the aforementioned works, a line of recent work also  
 608 propose practical implementation of the Mirror Dscnt (MD) algorithms, which can be used to learn  
 609 the Nash equilibrium via self-play. Rosset et al. [2024] propose Direct Nash Optimization (DNO),  
 610 where at each iteration, the model regresses the predicted preferences against the actual preferences  
 611 using cross-entropy loss. Similarly, Wu et al. [2024] introduce the Self-Play Preference Optimization  
 612 (SPPO) method, Gao et al. [2024] introduce Reinforcement Learning via Regressing Relative Rewards  
 613 (REBEL), and Richemond et al. [2024] introduce the Direct Reward Optimization (DRO) which  
 614 regresses the loss using the  $L_2$  distance at each iteration. Since these algorithms simulate the MD  
 615 update, when applied in a (unregularized) zero-sum game, they only have average-iterate convergence  
 616 but all *diverge in last iterate*. Moreover, all these methods require the estimation of the win rate,  
 617 which can be computationally intensive and may introduce estimation errors.

618 Most closely related to our work is Iterative Nash Policy Optimization (INPO) by Zhang et al. [2024],  
 619 which continues to use  $L_2$  distance regression. However, by further reformulating and simplifying  
 620 the objective similar to IPO, INPO eliminates the need to estimate the expected win rate. The primary

621 distinction between our approach and INPO is that INPO is designed for the KL-regularized game  
622 and is equivalent to MD; while our algorithm LINE-PO is inspired by the Conceptual Prox algorithm  
623 and guarantees last-iterate convergence in the unregularized game. This fundamental difference  
624 allows LINE-PO to achieve more favourable convergence properties with robust alignment (i.e., 50%  
625 against any other policy) for large language models.

626 **Last-Iterate Convergence on Games** It is well-established that Mirror Descent fails to converge  
627 in simple zero-sum games, often resulting in cycling behavior [Mertikopoulos et al., 2018]. In  
628 contrast, several prominent algorithms have been shown to achieve last-iterate convergence including  
629 the Proximal Point (PP) method [Rockafellar, 1976], Extra-Gradient (EG) [Korpelevich, 1976],  
630 Optimistic Gradient Descent (OGD) [Popov, 1980, Rakhlin and Sridharan, 2013], and the Conceptual  
631 Prox/Mirror Prox methods [Nemirovski, 2004]. The asymptotic convergence properties of these  
632 algorithms have been extensively studied [Popov, 1980, Facchinei and Pang, 2003, Iusem et al., 2003,  
633 Nemirovski, 2004, Daskalakis and Panageas, 2018]. Recently, there has been a growing focus on  
634 establishing finite-time convergence guarantees for these methods, addressing the practical necessity  
635 of understanding their performance within a limited number of iterations (see e.g. [Mokhtari et al.,  
636 2020b,a, Golowich et al., 2020b,a, Bauschke et al., 2021, Wei et al., 2021, Cai et al., 2022, Gorbunov  
637 et al., 2022] and references therein).

## 638 C Properties of the Prox Operator

639 Recall that  $\text{Prox}(z, g) = \operatorname{argmax}_{z' \in \mathcal{Z}} \langle g, z' \rangle - D_\varphi(z' \| z) = \operatorname{argmax}_{z' \in \mathcal{Z}} \langle g + \nabla \varphi(z), z' \rangle - \varphi(z')$ .  
640 The following properties of the prox operator are well-known in the literature (e.g., [Nemirovski,  
641 2004])

642 **Lemma 1.**  $\text{Prox}(z, g) = z'$  if and only if  $\langle g + \nabla \varphi(z) - \nabla \varphi(z'), z' - z^* \rangle \geq 0$  for all  $z^* \in \mathcal{Z}$ .

643 **Corollary 1.** Let  $\text{Prox}(z, g) = z'$ , then

$$\langle g, z^* - z' \rangle \leq D_\varphi(z^* \| z) - D_\varphi(z^* \| z') - D_\varphi(z' \| z), \quad \forall z^* \in \mathcal{Z}$$

## 644 D Proof of Theorem 1

645 The proof of Theorem 1 is relatively standard in the literature [Facchinei and Pang, 2003, Nemirovski,  
646 2004]. We include a formal proof here for completeness. In Theorem 1, we make the following  
647 assumption.

648 **Assumption 1.** We assume there exists a Nash equilibrium  $\pi^*$  such that  $\operatorname{supp}(\pi^*) = \operatorname{supp}(\pi_{\text{sft}})$ .

649 This assumption is mild and **much weaker** than the ‘‘Bounded Log Density’’ assumptions used in  
650 previous works [Rosset et al., 2024, Zhang et al., 2024], which requires  $|\log \frac{\pi^t}{\pi_{\text{sft}}}|$  is bounded.

651 Recall that  $\Pi := \{\pi : \operatorname{supp}(\pi) \subseteq \operatorname{supp}(\pi_{\text{sft}})\}$ . Then  $\text{KL}(\pi \| \pi_{\text{sft}}) \leq D :=$   
652  $\max_{y: \pi_{\text{sft}}(y) > 0} \log \pi_{\text{sft}}(y)$  is bounded for any  $\pi \in \Pi$ . We first prove  $\text{KL}(\pi^* \| \pi^{t+1}) \leq \text{KL}(\pi^* \| \pi^t)$   
653 for any  $t \geq 1$ .

654 In the proof, we assume that each step of LINE-PO,  $\pi^{t+1} \leftarrow \operatorname{argmax}_{\pi_1} \min_{\pi_2} J_\tau(\pi_1, \pi_2, \pi_{\text{ref}})$  can  
655 be solved exactly. Our proof extends to the case the optimization problem is solved approximately  
656 with sufficient accuracy.

657 **Lemma 2.** Let  $\pi^*$  be a Nash equilibrium of  $J(\pi_1, \pi_2)$ . Then for any  $\tau > 0$ , if

$$(\pi, \pi) = \operatorname{argmax}_{\pi_1} \operatorname{argmin}_{\pi_2} J_\tau(\pi_1, \pi_2, \pi_{\text{ref}}),$$

658 then

$$\text{KL}(\pi^* \| \pi) \leq \text{KL}(\pi^* \| \pi_{\text{ref}}) - \text{KL}(\pi \| \pi_{\text{ref}})$$



659 *Proof.* By definition of the prox operator, we have

$$\begin{aligned}
\pi &= \operatorname{argmax}_{\pi_1} J_\tau(\pi_1, \pi, \pi_{\text{ref}}) \\
&= \operatorname{argmax}_{\pi_1} \mathbb{P}(\pi_1 \succ \pi) - \tau \operatorname{KL}(\pi_1, \pi_{\text{ref}}) \\
&= \operatorname{Prox}(\pi_{\text{ref}}, \frac{1}{\tau} \mathbb{P}(\cdot \succ \pi)).
\end{aligned} \tag{6}$$

660 Using [Corollary 1](#), we have for any  $\pi' \in \Pi$ ,

$$\frac{1}{\tau} (\mathbb{P}(\pi' \succ \pi) - \mathbb{P}(\pi \succ \pi)) \leq \operatorname{KL}(\pi' || \pi_{\text{ref}}) - \operatorname{KL}(\pi' || \pi) - \operatorname{KL}(\pi || \pi_{\text{ref}}). \tag{7}$$

661 Plugging  $\pi' = \pi^*$  into the above inequality and noting that  $\mathbb{P}(\pi \succ \pi) = \frac{1}{2}$ , we get

$$\frac{1}{\tau} \left( \mathbb{P}(\pi^* \succ \pi) - \frac{1}{2} \right) \leq \operatorname{KL}(\pi^* || \pi_{\text{ref}}) - \operatorname{KL}(\pi^* || \pi) - \operatorname{KL}(\pi || \pi_{\text{ref}}).$$

662 Since  $\pi^*$  is a Nash equilibrium and thus  $\mathbb{P}(\pi^* \succ \pi) \geq \frac{1}{2}$ , the lefthand side of the above inequality is  
663  $\geq 0$ . The we have

$$\operatorname{KL}(\pi^* || \pi) \leq \operatorname{KL}(\pi^* || \pi_{\text{ref}}) - \operatorname{KL}(\pi || \pi_{\text{ref}}).$$

664

□

665 [Lemma 2](#) implies the following properties on the trajectory  $\{\pi^t\}$ .

666 **Corollary 2.** *In LINE-PO, we have*

- 667 1.  $\operatorname{KL}(\pi^* || \pi^{t+1}) \leq \operatorname{KL}(\pi^* || \pi^t)$  for all  $t \geq 1$ .
- 668 2.  $\sum_{t=1}^{\infty} \operatorname{KL}(\pi^{t+1} || \pi^t) \leq \operatorname{KL}(\pi^* || \pi_{\text{sft}}) < +\infty$ .
- 669 3.  $\operatorname{supp}(\pi^t) = \operatorname{supp}(\pi_{\text{sft}})$  for all  $t \geq 1$ .

670 *Proof.* The first item is direct from [Lemma 2](#). The second item is also direct by applying [Lemma 2](#)  
671 for  $t = 1, 2, \dots$ :

$$\sum_{t=1}^{\infty} \operatorname{KL}(\pi^{t+1} || \pi^t) \leq \sum_{t=1}^{\infty} \operatorname{KL}(\pi^* || \pi^t) - \operatorname{KL}(\pi^* || \pi^{t+1}) \leq \operatorname{KL}(\pi^* || \pi_{\text{sft}}) \leq D < \infty.$$

672 For the third item, let  $\pi^*$  be a Nash equilibrium such that  $\operatorname{supp}(\pi^*) = \operatorname{supp}(\pi_{\text{sft}})$  as guaranteed  
673 by [Assumption 1](#). On one hand, since  $\operatorname{KL}(\pi^t || \pi^{t-1}) < \infty$  for all  $t \geq 1$ , we have  $\operatorname{supp}(\pi^t) \subseteq$   
674  $\operatorname{supp}(\pi^{t-1}) \subseteq \dots \subseteq \operatorname{supp}(\pi_{\text{sft}})$ . On the other hand,  $\operatorname{KL}(\pi^* || \pi^t) < \infty$  implies  $\operatorname{supp}(\pi^*) \subseteq$   
675  $\operatorname{supp}(\pi^t)$ . Since  $\operatorname{supp}(\pi_{\text{sft}}) = \operatorname{supp}(\pi^*)$ , we have  $\operatorname{supp}(\pi^t) = \operatorname{supp}(\pi_{\text{sft}}) = \operatorname{supp}(\pi^*)$ . □

676 Since the sequence  $\{\pi^t\}$  is bounded (all lies in the simplex), it has at least one limit point  $\hat{\pi}$ . The  
677 next lemma shows that a limit point must be a Nash equilibrium.

678 **Lemma 3.** *If  $\hat{\pi}$  is a limit point of  $\{\pi^t\}$ , then  $\hat{\pi}$  is a Nash equilibrium of  $J(\pi_1, \pi_2)$ .*

679 *Proof.* By item 2 in [Corollary 2](#), we have  $\lim_{t \rightarrow \infty} \operatorname{KL}(\pi^{t+1} || \pi^t) = 0$ . This implies  
680  $\lim_{t \rightarrow \infty} \|\pi^{t+1} - \pi^t\| = 0$ . As  $\hat{\pi}$  is a limit point of  $\{\pi^t\}$ , we let  $\{\pi^{k_i} : k_i \in \kappa\}$  be the subsequence  
681 that converges to  $\hat{\pi}$ . Then by [Equation \(6\)](#), we have

$$\begin{aligned}
\lim_{k \in \kappa, k \rightarrow \infty} \pi^{k+1} &= \lim_{k \in \kappa, k \rightarrow \infty} \operatorname{Prox}(\pi^k, \frac{1}{\tau} \mathbb{P}(\cdot \succ \pi^{k+1})) \\
&\Rightarrow \hat{\pi} = \operatorname{Prox}(\hat{\pi}, \frac{1}{\tau} \mathbb{P}(\cdot \succ \hat{\pi})).
\end{aligned}$$

682 Thus  $\hat{\pi}$  is a fixed point of  $\operatorname{Prox}(\pi, \frac{1}{\tau} \mathbb{P}(\cdot \succ \pi))$ . Moreover, by item 3 in [Corollary 2](#), we have  
683  $\operatorname{supp}(\hat{\pi}) = \operatorname{supp}(\pi_{\text{sft}})$ . Now consider both the max and min player running MWU initialized with  
684  $\pi^1 = \hat{\pi}$ . Then we have  $\pi^t = \hat{\pi}$  for all  $t \geq 1$ . By [Equation \(7\)](#), we have for any  $\pi' \in \Pi$ ,

$$\frac{1}{\tau} \sum_{t=1}^{\infty} \left( \mathbb{P}(\pi' \succ \hat{\pi}) - \frac{1}{2} \right) \leq \operatorname{KL}(\pi' || \hat{\pi}) < \infty,$$

685 where the inequality holds since  $\text{supp}(\pi') \subseteq \text{supp}(\hat{\pi})$ . As a result, we get

$$\mathbb{P}(\pi' \succ \hat{\pi}) \leq \frac{1}{2}, \forall \pi' \in \Pi \Leftrightarrow \mathbb{P}(\hat{\pi} \succ \pi') \geq \frac{1}{2}, \forall \pi' \in \Pi$$

686 Thus  $\hat{\pi}$  is a Nash equilibrium of  $J(\pi_1, \pi_2)$ .  $\square$

687 **Proof of Theorem 1** Since  $\hat{\pi}$  is a Nash equilibrium, by [Corollary 2](#),  $\{\text{KL}(\hat{\pi}||\pi^t) \geq 0\}$  is a  
688 decreasing sequence. Thus  $\{\text{KL}(\hat{\pi}||\pi^t)\}$  converges. As a result,

$$\lim_{t \rightarrow \infty} \text{KL}(\hat{\pi}||\pi^t) = \lim_{k \in \kappa, k \rightarrow \infty} \text{KL}(\hat{\pi}||\pi^k) = \text{KL}(\hat{\pi}||\hat{\pi}) = 0.$$

689 Thus we have  $\lim_{t \rightarrow \infty} \pi^t = \hat{\pi}$  is a Nash equilibrium. This completed the proof of [Theorem 1](#).

## 690 E Proof of Theorem 2

691 We show that MWU has linear convergence to the unique Nash equilibrium of a KL-regularized  
692 zero-sum game  $J(\pi_1, \pi_2, \pi_{\text{ref}})$ .

693 We denote  $\mu^* = \pi_{\tau}^*$  the unique Nash equilibrium of the KL regularized game  $J_{\tau}(\pi_1, \pi_2, \pi_{\text{ref}})$ . We  
694 note that  $J(\pi_1, \pi_2)$  is 1-smooth. We then can adapt [[Abe et al., 2024, Lemma F.1](#)] to our setting.

695 **Lemma 4** (Adapted from Lemma F.1 in [Abe et al. \[2024\]](#)). *If we choose  $\eta \in (0, \frac{2\tau}{3\tau^2+8}]$ , then we  
696 have for every  $k \geq 1$*

$$\text{KL}(\mu^*, \mu^{k+1}) \leq (1 - \frac{\eta\tau}{2}) \text{KL}(\mu^*, \mu^k).$$

697 Applying the lemma recursively implies  $\text{KL}(\mu^*||\mu^{k+1}) \leq (1 - \frac{\eta\tau}{2})^k \text{KL}(\mu^*||\pi_{\text{ref}})$  and completes  
698 the proof.

## 699 F Computing the Prox Operator using Preference Learning Methods

700 **Reinforcement Learning algorithms** We can use the Proximal Policy Optimization (PPO) algo-  
701 rithm [[Schulman et al., 2017](#)] to solve  $\text{Prox}(\pi, \eta g)$ . Observe that

$$\begin{aligned} \text{Prox}(\pi, \eta g) &= \underset{\pi'}{\text{argmax}} \{ \langle \eta g, \pi' \rangle - \text{KL}(\pi' || \pi) \} \\ &= \underset{\pi'}{\text{argmax}} \mathbb{E}_{y \sim \pi'} [g[y] - \eta^{-1} \cdot \text{KL}(\pi' || \pi)] \end{aligned}$$

702 shares the same form as the objective in (4). Typically, we parameterize  $\pi' = \pi_{\theta}$  with neural network  
703 parameters  $\theta$  and optimize over  $\theta$ .

704 **Loss minimization algorithms** Let us denote  $\hat{\pi}$  the prox operator  $\text{Prox}(\pi, \eta g)$ , then we have

$$\hat{\pi}[y] = \frac{\pi(y) \exp(\eta g(y))}{Z} \Leftrightarrow \log \frac{\hat{\pi}(y)}{\pi(y)} - \eta g(y) + \log Z = 0,$$

705 where  $Z = \mathbb{E}_{y \sim \pi}[\exp(\eta g(y))]$  is the partition function. We can directly compute the partition  
706 function  $Z$  and thus  $\hat{\pi}$  in small tabular cases. However, the partition function is hard to compute in  
707 general large-scale applications. Several works have recently proposed to solve the above equality by  
708 optimizing the corresponding  $L_2$  loss. Specifically, the Self-Play Preference Optimization (SPPO)  
709 loss [[Wu et al., 2024](#)] assumes  $\log Z = \frac{\eta}{2}$  and optimizes

$$\ell_{\text{SPPO}}(\theta) = \left( \log \frac{\pi_{\theta}(y)}{\pi(y)} - \eta g(y) - \frac{\eta}{2} \right)^2.$$

710 The Direct Reward Optimization (DRO) loss [[Richemond et al., 2024](#)] parameterizes both  $\hat{\pi}$  and  
711  $\log Z$  with  $\theta$  and  $V_{\phi}$  respectively and optimize<sup>5</sup>

$$\ell_{\text{DRO}}(\theta, \phi) = \left( \log \frac{\pi_{\theta}(y)}{\pi(y)} - \eta g(y) - \eta V_{\phi} \right)^2.$$

<sup>5</sup>we modified some constants in the original DRO loss to make it consistent with our presentation. The modification has no other effects.

712 The REBEL loss [Gao et al., 2024] uses *differences in rewards* to eliminate the partition function  $Z$   
 713 and optimize the regression loss

$$\ell_{\text{REBEL}}(\theta) = \left( \eta^{-1} \left( \log \frac{\pi_{\theta}(y)}{\pi(y)} - \log \frac{\pi_{\theta}(y')}{\pi(y')} \right) - (g(y) - g(y')) \right)^2.$$

714 All the above approaches can be used to solve  $\text{Prox}(\pi, \eta g)$ . However, directly applying them  
 715 iteratively on  $J(\pi_1, \pi_2)$  is equivalent to running MWU, which provably diverges. In contrast, we can  
 716 apply them in Algorithm 2 and then apply our meta-algorithm LINE-PO to guarantee convergence to  
 717 a Nash equilibrium with robust alignment.

718 **Remark 1.** *The above approaches are versatile and work well for any  $g$  that can be evaluated*  
 719 *efficiently. In particular, we should consider using them when (1)  $g = r$  is a reward function and*  
 720 *we can efficiently query  $r$ ; (2)  $g = \mathbb{P}(\cdot | \mu)$  is the win rate against a reference policy  $\mu$ , and we can*  
 721 *efficiently sample from  $\mu$  and have oracle access to  $\mathbb{P}$ . These two settings are popular and practical in*  
 722 *the LLM alignment setting.*

723 Now we turn attention to the more specific setting where  $g$  corresponds to a preference model  $\mathbb{P}$   
 724 (could be a BT model or a general preference) and that we can collect a win-loss preference data set  
 725  $\mathcal{D} = \{(y_w, y_l)\}$ , which is standard for LLM alignment. Although the abovementioned algorithms  
 726 apply, they all require estimating  $g$  (the win rate) and may be inefficient in practice. In the following,  
 727 we present algorithms directly working on the sampled dataset  $\mathcal{D}$  without further estimation.

728 **Sampled loss based on the BT preference model** Assume  $g = r$  is the reward of the Bradley-  
 729 Terry model, and the dataset  $\{(y_w, y_l)\}$  consists of win-lose pairs of responses. Then we can solve  
 730  $\text{Prox}(\pi, \eta g)$  by optimize the DPO loss [Rafailov et al., 2024] defined as

$$\ell_{\text{DPO}}((y_w, y_l); \theta) = -\log \sigma \left( \eta^{-1} \log \frac{\pi_{\theta}(y_w)}{\pi(y_w)} - \eta^{-1} \log \frac{\pi_{\theta}(y_l)}{\pi(y_l)} \right).$$

731 **Sampled loss for general preference** The DPO loss inspires many other loss functions that work  
 732 under even weaker assumptions on the preference model. Now, we assume a general preference  
 733 model  $\mathbb{P}$  over  $\mathcal{Y}$  (not necessarily the BT model). We assume  $g$  is the win-rate against some policy  
 734  $\mu$  such that  $g_{\mu}(y) = \mathbb{P}[y \succ \mu] := \mathbb{E}_{y' \sim \mu} [\mathbb{P}[y \succ y']]$  (think of  $\mu$  as the reference policy  $\pi_{\text{ref}}$  or other  
 735 online policy  $\pi_t$ ). We assume the dataset contains win-lose pairs sampled from  $\mu$ :  $\{y_w, y_l \sim \mu\}$ . We  
 736 denote the preference distribution  $\lambda_{\mathbb{P}}(y, y')$  as a binary distribution:

$$\lambda_{\mathbb{P}}(y, y') = \begin{cases} (y, y') & \text{with probability } \mathbb{P}[y \succ y'] \\ (y', y) & \text{with probability } 1 - \mathbb{P}[y \succ y'] \end{cases}$$

737 The (population) IPO loss [Tang et al., 2024, Calandriello et al., 2024] is defined as

$$\ell_{\text{IPO}}(\theta, \mu) := \mathbb{E}_{(y_w, y_l) \sim \mu, (y^+, y^-) \sim \lambda_{\mathbb{P}}(y_w, y_l)} \left[ \left( \log \frac{\pi_{\theta}(y^+)}{\pi(y^+)} - \log \frac{\pi_{\theta}(y^-)}{\pi(y^-)} - \frac{\eta}{2} \right)^2 \right].$$

738 It has been proved that the minimizer of the  $\ell_{\text{IPO}}(\theta, \mu)$  satisfies

$$\pi_{\theta}(y) \propto \pi(y) \exp(-\eta \mathbb{P}[y \succ \mu]) \Leftrightarrow \pi_{\theta} = \text{Prox}(\pi, \eta g_{\mu}).$$

739 Thus we can compute the prox operator  $\text{Prox}(\pi, \eta g_{\mu})$  where  $g_{\mu} = \mathbb{P}(\cdot \succ \mu)$  by minimizing the IPO  
 740 loss against policy  $\mu$ .

741 A variant of the IPO loss applied to the regularized preference setting is the Iterative Nash Policy  
 742 Optimization (INPO) loss [Zhang et al., 2024]. Here, we define  $g_{\mu}^{\tau}$  the gradient  $\nabla_{\pi} \mathcal{J}_{\tau}(\pi, \mu, \pi_{\text{ref}}) =$   
 743  $\mathbb{P}(\cdot \succ \mu) - \tau \log \frac{\mu(\cdot)}{\pi_{\text{ref}}(\cdot)}$  of the regularized objective. The corresponding INPO loss is

$$\ell_{\text{INPO}} := \mathbb{E}_{(y_w, y_l) \sim \mu, (y^+, y^-) \sim \lambda_{\mathbb{P}}(y_w, y_l)} \left[ \left( \log \frac{\pi_{\theta}(y^+)}{\pi_{\theta}(y^-)} - \eta \tau \log \frac{\pi_{\text{ref}}(y^+)}{\pi_{\text{ref}}(y^-)} - (1 - \eta \tau) \log \frac{\mu(y^+)}{\mu(y^-)} - \frac{\eta}{2} \right)^2 \right].$$

744 Similarly, it has been shown that the INPO loss minimizer corresponds to the prox operator’s solution  
 745  $\text{Prox}(\pi, \eta g_{\mu}^{\tau})$ . Thus we can use the INPO in Algorithm 2 directly.

## 746 G Practical Implementation of Algorithms

747 We present an implementation of LINE-PO using the INPO [Zhang et al., 2024] as a subgame solver  
 748 here. We remark that LINE-PO can also be implemented using PPO or many other preference learning  
 749 algorithms, as we show in Section 3.3. Given the implementation of these existing methods, our  
 750 meta-algorithm requires minimal change but archives last-iterate convergence to Nash equilibrium  
 with robust alignment.

---

### Algorithm 3: Practical Implementation of LINE-PO integrated with INPO (Algorithm 4)

---

**Input:** Initial policy  $\pi_{\text{sft}}$ , regularization  $\{\tau_t > 0\}$ , step size  $\{\eta_t > 0\}$ , number of iterations  
 $T \geq 1$ , number of inner optimization steps  $\{K_t \geq 1\}$ , preference oracle  $\mathbb{P}$ .

```

1 Initialize  $\pi^1, \pi_{\text{ref}} \leftarrow \pi_{\text{sft}}$ 
2 for  $t = 1, 2, \dots, T - 1$  do
3    $\pi^{t+1} \leftarrow \text{INPO}(\pi_{\text{ref}}, \tau_t, \eta_t, K_t, \mathbb{P})$ 
4    $\pi_{\text{ref}} \leftarrow \pi^{t+1}$ 
5 return  $\pi^T$ 

```

---

751

---

### Algorithm 4: INPO [Zhang et al., 2024]

---

**Input:** Reference policy  $\pi_{\text{ref}}$ , regularization  $\tau > 0$ , step size  $\eta > 0$ , number of rounds  $K \geq 1$ ,  
 preference oracle  $\mathbb{P}$ .

```

1 Initialize  $\mu^1 \leftarrow \pi_{\text{ref}}$ 
2 for  $k = 1, 2, \dots, K - 1$  do
3   Generate response pairs  $\{y_1^{(i)}, y_2^{(i)}\}_{i=1}^n$  where  $y_1^{(i)}, y_2^{(i)} \sim \mu^k$ 
4   Query preference oracle  $\mathbb{P}$  to get preference data  $\mathcal{D}_k = \{y_w^{(i)}, y_l^{(i)}\}_{i=1}^n$ 
5   Compute  $\mu^{k+1}$  as
      
$$\mu^{k+1} = \underset{\pi \in \Pi}{\operatorname{argmin}} \mathbb{E}_{(y_w, y_l) \sim \mathcal{D}_k} \ell_{\text{INPO}}(\pi)$$

      
$$\ell_{\text{INPO}}(\pi) := \mathbb{E}_{(y^+, y^-) \sim \lambda_{\mathbb{P}}(y_w, y_l)} \left[ \left( \log \frac{\pi(y^+)}{\pi(y^-)} - \eta\tau \log \frac{\pi_{\text{ref}}(y^+)}{\pi_{\text{ref}}(y^-)} - (1 - \eta\tau) \log \frac{\mu^k(y^+)}{\mu^k(y^-)} - \frac{\eta}{2} \right)^2 \right]$$

6    $\mu^{k+1} = \underset{\pi \in \Pi}{\operatorname{argmin}} \mathbb{E}_{(y_w, y_l) \sim \mathcal{D}_k} \ell_{\text{INPO}}(\pi)$ 
7 return  $\mu^K$ 

```

---

## 752 H Implementation of Mirror-Prox and Optimistic Multiplicative Weights 753 Update

754 We note that there are other algorithms that has provable last-iterate convergence to Nash equilibrium  
 755 in (unregularized) zero-sum games, including the Mirror-Prox algorithm [Nemirovski, 2004] and  
 756 Optimistic Multiplicative Weights Update (OMWU) algorithm [Rakhlin and Sridharan, 2013, Syrgka-  
 757 nis et al., 2015, Hsieh et al., 2021]. We present practical implementations of these two algorithms  
 758 in the context of LLM alignment for solving  $J(\pi_1, \pi_2)$  (1), where we use preference optimization  
 759 algorithms to solve the prox operator as shown in Section 3.3.

760 We denote the gradient  $g(\pi) := \mathbb{P}(\cdot \succ \pi)$ .

761 **Mirror-Prox** The Mirror-Prox algorithm [Nemirovski, 2004] initialized  $\pi^1 = \pi_{\text{sft}}$  and updates in  
 762 each iteration  $t \geq 1$ :

$$\begin{aligned} \pi^{t+\frac{1}{2}} &= \text{Prox}(\pi^t, \eta g(\pi^t)) \\ \pi^{t+1} &= \text{Prox}(\pi^t, \eta g(\pi^{t+\frac{1}{2}})) \end{aligned}$$

763 As we have shown in Section 3.3, we can implement Mirror-Prox using  
 764 PPO/DPO/IPO/SPPO/DRO/REBEL to compute the prox operator. Specifically, we could

765 sample from  $\pi^t$  and construct a preference dataset  $D_t$  and optimize certain regression loss  
 766 (IPO/DRO/REBEL) to compute  $\pi^{t+\frac{1}{2}} = \text{Prox}(\pi^t, \eta g(\pi^t))$ . The procedure applies to the second step  
 767 in each iteration. Thus we require two sampling and two optimization procedure in each iteration.

768 **Optimistic Multiplicative Weights Update (OMWU)** The OMWU algorithm [Rakhlin and Srid-  
 769 haran, 2013] is an optimistic variant of the MWU algorithm. Although MWU diverges in zero-sum  
 770 games, it has been shown that OMWU has last-iterate convergence to Nash equilibrium [Wei et al.,  
 771 2021, Hsieh et al., 2021]. Initialized with  $\pi^1 = \pi^{\frac{1}{2}} = \pi_{\text{sft}}$ , OMWU updates in each iteration  $t \geq 1$ :

$$\begin{aligned}\pi^{t+\frac{1}{2}} &= \text{Prox}(\pi^t, \eta g(\pi^{t-\frac{1}{2}})) \\ \pi^{t+1} &= \text{Prox}(\pi^t, \eta g(\pi^{t+\frac{1}{2}}))\end{aligned}$$

772 Similarly, we can implement OMWU to solve  $J(\pi_1, \pi_2)$  using preference methods to compute the  
 773 prox operator as shown in Section 3.3. Moreover, OMWU has an equivalent update rule: initialize  
 774  $\pi^1 = \pi^0 = \pi_{\text{sft}}$

$$\pi^{t+1} = \text{Prox}(\pi^t, 2\eta g(\pi^t) - \eta g(\pi^{t-1})),$$

775 which requires computing only one prox operator in each iteration.

776 We leave testing the practical performance of Mirror-Prox and OMWU for large-scale applications  
 777 including LLM alignment as future works.

## 778 I Detailed Discussion on Synthetic Experiments

779 The sample-only setting is also more aligned with the practice. We use sufficient samples in each  
 780 iteration for every algorithm. As a result, the LINE-PO performs the same as in the noiseless gradient  
 781 setting, while the iterative IPO algorithm becomes equivalent to the MD algorithm. We present the  
 782 results in Figure 2 and noting that summarize the results below.

- 783 • Iterative DPO: We observe that iterative DPO diverges and cycles between extreme policies  
 784 (e.g., outputting  $y_a$  with probability close to 1). This is aligned with [Azar et al., 2024],  
 785 where they found DPO will converge to the deterministic policy regardless of the regulariza-  
 786 tion parameter in extreme preference settings. The cycling behavior of iterative DPO may  
 787 be explained as follows: in each iteration, DPO converges to a nearly deterministic policy  
 788 output  $y$ ; then the new preference data shows that  $y' \neq y$  is more preferred; finally, iterative  
 789 DPO cycles over  $\mathcal{Y}$  since the preference itself exhibits a cycle and there is no clear winner.
- 790 • Iterative IPO [Azar et al., 2024, Calandriello et al., 2024]: The IPO loss is a variant of  
 791 the DPO loss, but it does not rely on the BT model assumption and works for a general  
 792 preference model. However, as we have discussed before, (exactly) minimizing the IPO loss  
 793 is equivalent to performing one mirror descent step, and thus, iterative IPO is equivalent to  
 794 mirror descent up to sampling error. As a result, we observe that iterative IPO also exhibits  
 795 cycling behavior.
- 796 • SPPO [Wu et al., 2024]: The SPPO algorithm is not exactly the same as MWU since  
 797 SPPO assumes the partition function is always  $Z = \log \frac{\eta}{2}$  which may not be the case. We  
 798 observe that SPPO exhibits very similar cycling behavior as MD. We conclude that SPPO  
 799 approximates MD very well in this instance and exhibits similar behavior.
- 800 • INPO [Zhang et al., 2024]: The INPO algorithm is designed for finding the Nash equilibrium  
 801 of the KL regularized game  $J_\tau(\pi_1, \pi_2, \pi_{\text{ref}})$ . As we proved in Theorem 2, INPO does not  
 802 diverge but exhibits last-iterate convergence. However, it converges to a regularized Nash  
 803 equilibrium without the robust alignment property.