
Robust Preference Optimization through Reward Model Distillation

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Language model (LM) post-training (or alignment) involves maximizing a
2 reward function that is derived from preference annotations. Direct Preference
3 Optimization (DPO) is a popular offline alignment method that trains a policy
4 directly on preference data without the need to train a reward model or apply
5 reinforcement learning. However, typical preference datasets have only a single, or
6 at most a few, annotation per preference pair, which causes DPO to overconfidently
7 assign rewards that trend towards infinite magnitude. This frequently leads to
8 degenerate policies, sometimes causing even the probabilities of the *preferred*
9 generations to go to zero. In this work, we analyze this phenomenon and propose
10 *distillation* to get a better proxy for the true preference distribution over generation
11 pairs: we train the LM to produce probabilities that match the distribution induced
12 by a reward model trained on the preference data. Moreover, to account for
13 uncertainty in the reward model we are distilling from, we optimize against a
14 *family of reward models* that, as a whole, is likely to include at least one reasonable
15 proxy for the preference distribution. Our results show that distilling from such
16 a family of reward models leads to improved robustness to distribution shift in
17 preference annotations, while preserving the simple supervised nature of DPO.

18 1 Introduction

19 Language model (LM) post-training (or alignment) aims to steer language model policies towards
20 responses that agree with human preferences. Early state-of-the-art approaches have focused on
21 reward learning from human feedback. In this paradigm, preference annotations are used to train
22 reward models, which then guide the optimization of the language model policy through online
23 reinforcement learning (an approach broadly referred to as RLHF). Recent research on offline “Direct
24 Preference Optimization” [DPO; 23] and extensions thereof [3; 31], however, has demonstrated that
25 it is also possible to directly optimize policies on the preference data, which bypasses the need for a
26 separate reward model—and its offline nature also leads to faster, and simpler, training frameworks.

27 While this direct approach to preference optimization is attractive in terms of its simplicity and
28 efficiency, it also raises important questions about the effectiveness and robustness of the resulting
29 policies—as well as the broader utility of using an explicit reward model. In this paper, we argue that
30 explicit reward modeling can, in fact, offer substantial practical advantages that are not captured by
31 DPO’s formulation. In particular, we theoretically show that relying solely on the preference data
32 can be a precarious strategy, with few natural brakes in place to prevent policies trained under the
33 DPO objective from careening off towards degenerate policies when the preference data exhibits
34 certain idiosyncratic properties. On the other hand, explicit reward models can easily be regularized
35 and understood—regardless of whether they are Bradley-Terry models [4], margin-based ranking
36 models [40], or simply any other kind of function that correlates well with human preferences [31; 17].

37 Taking a step back from pure direct preference optimization, we propose a method that merges the
 38 best of both worlds: an efficient reward model distillation algorithm that (i) operates effectively in the
 39 offline setting, (ii) makes minimal assumptions about the true, optimal reward we aim to maximize,
 40 and (iii) demonstrates greater robustness to the specific distribution of prompt/response data used for
 41 policy alignment. Drawing inspiration from prior knowledge distillation techniques [14; 26; 35; 10],
 42 we leverage the same change of variables trick employed in DPO to express the language model
 43 policy in terms of its implicit reward model [23]. We then train the policy to match our desired,
 44 explicit reward via an L_2 loss that directly regresses the pairwise differences in target rewards for
 45 any two generation pairs (x, y_1) and (x, y_2) . We theoretically establish the equivalence between
 46 optimizing this distillation loss over a sufficiently diverse offline dataset of unlabeled examples and
 47 optimizing the traditional online RLHF objective.

48 Our reward model distillation approach, however, is not immune to some of the same challenges
 49 facing DPO-style learning of policies. In particular, reward model distillation requires having a
 50 reliable reward model—but having a reliable reward requires having a reliable method for extracting
 51 a reward model from a potentially noisy preference dataset. To address the uncertainty surrounding
 52 the “right” reward model, we introduce a pessimistic extension to our approach. This extension aims
 53 to maximize the worst-case improvement of our model across a plausible family of reward models
 54 (e.g., those sufficiently consistent with annotated preference data). This strategy aligns with that of
 55 existing work in conservative offline reinforcement learning [5; 16]. Interestingly, we derive that
 56 this pessimistic objective can be equivalently expressed and optimized by adding a simple additional
 57 KL-divergence regularization to the original distillation objective.

58 Empirically, we find that reward model distillation, particularly pessimistic reward model distillation,
 59 leads to similar performance to prior direct preference optimization methods in settings where the
 60 preference datasets used are unbiased, but significantly better performance in settings where the
 61 preference datasets are biased, when compared to DPO and the Identity Preference Optimization
 62 (IPO) framework of [3], which was introduced as a more robust alternative to DPO. To further support
 63 these empirical observations, we provide an extensive theoretical analysis that both (i) sheds more
 64 light on the degenerative tendencies of DPO and issues inherent to its objective, and (ii) highlights
 65 relative advantages of our explicitly regularized approaches.

66 2 Preliminaries

67 We begin with a brief review of Direct Preference Optimization (DPO) [23] and its analysis. Proofs
 68 of all theoretical results provided here, and in the rest of the paper, are deferred to Appendix A.

69 2.1 The preference alignment problem

70 Let x be an input prompt, and let $y \sim \pi_\theta(\cdot | x)$ be the language model policy π_θ ’s response to x .
 71 Given some reward function $r^*(x, y)$ and another reference policy $\pi_{\text{ref}}(y | x)$, the goal of alignment
 72 is to solve for the “aligned” policy $\pi_{\theta^*}(y | x)$ that maximizes the following RLHF objective, i.e.,

$$\pi_{\theta^*}(y | x) = \operatorname{argmax}_{\pi_\theta} \mathbb{E}_{\mu(x)} [\mathbb{E}_{\pi_\theta(y|x)} [r^*(x, y)] - \beta \mathbb{D}_{\text{KL}}[\pi_\theta(\cdot | x) \| \pi_{\text{ref}}(\cdot | x)]] , \quad (1)$$

73 where $\mu(x)$ is a fixed distribution over prompts, and the KL-divergence term prevents the aligned
 74 policy from being dramatically different from the anchoring reference policy, $\pi_{\text{ref}}(y | x)$. Here,
 75 the reward function r^* is typically not known in advance, but rather inferred from collected human
 76 preference data in the form of (x, y^w, y^ℓ) , where x is the prompt, y^w is the “winning”, or preferred,
 77 response, and y^ℓ is the “losing”, or dispreferred, response. A common approach is to assume that
 78 pairs (y_1, y_2) follow a Bradley-Terry model [4], under which the probability that y_1 is preferred to y_2
 79 given the reward function r^* and prompt x is $p^*(y_1 \succ y_2 | x) = \sigma(r^*(x, y_1) - r^*(x, y_2))$, where
 80 $\sigma(\cdot)$ is the sigmoid function and \succ denotes preference. Under this model, we can use the preference
 81 data $(x, y^w, y^\ell) \sim \mathcal{D}_{\text{pref}}$ to estimate r^* via maximum likelihood estimation, i.e.,

$$\hat{r} \in \operatorname{argmin}_r \mathbb{E}_{(y^w, y^\ell, x) \sim \mathcal{D}_{\text{pref}}} [-\log \sigma(r(x, y^w) - r_\phi(x, y^\ell))] . \quad (2)$$

82 With \hat{r} in hand, Eq. (1) can be optimized using standard reinforcement learning algorithms [27; 29; 6].

83 **2.2 Direct preference optimization**

84 DPO is a simple approach for offline policy optimization that uses preferences to directly align the
 85 language model policy, without training an intermediate reward model. Specifically, DPO leverages
 86 the fact that the optimal solution to the KL-constrained objective in (1) takes the form [15]

$$\pi_{\theta^*}(y | x) = \frac{1}{Z(x)} \pi_{\text{ref}}(y | x) \exp\left(\frac{1}{\beta} r^*(x, y)\right), \quad (3)$$

87 where $Z(x) = \sum_y \pi_{\text{ref}}(y | x) \exp(\frac{1}{\beta} r^*(x, y))$ is the partition function. DPO reparameterizes the
 88 true reward function r^* in terms of the optimal policy π_{θ^*} that it induces, i.e.,

$$r^*(x, y) = \beta \log\left(\frac{\pi_{\theta^*}(y | x)}{\pi_{\text{ref}}(y | x)}\right) + \beta \log Z(x). \quad (4)$$

89 Under the Bradley-Terry model, the likelihood that $y_1 \succ y_2$ can then be written as

$$p^*(y_1 \succ y_2 | x) = \sigma\left(\beta \log \frac{\pi_{\theta^*}(y_1) \pi_{\text{ref}}(y_2)}{\pi_{\theta^*}(y_2) \pi_{\text{ref}}(y_1)}\right), \quad (5)$$

90 where now π_{θ^*} can be directly estimated on $\mathcal{D}_{\text{pref}}$ following the objective in (2), in place of the
 91 intermediate reward model \hat{r} , i.e., $\pi_{\hat{\theta}}(y | x) \in \operatorname{argmin}_{\pi_{\theta}} \mathcal{L}_{\text{dpo}}(\pi_{\theta}; \mathcal{D}_{\text{pref}})$ where

$$\mathcal{L}_{\text{dpo}}(\pi_{\theta}; \mathcal{D}_{\text{pref}}) = \mathbb{E}_{(y^w, y^\ell, x) \sim \mathcal{D}_{\text{pref}}} \left[-\log \sigma\left(\beta \log \frac{\pi_{\theta^*}(y^w) \pi_{\text{ref}}(y^\ell)}{\pi_{\theta^*}(y^\ell) \pi_{\text{ref}}(y^w)}\right) \right]. \quad (6)$$

92 **2.3 Pitfalls of direct preference optimization**

93 As argued in [3], the Bradley-Terry assumption that DPO strongly relies on for maximum likelihood
 94 estimation is sensitive to the underlying preference data. Specifically, if we have any two responses y_1
 95 and y_2 where $p^*(y_1 \succ y_2 | x) = 1$, then the Bradley-Terry model dictates that $r^*(y_1) - r^*(y_2) = +\infty$,
 96 and therefore $\pi_{\theta^*}(y_2 | x) = 0$ for *any* finite KL-regularization strength β .

97 We can illustrate this phenomenon on a broader level with the following example.

98 **Assumption 1.** *Suppose we are given a preference dataset of (context-free) pairs $\mathcal{D}_{\text{pref}} =$
 99 $\{(y_i^w, y_i^\ell)\}_{i=1}^n$, the pairs (y_i^w, y_i^ℓ) are mutually disjoint in both the elements. Further suppose
 100 that we optimize the DPO objective on $\mathcal{D}_{\text{pref}}$ with a single parameter θ_y for each y .*

101 **Proposition 1.** *Under Assumption 1, for any (y, y') such that $y = y_i^w$ and $y' = y_i^\ell$ for some i , we
 102 have $\frac{\pi_{\theta^*}(y) \pi_{\text{ref}}(y')}{\pi_{\theta^*}(y') \pi_{\text{ref}}(y)} \rightarrow \infty$, for all global minimizers π_{θ^*} of the DPO objective in (6), for any $\beta > 0$.*

103 **Corollary 1.** *Under Assumption 1, further assume that $0 < \pi_{\text{ref}}(y) < 1$ for all y . Then π_{θ^*} is a
 104 global minimizer of the DPO objective in (6) iff $\pi_{\theta^*}(\mathcal{C}(y^\ell)^c) \rightarrow 1$ with $\pi_{\theta^*}(y_i^w) > 0 \forall i \in [n]$, where
 105 $\mathcal{C}(y^\ell)^c$ is the complement of the set of all responses y that appear as a dispreferred y_i^ℓ for any $i \in [n]$.*

106 Additional analysis of the training dynamics of DPO is also provided in §5. A significant, and non-
 107 obvious, implication of Corollary 1 is that the set of global optima of the DPO loss also includes poli-
 108 cies that can shift nearly all probability mass to responses that never even appear in the training set—
 109 and even assign near zero probability to all of the training data responses that do in fact correspond to
 110 winning generations, y^w , a phenomenon that has been observed empirically [e.g., 20]. Stated differ-
 111 ently, Corollary 1 implies that any θ^* merely satisfying $\pi_{\theta^*}(y_i^\ell) = 0$ with $\pi_{\theta^*}(y_i^w) > 0 \forall i \in [n]$ is a
 112 global minimizer of the DPO objective in this setting. Though simplistic, the scenario in Assumption 1
 113 is closer to reality than might first be appreciated: in many practical situations we can almost always
 114 expect the finite-sample preference data to contain one (or at most a few) preference annotations per
 115 example (x, y_1, y_2) , while the policies π_{θ} can have billions of parameters ($\gg n$). Of course, this issue
 116 can also be viewed as a classic instance of overfitting—with the additional caveat that as opposed to
 117 *overpredicting* responses within the training set, we might overfit to *almost never* producing anything
 118 like the “good” responses that do appear within the training set. Furthermore, without additional regu-
 119 larization (beyond β), we can expect this degeneration to easily happen in typical preference datasets.

120 3 Uncertainty-aware reward model distillation

121 As discussed in the previous section, a core issue in preference optimization is that the true preference
 122 distribution $p^*(y_1 \succ y_2 | x)$ is not known. Attempting to infer it from finite-sample preference data
 123 (that may further be biased or out-of-distribution with respect to the target domain) can then result
 124 in a failure to learn reasonable policies. In this section, we now propose an inherently regularized
 125 approach to direct preference optimization that uses uncertainty-aware reward model distillation.

126 3.1 Reward model distillation

127 Suppose for the moment that the reward function r^* was in fact known, and did not have to be
 128 inferred from sampled preference data. Under this setting, we can then define an efficient offline
 129 optimization procedure that is similar in spirit to DPO, but no longer relies directly on a preference
 130 dataset. Concretely, given unlabeled samples $(x, y_1, y_2) \sim \rho$ (where the number of samples can be
 131 potentially unlimited), we can define a simple ‘‘distillation’’ loss, $\mathcal{L}_{\text{distill}}(r^*, \pi_\theta)$, as follows:

$$\mathcal{L}_{\text{distill}}(r^*, \pi_\theta; \rho) = \mathbb{E}_{\rho(x, y_1, y_2)} \left[\left(r^*(x, y_1) - r^*(x, y_2) - \beta \log \frac{\pi_\theta(y_1 | x) \pi_{\text{ref}}(y_2 | x)}{\pi_\theta(y_2 | x) \pi_{\text{ref}}(y_1 | x)} \right)^2 \right]. \quad (7)$$

132 Intuitively, the distillation loss seeks to exactly match *differences* in reward model scores across
 133 all generation pairs (x, y_1, y_2) . It is then easy to see that under the Bradley-Terry model, this is
 134 equivalent to matching the strength of the preference relationship, $y_1 \succ y_2$. Furthermore, by only
 135 matching differences, we can still conveniently ignore the log partition term, $\log Z(x)$, in the implicit
 136 reward formulation for π_θ as shown in (4), as it is constant across different y for any given x . Finally,
 137 similar to the motivation in DPO, we can show that minimizing $\mathcal{L}_{\text{distill}}(r^*, \pi_\theta; \rho)$ indeed results in an
 138 optimally aligned policy π_{θ^*} , as long as the data distribution ρ has sufficient support.

139 **Theorem 1.** *Let \mathcal{Y} denote the set of all possible responses for any model π_θ . Assume that*
 140 *$\text{supp}(\pi_{\text{ref}}(y | x)) = \mathcal{Y}$, i.e., the reference policy may generate any outcome with non-zero probability.*
 141 *Further, let $\text{supp}(\rho(x, y_1, y_2)) = \text{supp}(\mu(x)) \times \mathcal{Y} \times \mathcal{Y}$. Let $\pi_{\theta^*}(y | x) \in \text{argmin}_{\pi_\theta} \mathcal{L}_{\text{distill}}(r^*, \pi_\theta; \rho)$*
 142 *be a minimizer over all possible policies, of the implicit reward distillation loss in (7), for which*
 143 *$r^*(x, y)$ is assumed to be deterministic, and finite everywhere. Then for any $\beta > 0$, π_{θ^*} also*
 144 *maximizes the alignment objective in (1).*

145 The above result holds for a broad class of data distributions $\rho(x, y_1, y_2)$, and makes no assumptions
 146 on r^* (e.g., it is no longer necessary for it to be defined using a Bradley-Terry model). In fact, this
 147 result can also be seen as strict generalization of the IPO framework of [3] when taking $r^*(x, y) \triangleq$
 148 $\mathbf{1}\{y = y_w\}$, if labeled pairs (x, y_w, y_l) are provided instead of the unlabeled pairs (x, y_1, y_2) .

149 Of course, the true reward r^* is usually not known in practice. Still, as in standard RLHF, we can
 150 go about constructing good proxies by using the preference data to identify plausible target reward
 151 models r_{tgt} —further guided by any amount of regularization and inductive bias that we desire. A
 152 natural choice is to first learn r_{tgt} on the preference data $\mathcal{D}_{\text{pref}}$ using standard methods, and then reuse
 153 $\mathcal{D}_{\text{pref}}$ to distill π_θ , which is similar to classical settings in teacher-based model distillation [14; 26].
 154 Furthermore, as r_{tgt} is a real-valued model, at a bare minimum it is guaranteed to induce a regularized
 155 Bradley-Terry preference distribution $p_{\text{tgt}}(y_1 \succ y_2 | x) > 0, \forall x, y_1, y_2 \in \mathcal{X} \times \mathcal{Y}$, and thereby avoid
 156 some of the degeneracies identified in §2.3 for the maximum likelihood estimate under DPO.

157 3.2 Pessimistic reward model distillation

158 Choosing a single reward model r_{tgt} for anchoring the LM policy can naturally still lead to degenerate
 159 behavior if r_{tgt} is a poor approximation of the true r^* that accurately reflects human preferences.
 160 However, we can easily extend our framework to handle uncertainty in the right target reward function
 161 by defining a confidence *set* of $k \geq 1$ plausible target reward models, $\mathcal{S} = \{r_{\text{tgt}}^1, \dots, r_{\text{tgt}}^k\}$, and
 162 training $\pi_{\theta^*}(y | x)$ to maximize the following ‘‘pessimistic’’ form of the objective in (1):

$$\max_{\pi_\theta} \min_{r_{\text{tgt}}^i \in \mathcal{S}} \mathbb{E}_{\mu(x)} \left[\underbrace{\mathbb{E}_{\pi_\theta(y|x)} [r_{\text{tgt}}^i(x, y)] - \mathbb{E}_{\pi_{\text{ref}}(y|x)} [r_{\text{tgt}}^i(x, y)]}_{\text{advantage over baseline policy}} - \beta \mathbb{D}_{\text{KL}}(\pi_\theta(\cdot | x) \| \pi_{\text{ref}}(\cdot | x)) \right]. \quad (8)$$

163 In this pessimistic objective we are no longer op-
 164 timizing π_θ for a single reward, but optimizing
 165 π_θ to produce generations that are scored favor-
 166 ably on average, even by the worst-case reward
 167 model in the set \mathcal{S} , relative to the generations of
 168 the baseline policy π_{ref} . When the set $\mathcal{S} = \{r^*\}$
 169 consists of only the ground-truth reward, the ob-
 170 jective (8) is equivalent to standard RLHF (1),
 171 up to a constant offset independent of θ . More
 172 generally, whenever \mathcal{S} includes a good proxy
 173 \tilde{r} for r^* , the pessimistic advantage evaluation
 174 ensures that the policy π_θ^* that maximizes
 175 eq. (8) still has a large advantage over π_{ref} under
 176 all $r \in \mathcal{S}$, including \tilde{r} . This use of pessimism
 177 to handle uncertainty in the knowledge of the
 178 true reward is related to similar techniques in
 179 the offline RL literature [16; 5].

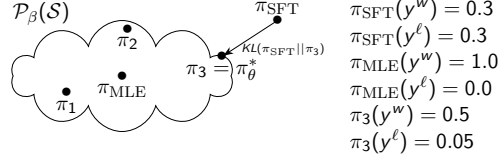


Figure 1: A toy illustration of Theorem 2, which states that the optimal π_{θ^*} for (8) is the policy in $\mathcal{P}_\beta(\mathcal{S})$ with the lowest forward-KL from π_{SFT} . The set $\mathcal{P}_\beta(\mathcal{S})$ contains a (potentially infinite) set of policies π_1, π_2, \dots corresponding to target reward models. Here, π_{SFT} assigns equal mass to target reward models y^w and y^l , π_{MLE} is the MLE solution for the DPO objective, which puts all probability mass on y^w , and π_3 is the policy in $\mathcal{P}_\beta(\mathcal{S})$ with lowest forward-KL.

180 For the objective to be meaningful, the set \mathcal{S} has to be chosen carefully. When \mathcal{S} is small, it might
 181 not include any good proxy for r^* . Conversely, if \mathcal{S} is too rich, it forces π_{θ^*} to be nearly identical to
 182 π_{ref} , since any deviations from π_{ref} might be penalized by some reward model in \mathcal{S} . Consequently,
 183 we want to design \mathcal{S} to be the smallest possible set which contains a reasonable approximation to r^* .

184 To optimize (8), it turns out that we can formulate it as an equivalent constrained offline optimization
 185 problem, that we will show to conveniently admit a similar loss form as (7).

186 **Theorem 2** (Pessimistic distillation). *Define the constrained minimizer*

$$\pi_{\theta^*}(y | x) \in \underset{\pi_\theta \in \mathcal{P}_\beta(\mathcal{S})}{\operatorname{argmin}} \beta \mathbb{E}_{\mu(x)} \mathbb{D}_{\text{KL}}(\pi_{\text{ref}}(\cdot | x) \| \pi_\theta(\cdot | x)), \quad (9)$$

187 where $\mathcal{P}_\beta(\mathcal{S})$ is the set of all possible policies with implicit reward models that are consistent with
 188 any target reward model $r_{\text{tgt}}^i \in \mathcal{S}$, i.e., $\mathcal{P}_\beta(\mathcal{S}) \triangleq \{\pi_{\theta_i}\}_{i=1}^{|\mathcal{S}|}$ where $\pi_{\theta_i} \propto \pi_{\text{ref}}(y | x) \exp \frac{1}{\beta} r_{\text{tgt}}^i(x, y)$.
 189 Then for any $\beta > 0$, π_{θ^*} also maximizes the pessimistic alignment objective in (8).

190 To unpack this result, Theorem 2 stipulates that the π_θ that maximizes the pessimistic objective in (8)
 191 is the policy in $\mathcal{P}_\beta(\mathcal{S})$ that is closest in *forward* KL-divergence to π_{ref} (see Figure 1).¹ In addition,
 192 this policy also maximizes the expected reward of one of the $r_{\text{tgt}}^i \in \mathcal{S}$ (minus the additional weighted
 193 reverse KL-divergence penalty term). Intuitively, the forward KL-divergence term serves the role of
 194 biasing the model towards optimizing for reward models that are similar to the implicit reward that
 195 π_{ref} already maximizes. Otherwise, there might exist a target reward model $r_{\text{tgt}}^i \in \mathcal{S}$ for which the
 196 advantage of π_θ relative to π_{ref} will be low, or even negative (a solution that we would like to avoid).

197 3.2.1 Optimization

198 The constraint in (9) can then be relaxed and approximately optimized by introducing an objective
 199 with a Lagrangian-style penalty with strength $\alpha > 0$ on a form of distillation loss as (7), i.e.,

$$\min_{\pi_\theta} \beta \mathbb{E}_{\mu(x)} \mathbb{D}_{\text{KL}}(\pi_{\text{ref}}(y | x) \| \pi_\theta(y | x)) + \alpha \min_{r_{\text{tgt}}^i \in \mathcal{S}} \mathcal{L}_{\text{distill}}(r_{\text{tgt}}^i, \pi_\theta; \rho), \quad (10)$$

200 where in practice we divide by α and instead optimize²

$$\mathcal{L}_{\text{pdistill}}(\mathcal{S}, \pi_\theta; \rho) = \min_{r_{\text{tgt}}^i \in \mathcal{S}} \mathcal{L}_{\text{distill}}(r_{\text{tgt}}^i, \pi_\theta; \rho) + \gamma \mathbb{E}_{\mu(x)} \mathbb{D}_{\text{KL}}(\pi_{\text{ref}}(\cdot | x) \| \pi_\theta(\cdot | x)), \quad (11)$$

201 where $\gamma = \beta\alpha^{-1}$. In reality, minimizing (11) for $\gamma > 0$ is equivalent to solving the constrained
 202 optimization problem in (9) with an implicitly larger set of possible reward models $\mathcal{S}_\gamma \supseteq \mathcal{S}$ indexed
 203 by γ . More specifically, \mathcal{S}_γ also contains all reward models \tilde{r} that are approximately consistent with
 204 the anchoring reward models r_{tgt}^i contained in \mathcal{S} , as the following result states.

¹Note that the objective in (9) minimizes the *forward* KL-divergence $\mathbb{D}_{\text{KL}}(\pi_{\text{ref}}(\cdot | x) \| \pi_\theta(\cdot | x))$ even though the pessimistic objective in (8) is regularized with *reverse* KL-divergence $\mathbb{D}_{\text{KL}}(\pi_\theta(\cdot | x) \| \pi_{\text{ref}}(\cdot | x))$.

²In practice, we compute and optimize the min over reward models per each mini-batch of examples.

205 **Proposition 2** (Soft pessimistic distillation). *Assume the same conditions as Theorem 1. Then for*
 206 *any $0 < \gamma < \infty$, there exists a $\lambda \geq 0$ such that $\pi_{\theta^*}(y | x) \in \operatorname{argmin}_{\pi_{\theta}} \mathcal{L}_{\text{pdistill}}(\mathcal{S}, \pi_{\theta}; \rho)$, where π_{θ^*}*
 207 *is a minimizer over all possible policies, is a solution to (9) for the effective reward model set*

$$\mathcal{S}_{\gamma} = \bigcup_{r_{\text{tgt}}^i \in \mathcal{S}} \left\{ \tilde{r}: \mathbb{E}_{\rho(x, y_1, y_2)} \left[(r_{\text{tgt}}^i(x, y_1) - r_{\text{tgt}}^i(x, y_2) - \tilde{r}(x, y_1) + \tilde{r}(x, y_2))^2 \right] \leq \lambda \right\}. \quad (12)$$

208 As a result, optimizing (11) even when using the singleton $\mathcal{S} = \{r_{\text{tgt}}\}$ yields an implicitly pessimistic
 209 objective, in which the pessimism is over all reward models \tilde{r} that are consistent up to λ with r_{tgt} .

210 3.3 Pessimistic DPO

211 We can also observe that Proposition 2 can be leveraged to obtain an alternative, implicitly pessimistic,
 212 objective that uses DPO directly instead of distillation. Consider the following regularized DPO loss:

$$\mathcal{L}_{\text{pdpo}}(\pi_{\theta}; \mathcal{D}_{\text{pref}}) = \mathcal{L}_{\text{dpo}}(\pi_{\theta}; \mathcal{D}_{\text{pref}}) + \gamma \mathbb{E}_{\mu(x)} \mathbb{D}_{\text{KL}}(\pi_{\text{ref}}(y | x) \| \pi_{\theta}(y | x)). \quad (13)$$

213 Following a similar analysis as in Proposition 2, we can derive that this implicitly corresponds to
 214 maximizing the pessimistic objective in (8) for the reward model set

$$\mathcal{S}_{\gamma} = \left\{ r_{\pi_{\theta}} : \mathcal{L}_{\text{dpo}}(\pi_{\theta}; \mathcal{D}_{\text{pref}}) \leq \min_{\pi'_{\theta}} \mathcal{L}_{\text{dpo}}(\pi'_{\theta}; \mathcal{D}_{\text{pref}}) + \lambda \right\}, \quad (14)$$

215 where $r_{\pi_{\theta}}(x, y) \triangleq \beta \log \pi_{\theta}(y | x) / \pi_{\text{ref}}(y | x) + \beta \log Z(x)$ is the implicit reward model defined by
 216 π_{θ} . \mathcal{S}_{γ} then corresponds to the set of reward models $r_{\pi_{\theta}}$ that are all approximate minimizers of the
 217 DPO loss. This not only includes the MLE, but also all other estimators that obtain nearly the same
 218 loss. In principle, this can be expected to help ameliorate some of the issues of §2.3: since driving the
 219 reward to $\pm\infty$ only marginally decreases the \mathcal{L}_{dpo} loss past a certain point, the set \mathcal{S} will also include
 220 finite reward functions $|r_{\pi_{\theta}}(x, y)| < \infty$ for any $\gamma > 0$. These rewards would then be preferred if they
 221 induce a policy with a smaller (forward) KL-divergence to π_{ref} than the degenerate, infinite rewards.

222 4 Experimental results

223 The main motivation for reward distillation and pessimism is to increase alignment robustness
 224 in challenging settings where it is difficult to learn good policies directly from the preference
 225 data. To demonstrate the effectiveness of our approach, we run experiments on the popular TL;DR
 226 summarization task [29; 32], in which we simulate a scenario where the preference data has a spurious
 227 correlation between the *length* of a summary and whether or not it is preferred.³

228 4.1 Experimental setup

229 We first train an “oracle” reward model on the TL;DR preference data training set [29] and relabel
 230 all preference pairs with this oracle. This enables us to use the oracle reward model for evaluation,
 231 without worrying about the gap to true human preferences. After relabeling, longer responses (where
 232 longer is defined as y_1 having at least 10% more tokens than y_2) are preferred in 61% of the examples.

233 To test the effect of a spurious correlation on preference-based policy optimization, we select as a
 234 training set 30K examples from the relabeled data such that the longer output is preferred in ρ fraction
 235 of examples, with $\rho \in \{0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8\}$. Each such training set is denoted \mathcal{D}_{ρ} . At
 236 each \mathcal{D}_{ρ} , we compare our approach to DPO [23] and IPO [3], which are currently the most commonly
 237 used offline alignment methods. We test the following variants of distillation and pessimism:

- 238 • **Distilled DPO** (d-DPO): Trains a reward model r_{ρ} on \mathcal{D}_{ρ} , and then optimizes $\mathcal{L}_{\text{distill}}(r_{\rho}, \pi_{\theta}; \rho)$.
- 239 • **Pessimistic DPO** (p-DPO): A pessimistic version of DPO as described in §3.3, trained on \mathcal{D}_{ρ} .
- 240 • **Pessimistic Distilled DPO** (pd-DPO): Combines the above two by training a reward model r_{ρ} on
 241 \mathcal{D}_{ρ} and optimizing the pessimistic distillation objective (Eq. (11)) with confidence set $\mathcal{S} = \{r_{\text{tgt}}\}$.
- 242 • **Pessimistic Ensemble DPO** (e-DPO): To create ensembles of reward models, we subsample from
 243 each \mathcal{D}_{ρ} five preference datasets, $\mathcal{D}_{\rho, b}$, at $b \in \mathcal{B} = \{0.2, 0.4, 0.5, 0.6, 0.8\}$, such that the fraction

³Length has been repeatedly shown in the past to correlate with reward [28; 21].

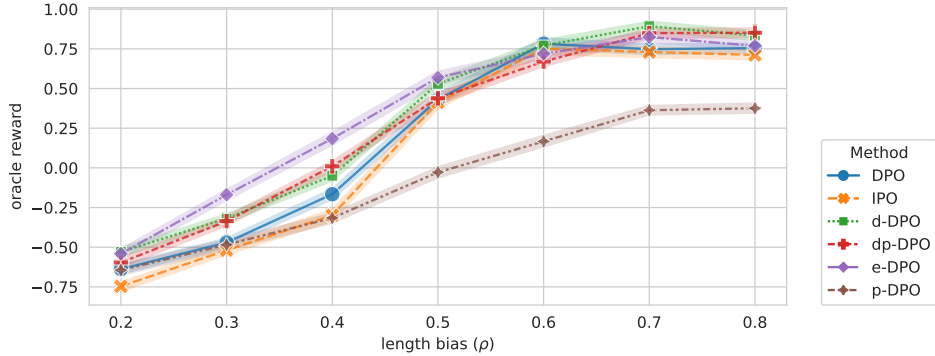


Figure 2: **Main results**, showing the advantage in oracle reward compared to the initial finetuned policy. Errorbars correspond to bootstrap 95% confidence intervals for finite sample variance. Ensemble DPO (e-DPO) is significantly better than DPO and IPO in the challenging setup where shorter responses are preferred ($\rho \leq 0.5$), and is generally the best-performing method overall in this regime. Distilled DPO (d-DPO) performs best when longer responses are preferred ($\rho > 0.6$).

244 of pairs where the longer response is preferred is b , and train reward models $r_{\rho,b}$ on those subsets.
 245 Consequently, sensitivity to length should vary across ensemble members. We then apply the
 246 same procedure as pd-DPO above, with a confidence set $\mathcal{S}_\rho = \{r_{\rho,b}\}_{b=1}^B$.

247 All reward models and policies are initialized from Palm-2-XS [2]. Policies also go through a
 248 supervised finetuning step on human-written summaries from the original TL;DR training set [32]
 249 prior to alignment, and we term this policy π_{SFT} . We evaluate performance by sampling summaries
 250 for test set prompts, evaluating the average reward according to the oracle reward model, and
 251 computing the advantage in average reward compared to π_{SFT} (before alignment). We train policies
 252 for 10^4 steps with batch size 16 and learning rate 10^{-6} , and reward models for $3k$ steps with
 253 batch size 64 and learning rate 4×10^{-6} . We use the validation set for model selection during
 254 policy training and to choose the following hyperparameters. For all DPO variants, we sweep over
 255 $\beta \in \{.01, .1, 1, 3, 10, 30, 100\}$. For IPO, we sweep over $\tau \in \{0.01, 0.1, 1, 3, 5, 10, 25\}$. For all
 256 pessimistic methods we anneal $\gamma = \alpha/\beta$ from 10^{-4} to 10^{-2} linearly during the $10k$ training steps.

257 4.2 Results

258 We present the results of our experiment in Figure 2. As can be seen in the plot, the more challenging
 259 setting is when $\rho < 0.5$, which corresponds to a sample of preference annotations in which shorter
 260 outputs are generally preferred. This distribution shift is more difficult because as mentioned the oracle
 261 reward model (trained on human annotations) has a bias in favor of longer outputs [28]. Nevertheless
 262 we get sizable improvements compared to the reference policy π_{SFT} for all length bias values.

263 All approaches that invoke distillation (d-DPO, e-DPO, dp-DPO) outperform IPO and DPO ($p < .01$
 264 by a Wald test) for $\rho \leq 0.5$, where shorter responses are preferred. Pessimistic ensemble DPO
 265 (e-DPO) performs particularly well in these settings, generally outperforming all methods that use
 266 a single reward model. When longer responses are preferred ($\rho > 0.6$), single reward distillation
 267 (d-DPO) leads to the highest performance, significantly outperforming both DPO and IPO ($p < .01$
 268 by a Wald test). Interestingly, p-DPO does not provide empirical benefits relative to the distillation
 269 based methods, indicating that the distillation loss itself is quite important. For the effect of
 270 hyper-parameter selection, see Figure D.1. In DPO-based methods, the optimal value of β is inversely
 271 correlated with the bias; in IPO the same holds for the τ hyperparameter.

272 To better understand the utility of reward ensembles in e-DPO, in particular when $\rho < 0.5$, we
 273 examine the role of each reward model in the ensemble across different biases. Specifically, given
 274 the final e-DPO policy per length bias, for each example we identify the reward model $r_{\rho,b}$ that
 275 best matches the implicit reward of this policy, i.e., for which reward model is $\mathcal{L}_{\text{distill}}$ minimized on
 276 that example (see Eq. (7) and (11)). We find that when the policy is trained on data where shorter
 277 preference are preferred ($\rho < .5$), the reward model that best matches the policy often has the opposite
 278 bias (b is high), and vice versa. Thus, the success of e-DPO may be explained by its ability to distill
 279 from reward models that do not suffer from the bias in the policy training data, which is particularly

280 helpful when $\rho \leq .5$ as this bias is also not shared by the oracle RM. We provide the full distribution
 281 over reward models for all ρ and β in App. C. Overall, these results demonstrate the efficacy of
 282 training a policy by distilling from a reward model in the presence of distribution shifts, and that a
 283 careful design of an ensemble to mitigate spurious correlations can lead to further performance gains.⁴

284 5 Theoretical analysis

285 This section characterizes problems with the DPO objective and solutions offered by pessimistic DPO
 286 and distillation, focusing on the simplified scenario in which we optimize with respect to a single
 287 preference pairs (y^w, y^ℓ) . Once again, all proofs are deferred to Appendix A.

288 In its Lagrangian formulation, pessimistic DPO adds a forward KL term to the DPO objective (§3.3).
 289 For the sake of analysis, we assume that the preference annotations are sampled from the reference
 290 distribution, $\mu(x) \times \pi_{\text{ref}}(y | x) \times \pi_{\text{ref}}(y | x)$. Then a finite-sample approximation of the forward
 291 KL term is $\hat{\Omega}(\Theta) := \sum_{(y^w, y^\ell) \in \mathcal{D}_{\text{Pref}}} -(\log \pi_\theta(y^\ell) + \log \pi_\theta(y^w))$. By applying this finite-sample
 292 approximation, *p-DPO has a finite optimum, unlike DPO*, as shown in Proposition 1. Note that this
 293 analysis is limited in two ways: (1) as mentioned, we compute the KL term over the completions
 294 in the preference data; (2) we directly optimize the probability ratios $\psi_w = \pi_\theta(y^w)/\pi_{\text{ref}}(y^w)$ and
 295 $\psi_\ell = \pi_\theta(y^\ell)/\pi_{\text{ref}}(y^\ell)$, rather than optimizing them jointly through the parameters. For sufficiently ex-
 296 pressive π_θ , however, this approximation captures the behavior of the two algorithms reasonably well.

297 **Proposition 3.** *Let $\hat{\mathcal{L}}_{\text{pdpo}}$ represent a finite-sample approximation to $\mathcal{L}_{\text{pdpo}}$ with the empir-*
 298 *ical forward KL term $\hat{\Omega}(\Theta)$. For a fixed $\hat{\pi}_\theta(y_i^w)$ and $\alpha > 1$, the $\text{argmin}_{\pi_\theta(y_i^\ell)} \hat{\mathcal{L}}_{\text{pdpo}}$ is*
 299 $\min(1 - \hat{\pi}_\theta(y_i^w), \hat{\pi}_\theta(y_i^\ell))$, *with $\log \hat{\pi}_\theta(y_i^\ell) = -\frac{1}{\beta} \log(\alpha - 1) + \log \hat{\pi}_\theta(y_i^w) + \log \frac{\pi_{\text{ref}}(y_i^\ell)}{\pi_{\text{ref}}(y_i^w)}$.*

300 The optimum in Proposition 3 corresponds to $\log \psi_w / \psi_\ell = \beta^{-1} \log(\alpha - 1)$. Recall that IPO seeks
 301 to assign a constant value to this ratio by minimizing $(\log \frac{\psi_w}{\psi_\ell} - \tau^{-1})^2$; the (unconstrained) optima
 302 are identical for $\tau^{-1} := \beta^{-1} \log(\alpha - 1)$, but the loss surfaces are different (see Appendix B). DPO
 303 sets $\pi_\theta(y_i^\ell) \rightarrow 0$, as shown in Corollary 1; this is due not only to competition from $\pi_\theta(y_i^w)$ but from
 304 DPO penalizing positive probability on y_i^ℓ . Analysis of the distilled loss gives a similar result:

305 **Proposition 4.** *For any fixed $\hat{\pi}_\theta(y_i^w)$ and $\beta > 0$, the argmin of the distilled DPO objective (eq. (7))*
 306 *is $\min(1 - \hat{\pi}_\theta(y_i^w), \hat{\pi}_\theta(y_i^\ell))$, with $\log \hat{\pi}_\theta(y_i^\ell) = \frac{1}{\beta}(r_t(x, y_i^\ell) - r_t(x, y_i^w)) + \log \hat{\pi}_\theta(y_i^w) + \log \frac{\pi_{\text{ref}}(y_i^\ell)}{\pi_{\text{ref}}(y_i^w)}$.*

307 While the setting is simplistic, the results are comforting: here the additional regularization effects of
 308 both distillation and pessimism (in the case of p-DPO) clearly help to avoid degenerate optima.

309 **Why DPO can drive $\pi(y^w)$ to zero.** In §2.3 we pointed out a peculiarity of the DPO global optima:
 310 in certain cases, it can include policies where $\pi(y^w)$ may be nearly 0 for all y^w in the training set. This
 311 undesirable behavior has also been observed in practice [20; 22; 30]. For intuition on why this may
 312 happen, consider the simplified case where the policy is a bag-of-words model, $\pi_\theta(y) \propto \exp(c(y) \cdot \theta)$
 313 for $c(y)$ representing a vector of counts in y and θ_i representing the unnormalized log-probability of
 314 token i . Then we can formally show that DPO optimization monotonically decreases an upper bound
 315 on the probability of the *preferred* completion, $\tilde{\pi}_{\theta^{(t-1)}}(y^w) \geq \tilde{\pi}_{\theta^{(t)}}(y^w) \geq \pi_{\theta^{(t)}}(y^w)$.

316 **Proposition 5.** *Let $y^w, y^\ell \in \mathcal{V}^n$ be preferred vs. dispreferred outputs of length n , with*
 317 $\pi_{\text{ref}}(y^w), \pi_{\text{ref}}(y^\ell) > 0$ *and corresponding count vectors $c(y^w), c(y^\ell)$. Let $\log \pi_\theta(y) = c(y) \cdot \theta -$*
 318 $nZ(\theta)$ *for $Z(\theta) = \log \sum_i e^{\theta_i}$, with upper bound $\log \tilde{\pi}_\theta(y) = c(y) \cdot \theta - n \max_j \theta_j$. Let $\theta^{(t)}$ represent*
 319 *the parameters of π after t steps of gradient descent on $\mathcal{L}_{\text{dpo}}(\{y^\ell, y^w, x\})$, with $\theta^{(0)} = 0$. Then*
 320 $\pi_{\theta^{(t)}}(y^w) \leq \tilde{\pi}_{\theta^{(t)}}(y^w) \leq \tilde{\pi}_{\theta^{(t-1)}}(y^w)$ *for all t .*

321 **Where does the probability mass go?** If $\pi_{\theta^{(t)}}(y^w)$ decreases in t , what other strings become
 322 more probable? In the following proposition, we show that under the bag-of-words model, DPO
 323 optimization moves probability mass away from y^w to sequences that contain only the tokens that
 324 maximize the difference between y^w and y^ℓ . This is a concrete example of the type of undesirable
 325 optima described in §2.3, now shown here to be realizable.

⁴We also experimented with an ensemble where members are different checkpoints across training of a
 reward model on the preference data and did not observe any empirical gains from this form of ensemble.

326 **Proposition 6.** *Let y^w and y^ℓ be preferred / dispreferred outputs of length n . Let $\Delta = c(y^w) - c(y^\ell)$
 327 *be the difference in unigram counts. Let $\hat{y} = [i, i, \dots, i]$, for $i \in \arg \max \Delta$, with $\|c(\hat{y})\|_1 = n$.
 328 *Then $\pi_{\theta(t)}(y^w) - \pi_{\theta(t)}(\hat{y}) = \tau(t)k$ for some $k \leq 0$ and some non-decreasing $\tau : \mathbb{Z}_+ \rightarrow \mathbb{R}_+$.***

329 We have $k = 0$ when $c(y^w) = c(\hat{y})$, and $k \ll 0$ when $\|c(y^w)\|_2 \ll \|c(\hat{y})\|_2 = n$ (dense $c(y^w)$) and
 330 $\|\Delta\|_2 = \|\Delta\|_\infty$ (sparse Δ). This implies that when y^w and y^ℓ are similar, $\pi_\theta(y^w)$ will degrade more
 331 rapidly. Early stopping will therefore tradeoff between reaching the degenerate solution on such
 332 cases, and underfitting other cases in which y^w and y^ℓ are more distinct.

333 6 Related work

334 Recent work in offline alignment has focused on DPO [23] as a simpler alternative for aligning
 335 language models from preference data. Subsequent work has identified issues with DPO, including
 336 weak regularization [3] and a tendency to decrease the probability of winning generations during
 337 training [20]. Other methods have explored various avenues for improvement. These include
 338 analyzing the impact of noise on DPO alignment [11], proposing to update the reference policy
 339 during training [12], and suggesting a variant of IPO with a per-context margin [1]. Additional
 340 research has focused on token-level alignment methods [38; 22] and on developing a unified view of
 341 various offline alignment methods [31]. This work builds upon several these findings, and provides
 342 further analysis, as well as a solution based on pessimism and reward distillation.

343 While offline alignment methods are popular, recent evidence suggests that online alignment methods
 344 such as RLHF [6; 29], may lead to more favorable outcomes [13; 30; 8; 34]. Notably, Zhu et al. [41]
 345 proposed iterative data smoothing, which uses a trained model to softly label data during RLHF.
 346 Whether online or offline, however, policies are still susceptible to overfitting to certain degenerate
 347 phenomena. To this end, reward ensembles have been widely investigated recently as a mechanism
 348 for tackling reward hacking in RLHF [9; 7; 39; 25], and in the context of multi-objective optimization
 349 [19; 24]. We use an ensemble of rewards to represent the uncertainty with respect to reward models
 350 that are suitable given preference data. Moskovitz et al. [19] focus on “composite” rewards, with the
 351 goal of achieving high task reward while ensuring that every individual component is above some
 352 threshold—also by applying a Lagrangian relaxation. In this work, we also consider multiple reward
 353 models, but we only focus on cases where there is no known, obvious reward decomposition.

354 Finally, the question of using a small amount of offline data to learn high-quality policies, instead
 355 of online access to reward feedback, has been widely studied in the offline reinforcement learning
 356 (RL) literature. The predominant approach here is to use pessimism, that is, to learn a policy with
 357 the highest reward under all plausible environment models consistent with the data, with an extensive
 358 theoretical [18; 37; 33] and empirical [16; 5; 36] body of supporting work. The key insight in this
 359 literature is that without pessimism, the RL algorithm learns undesirable behaviors which are not
 360 explicitly ruled out in the training data, and pessimism provides a robust way of preventing such
 361 undesirable extrapolations, while still preserving generalization within the support of the data.

362 7 Conclusion

363 LM alignment is crucial for deploying safe and helpful assistants, but is difficult due to lack of
 364 access to perfect preference oracles. We presented a thorough theoretical analysis of some of
 365 the degeneracies that DPO is susceptible to when learning from sampled human preference data.
 366 Furthermore, our findings suggest that explicit reward modeling remains a powerful vehicle for
 367 introducing regularization into post-training. By distilling the reward assigned by a single, explicit
 368 reward model—or a family of explicit reward models—directly into the implicit reward maximized
 369 by our policies using offline data, we demonstrated that we can achieve improved robustness to
 370 variations in preference dataset quality, while maintaining the simplicity of the DPO framework.

371 **Limitations.** The empirical results in the paper are based on one dataset and form of distribution shift.
 372 For deeper understanding of pessimism and ensembling, additional settings should be explored. The
 373 theoretical aspects of the paper are sometimes based on restrictive assumptions and simplifications.
 374 Nonetheless, they provide potential explanations for phenomena observed in real-world settings.

375 **Broader impact.** We introduce new ideas to the active field of research on preference-based post-
 376 training, which we hope will help facilitate the alignment of large models, and improve understanding
 377 of current approaches—ultimately supporting the development of capable and reliable AI systems.

378 **References**

- 379 [1] Afra Amini, Tim Vieira, and Ryan Cotterell. Direct preference optimization with an offset.
380 *arXiv preprint arXiv:2402.10571*, 2024.
- 381 [2] Rohan Anil, Andrew M. Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos,
382 Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, Eric Chu, Jonathan H. Clark,
383 Laurent El Shafey, Yanping Huang, Kathy Meier-Hellstern, Gaurav Mishra, Erica Moreira,
384 Mark Omernick, Kevin Robinson, Sebastian Ruder, Yi Tay, Kefan Xiao, Yuanzhong Xu, Yujing
385 Zhang, Gustavo Hernandez Abrego, Junwhan Ahn, Jacob Austin, Paul Barham, Jan Botha,
386 James Bradbury, Siddhartha Brahma, Kevin Brooks, Michele Catasta, Yong Cheng, Colin
387 Cherry, Christopher A. Choquette-Choo, Aakanksha Chowdhery, Clément Crepy, Shachi Dave,
388 Mostafa Dehghani, Sunipa Dev, Jacob Devlin, Mark Díaz, Nan Du, Ethan Dyer, Vlad Feinberg,
389 Fangxiaoyu Feng, Vlad Fienber, Markus Freitag, Xavier Garcia, Sebastian Gehrmann, Lucas
390 Gonzalez, Guy Gur-Ari, Steven Hand, Hadi Hashemi, Le Hou, Joshua Howland, Andrea Hu,
391 Jeffrey Hui, Jeremy Hurwitz, Michael Isard, Abe Ittycheriah, Matthew Jagielski, Wenhao Jia,
392 Kathleen Kenealy, Maxim Krikun, Sneha Kudugunta, Chang Lan, Katherine Lee, Benjamin
393 Lee, Eric Li, Music Li, Wei Li, YaGuang Li, Jian Li, Hyeontaek Lim, Hanzhao Lin, Zhongtao
394 Liu, Frederick Liu, Marcello Maggioni, Aroma Mahendru, Joshua Maynez, Vedant Misra,
395 Maysam Moussalem, Zachary Nado, John Nham, Eric Ni, Andrew Nystrom, Alicia Parrish,
396 Marie Pellat, Martin Polacek, Alex Polozov, Reiner Pope, Siyuan Qiao, Emily Reif, Bryan
397 Richter, Parker Riley, Alex Castro Ros, Aurko Roy, Brennan Saeta, Rajkumar Samuel, Renee
398 Shelby, Ambrose Slone, Daniel Smilkov, David R. So, Daniel Sohn, Simon Tokumine, Dasha
399 Valter, Vijay Vasudevan, Kiran Vodrahalli, Xuezhi Wang, Pidong Wang, Zirui Wang, Tao Wang,
400 John Wieting, Yuhuai Wu, Kelvin Xu, Yunhan Xu, Linting Xue, Pengcheng Yin, Jiahui Yu,
401 Qiao Zhang, Steven Zheng, Ce Zheng, Weikang Zhou, Denny Zhou, Slav Petrov, and Yonghui
402 Wu. Palm 2 technical report, 2023.
- 403 [3] Mohammad Gheshlaghi Azar, Zhaohan Daniel Guo, Bilal Piot, Remi Munos, Mark Rowland,
404 Michal Valko, and Daniele Calandriello. A general theoretical paradigm to understand learning
405 from human preferences. In *International Conference on Artificial Intelligence and Statistics*,
406 pages 4447–4455. PMLR, 2024.
- 407 [4] Ralph Allan Bradley and Milton E. Terry. Rank analysis of incomplete block designs: I. the
408 method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952. ISSN 00063444. URL
409 <http://www.jstor.org/stable/2334029>.
- 410 [5] Ching-An Cheng, Tengyang Xie, Nan Jiang, and Alekh Agarwal. Adversarially trained actor
411 critic for offline reinforcement learning. In *International Conference on Machine Learning*,
412 pages 3852–3878. PMLR, 2022.
- 413 [6] Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep
414 reinforcement learning from human preferences. *Advances in neural information processing*
415 *systems*, 30, 2017.
- 416 [7] Thomas Coste, Usman Anwar, Robert Kirk, and David Krueger. Reward model ensembles help
417 mitigate overoptimization. *arXiv preprint arXiv:2310.02743*, 2023.
- 418 [8] Hanze Dong, Wei Xiong, Bo Pang, Haoxiang Wang, Han Zhao, Yingbo Zhou, Nan Jiang, Doyen
419 Sahoo, Caiming Xiong, and Tong Zhang. Rlhf workflow: From reward modeling to online rlhf.
420 2024.
- 421 [9] Jacob Eisenstein, Chirag Nagpal, Alekh Agarwal, Ahmad Beirami, Alex D’Amour, DJ Dvi-
422 jotham, Adam Fisch, Katherine Heller, Stephen Pfohl, Deepak Ramachandran, et al. Helping
423 or herding? Reward model ensembles mitigate but do not eliminate reward hacking. *arXiv*
424 *preprint arXiv:2312.09244*, 2023.
- 425 [10] Tommaso Furlanello, Zachary Lipton, Michael Tschannen, Laurent Itti, and Anima Anandkumar.
426 Born again neural networks. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the*
427 *35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine*
428 *Learning Research*, pages 1607–1616. PMLR, 10–15 Jul 2018. URL [https://proceedings.](https://proceedings.mlr.press/v80/furlanello18a.html)
429 [mlr.press/v80/furlanello18a.html](https://proceedings.mlr.press/v80/furlanello18a.html).

- 430 [11] Yang Gao, Dana Alon, and Donald Metzler. Impact of preference noise on the alignment
431 performance of generative language models. *arXiv preprint arXiv:2404.09824*, 2024.
- 432 [12] Alexey Gorbatovski, Boris Shaposhnikov, Alexey Malakhov, Nikita Surnachev, Yaroslav Ak-
433 senov, Ian Maksimov, Nikita Balagansky, and Daniil Gavrilov. Learn your reference model for
434 real good alignment. *arXiv preprint arXiv:2404.09656*, 2024.
- 435 [13] Shangmin Guo, Biao Zhang, Tianlin Liu, Tianqi Liu, Misha Khalman, Felipe Llinares, Alexan-
436 dre Rame, Thomas Mesnard, Yao Zhao, Bilal Piot, et al. Direct language model alignment from
437 online ai feedback. *arXiv preprint arXiv:2402.04792*, 2024.
- 438 [14] Geoffrey Hinton, Oriol Vinyals, and Jeffrey Dean. Distilling the knowledge in a neural network.
439 In *NIPS Deep Learning and Representation Learning Workshop*, 2015. URL [http://arxiv.
440 org/abs/1503.02531](http://arxiv.org/abs/1503.02531).
- 441 [15] Tomasz Korbak, Ethan Perez, and Christopher Buckley. RL with KL penalties is better viewed
442 as bayesian inference. In *Findings of the Association for Computational Linguistics: EMNLP
443 2022*, pages 1083–1091, 2022.
- 444 [16] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning
445 for offline reinforcement learning. *Advances in Neural Information Processing Systems*, 33:
446 1179–1191, 2020.
- 447 [17] Harrison Lee, Samrat Phatale, Hassan Mansoor, Kellie Lu, Thomas Mesnard, Colton Bishop,
448 Victor Carbune, and Abhinav Rastogi. RLHF: Scaling reinforcement learning from human
449 feedback with ai feedback. *arXiv preprint arXiv:2309.00267*, 2023.
- 450 [18] Yao Liu, Adith Swaminathan, Alekh Agarwal, and Emma Brunskill. Provably good batch
451 off-policy reinforcement learning without great exploration. *Advances in neural information
452 processing systems*, 33:1264–1274, 2020.
- 453 [19] Ted Moskovitz, Aaditya K Singh, DJ Strouse, Tuomas Sandholm, Ruslan Salakhutdinov, Anca D
454 Dragan, and Stephen McAleer. Confronting reward model overoptimization with constrained
455 rlhf. *arXiv preprint arXiv:2310.04373*, 2023.
- 456 [20] Arka Pal, Deep Karkhanis, Samuel Dooley, Manley Roberts, Siddhartha Naidu, and Colin White.
457 Smaug: Fixing failure modes of preference optimisation with dpo-positive. *arXiv preprint
458 arXiv:2402.13228*, 2024.
- 459 [21] Ryan Park, Rafael Rafailov, Stefano Ermon, and Chelsea Finn. Disentangling length from
460 quality in direct preference optimization. *arXiv preprint arXiv:2403.19159*, 2024.
- 461 [22] Rafael Rafailov, Joey Hejna, Ryan Park, and Chelsea Finn. From r to q^* : Your language model
462 is secretly a q-function. 2024.
- 463 [23] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and
464 Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model.
465 *Advances in Neural Information Processing Systems*, 36, 2024.
- 466 [24] Alexandre Ramé, Guillaume Couairon, Mustafa Shukor, Corentin Dancette, Jean-Baptiste Gaya,
467 Laure Soulier, and Matthieu Cord. Rewarded soups: towards pareto-optimal alignment by
468 interpolating weights fine-tuned on diverse rewards, 2023.
- 469 [25] Alexandre Ramé, Nino Vieillard, Léonard Hussenot, Robert Dadashi, Geoffrey Cideron, Olivier
470 Bachem, and Johan Ferret. Warm: On the benefits of weight averaged reward models. 2024.
- 471 [26] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and
472 Yoshua Bengio. Fitnets: Hints for thin deep nets. In *In Proceedings of ICLR*, 2015.
- 473 [27] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal
474 policy optimization algorithms. 2017.
- 475 [28] Prasann Singhal, Tanya Goyal, Jiacheng Xu, and Greg Durrett. A long way to go: Investigating
476 length correlations in rlhf. 2023.

- 477 [29] Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec
478 Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback.
479 *Advances in Neural Information Processing Systems*, 33:3008–3021, 2020.
- 480 [30] Fahim Tajwar, Anikait Singh, Archit Sharma, Rafael Rafailov, Jeff Schneider, Tengyang Xie,
481 Stefano Ermon, Chelsea Finn, and Aviral Kumar. Preference fine-tuning of llms should leverage
482 suboptimal, on-policy data, 2024.
- 483 [31] Yunhao Tang, Zhaohan Daniel Guo, Zeyu Zheng, Daniele Calandriello, Rémi Munos, Mark
484 Rowland, Pierre Harvey Richemond, Michal Valko, Bernardo Ávila Pires, and Bilal Piot.
485 Generalized preference optimization: A unified approach to offline alignment. 2024.
- 486 [32] Michael Völske, Martin Potthast, Shahbaz Syed, and Benno Stein. TL;DR: Mining Reddit
487 to learn automatic summarization. In Lu Wang, Jackie Chi Kit Cheung, Giuseppe Carenini,
488 and Fei Liu, editors, *Proceedings of the Workshop on New Frontiers in Summarization*, pages
489 59–63, Copenhagen, Denmark, September 2017. Association for Computational Linguistics.
490 doi: 10.18653/v1/W17-4508. URL <https://aclanthology.org/W17-4508>.
- 491 [33] Tengyang Xie, Ching-An Cheng, Nan Jiang, Paul Mineiro, and Alekh Agarwal. Bellman-
492 consistent pessimism for offline reinforcement learning. *Advances in neural information*
493 *processing systems*, 34:6683–6694, 2021.
- 494 [34] Shusheng Xu, Wei Fu, Jiaxuan Gao, Wenjie Ye, Weilin Liu, Zhiyu Mei, Guangju Wang, Chao
495 Yu, and Yi Wu. Is dpo superior to ppo for llm alignment? a comprehensive study. 2024.
- 496 [35] Chenglin Yang, Lingxi Xie, Siyuan Qiao, and Alan L. Yuille. Training deep neural networks
497 in generations: a more tolerant teacher educates better students. In *Proceedings of the Thirty-*
498 *Third AAAI Conference on Artificial Intelligence and Thirty-First Innovative Applications of*
499 *Artificial Intelligence Conference and Ninth AAAI Symposium on Educational Advances in*
500 *Artificial Intelligence*, AAAI’19/IAAI’19/EAAI’19. AAAI Press, 2019. ISBN 978-1-57735-809-
501 1. doi: 10.1609/aaai.v33i01.33015628. URL [https://doi.org/10.1609/aaai.v33i01.](https://doi.org/10.1609/aaai.v33i01.33015628)
502 [33015628](https://doi.org/10.1609/aaai.v33i01.33015628).
- 503 [36] Tianhe Yu, Aviral Kumar, Rafael Rafailov, Aravind Rajeswaran, Sergey Levine, and Chelsea
504 Finn. Combo: Conservative offline model-based policy optimization. *Advances in neural*
505 *information processing systems*, 34:28954–28967, 2021.
- 506 [37] Andrea Zanette, Martin J Wainwright, and Emma Brunskill. Provable benefits of actor-critic
507 methods for offline reinforcement learning. *Advances in neural information processing systems*,
508 34:13626–13640, 2021.
- 509 [38] Yongcheng Zeng, Guoqing Liu, Weiyu Ma, Ning Yang, Haifeng Zhang, and Jun Wang. Token-
510 level direct preference optimization. 2024.
- 511 [39] Yuanzhao Zhai, Han Zhang, Yu Lei, Yue Yu, Kele Xu, Dawei Feng, Bo Ding, and Huaimin
512 Wang. Uncertainty-penalized reinforcement learning from human feedback with diverse reward
513 lora ensembles. 2023.
- 514 [40] Yao Zhao, Rishabh Joshi, Tianqi Liu, Misha Khalman, Mohammad Saleh, and Peter J. Liu.
515 Slic-hf: Sequence likelihood calibration with human feedback. 2023.
- 516 [41] Banghua Zhu, Michael I. Jordan, and Jiantao Jiao. Iterative data smoothing: Mitigating reward
517 overfitting and overoptimization in rlhf. 2024.

518 **A Proofs**

519 **A.1 Proof of Proposition 1**

520 *Proof.* Since all the preference pairs (y, y') are mutually disjoint, and θ_y is specific to each y , the
 521 DPO objective over $\mathcal{D}_{\text{pref}}$ is convex in $\Delta = \{\Delta_1, \dots, \Delta_n\}$, where

$$\Delta_i = \beta \log \frac{\pi_\theta(y_i^w) \pi_{\text{ref}}(y_i^\ell)}{\pi_\theta(y_i^\ell) \pi_{\text{ref}}(y_i^w)}. \quad (15)$$

522 Furthermore, the different Δ_i are completely independent from each other due to the preference pairs
 523 being disjoint, so they can be optimized over separately.

524 In particular, for every i we have that

$$\lim_{\Delta_i \rightarrow \infty} -\log(\sigma(\Delta_i)) = 0, \quad (16)$$

525 which implies that $\Delta^* = \{\infty\}^n$ is the unique global minimizer of the DPO loss over $\mathcal{D}_{\text{pref}}$ in the
 526 space of Δ 's, and any θ^* that is a global minimizer must therefore satisfy

$$\log \frac{\pi_{\theta^*}(y_i^w) \pi_{\text{ref}}(y_i^\ell)}{\pi_{\theta^*}(y_i^\ell) \pi_{\text{ref}}(y_i^w)} = \infty. \quad (17)$$

527

□

528 **A.2 Proof of Corollary 1**

529 *Proof.* Following the same argument of the proof of Proposition 1, we have that all global minimizers
 530 θ^* of the DPO satisfy $\Delta_i^* = \infty$, which in turn implies that

$$\frac{\pi_{\theta^*}(y_i^w) \pi_{\text{ref}}(y_i^\ell)}{\pi_{\theta^*}(y_i^\ell) \pi_{\text{ref}}(y_i^w)} = \infty. \quad (18)$$

531 Since $\pi_{\text{ref}}(y)$ is assumed to satisfy $0 < \pi_{\text{ref}}(y) < 1$ for all y , this implies that all θ^* satisfy

$$\frac{\pi_{\theta^*}(y_i^w)}{\pi_{\theta^*}(y_i^\ell)} = \infty, \quad (19)$$

532 which further implies that $\pi_{\theta^*}(y_i^\ell) = 0$ and $\pi_{\theta^*}(y_i^w) > 0$ for all $i \in [n]$, as $\pi_{\theta^*}(y_i^w) \leq 1$ for any y_i^w .
 533 **Aggregating**

$$\mathcal{C}(y_\ell) = \{y: \exists i \in [n] \text{ s.t } y_i^\ell = y\} \quad (20)$$

534 then gives that

$$\pi_{\theta^*}(\mathcal{C}(y_\ell)) = \sum_{y \in \mathcal{C}(y_\ell)} \pi_{\theta^*}(y) = 0 \implies \pi_{\theta^*}(\mathcal{C}(y_\ell)^c) = 1. \quad (21)$$

535

□

536 To prove the converse, let $\pi_{\theta'}$ be a policy that satisfies $\pi_{\theta'}(\mathcal{C}(y_\ell)^c) = 1$, with $\pi_{\theta'}(y_i^w) > 0, \forall i \in [n]$.
 537 As $\pi_{\theta'}(y) \geq 0$ for all y , this implies that $\pi_{\theta'}(y_i^\ell) = 0 \forall i \in [n]$. Then, we have

$$\frac{\pi_{\theta'}(y_i^w)}{\pi_{\theta'}(y_i^\ell)} = \infty, \quad (22)$$

538 which by Proposition 1 implies that $\pi_{\theta'}$ is a global optimum.

539 **A.3 Proof of Theorem 1**

540 *Proof.* We know that the optimal policy for the RLHF objective (1) is given by $\pi_{\theta^*}(y|x) \propto$
 541 $\pi_{\text{ref}}(y|x) \exp(r^*(x, y)/\beta)$. Plugging this policy into the distillation objective (7), we see that
 542 $\mathcal{L}_{\text{distill}}(r^*, \pi_{\theta^*}, \rho) = 0$ for all ρ . In fact, the loss is equal to 0 pointwise, meaning that π_{θ^*} is

543 a global minimizer of the distillation objective (7). Further, let π be some other minimizer of
 544 $\mathcal{L}_{\text{distill}}(r^*, \cdot, \rho)$. Then π also has to attain a loss of 0 at all (x, y, y') in the support of ρ , meaning
 545 that $\log \pi(y|x) - \log \pi(y'|x) = \log \pi_{\theta^*}(y|x) - \log \pi_{\theta^*}(y'|x)$ for all (x, y, y') in the support of ρ .
 546 Consequently, the two policies coincide in the support of ρ (due to the normalization constraint, there
 547 is no additional offset term allowed as the support of ρ covers all of \mathcal{Y}). Finally, noting that the
 548 support of the chosen ρ is such that π_{θ^*} puts no mass outside its support due to the KL constraint
 549 in (1), we complete the proof. \square

550 A.4 Proof of Theorem 2

551 *Proof.* Consider the pessimistic objective:

$$\max_{\pi_{\theta}} \min_{r_{\text{tgt}} \in \mathcal{S}} \mathbb{E}_{\mu(x)} \left[\mathbb{E}_{\pi_{\theta}(y|x)} [r_{\text{tgt}}(x, y)] - \mathbb{E}_{\pi_{\text{ref}}(y|x)} [r_{\text{tgt}}(x, y)] \right] - \beta \mathbb{D}_{\text{KL}}(\pi_{\theta} \| \pi_{\text{ref}}). \quad (23)$$

552 As it is linear in r_{tgt} and convex in π , we can switch the order of min and max:

$$\min_{r_{\text{tgt}} \in \mathcal{S}} \left[\max_{\pi \in \Pi} \mathbb{E}_{\mu(x)} \left[\mathbb{E}_{\pi(y|x)} [r_{\text{tgt}}(x, y)] - \mathbb{E}_{\pi_{\text{ref}}(y|x)} [r_{\text{tgt}}(x, y)] \right] - \beta \mathbb{D}_{\text{KL}}(\pi \| \pi_{\text{ref}}) \right]. \quad (24)$$

553 Note that every $r_{\text{tgt}} \in \mathcal{S}$ can be written in terms of the KL-constrained policy $\pi_{r_{\text{tgt}}}^*$ it induces, i.e.,

$$r_{\text{tgt}}(x, y) = \beta \log \frac{\pi_{r_{\text{tgt}}}^*(y | x)}{\pi_{\text{ref}}(y | x)} + \beta \log Z(x, r_{\text{tgt}}), \quad (25)$$

554 where

$$\pi_{r_{\text{tgt}}}^* = \operatorname{argmax}_{\pi_{\theta}} \mathbb{E}_{\mu(x)} \mathbb{E}_{\pi_{\theta}(y|x)} [r_{\text{tgt}}(x, y)] - \beta \mathbb{D}_{\text{KL}}(\pi_{\theta} \| \pi_{\text{ref}}) \quad (26)$$

555 which has the form

$$\pi_{r_{\text{tgt}}}^*(y | x) = \frac{1}{Z(x, r_{\text{tgt}})} \pi_{\text{ref}}(y | x) \exp \left(\frac{1}{\beta} r_{\text{tgt}}(x, y) \right) \quad (27)$$

556 where $Z(x, r_{\text{tgt}})$ is the partition function:

$$Z(x, r_{\text{tgt}}) = \sum_{y \in \mathcal{Y}} \pi_{\text{ref}}(y | x) \exp \left(\frac{1}{\beta} r_{\text{tgt}}(x, y) \right). \quad (28)$$

557 Substituting $\pi_{r_{\text{tgt}}}^*$ in for \max_{π} and writing r_{tgt} in terms of $\pi_{r_{\text{tgt}}}^*$, we get the simplified objective

$$\begin{aligned} & \min_{r_{\text{tgt}} \in \mathcal{S}} \left[\max_{\pi \in \Pi} \mathbb{E}_{\mu(x)} \left[\mathbb{E}_{\pi(y|x)} [r_{\text{tgt}}(x, y)] - \mathbb{E}_{\pi_{\text{ref}}(y|x)} [r_{\text{tgt}}(x, y)] \right] - \beta \mathbb{D}_{\text{KL}}(\pi \| \pi_{\text{ref}}) \right] \\ &= \min_{r_{\text{tgt}} \in \mathcal{S}} \left[\mathbb{E}_{\mu(x)} \left[\mathbb{E}_{\pi_{r_{\text{tgt}}}^*(y|x)} \left[\beta \log \frac{\pi_{r_{\text{tgt}}}^*(y | x)}{\pi_{\text{ref}}(y | x)} + \beta \log Z(x, r_{\text{tgt}}) \right] \right. \right. \\ & \quad \left. \left. - \mathbb{E}_{\pi_{\text{ref}}(y|x)} \left[\beta \log \frac{\pi_{r_{\text{tgt}}}^*(y | x)}{\pi_{\text{ref}}(y | x)} + \beta \log Z(x, r_{\text{tgt}}) \right] \right. \right. \\ & \quad \left. \left. - \beta \mathbb{D}_{\text{KL}}(\pi_{r_{\text{tgt}}}^* \| \pi_{\text{ref}} | x) \right] \right] \\ &= \min_{r_{\text{tgt}} \in \mathcal{S}} \beta \left[\mathbb{E}_{\mu(x)} \left[\mathbb{D}_{\text{KL}}(\pi_{r_{\text{tgt}}}^* \| \pi_{\text{ref}} | x) + \mathbb{D}_{\text{KL}}(\pi_{\text{ref}} \| \pi_{r_{\text{tgt}}}^* | x) - \mathbb{D}_{\text{KL}}(\pi_{r_{\text{tgt}}}^* \| \pi_{\text{ref}} | x) \right] \right] \\ &= \min_{r_{\text{tgt}} \in \mathcal{S}} \beta \mathbb{E}_{\mu(x)} \left[\mathbb{D}_{\text{KL}}(\pi_{\text{ref}} \| \pi_{r_{\text{tgt}}}^* | x) \right]. \end{aligned} \quad (29)$$

558 \square

559 **A.5 Proof of Proposition 2**

560 *Proof.* The proof is a standard Lagrangian duality argument, which we reproduce here for complete-
561 ness. For two functions $f(z)$ and $g(z)$, let us define

$$z^* = \operatorname{argmin}_z f(z) + \alpha g(z). \quad (30)$$

562 Let us also consider the constrained problem

$$z' = \operatorname{argmin}_z f(z) \quad \text{s.t.} \quad g(z) \leq g(z^*). \quad (31)$$

563 Suppose by contradiction that z^* is not a minimizer of (31). Since z^* is feasible for the constraint by
564 construction, we get that $f(z') < f(z^*)$. Consequently, we further have

$$f(z') + \alpha g(z') < f(z^*) + \alpha g(z^*),$$

565 where the inequality follows from the feasibility of z' in (31). This contradicts the optimality
566 of z^* in (30), meaning that z^* must be a minimizer of (31). Applying this general result with
567 $f = \beta \mathbb{E}_{\mu(x)} \mathbb{D}_{\text{KL}}(\pi_{\text{ref}}(y | x) \| \pi_{\theta}(y | x))$, $g = \min_{r_{\text{tgt}}^i \in \mathcal{S}} \mathcal{L}_{\text{distill}}(r_{\text{tgt}}^i, \pi_{\theta}; \rho)$, and $z = \pi_{\theta}$ completes
568 the proof, since we recognize the set \mathcal{S}_{γ} in (12) to be equivalent to $\bigcup_{r_{\text{tgt}}^i \in \mathcal{S}} \mathcal{L}_{\text{distill}}(r_{\text{tgt}}^i, \pi_{\theta}; \rho) \leq \lambda$.

569 \square

570 **A.6 Proof of Proposition 3**

571 *Proof.* We differentiate $\mathcal{L}_{\text{pdpo}}$ with respect to $\psi_{\ell} = \pi_{\theta}(y^{\ell}) / \pi_{\text{ref}}(y^{\ell})$ with i implicit, obtaining,

$$\frac{\partial \mathcal{L}_{\text{pdpo}}}{\partial \psi_{\ell}} = \beta \frac{\psi_{\ell}^{\beta}}{\psi_w^{\beta} + \psi_{\ell}^{\beta}} \psi_{\ell}^{-1} - \frac{\beta}{\alpha} \psi_{\ell}^{-1} = \beta \psi_{\ell}^{-1} \left(\frac{\psi_{\ell}^{\beta}}{\psi_w^{\beta} + \psi_{\ell}^{\beta}} - \alpha^{-1} \right) \quad (32)$$

572 which is zero when,

$$\alpha \psi_{\ell}^{\beta} = \psi_w^{\beta} + \psi_{\ell}^{\beta} \quad (33)$$

$$\psi_{\ell} = \left(\frac{1}{\alpha - 1} \right)^{1/\beta} \psi_w \quad (34)$$

$$\log \psi_{\ell} = -\frac{1}{\beta} \log(\alpha - 1) + \log \psi_w \quad (35)$$

$$\log \pi_{\theta}(y^{\ell}) = \log \pi_{\text{ref}}(y^{\ell}) - \frac{1}{\beta} \log(\alpha - 1) + \log \pi_{\theta}(y^w) - \log \pi_{\text{ref}}(y^w). \quad (36)$$

573 By the second-order condition, the critical point is a minimum. The objective $\mathcal{L}_{\text{pdpo}}$ is the sum of two
574 components: the negative log sigmoid term for \mathcal{L}_i and the negative log probability for $\hat{\Omega}$. Because
575 each component is a convex function of ψ_i , so is $\mathcal{L}_{\text{pdpo}}$. As a result, the local minimum $\log \hat{\pi}_{\theta}(y^{\ell})$ is
576 also a global minimum. \square

577 **A.7 Proof of Proposition 4**

578 *Proof.* This follows directly from differentiating eq. (7) with respect to $\pi_{\theta}(y_2)$. \square

579 **A.8 Proof of Proposition 5**

580 *Proof.* Let $\Delta = [c(y^w) - c(y^{\ell})]$ and $\rho = \pi_{\text{ref}}(y^w) / \pi_{\text{ref}}(y^{\ell})$. The theorem assumes $|y^w| = |y^{\ell}|$.
581 Then $\mathcal{L}_{\text{dpo}} = -\log \sigma(\beta(\Delta \cdot \theta) + \beta \log \rho)$. The derivative with respect to θ is,

$$\frac{\partial \mathcal{L}_{\beta}(\theta)}{\partial \theta} = -(1 - \sigma(\beta(\Delta \cdot \theta) + \beta \log \rho)) \beta \Delta = -\Pr(y^{\ell} \succ y^w; \theta) \beta \Delta \prec 0. \quad (37)$$

582 Let $\delta_t = \beta \Pr(y^\ell \succ y^w; \theta^{(t)})$. Then,

$$\tilde{\pi}_{\theta^{(t)}} = \theta^{(t)} \cdot c(y^w) - n \max_j \theta_j^{(t)} \quad (38)$$

$$= (\theta^{(t-1)} + \delta_t \Delta) \cdot c(y^w) - n \max_j (\theta_j^{(t-1)} + \delta_t \Delta_j) \quad (39)$$

$$= \theta^{(t-1)} \cdot c(y^w) - n \max_j \theta_j^{(t-1)} + \delta_t \Delta \cdot c(y^w) - n \delta_t \max_j \Delta_j \quad (40)$$

$$= \tilde{\pi}_{\theta^{(t-1)}} + \delta_t \left(\Delta \cdot c(y^w) - n \max_j \Delta_j \right) \quad (41)$$

$$= \tilde{\pi}_{\theta^{(t-1)}} + \delta_t \sum_j c_j(y^w) (\Delta_j - \max_{j'} \Delta_{j'}) \leq \tilde{\pi}_{\theta^{(t-1)}}. \quad (42)$$

583 We obtain $\max_j (\theta_j^{(t-1)} + \delta_t \Delta_j) = \max_j \theta_j^{(t-1)} + \max_j \delta_t \Delta_j$ from the fact that $\theta^{(0)} = 0$ and
 584 therefore $j \in \arg \max \Delta$ implies $j \in \arg \max \theta^{(t')}$ for all $t' > 0$. The second-to-last step uses
 585 $n = \sum_j c_j(y^w)$ and the final step uses $\Delta_j \leq \max_{j'} \Delta_{j'}$. Finally, we have $\pi_{\theta^{(t)}}(y) \leq \tilde{\pi}_{\theta^{(t)}}(y^w)$
 586 because $Z(\theta) = \log \sum_j \exp \theta_j \geq \log \max_j \exp \theta_j = \max_j \theta_j$. \square

587 A.9 Proof of Proposition 6

588 *Proof.* Applying gradient descent with learning rate η to the gradient from Equation (37), at each
 589 step t the parameters are,

$$\theta^{(t)} = \theta^{(t-1)} + \eta \beta \Pr(y^\ell \succ y^w; \theta^{(t-1)}) \Delta = \left(\sum_{t'=1}^t \eta \beta \Pr(y^\ell \succ y^w; \theta^{(t')}) \right) \Delta = \tau(t) \Delta. \quad (43)$$

590 Plugging these parameters into the likelihoods,

$$\ell_{\theta^{(t)}}(c(y^w)) - \ell_{\theta^{(t)}}(\hat{y}) = c(y^w) \cdot \theta^{(t)} - n Z(\theta^{(t)}) - c(\hat{y}) \cdot \theta^{(t)} + n Z(\theta^{(t)}) \quad (44)$$

$$= (c(y^w) - c(\hat{y})) \cdot \theta^{(t)} = (c(y^w) - c(\hat{y})) \cdot (\tau(t) \Delta) \quad (45)$$

$$= \tau(t) (c(y^w) \cdot \Delta - n \max \Delta) = \tau(t) k, \quad (46)$$

591 with $k \leq 0$ by $c(y^w) \cdot \Delta \leq \|c(y^w)\|_1 \times \|\Delta\|_\infty = n \max \Delta$. \square

592 B Transitive closure

593 Both p-DPO and IPO target a constant ratio for $\log \psi_w / \psi_l$. However, the loss surfaces are different.
 594 To see this, we consider a simplified setting with three possible outputs, y_1, y_2, y_3 . We observe either
 595 $\mathcal{D} = \{(y_1 \prec y_2), (y_2 \prec y_3)\}$ or $\bar{\mathcal{D}} = \mathcal{D} \cup \{(y_1 \prec y_3)\}$. If we treat this problem as a multi-arm
 596 bandit, the goal is to assign a weight to each arm, which we denote $\psi_i = \log \pi_\theta(y_i | x) + Z_x$, with Z_x
 597 an underdetermined log-partition function.

598 **Proposition 7.** Let $\mathcal{D} = \{(i, i+1) : i \in 1, 2, \dots, n\}$ for $n > 2$. Let $\bar{\mathcal{D}}$ be the dataset arising from the
 599 transitive closure of \mathcal{D} . Assume π_{ref} is indifferent to all (y_i, y_j) . Let $\psi_\infty^{(\mathcal{D})} = \max_i \psi_i^{(\mathcal{D})} - \min_i \psi_i^{(\mathcal{D})}$.
 600 Then $\psi_\infty^{(\mathcal{D})} = (n-1)\tau^{-1} > \psi_\infty^{(\bar{\mathcal{D}})} = 2\frac{n-1}{n}\tau^{-1}$.

601 *Proof.* For \mathcal{D} , the IPO objective can be minimized at zero, so that $\psi_\infty^{(\mathcal{D})} = (n-1)\tau^{-1}$. For $\bar{\mathcal{D}}$,
 602 each adjacent pair of completions is separated by γ , and the objective is $\sum_{i=1}^{n-1} (n-i)(i\gamma - \tau^{-1})^2$.
 603 The minimum is $\gamma = \frac{n(n+1)(n-1)/6}{n^2(n+1)(n-1)/12} \tau^{-1} = \frac{2}{n} \tau^{-1}$, so that $\psi_\infty^{(\bar{\mathcal{D}})} = (n-1)\gamma = 2\frac{n-1}{n}\tau^{-1} <$
 604 $(n-1)\tau^{-1} = \psi_\infty^{(\mathcal{D})}$ for $n > 2$. \square

605 Intuitively, the observation of $(y_1 \prec y_3)$ should increase confidence that y_3 is superior to y_1 , but
 606 in IPO it has the opposite effect, drawing their scores closer together. While pessimistic DPO also
 607 has a target ratio between each preference pair, its loss surface is different: in particular, it does not

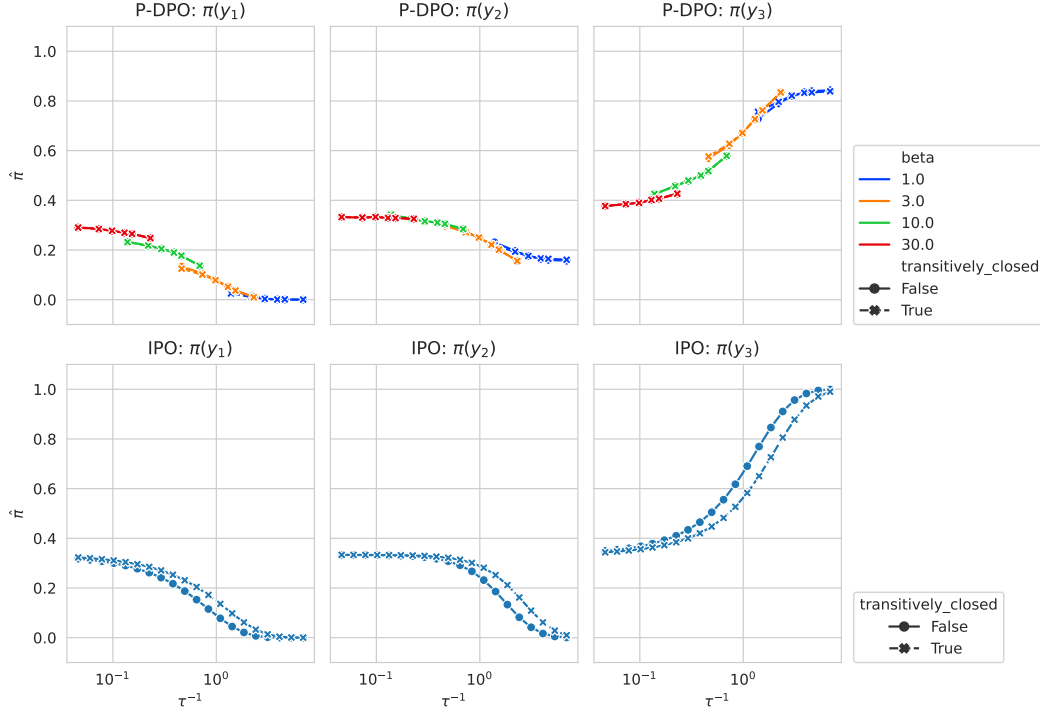


Figure B.1: **Effect of transitive closure on p-DPO and IPO solutions to preference learning in a multi-arm bandit.** Each column shows the learned policy probability for a given arm, based on the preferences $y_1 \prec y_2 \prec y_3$. The top row shows that in p-DPO, the probabilities are not materially affected by the transitive closure $y_1 \prec y_3$. The bottom row shows that in IPO, transitive closure causes the probabilities to be compressed. In each subfigure, we sweep a range of effective values of τ^{-1} , shown on the x-axis.

608 increase quadratically as we move away from the target. We find empirically that pessimistic DPO is
 609 robust to the transitive closure of preference annotations in the multi-arm bandit setting, as shown in
 610 [Figure B.1](#). As discussed above, DPO will set $\psi_1 \rightarrow -\infty$ because y_1 is never preferred.

611 In our empirical experiments we solve the p-DPO and IPO objectives for both $\mathcal{D} =$
 612 $\{(y_1, y_2), (y_2, y_3)\}$ and $\overline{\mathcal{D}} = \mathcal{D} \cup \{(y_1, y_3)\}$, solving with respect to $\{\pi_\theta(y_i)\}$. IPO is solved analytically
 613 as a quadratic program; for pessimistic DPO we used projected gradient descent. We consider
 614 $\beta \in (1, 3, 10, 30)$ and $\alpha \in (5, 10, 20, 50, 100, 1000)$. As shown in [Figure B.1](#), there are significant
 615 differences in the IPO solutions with and without transitive closure, while for p-DPO these differences
 616 are imperceptible.

617 C Distribution over reward models for e-DPO

618 [Figure C.1](#) investigates the reason for the success of e-DPO, especially when $\rho < .5$. For every length
 619 bias, we show across all training examples the fraction of cases where a certain reward model, $r_{\rho,b}$,
 620 best matched the implicit reward of the final e-DPO policy. The policy matches different reward
 621 models in different examples. Moreover, there is inverse correlation between the data bias for policy
 622 training (ρ) and the data bias for training the reward models (b). This suggests that the ensemble
 623 in e-DPO helps as the policy is distilling from reward models that do not share the data bias of the
 624 policy training set.

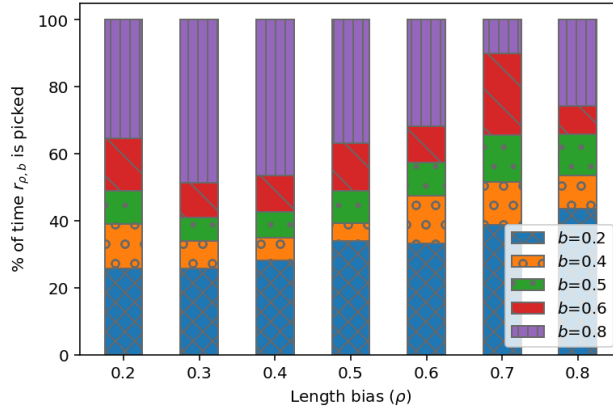


Figure C.1: We show for every length bias, ρ , the distribution over reward models that best match the final policy trained by e-DPO across all training examples. We observe that the e-DPO policy matches different reward models across examples. Moreover, when the policy is trained with data biased towards preferring short responses, the reward model that was trained on longer responses is often preferred and vice versa.

625 D Hyperparameters

626 Validation set performance across the range of hyperparameter settings is shown in Figure D.1. In
 627 pilot studies we found that these results were relatively robust to variation in the random seed, but did
 628 not conduct extensive investigation of this effect across all methods and hyperparameters due to cost.

629 E Compute resources

630 We train policies on 32 TPU v3 chips and reward models on 16 TPU v3 chips. We obtain roughly 0.1
 631 steps per second when training, for both the policy and reward models.

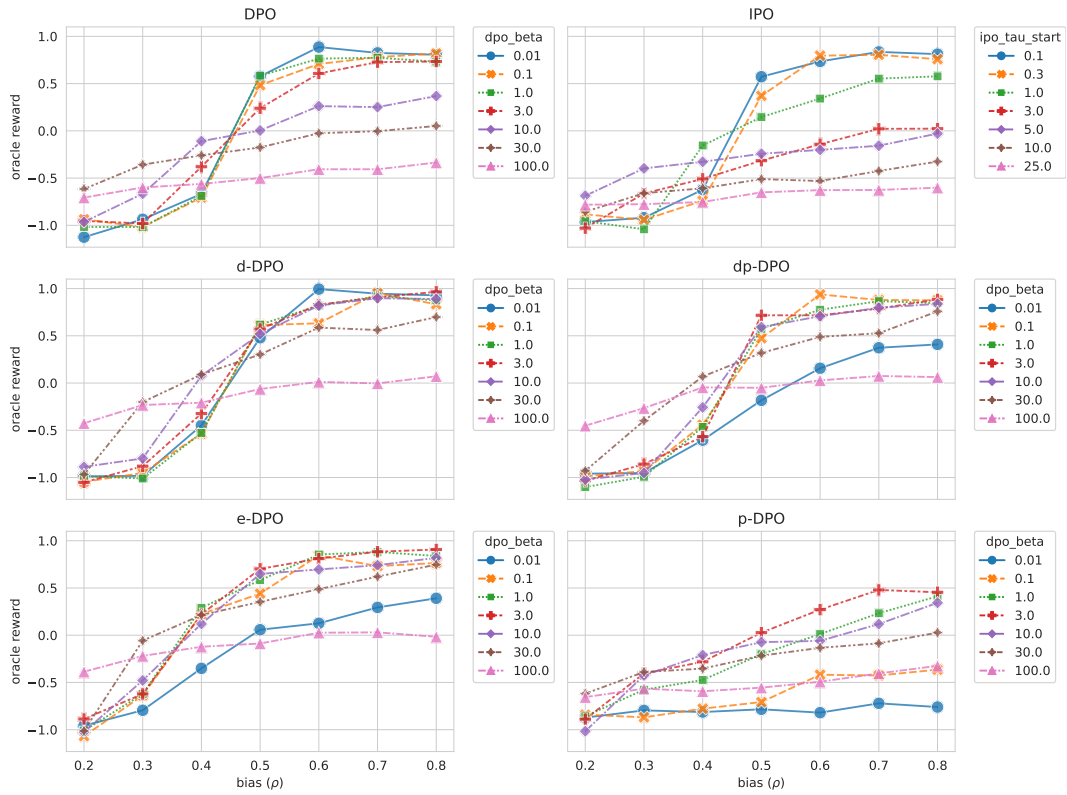


Figure D.1: **Validation set results** across hyperparameters for each method. For all methods, different values of ρ induce different optimal hyperparameters β and τ^{-1} .

632 **NeurIPS Paper Checklist**

633 The checklist is designed to encourage best practices for responsible machine learning research,
634 addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove
635 the checklist: **The papers not including the checklist will be desk rejected.** The checklist should
636 follow the references and precede the (optional) supplemental material. The checklist does NOT
637 count towards the page limit.

638 Please read the checklist guidelines carefully for information on how to answer these questions. For
639 each question in the checklist:

- 640 • You should answer [Yes] , [No] , or [NA] .
- 641 • [NA] means either that the question is Not Applicable for that particular paper or the
642 relevant information is Not Available.
- 643 • Please provide a short (1–2 sentence) justification right after your answer (even for NA).

644 **The checklist answers are an integral part of your paper submission.** They are visible to the
645 reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it
646 (after eventual revisions) with the final version of your paper, and its final version will be published
647 with the paper.

648 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.
649 While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a
650 proper justification is given (e.g., "error bars are not reported because it would be too computationally
651 expensive" or "we were unable to find the license for the dataset we used"). In general, answering
652 "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased in a binary way, we
653 acknowledge that the true answer is often more nuanced, so please just use your best judgment and
654 write a justification to elaborate. All supporting evidence can appear either in the main paper or the
655 supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification
656 please point to the section(s) where related material for the question can be found.

657 **IMPORTANT, please:**

- 658 • **Delete this instruction block, but keep the section heading “NeurIPS paper checklist”.**
- 659 • **Keep the checklist subsection headings, questions/answers and guidelines below.**
- 660 • **Do not modify the questions and only use the provided macros for your answers.**

661 **1. Claims**

662 Question: Do the main claims made in the abstract and introduction accurately reflect the
663 paper’s contributions and scope?

664 Answer: [Yes]

665 Justification: In our view, the abstract and introduction accurately summarize the contribu-
666 tions of the paper.

667 Guidelines:

- 668 • The answer NA means that the abstract and introduction do not include the claims
669 made in the paper.
- 670 • The abstract and/or introduction should clearly state the claims made, including the
671 contributions made in the paper and important assumptions and limitations. A No or
672 NA answer to this question will not be perceived well by the reviewers.
- 673 • The claims made should match theoretical and experimental results, and reflect how
674 much the results can be expected to generalize to other settings.
- 675 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
676 are not attained by the paper.

677 **2. Limitations**

678 Question: Does the paper discuss the limitations of the work performed by the authors?

679 Answer: [Yes]

680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731

Justification: See Section 7

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: See Appendix A

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: Details are provided in Section 4.1 and Appendix D.

Guidelines:

- The answer NA means that the paper does not include experiments.

- 732 • If the paper includes experiments, a No answer to this question will not be perceived
733 well by the reviewers: Making the paper reproducible is important, regardless of
734 whether the code and data are provided or not.
- 735 • If the contribution is a dataset and/or model, the authors should describe the steps taken
736 to make their results reproducible or verifiable.
- 737 • Depending on the contribution, reproducibility can be accomplished in various ways.
738 For example, if the contribution is a novel architecture, describing the architecture fully
739 might suffice, or if the contribution is a specific model and empirical evaluation, it may
740 be necessary to either make it possible for others to replicate the model with the same
741 dataset, or provide access to the model. In general, releasing code and data is often
742 one good way to accomplish this, but reproducibility can also be provided via detailed
743 instructions for how to replicate the results, access to a hosted model (e.g., in the case
744 of a large language model), releasing of a model checkpoint, or other means that are
745 appropriate to the research performed.
- 746 • While NeurIPS does not require releasing code, the conference does require all submis-
747 sions to provide some reasonable avenue for reproducibility, which may depend on the
748 nature of the contribution. For example
 - 749 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
750 to reproduce that algorithm.
 - 751 (b) If the contribution is primarily a new model architecture, the paper should describe
752 the architecture clearly and fully.
 - 753 (c) If the contribution is a new model (e.g., a large language model), then there should
754 either be a way to access this model for reproducing the results or a way to reproduce
755 the model (e.g., with an open-source dataset or instructions for how to construct
756 the dataset).
 - 757 (d) We recognize that reproducibility may be tricky in some cases, in which case
758 authors are welcome to describe the particular way they provide for reproducibility.
759 In the case of closed-source models, it may be that access to the model is limited in
760 some way (e.g., to registered users), but it should be possible for other researchers
761 to have some path to reproducing or verifying the results.

762 5. Open access to data and code

763 Question: Does the paper provide open access to the data and code, with sufficient instruc-
764 tions to faithfully reproduce the main experimental results, as described in supplemental
765 material?

766 Answer: [No]

767 Justification: Experiments are on publicly-available data, but it is not possible for us to share
768 code. We believe that the implementation should be relatively straightforward, given the
769 mathematical descriptions presented here.

770 Guidelines:

- 771 • The answer NA means that paper does not include experiments requiring code.
- 772 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
773 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 774 • While we encourage the release of code and data, we understand that this might not be
775 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
776 including code, unless this is central to the contribution (e.g., for a new open-source
777 benchmark).
- 778 • The instructions should contain the exact command and environment needed to run to
779 reproduce the results. See the NeurIPS code and data submission guidelines ([https:
780 //nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 781 • The authors should provide instructions on data access and preparation, including how
782 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 783 • The authors should provide scripts to reproduce all experimental results for the new
784 proposed method and baselines. If only a subset of experiments are reproducible, they
785 should state which ones are omitted from the script and why.
- 786 • At submission time, to preserve anonymity, the authors should release anonymized
787 versions (if applicable).

- 788 • Providing as much information as possible in supplemental material (appended to the
789 paper) is recommended, but including URLs to data and code is permitted.

790 6. Experimental Setting/Details

791 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
792 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
793 results?

794 Answer: [Yes]

795 Justification: These details are provided in Section 4.1.

796 Guidelines:

- 797 • The answer NA means that the paper does not include experiments.
- 798 • The experimental setting should be presented in the core of the paper to a level of detail
799 that is necessary to appreciate the results and make sense of them.
- 800 • The full details can be provided either with the code, in appendix, or as supplemental
801 material.

802 7. Experiment Statistical Significance

803 Question: Does the paper report error bars suitably and correctly defined or other appropriate
804 information about the statistical significance of the experiments?

805 Answer: [Yes]

806 Justification: Section 4.2 includes bootstrap 95% confidence intervals on the main figure
807 and hypothesis tests for specific comparisons between methods.

808 Guidelines:

- 809 • The answer NA means that the paper does not include experiments.
- 810 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
811 dence intervals, or statistical significance tests, at least for the experiments that support
812 the main claims of the paper.
- 813 • The factors of variability that the error bars are capturing should be clearly stated (for
814 example, train/test split, initialization, random drawing of some parameter, or overall
815 run with given experimental conditions).
- 816 • The method for calculating the error bars should be explained (closed form formula,
817 call to a library function, bootstrap, etc.)
- 818 • The assumptions made should be given (e.g., Normally distributed errors).
- 819 • It should be clear whether the error bar is the standard deviation or the standard error
820 of the mean.
- 821 • It is OK to report 1-sigma error bars, but one should state it. The authors should
822 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
823 of Normality of errors is not verified.
- 824 • For asymmetric distributions, the authors should be careful not to show in tables or
825 figures symmetric error bars that would yield results that are out of range (e.g. negative
826 error rates).
- 827 • If error bars are reported in tables or plots, The authors should explain in the text how
828 they were calculated and reference the corresponding figures or tables in the text.

829 8. Experiments Compute Resources

830 Question: For each experiment, does the paper provide sufficient information on the com-
831 puter resources (type of compute workers, memory, time of execution) needed to reproduce
832 the experiments?

833 Answer: [Yes]

834 Justification: Please see Section 4.1.

835 Guidelines:

- 836 • The answer NA means that the paper does not include experiments.
- 837 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
838 or cloud provider, including relevant memory and storage.

- 839
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- 840
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).
- 841
- 842
- 843

844 9. Code Of Ethics

845 Question: Does the research conducted in the paper conform, in every respect, with the
846 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

847 Answer: [Yes]

848 Justification: The research does not involve human subjects and does not introduce new data.
849 Its main impact should be to improve effectiveness and understanding of preference-based
850 post-training.

851 Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
 - If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
 - The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).
- 852
- 853
- 854
- 855
- 856

857 10. Broader Impacts

858 Question: Does the paper discuss both potential positive societal impacts and negative
859 societal impacts of the work performed?

860 Answer: [Yes]

861 Justification: See Section 7

862 Guidelines:

- The answer NA means that there is no societal impact of the work performed.
 - If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
 - Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
 - The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
 - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
 - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).
- 863
- 864
- 865
- 866
- 867
- 868
- 869
- 870
- 871
- 872
- 873
- 874
- 875
- 876
- 877
- 878
- 879
- 880
- 881
- 882
- 883
- 884

885 11. Safeguards

886 Question: Does the paper describe safeguards that have been put in place for responsible
887 release of data or models that have a high risk for misuse (e.g., pretrained language models,
888 image generators, or scraped datasets)?

889 Answer: [NA]

890 Justification: No data or models are released.

891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: The main external resource is the TLDR dataset, which we cite. Its license is CC BY 4.0.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: No new assets are introduced.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and Research with Human Subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972

Answer: [NA]

Justification: The paper does not involve crowdsourcing or research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing or research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.