
Scaling Laws for Reward Model Overoptimization in Direct Alignment Algorithms

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Reinforcement Learning from Human Feedback (RLHF) has been crucial to the
2 recent success of Large Language Models (LLMs), however, it is often a complex
3 and brittle process. In the classical RLHF framework, a reward model is first trained
4 to represent human preferences, which is in turn used by an online reinforcement
5 learning (RL) algorithm to optimize the LLM. A prominent issue with such meth-
6 ods is *reward over-optimization* or *reward hacking*, where performance as measured
7 by the learned proxy reward model increases, but true quality plateaus or even dete-
8riorates. Direct Alignment Algorithms (DAAs) like Direct Preference Optimization
9 have emerged as alternatives to the classical RLHF pipeline by circumventing
10 the reward modeling phase. However, although DAAs do not use separate proxy
11 reward model, they still commonly deteriorate from over-optimization. While
12 the so-called reward hacking phenomenon is not well-defined for DAAs, we still
13 uncover similar trends: at higher KL-budgets, DAA algorithms exhibit similar
14 degradation patterns to their classic RLHF counterparts. In particular, we find that
15 DAA methods deteriorate not only across a wide range of KL-budgets, but also
16 often before even a single epoch of the dataset is completed. Through extensive
17 empirical experimentation, this work formulates and formalizes the reward over-
18 optimization or hacking problem for DAAs and explores its consequences across
19 objectives, training regimes, and model scales.

20 1 Introduction

21 Recent advancements in Large Language Models (LLMs) have broadened their capabilities signifi-
22 cantly, enabling applications in code generation, mathematical reasoning, tool use, and interactive
23 communication. These improvements have popularized LLMs across various domains. Reinforce-
24 ment Learning from Human Feedback (RLHF) has been instrumental in these advances and is now
25 integral to sophisticated LLM training regimes [10, 55]. Before alignment, LLMs, trained on vast text
26 corpora to predict subsequent tokens [45, 8] are often unwieldy and hard to use. Today, leading LLMs
27 incorporate variants of the RLHF framework [14, 68, 36] to align them with human intent, which
28 generally involves a multi-stage process. Specifically, users evaluate model responses to assorted
29 prompts in order to train a reward model that encapsulates human preferences [10, 55, 71, 5, 61].
30 Then, the refined LLM maximizes the expected learned reward function using a reinforcement learn-
31 ing (RL) algorithm [50, 1, 64]. Despite its efficacy, this procedure is complex and computationally
32 intensive, particularly in its latter stages.

33 Goodhart’s Law [25, 11], that “when a measure becomes a target, it ceases to be a good measure”,
34 has often been cited as a core shortcoming of RLHF. Standard RLHF methods optimize a learned, but
35 imperfect reward function which ends up amplifying the reward model’s shortcomings. Empirically,
36 this phenomena was first extensively characterized by Gao et al. [21], who coined the term “reward

37 over-optimization”, and has been seen consistently in recent findings [61, 16, 14]. While reward
 38 over-optimization has been studied in the context of the aforementioned RLHF procedure, recent
 39 contemporary methods for aligning LLMs circumvent the reward learning procedure, necessitating a
 40 new characterization of the over-optimization phenomena.

41 This new broad class of algorithms, which we refer to as Direct Alignment Algorithms (DAAs),
 42 bypass the traditional RLHF pipeline by re-parameterizing the reward model directly through the
 43 optimal policy derived during the reinforcement learning phase. DAA methods, like Direct Preference
 44 Optimization [46], have gained popularity [14, 28] as they often reduce computational demands. Yet,
 45 despite not fitting a reward function, DAAs still exhibit over-optimization trends similar to those of
 46 traditional RLHF methods using a learned reward function. In some sense, this is puzzling: DAAs
 47 can be viewed as simply learning a reward function with supervised learning from which the optimal
 48 policy is deterministically mapped, however more seems to be at play than simple supervised learning.

49 In this work we investigate the over-fitting phenomena present in DAA algorithms through extensive
 50 experimentation. First, we unify a number of different recent methods [46, 67, 4] under the DAA
 51 framework. Then, across different model scales and hyper-parameters, we show that DAAs exhibit a
 52 type of reward over-optimization consistent with that previously observed in RLHF [21]. Specifically,
 53 we find that at different KL-divergence budgets DAAs exhibit degradation patterns similar to those
 54 found in RLHF. Interestingly, we also find that performance within a single epoch is not always
 55 consistent as expected for DAAs. Finally, we explain why this happens by appealing to the under-
 56 constrained nature of the optimization problem used in DAAs.

57 2 Preliminaries

58 In this section, we first outline the core components of the standard RLHF pipeline [71, 55, 5, 41]).
 59 Then, we examine prior literature to characterize the reward over-optimization exhibited by standard
 60 RLHF methods. Finally, we provide a unifying view of direct alignment algorithms (DAAs) which
 61 will guide our analysis of their training dynamics in the next section.

62 2.1 Reinforcement Learning From Human Feedback

63 The standard RLHF pipeline consists of three distinct stages with the goal of aligning the LLM with
 64 human preferences.

65 **Supervised Fine Tuning (SFT):** First, a dataset of prompts x and high-quality answers y are used to
 66 train an LLM for instruction following via maximum likelihood estimation over next-tokens. We
 67 refer to the resultant model as $\pi_{\text{SFT}}(y|x)$ and consider the entire prompt and answer strings to be
 68 single variables.

69 **Reward Modeling:** Second, the SFT model $\pi_{\text{SFT}}(y|x)$ is used to learn a reward function over human
 70 preferences. Specifically, the SFT model is queried to produce pairs of answers $(y_1, y_2) \sim \pi_{\text{SFT}}(y|x)$,
 71 for every prompt x in a dataset. Then, users select their preferred answers, resulting in ranking
 72 $y_w \succ y_l \mid x$ where y_w and y_l are the preferred and dispreferred answers respectively. Typically, user
 73 rankings are assumed to be distributed according to the Bradley-Terry (BT) model [7]

$$p(y_1 \succ y_2 \mid x) = \frac{\exp(r(x, y_1))}{\exp(r(x, y_1)) + \exp(r(x, y_2))} = \sigma(r(x, y_1) - r(x, y_2)) \quad (1)$$

74 where the preference distribution p results from an unobserved latent reward $r(x, y)$, and σ is the
 75 logistic function. Given this model and a dataset of rankings, denoted $\mathcal{D} = \{x^{(i)}, y_w^{(i)}, y_l^{(i)}\}_{i=1}^N$, we
 76 can train a parameterized model $r_\phi(x, y)$ to predict the unobserved reward using maximum likelihood
 77 estimation. This yields the following loss function,

$$\mathcal{L}_{\text{rew}}(r_\phi) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} [\log \sigma(r_\phi(x, y_w) - r_\phi(x, y_l))]. \quad (2)$$

78 **Reinforcement Learning (RL):** The final stage of the standard RLHF pipeline uses the learned reward
 79 model $r_\phi(x, y)$ to update the LLM π_θ with an on-policy RL algorithm like PPO [50], optimizing the
 80 model to provide responses more preferred by human raters. The most common objective is

$$\max_{\pi_\theta} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta(\cdot|x)} [r_\phi(x, y)] - \beta \mathbb{D}_{\text{KL}}[\pi_\theta(y \mid x) \parallel \pi_{\text{ref}}(y|x)] \quad (3)$$

81 which enforces a Kullback-Leibler (KL) divergence penalty with a reference distribution $\pi_{\text{ref}}(y|x)$
 82 (usually taken to be $\pi_{\text{SFT}}(y|x)$) to prevent the LLM π_{θ} from straying too far from its initialization.
 83 Thus, the hyper-parameter β directly trades off exploiting the reward function and deviating from
 84 $\pi_{\text{ref}}(y|x)$.

85 2.2 Reward Exploitation in RLHF

86 Unfortunately, repeating the above procedure without careful tuning of the RL phase can lead to
 87 disastrous performance. This is because in the context of RLHF the LLM policy is optimizing the
 88 surrogate reward estimate $r_{\phi}(x, y)$ and not the true reward function as is often the case in other
 89 domains. Thus, prior works have observed that while the LLM’s expected reward according to
 90 eq. (3) increases the actual quality of the model’s outputs can decrease [54, 43, 9, 34]. This particular
 91 instantiation of the reward exploitation or hacking problem [3] is often referred to as reward “over-
 92 optimization” in RLHF literature and has been studied empirically in both controlled experiments
 93 [21] and user studies [14]. There are two prevailing explanations for why this phenomena occurs.

94 **1. OOD Robustness:** In the classical RLHF pipeline, the RL objective (eq. (3)) is optimized using
 95 on-policy samples from π_{θ} . This means that the reward function is continuously queried using unseen
 96 model samples which are potentially out-of-distribution. Beyond the support of the reward modeling
 97 distribution, r_{ϕ} may assign high reward to sub-par responses, leading the policy to believe it is doing
 98 well when it may not be. While the KL-regularization term is designed to prevent the model from drift-
 99 ing too far out of distribution, this term alone has proven inadequate to prevent reward hacking [21].

100 **2. Reward Mis-specification.** Learned reward functions may exhibit spurious correlations that cause
 101 them to prefer unintended behaviors. While this issue is not at the forefront of LLM research, it is
 102 known to be pervasive in RL [43, 34]. Most efforts to address these problems exist at the intersection
 103 of robustness and offline RL literature [13, 66, 16] and use measures of epistemic uncertainty to
 104 penalize the predicted reward.

105 2.3 Direct Alignment Algorithms

106 Due to its complex multi-step nature, recent works have sought alternatives to the classic RLHF
 107 pipeline. A new class of algorithms, which we broadly classify as Direct Alignment Algorithms
 108 (DAAs), directly update the LLM’s policy π_{θ} using user feedback instead of fitting a reward function
 109 to it and then employing an RL algorithm. Perhaps the most known example is Direct Preference Op-
 110 timization (DPO). DPO as well as other DAAs are derived using the closed form solution to the RLHF
 111 objective in eq. (3) [70], $\pi^*(y|x) \propto \pi_{\text{ref}}(y|x)e^{r(x,y)/\beta}$, where $r(x, y)$ is the ground-truth reward.
 112 By isolating $r(x, y)$ in this relationship and substituting it into the reward optimization objective in
 113 eq. (2), we arrive at a general objective that allows us to train the LLM directly using feedback data:

$$\mathcal{L}_{\text{DAA}}(\pi_{\theta}; \pi_{\text{ref}}) = \mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[g \left(\beta \log \frac{\pi_{\theta}(y_w | x)}{\pi_{\text{ref}}(y_w | x)} - \beta \log \frac{\pi_{\theta}(y_l | x)}{\pi_{\text{ref}}(y_l | x)} \right) \right] \quad (4)$$

114 where g is a convex loss function. Using $g(x) = -\log \sigma(x)$ coincides with the standard
 115 Bradley-Terry model and the original DPO objective. Other methods choose different loss functions:
 116 IPO [4] uses the quadratic objective $g(x) = (x - 1)^2$ and SLiC-HF [67, 38] uses the hinge loss
 117 $g(x) = \max(0, 1 - x)$. Additional objectives were also considered in [59], but due to limited
 118 computational resources, we focus on the three objectives outlined above.

119 Crucially, the DAA approach allows us to recover the optimal policy using a straightforward classifi-
 120 cation loss without the need for learning a reward function, on-policy sampling, or RL, which can be
 121 notoriously difficult to tune and computationally expensive. Because of this, DAAs have emerged as
 122 a popular alternative. However, just like classical RLHF methods, DAAs exhibit strong over-fitting
 123 and even reward-hacking like behaviors. For example, Park et al. [44] show that LLMs trained with
 124 DPO generate responses with increasing length throughout the course of training, but do not improve
 125 in ground-truth win-rate after a certain point. Since DAAs do not explicitly learn a reward function, it
 126 is unclear how “reward-overoptimization” fits into the picture. In this work, we aim to shed some
 127 light on this phenomena in DAAs.

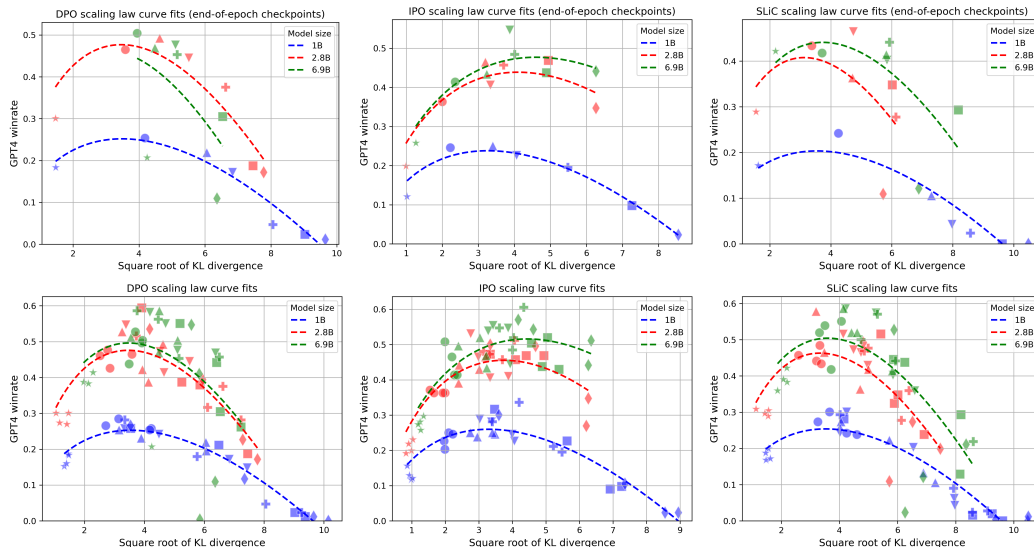


Figure 1: Results on over-optimization in Direct Alignment Algorithms for DPO, IPO and SLiC. Results shows model win-rates over the dataset summary on an evaluation set of prompts as judged by GPT-4. The top row shows final performance after 1 epoch of training, while the second row also includes 4 intermediate checkpoints as well. The fitted dotted curves utilize scaling laws from [21] applied to direct alignment, with GPT4 winrates taking the place of the gold reward model score.

128 3 Empirical Analysis of Overoptimization in DAAs

129 First, we examine the over-optimization problem in DAAs and compare it to those observed in
 130 traditional RLHF methods. All our experiments are carried using the Reddit TL;DR summarization
 131 dataset [55] and the Pythia family of Large Language Models [6].

132 3.1 Evaluating Model-Overoptimization

133 In our first set of experiments we evaluate the reward model over-optimization phenomenon. We
 134 evaluate three training objectives DPO, IPO and SLiC using seven β parameters, representing different
 135 KL budgets at three model size - 1B, 2.8B and 6.9B. Our main results are shown in Fig. 1 which
 136 presents results for different configurations after 1 epoch of training (row 1) and including 4 uniform
 137 intermediate checkpoints (row 2). We include additional results on the training dynamics in Fig. 2,
 138 which shows win rates and KL bounds for intra-epoch training. We present our findings below.

139 **Model Over-Optimization:** We see clear over-optimization for all objectives as performance exhibits
 140 a hump-shaped pattern, where an additional increase in the KL budget leads to decreasing model
 141 performance. Moreover in Fig. 2 we observe similar intra-epoch training dynamics patterns as
 142 configurations with wider KL budgets achieve their best performance after training on only 25% of
 143 the data, after which performance starts decreasing in conjunction with increasing KL divergence
 144 metrics.

145 **Effect of Training Objective:** In the IPO work [4] the authors present theoretical arguments that
 146 due to the monotone sigmoid objective in the DPO formulation, the KL constraint is not effectively
 147 enforced and propose the quadratic fixed-margin loss as an alternative. Across all objectives, there
 148 are clear dependencies between the β parameter and the corresponding KL achieved at the end of
 149 training. While DPO and SLiC exhibit similar performance, IPO indeed seems to be less prone to
 150 over-optimization and in general achieve lower KLs under the same constraint. Our observations
 151 with IPO also align with prior works in preference-based RL and imitation learning where imposing
 152 a fixed margin led to more stable and performant methods [48, 51].

153 **Effect of Model Size:** The results also show strong parameter count scaling effect. The Pythia 1B
 154 model achieves low performance under the same set of constraints it reaches much higher KL values,
 155 while almost immediately exhibiting signs of over-optimization. This behavior holds under all three

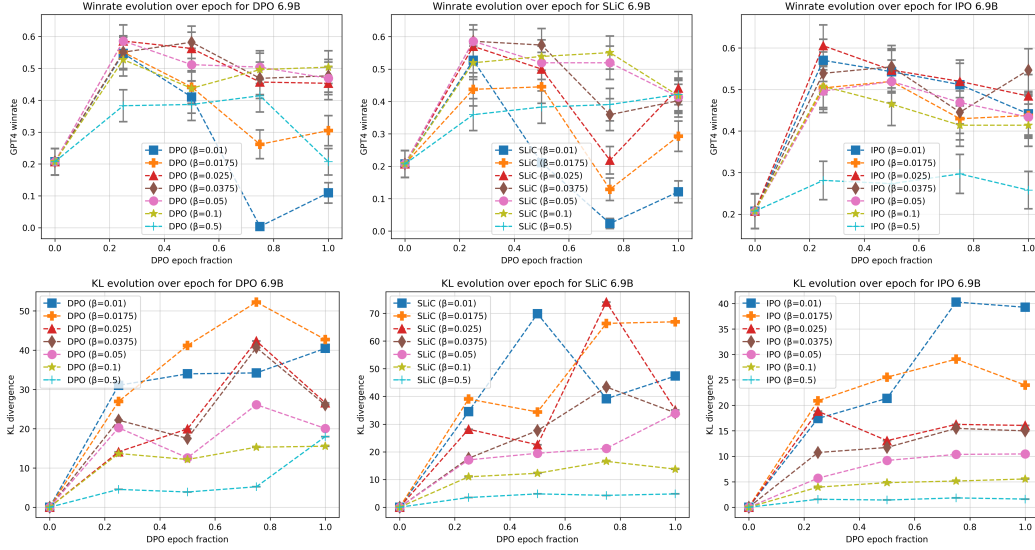


Figure 2: Results on intra-epoch optimization dynamics. The top row shows win-rates against fraction of an epoch so far, while the bottom row shows the corresponding KL values. Under a lower KL constraint most experiments reach their best performance in the first 25% of the epoch and degrade over the remaining of training, while the model deviates from the reference under increasing KL. All models are 6.9B and vary across DPO, SLiC, and IPO loss formulations.

156 objectives. At larger scales, the 6.9B Pythia model tends to exhibit more win-rate - KL trade-offs and
 157 be less prone to over-optimization, with both models significantly outperforming the 1B model. In
 158 the case of the IPO objective, the 6.9B also exhibits significantly better control over the KL objective
 159 and shows little to no over-optimization behavior.

160 3.2 Scaling Law Fits

161 Given we have established a framework for evaluating over-optimization in DAAs and empirically
 162 validated it (section 3.1), we now develop scaling laws for this phenomenon. Previous work in
 163 classical RLHF have established such scaling laws for reward model scores as a function of the KL
 164 divergence between initial and optimized policies [21]. The relevant functional of the reward $R(d)$ is

$$R(d) = d(\alpha - \beta \log d) \quad (5)$$

165 where α, β are constants dependent on the size of the reward model dataset and parameter count,
 166 and $d = \sqrt{D_{\text{KL}}(\pi || \pi_{\text{ref}})}$. As DAAs do not train a proxy reward model, we treat GPT4 winrates over
 167 dataset completions as a proxy for gold reward. Somewhat surprisingly, we find that this scaling law
 168 accurately relates d and winrates for DAAs. Compared to a quadratic fit between $D_{\text{KL}}(\pi || \pi_{\text{ref}})$ and
 169 winrates, this scaling law halves the RMSE. It is worth noting, however, that a quadratic fit between
 170 d and winrates yields similar error compared to Equation 5.

171 3.3 Length Correlations

172 Prior work [44] has shown that the DPO algorithm is prone to length exploitation as it amplifies
 173 verbosity biases in preference datasets. Here we show that length is not the only dimension on which
 174 exploitation can occur. Our experimental results are shown in Fig. 3. On the left we show results
 175 for the 2.8B Pythia model with standard training plus the length-regularization approach from [44].
 176 Both approaches suffer from over-optimization, but the dynamics differ depending on the KL budget.
 177 Moreover, even though the regularized model achieves higher win rates on a length-correct basis,
 178 it under-performs the model trained with the standard objective in the lower KL constraint region.

179 Recent work [27] has also shown that DAAs prioritize features of the data based on their complexity
 180 and prevalence (with length a clear example of human datasets). [44] further showed that models
 181 trained with the DPO algorithm extrapolate significantly based on length. We extend this analysis in

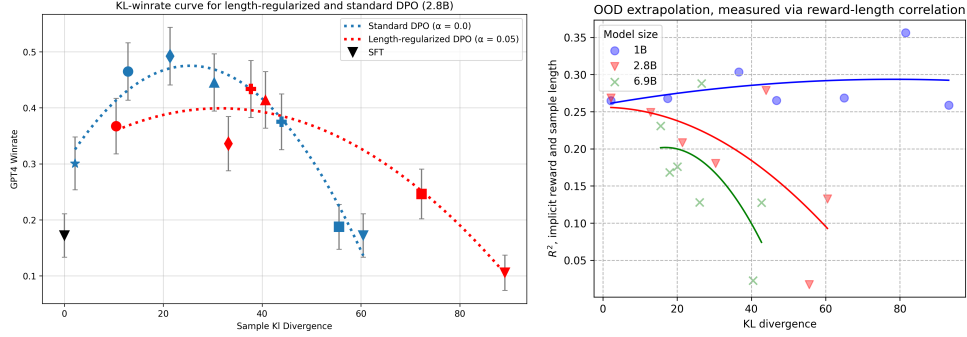


Figure 3: **Left:** KL budget versus win-rates (over dataset human answer) with and without length-regularization [44]. While including a length-correction in the optimization objective changes the KL-win-rate Pareto Frontier, it does not alleviate reward over-optimization and might even exacerbate it. **Right:** Scaling behaviour for length extrapolation - smaller capacity models (either by size or KL budget) extrapolate more strongly on simpler features such as length.

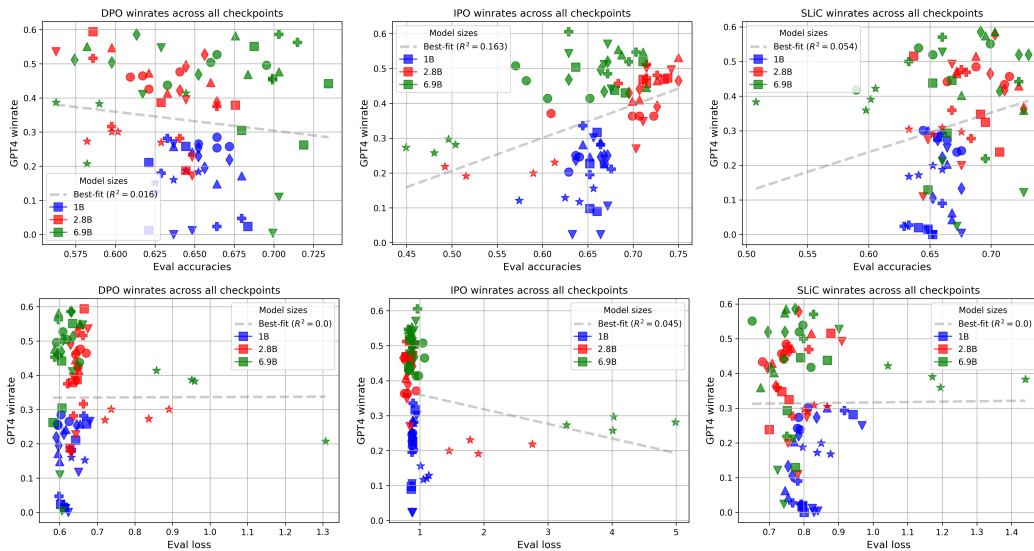


Figure 4: **Top:** We plot the DAA implicit reward accuracy in preference classification versus win rates. **Bottom:** DAA optimization loss versus checkpoint win rates. Model training statistics, do not exhibit strong relationship with downstream performance.

182 Fig. 3 (right). We consider a linear regression of the form

$$\log \frac{\pi_{\theta}(y^{(i)}|x^{(i)})}{\pi_{ref}(y^{(i)}|x^{(i)})} = \hat{\gamma}|y^{(i)}| + \epsilon^{(i)} \quad (6)$$

183 where $x^{(i)}$ are held-out prompts and $y^{(i)}$ are samples from the corresponding model between the
 184 DPO implicit reward and length. We fit a different regression for each model size and checkpoint
 185 and plot the corresponding R^2 values. We observe two main effects; first there is a clear scaling
 186 law behaviour. Weaker models extrapolate across the simple length feature to a much higher degree
 187 than stronger ones. This is especially clear comparing the behaviour of the Pythia 1B versus the
 188 2.8B and 6.9B models. Second, we see significant effects based on the KL budget - under a smaller
 189 budget all model sizes exhibit higher extrapolation behaviour. Based on these results we formulate
 190 the hypothesis that under limited capacity, either from model capability or limited KL budgets the
 191 model will extrapolate more strongly based on simpler features, which can lead to OOD issues.

192 3.4 Reward Metrics Correlations

193 Prior works have measured reward model quality in ranking settings by classification accuracy. We
 194 evaluate the relationship between the DAA implicit reward model accuracy and policy performance

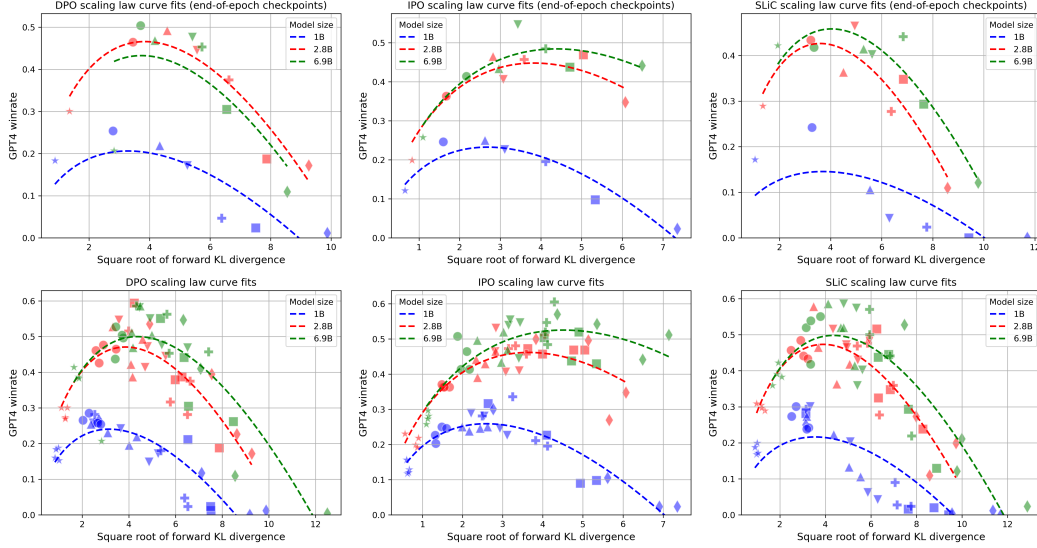


Figure 5: Over-optimization results for $\sqrt{\text{Forward KL}}$ vs. winrates. The top row shows final performance after 1 epoch of training, while the second row also includes 4 intermediate checkpoints. The fitted dotted curves are scaling laws from [21] applied to DAAs, with GPT4 winrates taking the place of the gold reward model score.

195 in Figure 4. The DPO and SLiC algorithms exhibit little to no correlation between reward model
 196 accuracy and downstream model performance. The IPO model shows a weak positive relationship,
 197 but upon further examinations, this is entirely due to model size scaling - stronger models both
 198 fit the data better and produce better generations as well, however within each particular model
 199 size, there is no discernible relationship between the DAA implicit reward accuracy and the actual
 200 policy performance. Similar observations hold when comparing the empirical DAA loss with model
 201 performance, which is contrary to observations in supervised pre-training and instruction tuning [30].

202 3.5 Decreasing Likelihoods and Model Performance

203 A number of recent works have observed that the implicit DAA rewards of both preferred and
 204 dis-preferred responses decrease doing training, which may be counter-intuitive. In [47] the authors
 205 make a counter-point that in offline training of DAAs π_{ref} is usually pre-trained with SFT on the
 206 preferred response and thus

$$\mathbb{E}_{p_{\mathcal{D}}(y_w|x)} \left[\log \frac{\pi_{\theta}(y_w|x)}{\pi_{\text{ref}}(y_w|x)} \right] \approx \mathbb{E}_{\pi_{\text{ref}}(y_w|x)} \left[\log \frac{\pi_{\theta}(y_w|x)}{\pi_{\text{ref}}(y_w|x)} \right] = -\mathbb{D}_{\text{KL}}[\pi_{\text{ref}}(y|x) \parallel \pi_{\theta}(y|x)] \quad (7)$$

207 where $p_{\mathcal{D}}(y^w|x)$ is the dataset distribution of preferred answers. That is the expected implicit reward
 208 represent a forward KL divergence between the reference policy and the optimisation policy, thus it
 209 is expected to be negative and decrease with training as the optimisation model moves away from
 210 the reference. In this section we study whether this empirical phenomenon presents a challenge for
 211 DAA learning. Similar to Fig. 1 we plot the win rates against the square-root-transformed (negative)
 212 expected implicit reward of the preferred response (evaluated on a held-out evaluation dataset), which
 213 as stated above approximates the (square-root-transformed) forward KL $\mathbb{D}_{\text{KL}}[\pi_{\text{ref}}(y|x) \parallel \pi_{\theta}(y|x)]$.
 214 Results are included in Fig. 5, which follow closely the pattern in Fig. 1 with performance initially
 215 increasing before it starts dipping down after a certain threshold. This indicates that under the
 216 standard DAA training pipeline decreasing likelihoods are not necessary an issue for performance,
 217 and are even necessary for improvement, but exhibit a non-linear over-optimization dynamics.

218 4 Reward Exploitation in Direct Alignment Algorithms

219 While the phenomena observed in the previous section echo those observed in classical RLHF, their
 220 underlying causes may be distinct. Reward over-optimization in classical RLHF is largely attributed to
 221 querying a proxy reward function that is potentially OOD, while DAAs do not train a separate reward

222 model. Instead, DAAs are generally understood as fitting an “implicit” reward model to preference
 223 data with the parameterization $r_\theta(x, y) = \beta \log \frac{\pi_\theta(y|x)}{\pi_{\text{ref}}(y|x)}$ using the objective in eq. (2). Therefore, the
 224 OOD behavior of the policy is inextricably linked to the OOD behavior of the implicit reward model.
 225 Here we demonstrate that the reward modeling objective used is heavily under-constrained, allowing
 226 for a potentially large number of solutions that can place weight on OOD responses. This is especially
 227 problematic for DAAs which deterministically map the optimal policy from the “implicit” reward.

228 **Rank Deficiency with Finite Preferences.** In DAAs, the language modeling problem is treated
 229 as contextual bandit. However, the space of possible prompts $x \in \mathcal{X}$ and answers $y \in \mathcal{Y}$ are both
 230 exponentially large in sequence length. However, as highlighted by Tang et al. [59], DAAs often
 231 assume full support of the reference distribution when mapping from the implicit reward to the optimal
 232 policy π by eq. (10). However, in practice such coverage is impossible. Instead, preference datasets
 233 cover a minuscule portion of the prompt-response space. Unfortunately, as DAA objectives are not
 234 strictly convex, this means that many optima of eq. (4) can place a high weight on OOD responses.

235 We can demonstrate this using the regression interpretation from Hejna et al. [23]. Consider re-
 236 writing the DAA objective from eq. (4) in terms of preference query vectors q which select the win
 237 response pair (x, y^w) and the loss response pair (x, y^l) from the prompt-response space. Each vector
 238 q represents both the preferred and dis-preferred responses, with the entree corresponding to (x, y^w)
 239 being +1 and the entree corresponding to (x, y^l) being -1. We can then write the generalized DAA
 240 loss function with finite preference data as

$$\mathcal{L}_{\text{DAA}}(\pi_\theta, \mathcal{D}) = \sum_{i=1}^{|\mathcal{D}|} g(\beta q_i^\top (\log \pi(y|x) - \log \pi_{\text{ref}}(y|x))), \text{ where } q_i[x, y] = \begin{cases} 1 & \text{if } (x^{(i)}, y_w^{(i)}) = (x, y) \\ -1 & \text{if } (x^{(i)}, y_l^{(i)}) = (x, y) \\ 0 & \text{otherwise} \end{cases}$$

241 where the policy π is a single vector of size $|\mathcal{X} \times \mathcal{Y}|$. In practice, the distributional constraint of π
 242 also applies. Choosing g to be the negative log sigmoid above recovers DPO with finite preferences,
 243 but also logistic regression with a data matrix Q of shape $|\mathcal{D}|$ by $|\mathcal{X} \times \mathcal{Y}|$ constructed by stacking
 244 the aforementioned query vectors q . As $|\mathcal{X} \times \mathcal{Y}| \gg |\mathcal{D}|$, this matrix is likely to have a non-trivial
 245 null-space, making the problem not strictly convex. Thus, there are many possible policies π that
 246 can thus achieve the same optima, some of which can place a high weight on out-of-distribution
 247 responses [23, 69]. For more details and constructions, we defer to Hejna et al. [23].

248 **Understanding OOD behavior for DAA algorithms with a Toy MDP:** To illustrate that DAA
 249 algorithms, in general and not an artifact of training LLM’s, end up placing probability mass on
 250 OOD sequences during training we design a simple Tree MDP (shown in Figure 6) to mimic the
 251 token-level MDP in LLMs. We use a dataset containing a single preference between two trajectories
 252 and follow the standard procedure of running SFT on preferred responses before updating an RNN
 253 policy using a DAA. Figure 7 shows that even in this simple setup, popular DAAs (DPO/IPO/SLiC)
 254 end up extrapolating incorrectly out of distribution revealing a fundamental shortcoming. Unlike in
 255 standard RLHF, the non-strict convexity of the reward function in DAAs ends up directly affecting
 256 the policy. Detailed experimental details can be found in Appendix E.

257 5 Related Work

258 Broadly, over-optimization has been a widely studied phenomena across different settings [60, 18].
 259 Over-fitting can be characterized as over-optimization in the supervised learning setting [39, 32],
 260 which can harm generalization [19, 12, 24] or lead to susceptibility to adversarial attacks [56, 37, 15].
 261 Reward hacking in reinforcement learning (RL) [54], where an agent maximizes its reward through
 262 behavior that deviates from the intended goal, can be viewed as a different type of over-optimization,
 263 commonly observed in prior work [43, 3, 22].

264 We study over-optimization in the context of aligning LLMs with human feedback, for which the
 265 most common approach is RLHF as outlined in section 2.1. Similar RLHF techniques were originally
 266 pioneered for control [31, 2, 10]. Standard RLHF methods suffer from both potential over-fitting
 267 of the reward function and reward exploitation by the RL algorithm. Several works have considered
 268 how to reduce over-fitting or increase the robustness of learned reward functions using ensembles
 269 [13, 66, 16] or data smoothing [69]. Other approaches, like Moskovitz et al. [40] consider how
 270 reward exploitation can be reduced by using different optimization techniques in the RL stage. Much

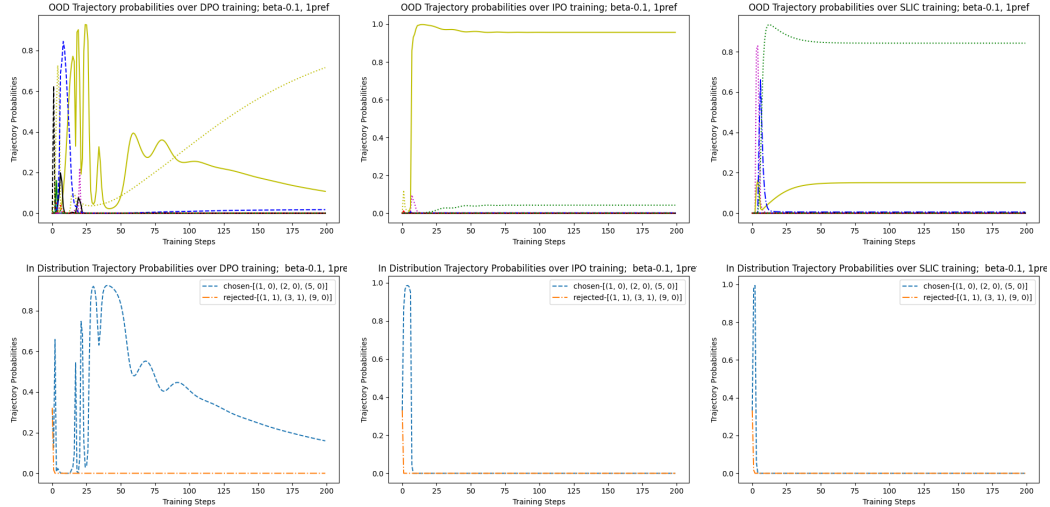


Figure 7: (Top row) Probability of OOD trajectories. DAA algorithms end up placing a substantial probability mass of some of the OOD trajectories during training. (Bottom row) Probability of in-distribution (preference-pair) trajectories decrease during training.

271 of this work is motivated by Gao et al. [21], which first characterized and provided scaling laws
 272 for over-optimization in RLHF.

273 Unlike Gao et al. [21], we consider the over-
 274 optimization problem in DAAs, which differ
 275 significantly from the standard RLHF pipeline.
 276 Different DAAs have been derived theoretically
 277 [47, 46, 67, 4, 63], and applied to problems be-
 278 yond language modeling like image generation
 279 [62] and control [23]. In all of these scenarios,
 280 over-optimization problems have persisted. Park
 281 et al. [44] show that DAAs commonly over-fit to
 282 length and the expense of performance, which
 283 has been linked to inherent bias in training data
 284 [53, 29]. Other works have tried to allow DAAs
 285 to use more types of data like demonstrations
 286 [49] or ratings [17] to get better performance.
 287 Recently, incorporating online data has proven
 288 critical to improving performance [65, 26, 57].
 289 Concurrent to our work, Tang et al. [58] study
 290 the differences between offline DAAs and stan-
 291 dard RLHF methods. Unlike us, they focus on
 292 comparisons with online sampling whereas we
 293 focus on the purely offline setting.

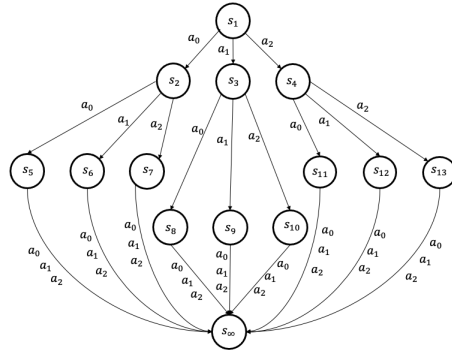


Figure 6: An illustration of the Tree MDP. At each state, we can choose one of 3 actions (a_0, a_1, a_2), which deterministically maps to the next state. Furthermore, all the leaf nodes in this tree MDP, transition to the terminal absorbing state s_∞ , irrespective of the chosen action

294 6 Conclusion

295 In this work we present an analysis of the over-optimization problem in Direct Alignment Algorithms.
 296 Through extensive experimentation on different algorithms (DPO, IPO, SLIC) and at different model
 297 scales (1B, 2.8B, 6.9B), we observe consistent over-optimization trends at different KL-divergence
 298 budgets. While our analysis is a first step, it is not a complete picture of understanding the over-
 299 optimization phenomena. More work can be done characterizing this effect at larger model scales,
 300 which we were unable to do due to computational limitations. Nevertheless, we believe our work
 301 sheds light on important problems in Direct Alignment Algorithms that can spur future research.

302 **References**

- 303 [1] A. Ahmadian, C. Cremer, M. Gallé, M. Fadaee, J. Kreutzer, A. Üstün, and S. Hooker. Back to
304 basics: Revisiting reinforce style optimization for learning from human feedback in llms. *arXiv*
305 *preprint arXiv:2402.14740*, 2024.
- 306 [2] R. Akrou, M. Schoenauer, and M. Sebag. Preference-based policy learning. In *Joint European*
307 *Conference on Machine Learning and Knowledge Discovery in Databases*, 2011.
- 308 [3] D. Amodei, C. Olah, J. Steinhardt, P. Christiano, J. Schulman, and D. Mané. Concrete problems
309 in ai safety. *arXiv preprint arXiv:1606.06565*, 2016.
- 310 [4] M. G. Azar, M. Rowland, B. Piot, D. Guo, D. Calandriello, M. Valko, and R. Munos. A general
311 theoretical paradigm to understand learning from human preferences, 2023.
- 312 [5] Y. Bai, A. Jones, K. Ndousse, A. Askell, A. Chen, N. DasSarma, D. Drain, S. Fort, D. Ganguli,
313 T. Henighan, N. Joseph, S. Kadavath, J. Kernion, T. Conerly, S. El-Showk, N. Elhage, Z. Hatfield-
314 Dodds, D. Hernandez, T. Hume, S. Johnston, S. Kravec, L. Lovitt, N. Nanda, C. Olsson,
315 D. Amodei, T. Brown, J. Clark, S. McCandlish, C. Olah, B. Mann, and J. Kaplan. Training a
316 helpful and harmless assistant with reinforcement learning from human feedback, 2022.
- 317 [6] S. Biderman, H. Schoelkopf, Q. Anthony, H. Bradley, K. O’Brien, E. Hallahan, M. A. Khan,
318 S. Purohit, U. S. Prashanth, E. Raff, A. Skowron, L. Sutawika, and O. van der Wal. Pythia: A
319 suite for analyzing large language models across training and scaling, 2023.
- 320 [7] R. A. Bradley and M. E. Terry. Rank analysis of incomplete block designs: I. the method of
321 paired comparisons. *Biometrika*, 39(3/4):324–345, 1952. doi: <https://doi.org/10.2307/2334029>.
- 322 [8] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam,
323 G. Sastry, A. Askell, et al. Language models are few-shot learners. *Advances in neural*
324 *information processing systems*, 33:1877–1901, 2020.
- 325 [9] S. Casper, X. Davies, C. Shi, T. K. Gilbert, J. Scheurer, J. Rando, R. Freedman, T. Korbak,
326 D. Lindner, P. Freire, T. Wang, S. Marks, C.-R. Segerie, M. Carroll, A. Peng, P. Christoffersen,
327 M. Damani, S. Slocum, U. Anwar, A. Siththaranjan, M. Nadeau, E. J. Michaud, J. Pfau,
328 D. Krasheninnikov, X. Chen, L. Langosco, P. Hase, E. Bıyık, A. Dragan, D. Krueger, D. Sadigh,
329 and D. Hadfield-Menell. Open problems and fundamental limitations of reinforcement learning
330 from human feedback, 2023.
- 331 [10] P. F. Christiano, J. Leike, T. Brown, M. Martic, S. Legg, and D. Amodei. Deep reinforcement
332 learning from human preferences. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus,
333 S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Sys-*
334 *tems*, volume 30. Curran Associates, Inc., 2017. URL [https://proceedings.neurips.cc/
335 paper_files/paper/2017/file/d5e2c0adad503c91f91df240d0cd4e49-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/d5e2c0adad503c91f91df240d0cd4e49-Paper.pdf).
- 336 [11] J. Clark and D. Amodei. Faulty reward functions in the wild, 2016. URL [https://openai.
337 com/research/faulty-reward-functions](https://openai.com/research/faulty-reward-functions).
- 338 [12] K. Cobbe, O. Klimov, C. Hesse, T. Kim, and J. Schulman. Quantifying generalization in
339 reinforcement learning. In K. Chaudhuri and R. Salakhutdinov, editors, *Proceedings of the*
340 *36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine*
341 *Learning Research*, pages 1282–1289. PMLR, 09–15 Jun 2019. URL [https://proceedings.
342 mlr.press/v97/cobbe19a.html](https://proceedings.mlr.press/v97/cobbe19a.html).
- 343 [13] T. Coste, U. Anwar, R. Kirk, and D. Krueger. Reward model ensembles help mitigate overopti-
344 mization, 2023.
- 345 [14] Y. Dubois, X. Li, R. Taori, T. Zhang, I. Gulrajani, J. Ba, C. Guestrin, P. Liang, and T. B.
346 Hashimoto. AlpacaFarm: A simulation framework for methods that learn from human feedback,
347 2024.
- 348 [15] J. Ebrahimi, D. Lowd, and D. Dou. On adversarial examples for character-level neural machine
349 translation. *arXiv preprint arXiv:1806.09030*, 2018.

- 350 [16] J. Eisenstein, C. Nagpal, A. Agarwal, A. Beirami, A. D’Amour, D. Dvijotham, A. Fisch,
351 K. Heller, S. Pfohl, D. Ramachandran, P. Shaw, and J. Berant. Helping or herding? reward
352 model ensembles mitigate but do not eliminate reward hacking, 2023.
- 353 [17] K. Ethayarajh, W. Xu, N. Muennighoff, D. Jurafsky, and D. Kiela. Kto: Model alignment as
354 prospect theoretic optimization. *arXiv preprint arXiv:2402.01306*, 2024.
- 355 [18] T. Everitt, V. Krakovna, L. Orseau, M. Hutter, and S. Legg. Reinforcement learning with a
356 corrupted reward channel. *arXiv preprint arXiv:1705.08417*, 2017.
- 357 [19] J. Farebrother, M. C. Machado, and M. Bowling. Generalization and regularization in dqn.
358 *arXiv preprint arXiv:1810.00123*, 2018.
- 359 [20] S. Fujimoto, D. Meger, and D. Precup. Off-policy deep reinforcement learning without explo-
360 ration, 2019.
- 361 [21] L. Gao, J. Schulman, and J. Hilton. Scaling laws for reward model overoptimization. *Interna-
362 tional Conference on machine Learning*, 2023.
- 363 [22] D. Hadfield-Menell, S. Milli, P. Abbeel, S. J. Russell, and A. Dragan. Inverse reward design. In
364 I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett,
365 editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates,
366 Inc., 2017. URL [https://proceedings.neurips.cc/paper_files/paper/2017/file/
367 32fdab6559cdfa4f167f8c31b9199643-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/32fdab6559cdfa4f167f8c31b9199643-Paper.pdf).
- 368 [23] J. Hejna, R. Rafailov, H. Sikchi, C. Finn, S. Niekum, W. B. Knox, and D. Sadigh. Contrastive
369 preference learning: Learning from human feedback without reinforcement learning. In
370 *The Twelfth International Conference on Learning Representations*, 2024. URL [https://
371 openreview.net/forum?id=iX1RjVQDj](https://openreview.net/forum?id=iX1RjVQDj).
- 372 [24] D. Hernandez, J. Kaplan, T. Henighan, and S. McCandlish. Scaling laws for transfer. *arXiv
373 preprint arXiv:2102.01293*, 2021.
- 374 [25] K. Hoskin. The ‘awful idea of accountability’: inscribing people into the measurement of
375 objects. *Accountability: Power, ethos and the technologies of managing*, 265, 1996.
- 376 [26] A. Hosseini, X. Yuan, N. Malkin, A. Courville, A. Sordoni, and R. Agarwal. V-star: Training
377 verifiers for self-taught reasoners. *arXiv preprint arXiv:2402.06457*, 2024.
- 378 [27] S. Im and Y. Li. Understanding the learning dynamics of alignment with human feedback, 2024.
- 379 [28] A. Q. Jiang, A. Sablayrolles, A. Roux, A. Mensch, B. Savary, C. Bamford, D. S. Chaplot,
380 D. de las Casas, E. B. Hanna, F. Bressand, G. Lengyel, G. Bour, G. Lample, L. R. Lavaud,
381 L. Saulnier, M.-A. Lachaux, P. Stock, S. Subramanian, S. Yang, S. Antoniak, T. L. Scao,
382 T. Gervet, T. Lavril, T. Wang, T. Lacroix, and W. E. Sayed. Mixtral of experts, 2024.
- 383 [29] S. Kabir, D. N. Udo-Imeh, B. Kou, and T. Zhang. Who answers it better? an in-depth analysis
384 of chatgpt and stack overflow answers to software engineering questions, 2023.
- 385 [30] J. Kaplan, S. McCandlish, T. Henighan, T. B. Brown, B. Chess, R. Child, S. Gray, A. Radford,
386 J. Wu, and D. Amodei. Scaling laws for neural language models, 2020.
- 387 [31] W. B. Knox and P. Stone. Tamer: Training an agent manually via evaluative reinforcement. In
388 *2008 7th IEEE international conference on development and learning*, pages 292–297. IEEE,
389 2008.
- 390 [32] V. Krakovna and R. Kumar. Classifying specification problems as variants of
391 goodhart’s law, 8 2019. URL [https://vkrakovna.wordpress.com/2019/08/19/
392 classifying-specification-problems-as-variants-of-goodharts-law/](https://vkrakovna.wordpress.com/2019/08/19/classifying-specification-problems-as-variants-of-goodharts-law/).
- 393 [33] A. Kumar, A. Zhou, G. Tucker, and S. Levine. Conservative q-learning for offline reinforcement
394 learning. *Advances in Neural Information Processing Systems*, 33:1179–1191, 2020.
- 395 [34] N. Lambert and R. Calandra. The alignment ceiling: Objective mismatch in reinforcement
396 learning from human feedback, 2023.

- 397 [35] S. Levine, A. Kumar, G. Tucker, and J. Fu. Offline reinforcement learning: Tutorial, review,
398 and perspectives on open problems, 2020.
- 399 [36] P. Liang, R. Bommasani, T. Lee, D. Tsipras, D. Soylu, M. Yasunaga, Y. Zhang, D. Narayanan,
400 Y. Wu, A. Kumar, B. Newman, B. Yuan, B. Yan, C. Zhang, C. Cosgrove, C. D. Manning,
401 C. Ré, D. Acosta-Navas, D. A. Hudson, E. Zelikman, E. Durmus, F. Ladhak, F. Rong, H. Ren,
402 H. Yao, J. Wang, K. Santhanam, L. Orr, L. Zheng, M. Yuksekogul, M. Suzgun, N. Kim,
403 N. Guha, N. Chatterji, O. Khattab, P. Henderson, Q. Huang, R. Chi, S. M. Xie, S. Santurkar,
404 S. Ganguli, T. Hashimoto, T. Icard, T. Zhang, V. Chaudhary, W. Wang, X. Li, Y. Mai, Y. Zhang,
405 and Y. Koreeda. Holistic evaluation of language models, 2023.
- 406 [37] Y.-C. Lin, Z.-W. Hong, Y.-H. Liao, M.-L. Shih, M.-Y. Liu, and M. Sun. Tactics of adversarial
407 attack on deep reinforcement learning agents. *arXiv preprint arXiv:1703.06748*, 2017.
- 408 [38] T. Liu, Y. Zhao, R. Joshi, M. Khalman, M. Saleh, P. J. Liu, and J. Liu. Statistical rejection
409 sampling improves preference optimization, 2024.
- 410 [39] D. Manheim and S. Garrabrant. Categorizing variants of goodhart’s law, 2019.
- 411 [40] T. Moskovitz, A. K. Singh, D. Strouse, T. Sandholm, R. Salakhutdinov, A. Dragan, and S. M.
412 McAleer. Confronting reward model overoptimization with constrained RLHF. In *The Twelfth
413 International Conference on Learning Representations*, 2024. URL [https://openreview.
414 net/forum?id=gkfUvn0fLU](https://openreview.net/forum?id=gkfUvn0fLU).
- 415 [41] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal,
416 K. Slama, A. Ray, J. Schulman, J. Hilton, F. Kelton, L. Miller, M. Simens, A. Askell, P. Welinder,
417 P. F. Christiano, J. Leike, and R. Lowe. Training language models to follow instructions with
418 human feedback. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh,
419 editors, *Advances in Neural Information Processing Systems*, volume 35, pages 27730–27744.
420 Curran Associates, Inc., 2022. URL [https://proceedings.neurips.cc/paper_files/
421 paper/2022/file/b1efde53be364a73914f58805a001731-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/b1efde53be364a73914f58805a001731-Paper-Conference.pdf).
- 422 [42] A. Pal, D. Karkhanis, S. Dooley, M. Roberts, S. Naidu, and C. White. Smaug: Fixing failure
423 modes of preference optimisation with dpo-positive. *arXiv preprint arXiv:2402.13228*, 2024.
- 424 [43] A. Pan, K. Bhatia, and J. Steinhardt. The effects of reward misspecification: Mapping and
425 mitigating misaligned models. *International Conference on Learning Representations*, 2022.
- 426 [44] R. Park, R. Rafailov, S. Ermon, and C. Finn. Disentangling length from quality in direct
427 preference optimization, 2024.
- 428 [45] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever. Language models are
429 unsupervised multitask learners, 2019. OpenAI.
- 430 [46] R. Rafailov, A. Sharma, E. Mitchell, C. D. Manning, S. Ermon, and C. Finn. Direct preference
431 optimization: Your language model is secretly a reward model. In *Thirty-seventh Conference on
432 Neural Information Processing Systems*, 2023. URL <https://arxiv.org/abs/2305.18290>.
- 433 [47] R. Rafailov, J. Hejna, R. Park, and C. Finn. From r to q^* : Your language model is secretly a
434 q -function, 2024.
- 435 [48] N. D. Ratliff, J. A. Bagnell, and M. A. Zinkevich. Maximum margin planning. In *Proceedings
436 of the 23rd international conference on Machine learning*, pages 729–736, 2006.
- 437 [49] M. Rita, F. Strub, R. Chaabouni, P. Michel, E. Dupoux, and O. Pietquin. Countering reward
438 over-optimization in llm with demonstration-guided reinforcement learning. *arXiv preprint
439 arXiv:2404.19409*, 2024.
- 440 [50] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization
441 algorithms, 2017.
- 442 [51] H. Sikchi, A. Saran, W. Goo, and S. Niekum. A ranking game for imitation learning. *arXiv
443 preprint arXiv:2202.03481*, 2022.

- 444 [52] H. Sikchi, Q. Zheng, A. Zhang, and S. Niekum. Dual rl: Unification and new methods for
445 reinforcement and imitation learning. *arXiv preprint arXiv:2302.08560*, 2023.
- 446 [53] P. Singhal, T. Goyal, J. Xu, and G. Durrett. A long way to go: Investigating length correlations
447 in rlhf, 2023.
- 448 [54] J. Skalse, N. H. R. Howe, D. Krasheninnikov, and D. Krueger. Defining and characterizing
449 reward hacking, 2022.
- 450 [55] N. Stiennon, L. Ouyang, J. Wu, D. M. Ziegler, R. Lowe, C. Voss, A. Radford, D. Amodei, and
451 P. Christiano. Learning to summarize from human feedback, 2022.
- 452 [56] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus.
453 Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.
- 454 [57] F. Tajwar, A. Singh, A. Sharma, R. Rafailov, J. Schneider, T. Xie, S. Ermon, C. Finn, and
455 A. Kumar. Preference fine-tuning of llms should leverage suboptimal, on-policy data. *arXiv
456 preprint arXiv:2404.14367*, 2024.
- 457 [58] Y. Tang, D. Z. Guo, Z. Zheng, D. Calandriello, Y. Cao, E. Tarassov, R. Munos, B. Á. Pires,
458 M. Valko, Y. Cheng, et al. Understanding the performance gap between online and offline
459 alignment algorithms. *arXiv preprint arXiv:2405.08448*, 2024.
- 460 [59] Y. Tang, Z. D. Guo, Z. Zheng, D. Calandriello, R. Munos, M. Rowland, P. H. Richemond,
461 M. Valko, B. Ávila Pires, and B. Piot. Generalized preference optimization: A unified approach
462 to offline alignment, 2024.
- 463 [60] J. Taylor. Quantilizers: A safer alternative to maximizers for limited optimization. In *Workshops
464 at the Thirtieth AAAI Conference on Artificial Intelligence*, 2016.
- 465 [61] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal,
466 E. Hambro, F. Azhar, et al. Llama: Open and efficient foundation language models. *arXiv
467 preprint arXiv:2302.13971*, 2023.
- 468 [62] B. Wallace, M. Dang, R. Rafailov, L. Zhou, A. Lou, S. Purushwalkam, S. Ermon, C. Xiong,
469 S. Joty, and N. Naik. Diffusion model alignment using direct preference optimization, 2023.
- 470 [63] J. Watson, S. Huang, and N. Heess. Coherent soft imitation learning. In *Thirty-seventh
471 Conference on Neural Information Processing Systems*, 2023. URL [https://openreview.
472 net/forum?id=kCCD8d2aEu](https://openreview.net/forum?id=kCCD8d2aEu).
- 473 [64] R. J. Williams. Simple statistical gradient-following algorithms for connectionist rein-
474 forcement learning. *Mach. Learn.*, 8(3–4):229–256, may 1992. ISSN 0885-6125. doi:
475 10.1007/BF00992696. URL <https://doi.org/10.1007/BF00992696>.
- 476 [65] R. Yuanzhe Pang, W. Yuan, K. Cho, H. He, S. Sukhbaatar, and J. Weston. Iterative reasoning
477 preference optimization. *arXiv e-prints*, pages arXiv–2404, 2024.
- 478 [66] Y. Zhai, H. Zhang, Y. Lei, Y. Yu, K. Xu, D. Feng, B. Ding, and H. Wang. Uncertainty-penalized
479 reinforcement learning from human feedback with diverse reward lora ensembles, 2023.
- 480 [67] Y. Zhao, R. Joshi, T. Liu, M. Khalman, M. Saleh, and P. J. Liu. Slic-hf: Sequence likelihood
481 calibration with human feedback. *arXiv preprint arXiv:2305.10425*, 2023.
- 482 [68] L. Zheng, W.-L. Chiang, Y. Sheng, S. Zhuang, Z. Wu, Y. Zhuang, Z. Lin, Z. Li, D. Li, E. P.
483 Xing, H. Zhang, J. E. Gonzalez, and I. Stoica. Judging llm-as-a-judge with mt-bench and
484 chatbot arena. *Conference on Neural Information Processing Systems Track on Datasets and
485 Benchmarks.*, 2023.
- 486 [69] B. Zhu, M. I. Jordan, and J. Jiao. Iterative data smoothing: Mitigating reward overfitting and
487 overoptimization in rlhf. *arXiv preprint arXiv:2401.16335*, 2024.
- 488 [70] B. D. Ziebart. *Modeling purposeful adaptive behavior with the principle of maximum causal
489 entropy*. Carnegie Mellon University, 2010.
- 490 [71] D. M. Ziegler, N. Stiennon, J. Wu, T. B. Brown, A. Radford, D. Amodei, P. Christiano, and
491 G. Irving. Fine-tuning language models from human preferences, 2020.

492 **A Limitations and Societal Impacts**

493 Our discussion highlights a number of issues with direct alignment algorithms used widely as means
494 to align to human values. This work has mostly focused on pointing out those issues along with a
495 theoretical underpinning of the issue, but does not provide a way to resolve these issues. We still
496 assume an underlying model of human preferences, which is an ongoing research area as no model is
497 perfect in explaining the ways humans give preferences. Our work aims to drive the push towards
498 better alignment algorithms that do not overoptimize and generate models that safe to be deployed in
499 our society. We believe only through understanding and demonstrating the shortcomings of current
500 methods we can develop better alignment methods.

501 **B Experiment Details**

502 We largely follow the DPO setup unless otherwise mentioned and build on their code
503 (<https://github.com/eric-mitchell/direct-preference-optimization>) without changing any hyperparame-
504 ters unless otherwise mentioned.

505 For all DAA experiments, we used the curated OpenAI TL;DR dataset with 92K preferred-dispreferred
506 summary completions [55]. Each prompt is a Reddit post belonging to one of several topic forums,
507 with title/post metadata included. 256 prompts sampled from the held-out set are used for all
508 evaluations (e.g. loss, accuracy, KL, winrates, length), with temperature 1.0 and max length 512.

509 Model sizes include 1B, 2.8B, and 6.9B and were initialized from the base Pythia pre-trained weights.
510 All models underwent supervised fine-tuning on TL;DR prior to direct alignment. Across all SFT
511 and DAA runs, we used a batch size of 128 (8 gradient accumulation steps), and RMSProp with a
512 learning rate of 0.5×10^{-6} (linear warmup for 150 steps) for 1 epoch. 1B models were trained on 2
513 NVIDIA A40 GPUs, 2.8B models were trained on 4 NVIDIA A40 GPUs, and 6.9B models were
514 trained on 4 NVIDIA A100 GPUs. All evaluations were computed with "gpt-4-turbo-2024-04-09" as
515 judge, with random positional flips to avoid known bias.

516 **C Appendix A: Complete Intra-Epoch Training Dynamics**

517 This appendix contains similar intra-epoch KL divergence and winrate evolution results as in Fig. 2,
518 across all model sizes.

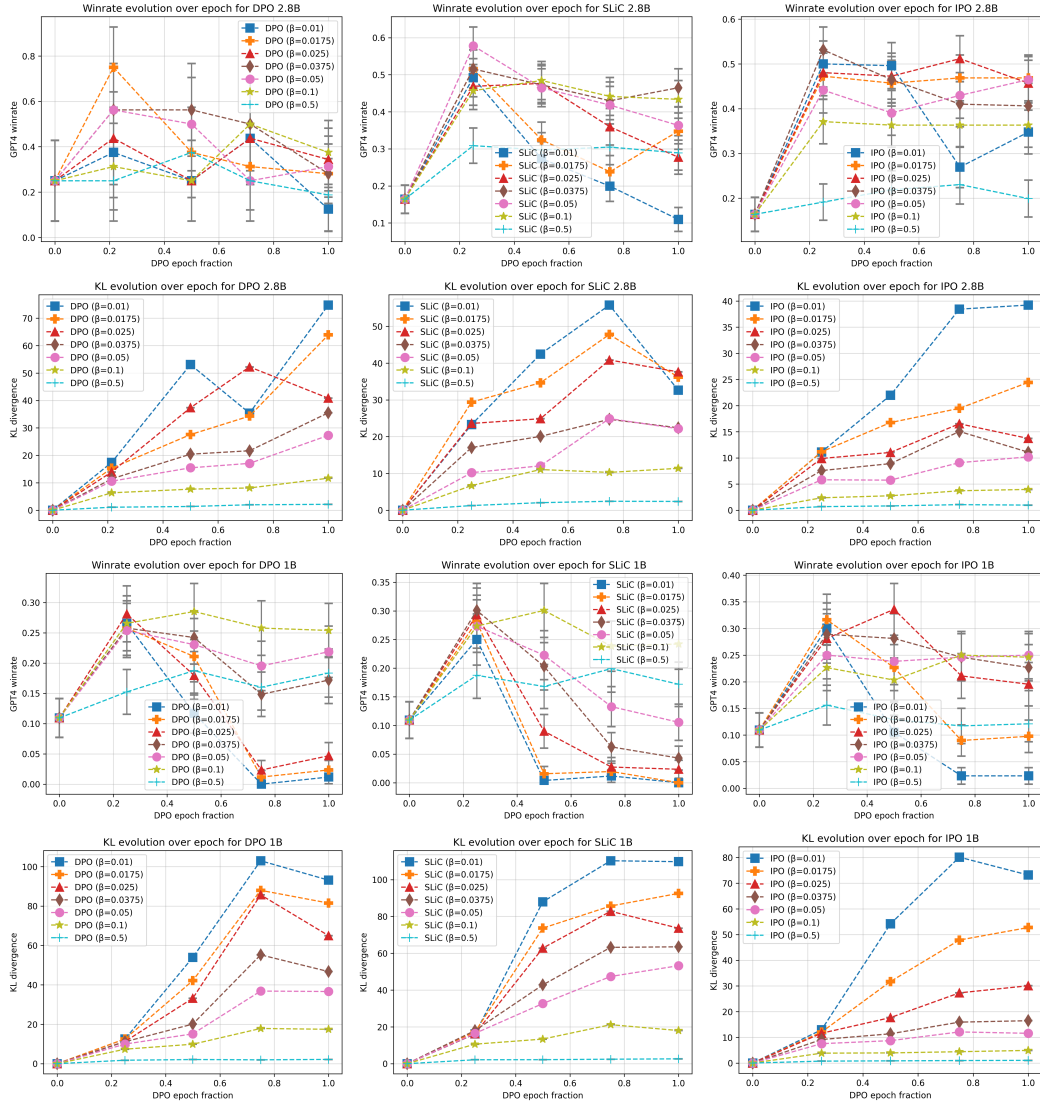


Figure 8: KL divergence and GPT4 winrate evolution for 2.8B and 1B models across DPO, SLiC, and IPO losses. Similar to the 6.9B models, performance tends to degrade after the first quarter epoch, particularly under a low KL budget, while KL increases almost monotonically.

519 **D Overoptimization from the lens of Implicit Bootstrapping**

520 Reward over-optimization is well understood in the classical RLHF setting, with a consensus that
 521 is driven by two main components - using a proxy reward function that is trained on limited data
 522 and continuous querying with new, potentially OOD samples during PPO training. At first glance
 523 none of these conditions hold in DAAs as we do not train a separate proxy reward model or generate
 524 new data during training. Therefore, understanding reward over-optimization in DAAs requires a
 525 new theory. We will base our analysis on [47] using the token-level MDP and corresponding (soft)
 526 Q-learning formulation. Consider the class of dense per-token reward functions $r_\theta(x, y_{\leq i})$, where
 527 $y_{\leq i}$ denotes the first i tokens of y , with sequence level-reward $r_\theta(x, y) = \sum_{i=1}^{|y|} r_\theta(x, y_{\leq i})$. This
 528 is a strictly more general class than the sparse reward function which returns a single score at the
 529 end of the sequence, since we can set all intermediate rewards as 0. Within the framework of [47]
 530 given a DAA-trained policy π_θ , there exists a dense per-token reward r_θ , that minimizes the reward
 531 modelling objective in Eq. 2 and satisfy the below.

532 The (soft) Bellman Equation holds:

$$Q^*(y_i, (x, y_{<i})) = \begin{cases} r(x, y_{\leq i}) + \beta \log \pi_{\text{ref}}(y_i | (x, y_{<i})) + V^*((x, y_{\leq i})), & \text{if } y_i \text{ is not EOS} \\ r(x, y_{\leq i}) + \beta \log \pi_{\text{ref}}(y_i | (x, y_{<i})), & \text{if } y_i \text{ is EOS} \end{cases} \quad (8)$$

533 where V^* is the corresponding soft-value function:

$$V^*((x, y_{<i})) = \beta \log \sum_{y \in |V|} e^{Q^*(y, (x, y_{<i}))/\beta} \quad (9)$$

534 then the DAA policy π_θ satisfies:

$$\pi_\theta(y_i | (x, y_{<i})) = \exp\left(\frac{1}{\beta} Q^*(y_i, (x, y_{<i})) - V^*((x, y_{<i}))\right) \quad (10)$$

535 in this interpretation, the LLM logits $l_\theta[i] = Q^*(y_i, (x, y_{<i}))/\beta$ represent Q-values. With a direct
 536 substitution we then have

$$Q^*(y_i, (x, y_{<i})) = r(x, y_{\leq i}) + \beta \log \pi_{\text{ref}}(y_i | (x, y_{<i})) + \underbrace{\beta \log \sum_{y_i \in |V|} e^{Q^*(y, (x, y_{<i}))/\beta}}_{\text{OOD bootstrapping}} \quad (11)$$

537 That is in this framework DAAs may suffer from the classical OOD bootstrapping issue in offline
 538 RL [20, 35, 33, 52]. In this case even though the objective is trained fully offline we still effectively
 539 query the model on the values of unseen tokens. This interpretation also provides further insight into
 540 the effect of the β coefficient and the training dynamics. For small values of beta the estimate

$$\beta \log \sum_{y_i \in |V|} e^{Q^*(y, (x, y_{<i}))/\beta} \approx \max_{y \in |V|} Q^*(y, (x, y_{<i})) \quad (12)$$

541 that is smaller parameter values yield a more optimistic estimate, which results in higher level of
 542 OOD bootstrapping. This interpretation would also explain the somewhat counter-intuitive results of
 543 section 3.4. While the implicit reward function can adequately fit and model the data, the resulting
 544 LLM might behave sub-optimally, due to OOD bootstrapping in the corresponding Q-value estimate.

545 E Understanding Behavior of DAAs on OOD sequences

546 We have established that common DAA objectives allow for placing a high-likelihood on OOD
547 data. In practice, while one might expect the likelihood of preferred responses to increase during
548 training, it has been observed that algorithms like DPO decrease the likelihood of both the preferred
549 and dis-preferred responses [42]. In fact, this is expected from a max-entropy RL perspective [47].
550 Since the total probability mass must sum to one, the probability of OOD responses must increase
551 during the course of training. A small amount of extrapolation may be necessary to reach the optimal
552 policy, however, too much is potentially detrimental to performance. Because they are not adequately
553 constrained to the reference distribution, current DAA objectives allow this to happen.

554 To understand how DAAs allocate probability mass out of distribution, we use a toy Markov Decision
555 Process (MDP), that mimics the LLM setting. The MDP is modeled as a tree, originating from a
556 single start state, featuring deterministic transitions. The Toy MDP is illustrated in fig. 6.

557 E.1 Designing a toy LLM MDP

558 The MDP is modeled as a tree, originating from a single start state. This configuration mirrors the
559 token-level MDP in Direct Preference Optimization (DPO) [47], or the scenario where both preferred
560 and dispreferred responses are conditioned on the same prompt in the broader Large Language Model
561 alignment context. Each leaf node in the MDP transitions deterministically to a terminal absorbing
562 state, regardless of the action taken. The deterministic transitions resemble the LLM setting, where
563 the current state is represented by the sequence of encountered tokens (s_1, s_2, \dots, s_i) , and the action
564 corresponds to predicting the next word s_{i+1} from the vocabulary, given the context. In this simplified
565 MDP, the deterministic transition is akin to a concatenation function, advancing the state to the next
566 step $(s_1, s_2, \dots, s_i, s_{i+1})$. Employing a toy MDP enables us to systematically evaluate the trajectory
567 probabilities for all feasible paths within the MDP, shedding light on the allocation of probability
568 mass by Direct Alignment Algorithms (DAAs) towards out-of-distribution (OOD) trajectories.

569 **The Experimental Setup.** We adhere to the standard direct alignment protocol [46][41], encompass-
570 ing two key stages:

- 571 1. **Supervised Fine-tuning (SFT) / Behavioral Cloning (BC):** This phase involves fine-
572 tuning the policy based on a limited number of trajectories. Specifically, we utilize
573 three demonstrations for SFT: $(s_1, a_0, s_2, a_0, s_5, a_0, s_\infty)$, $(s_1, a_1, s_3, a_1, s_9, a_0, s_\infty)$, and
574 $(s_1, a_2, s_4, a_2, s_{13}, a_2, s_\infty)$.
- 575 2. **Alignment with Preferences:** In this stage, preferences extracted from trajectories
576 are employed to align the policy. Notably, we have only one preference available:
577 $(s_1, a_1, s_3, a_1, s_9, a_0, s_\infty) \succ (s_1, a_0, s_2, a_0, s_5, a_0, s_\infty)$. This deliberate constraint exag-
578 gerates a scenario with limited data, enabling us to gauge the probability mass allocated
579 to out-of-distribution (OOD) trajectories under such conditions. Insights garnered from
580 this exaggerated low-data scenario hold relevance for Large Language Model (LLM) set-
581 tings where preference datasets are notably smaller compared to the scale of LLM models
582 deployed.

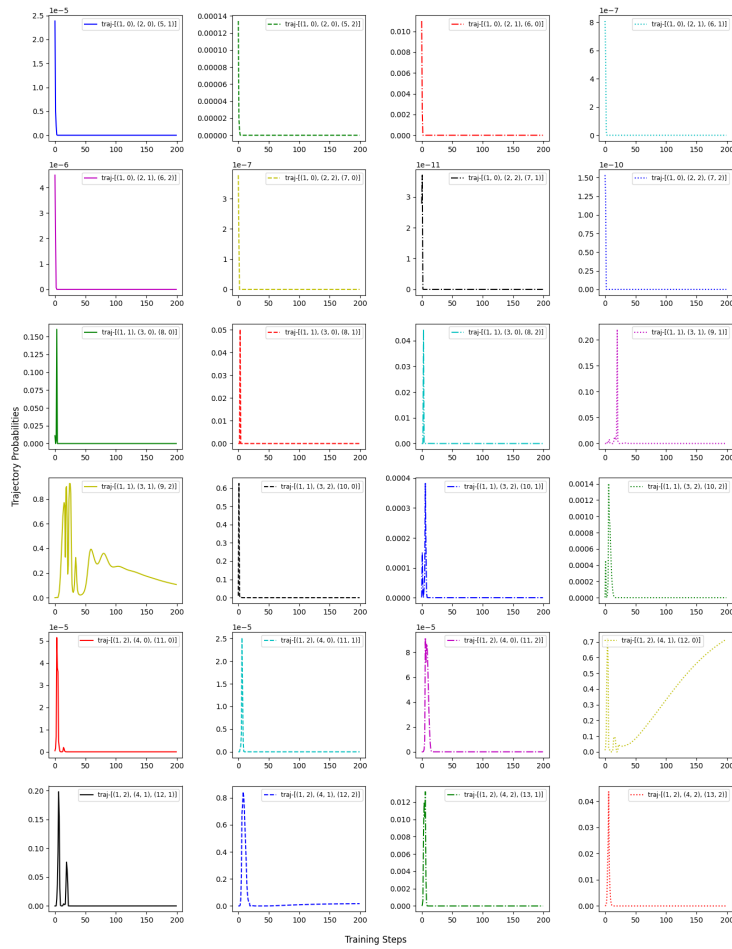
583 We utilize a Recurrent Neural Network (RNN) policy to navigate through the MDP, facilitating a
584 closer resemblance to real-world language modeling scenarios.

585 Subsequently, we explore three distinct direct alignment loss functions: Direct Preference Optimiza-
586 tion (DPO) [46], Identity Preference Optimization (IPO) [4], and Sequence Likelihood Calibration
587 (SLiC) [67]. Additionally, we investigate how the selection of the KL penalty coefficient β influences
588 the distribution of probability mass on OOD trajectories. This exploration encompasses three values
589 of β : (0.01, 0.1, 0.5).

590 In general, the plots illustrate that Direct Alignment Algorithms (DAAs) tend to allocate a significant
591 proportion of the probability mass to out-of-distribution (OOD) trajectories during the alignment
592 process. While Figure ?? may suggest that Direct Preference Optimization (DPO) can retain a
593 substantial amount of probability mass on the selected trajectory in the preference dataset, it's
594 noteworthy that the plots for DPO exhibit considerable noise. To provide further insight, Figure 18
595 displays the plots resulting from three additional repetitions of the DPO experiment. This elucidates
596 the unconstrained nature of the DPO problem: multiple solutions exist for the DPO loss, each

597 distributing varying amounts of probability mass to OOD trajectories. In the experiments with IPO
598 and SLiC, it's observed that the probability mass allocated to in-distribution trajectories diminishes
599 substantially over the course of training. Notably, the probability mass becomes concentrated on a
600 select few out-of-distribution trajectories. Moreover, consistent trends are discernible across various
601 values of β . All our experiments with Toy-MDP can be found in the following figures 12, 9, 15, 13,
602 10, 16, 14, 11, 17.

OOD Trajectory probabilities over DPO training; beta=0.1, 1pref



In Distribution Trajectory Probabilities over DPO training; beta=0.1, 1pref

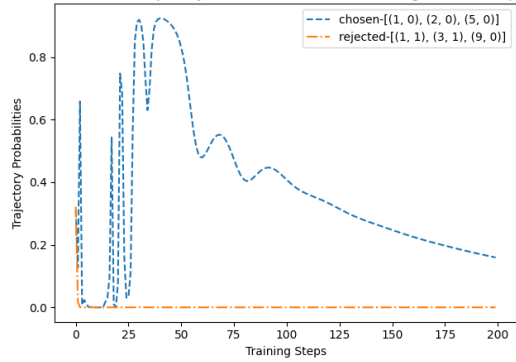
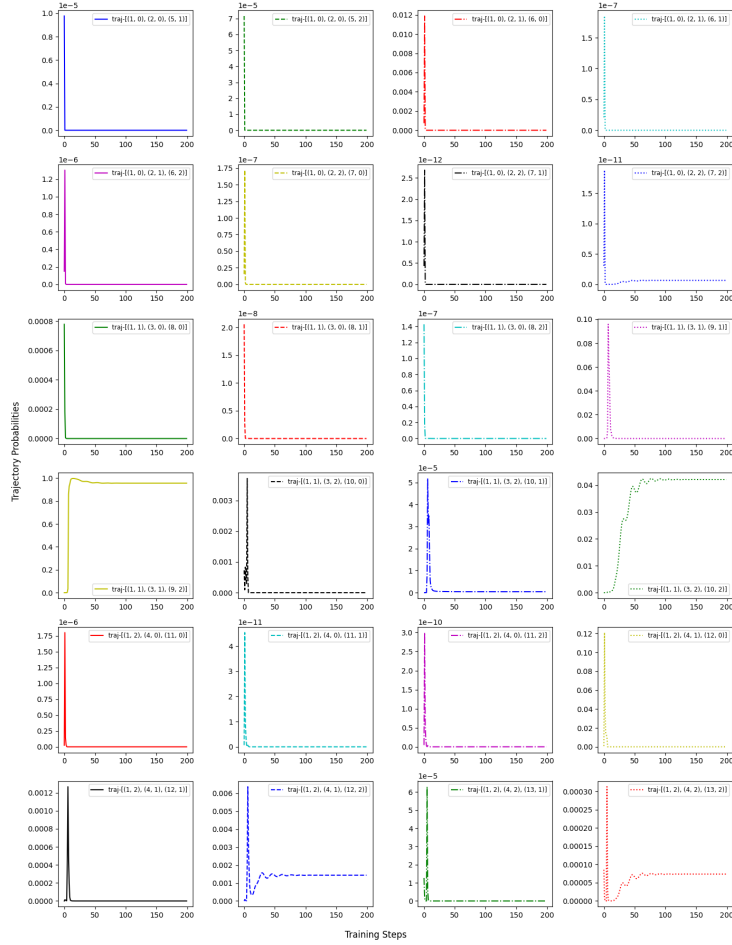


Figure 9: Trajectory probabilities throughout DPO training. The top plot shows how the probability mass of different OOD trajectories, changes throughout training. The bottom plot shows how the probability mass of the trajectories in our preference dataset (size 1) changes over training

OOD Trajectory probabilities over IPO training; beta-0.1, 1pref



In Distribution Trajectory Probabilities over IPO training; beta-0.1, 1pref

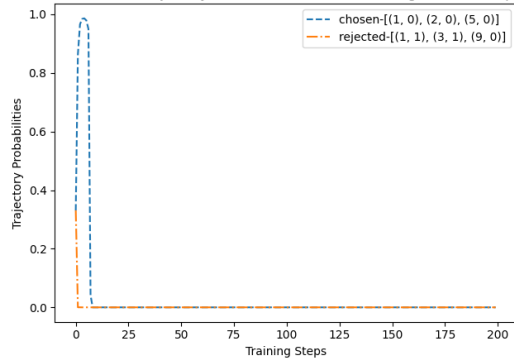
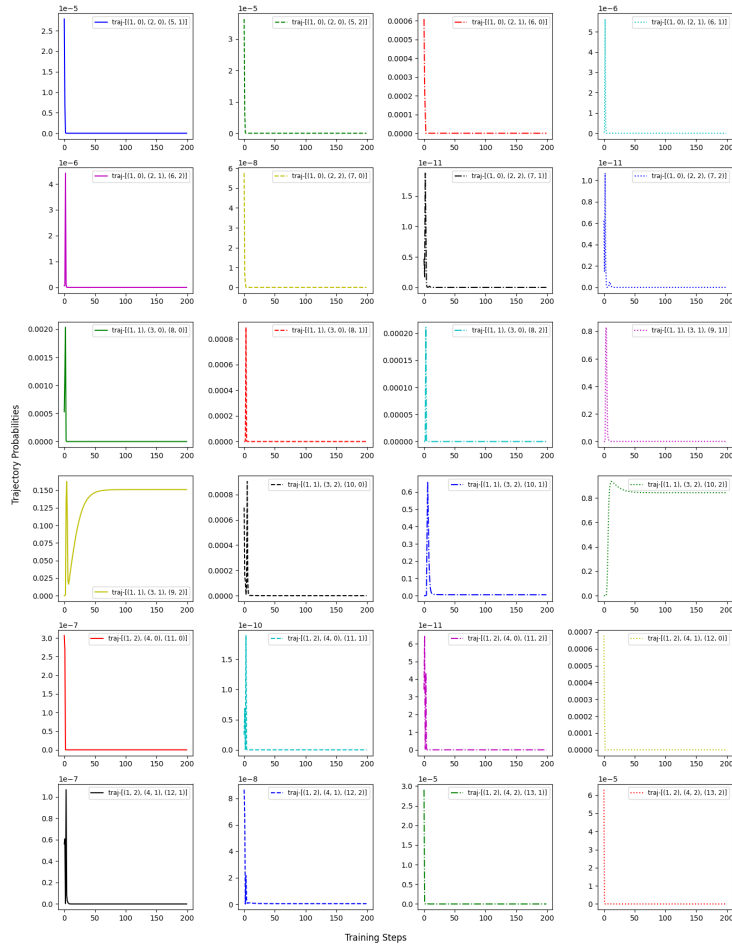


Figure 10: Trajectory probabilities throughout IPO training. The top plot shows how the probability mass of different OOD trajectories, changes throughout training. The bottom plot shows how the probability mass of the trajectories in our preference dataset (size 1) changes over training

OOD Trajectory probabilities over SLiC training: beta-0.1, 1pref



In Distribution Trajectory Probabilities over SLiC training: beta-0.1, 1pref

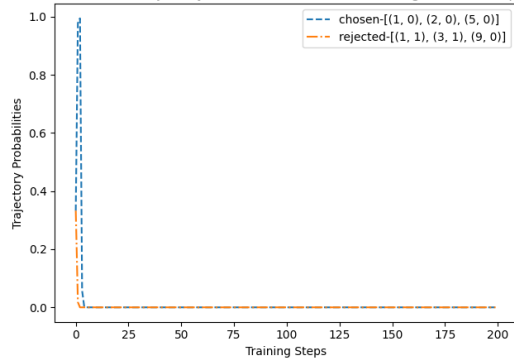
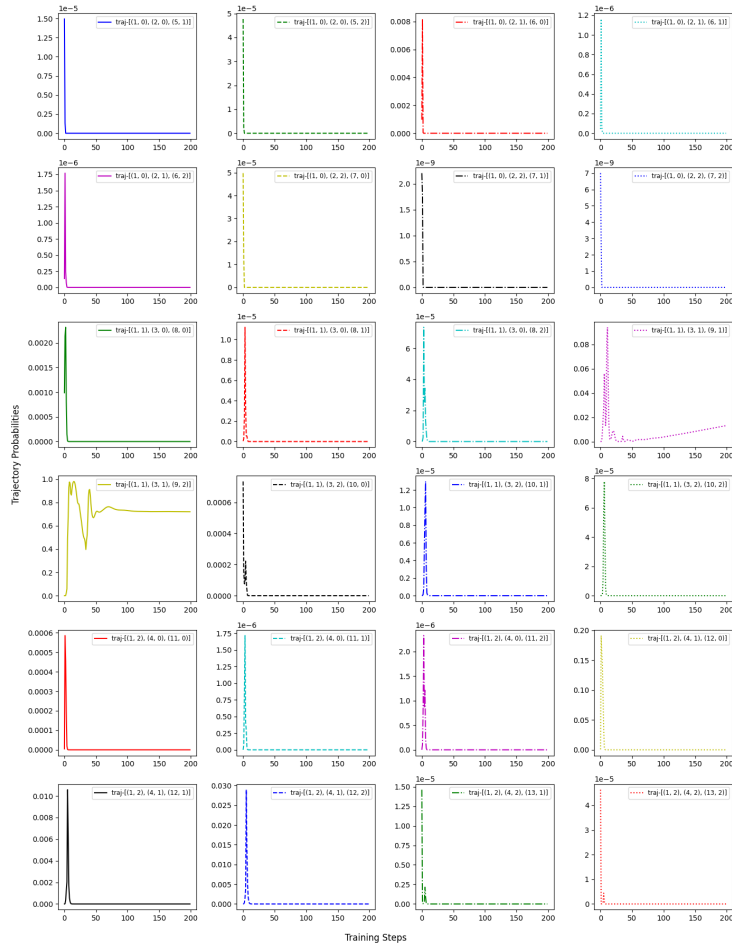


Figure 11: Trajectory probabilities throughout SLiC training. The top plot shows how the probability mass of different OOD trajectories, changes throughout training. The bottom plot shows how the probability mass of the trajectories in our preference dataset (size 1) changes over training

OOD Trajectory probabilities over DPO training; beta=0.01, 1prf



In Distribution Trajectory Probabilities over DPO training; beta=0.01, 1prf

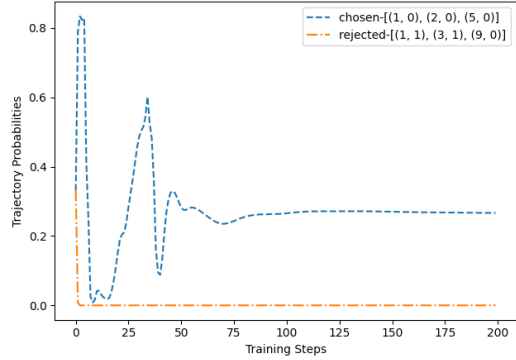
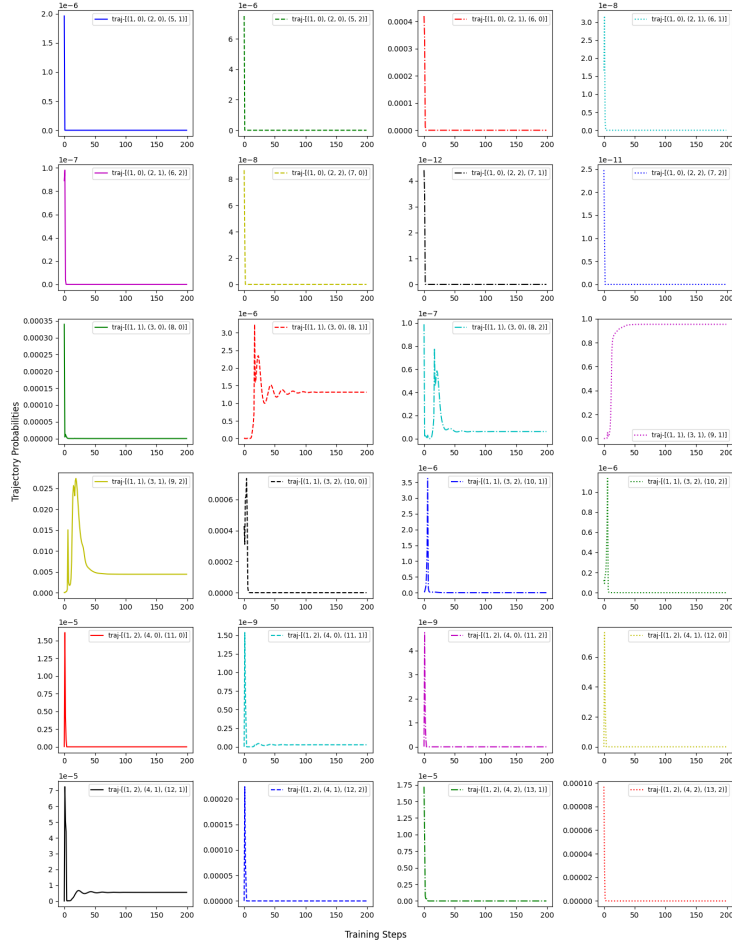


Figure 12: Trajectory probabilities throughout DPO training. The top plot shows how the probability mass of different OOD trajectories, changes throughout training. The bottom plot shows how the probability mass of the trajectories in our preference dataset (size 1) changes over training

OOD Trajectory probabilities over IPO training; beta=0.01, 1pref



In Distribution Trajectory Probabilities over IPO training; beta=0.01, 1pref

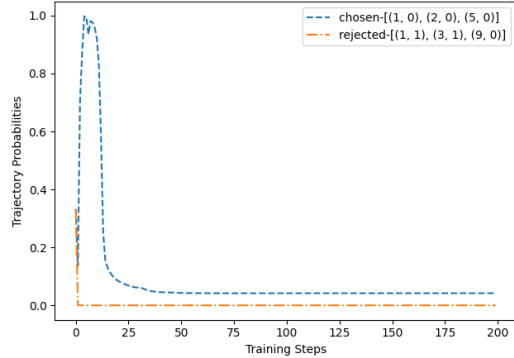
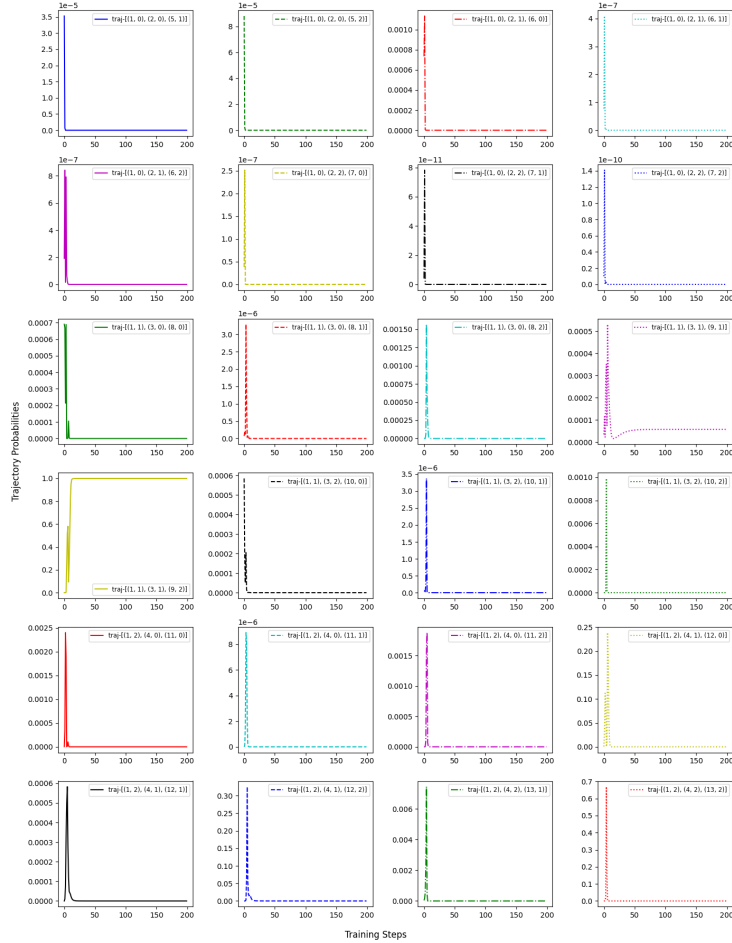


Figure 13: Trajectory probabilities throughout IPO training. The top plot shows how the probability mass of different OOD trajectories, changes throughout training. The bottom plot shows how the probability mass of the trajectories in our preference dataset (size 1) changes over training

OOD Trajectory probabilities over SLiC training; beta=0.01, 1pref



In Distribution Trajectory Probabilities over SLiC training; beta=0.01, 1pr

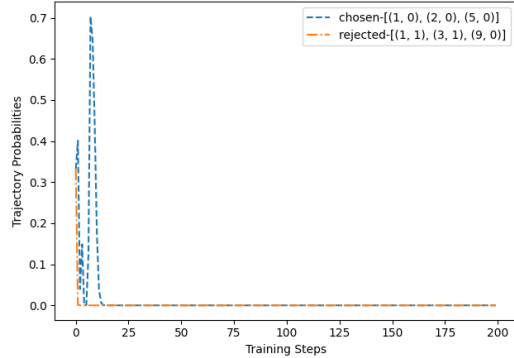
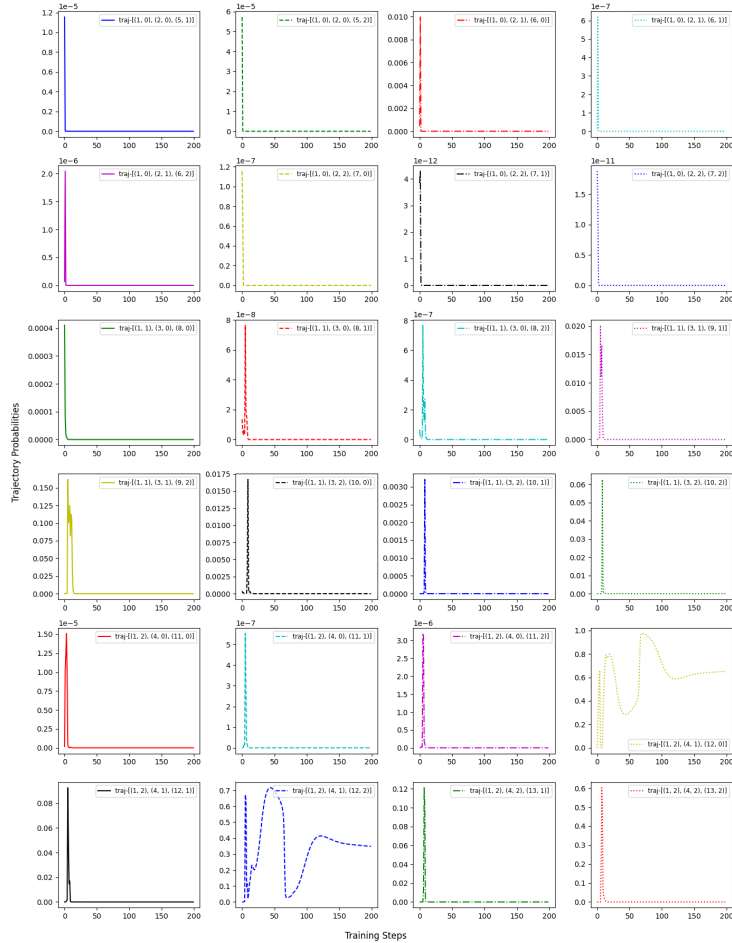


Figure 14: Trajectory probabilities throughout SLiC training. The top plot shows how the probability mass of different OOD trajectories, changes throughout training. The bottom plot shows how the probability mass of the trajectories in our preference dataset (size 1) changes over training

OOD Trajectory probabilities over DPO training; beta=0.5, 1pref



In Distribution Trajectory Probabilities over DPO training; beta=0.5, 1pre

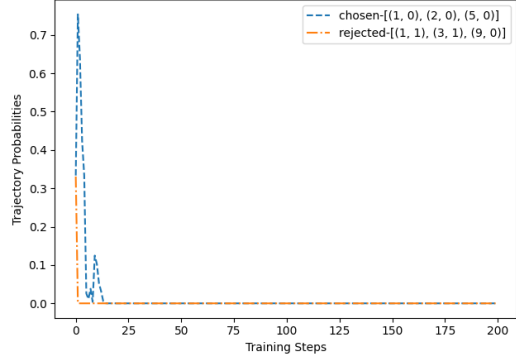
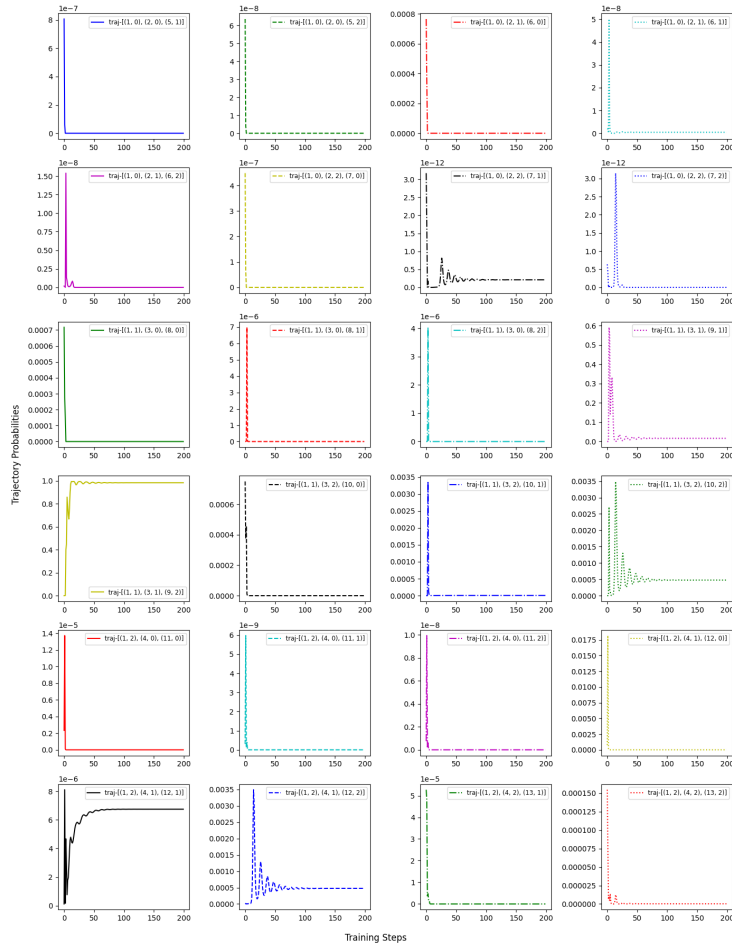


Figure 15: Trajectory probabilities throughout DPO training. The top plot shows how the probability mass of different OOD trajectories, changes throughout training. The bottom plot shows how the probability mass of the trajectories in our preference dataset (size 1) changes over training

OOD Trajectory probabilities over IPO training; beta-0.5, 1pref



In Distribution Trajectory Probabilities over IPO training; beta-0.5, 1pref

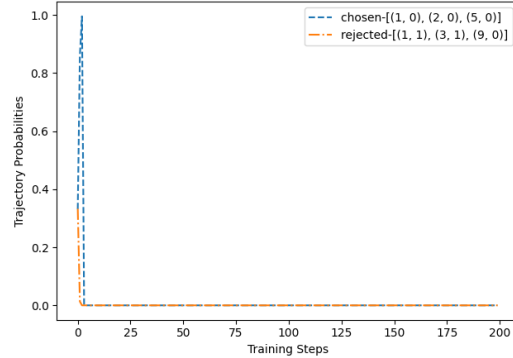


Figure 16: Trajectory probabilities throughout IPO training. The top plot shows how the probability mass of different OOD trajectories, changes throughout training. The bottom plot shows how the probability mass of the trajectories in our preference dataset (size 1) changes over training

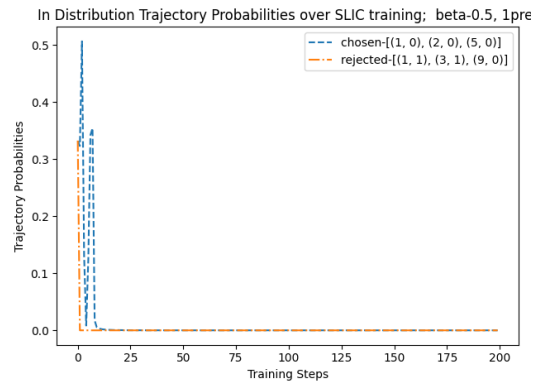
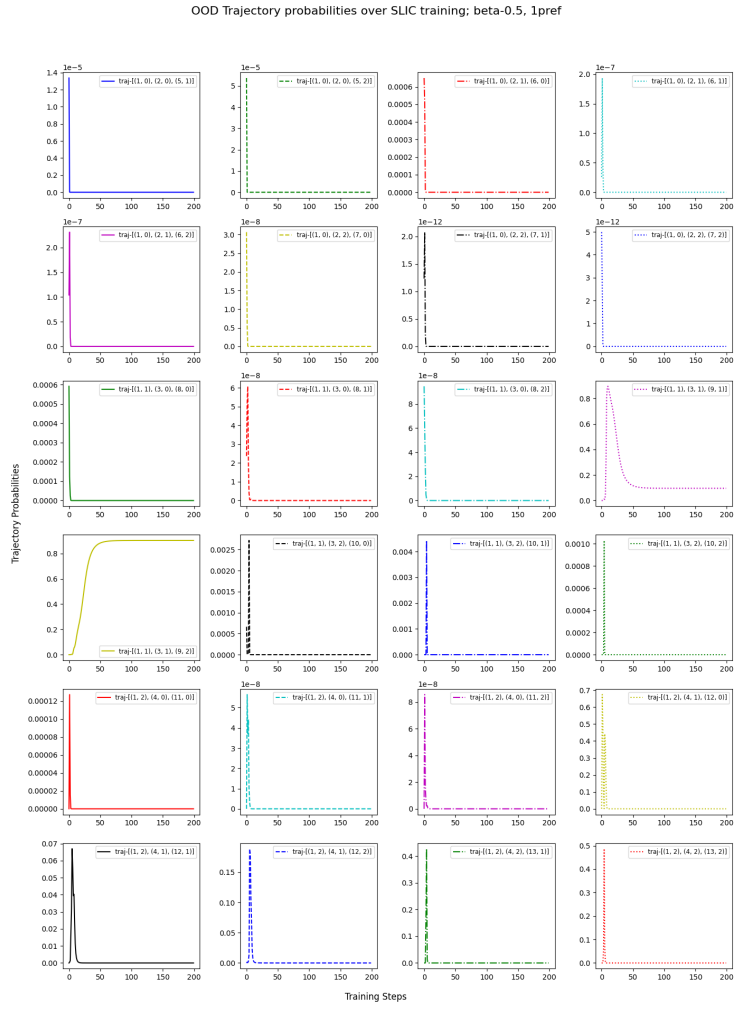


Figure 17: Trajectory probabilities throughout SLiC training. The top plot shows how the probability mass of different OOD trajectories, changes throughout training. The bottom plot shows how the probability mass of the trajectories in our preference dataset (size 1) changes over training

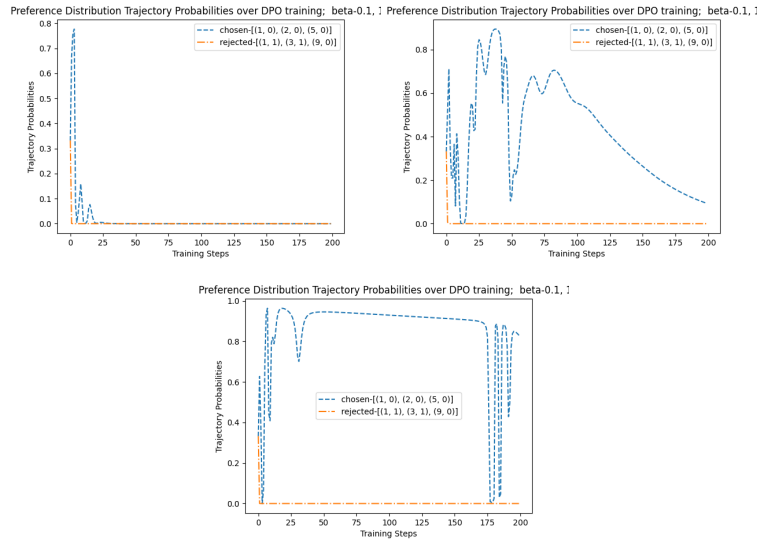


Figure 18: Trajectory probabilities throughout DPO training, over multiple runs

603 NeurIPS Paper Checklist

604 1. Claims

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

607 Answer: [Yes]

608 Justification: The paper faithfully adheres to the claims and motivation in the abstract and
609 provides proof and detailed empirical studies in support.

610 2. Limitations

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

612 Answer: [Yes]

613 Justification: A discussion of our limitations can be found as a separate section at the
614 beginning of the appendix.

615 Guidelines:

- 616 • The answer NA means that the paper has no limitation while the answer No means that
617 the paper has limitations, but those are not discussed in the paper.
- 618 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 619 • The paper should point out any strong assumptions and how robust the results are to
620 violations of these assumptions (e.g., independence assumptions, noiseless settings,
621 model well-specification, asymptotic approximations only holding locally). The authors
622 should reflect on how these assumptions might be violated in practice and what the
623 implications would be.
- 624 • The authors should reflect on the scope of the claims made, e.g., if the approach was
625 only tested on a few datasets or with a few runs. In general, empirical results often
626 depend on implicit assumptions, which should be articulated.
- 627 • The authors should reflect on the factors that influence the performance of the approach.
628 For example, a facial recognition algorithm may perform poorly when image resolution
629 is low or images are taken in low lighting. Or a speech-to-text system might not be
630 used reliably to provide closed captions for online lectures because it fails to handle
631 technical jargon.
- 632 • The authors should discuss the computational efficiency of the proposed algorithms
633 and how they scale with dataset size.

- 634
- 635
- 636
- 637
- 638
- 639
- 640
- 641
- 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.

642 3. Theory Assumptions and Proofs

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

645 Answer: [Yes]

646 Justification: We provide proofs and empirical evidence to support all our theoretical results.

647 Guidelines:

- 648
- 649
- 650
- 651
- 652
- 653
- 654
- 655
- 656
- 657
- 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.

658 4. Experimental Result Reproducibility

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

662 Answer: [Yes]

663 Justification: We provide detailed guidance on reproducibility by specifying all datasets,
664 code, and hyperparameters used in this work.

665 Guidelines:

- 666
- 667
- 668
- 669
- 670
- 671
- 672
- 673
- 674
- 675
- 676
- 677
- 678
- 679
- 680
- 681
- 682
- 683
- 684
- 685
- The answer NA means that the paper does not include experiments.
 - If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
 - If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
 - Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
 - While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.

- 686 (b) If the contribution is primarily a new model architecture, the paper should describe
687 the architecture clearly and fully.
- 688 (c) If the contribution is a new model (e.g., a large language model), then there should
689 either be a way to access this model for reproducing the results or a way to reproduce
690 the model (e.g., with an open-source dataset or instructions for how to construct
691 the dataset).
- 692 (d) We recognize that reproducibility may be tricky in some cases, in which case
693 authors are welcome to describe the particular way they provide for reproducibility.
694 In the case of closed-source models, it may be that access to the model is limited in
695 some way (e.g., to registered users), but it should be possible for other researchers
696 to have some path to reproducing or verifying the results.

697 5. Open access to data and code

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

701 Answer: [Yes]

702 Justification: We have only used open-source models with open-source datasets for all
703 aspects of the work. Please refer to section B for details on reproducing the results.

704 Guidelines:

- 705 • The answer NA means that paper does not include experiments requiring code.
- 706 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
707 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 708 • While we encourage the release of code and data, we understand that this might not be
709 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
710 including code, unless this is central to the contribution (e.g., for a new open-source
711 benchmark).
- 712 • The instructions should contain the exact command and environment needed to run to
713 reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 714 • The authors should provide instructions on data access and preparation, including how
715 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 716 • The authors should provide scripts to reproduce all experimental results for the new
717 proposed method and baselines. If only a subset of experiments are reproducible, they
718 should state which ones are omitted from the script and why.
- 719 • At submission time, to preserve anonymity, the authors should release anonymized
720 versions (if applicable).
- 721 • Providing as much information as possible in supplemental material (appended to the
722 paper) is recommended, but including URLs to data and code is permitted.

723 6. Experimental Setting/Details

724 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
725 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
726 results?
727

728 Answer: [Yes]

729 Justification: We list detailed information about the training and test details in Section ??.
730 Our experiments use open-source datasets and models.

731 Guidelines:

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

737 7. Experiment Statistical Significance

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

740 Answer: [No]

741 Justification: Training Large Language models is time-consuming and compute-intensive.
742 Our experiments do not run multiple seeds on one configuration due to limited computing
743 and financial budget. Instead, the focus of this work is extensive evaluation across multiple
744 configurations which we spent all our compute resources into. Our evaluation protocol is
745 similar to prior influential works in RLHF [46, 21].

746 Guidelines:

- 747 • The answer NA means that the paper does not include experiments.
- 748 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
749 dence intervals, or statistical significance tests, at least for the experiments that support
750 the main claims of the paper.
- 751 • The factors of variability that the error bars are capturing should be clearly stated (for
752 example, train/test split, initialization, random drawing of some parameter, or overall
753 run with given experimental conditions).
- 754 • The method for calculating the error bars should be explained (closed form formula,
755 call to a library function, bootstrap, etc.)
- 756 • The assumptions made should be given (e.g., Normally distributed errors).
- 757 • It should be clear whether the error bar is the standard deviation or the standard error
758 of the mean.
- 759 • It is OK to report 1-sigma error bars, but one should state it. The authors should
760 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
761 of Normality of errors is not verified.
- 762 • For asymmetric distributions, the authors should be careful not to show in tables or
763 figures symmetric error bars that would yield results that are out of range (e.g. negative
764 error rates).
- 765 • If error bars are reported in tables or plots, The authors should explain in the text how
766 they were calculated and reference the corresponding figures or tables in the text.

767 8. Experiments Compute Resources

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

771 Answer: [Yes]

772 Justification: We provide information on compute resources in the experimental details
773 section B in the appendix.

774 Guidelines:

- 775 • The answer NA means that the paper does not include experiments.
- 776 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
777 or cloud provider, including relevant memory and storage.
- 778 • The paper should provide the amount of compute required for each of the individual
779 experimental runs as well as estimate the total compute.
- 780 • The paper should disclose whether the full research project required more compute
781 than the experiments reported in the paper (e.g., preliminary or failed experiments that
782 didn't make it into the paper).

783 9. Code Of Ethics

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

786 Answer: [Yes]

787 Justification: We abide by the code of ethics in every respect.

788 Guidelines:

- 789 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.

- 790
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- 791
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).
- 792
- 793

794 **10. Broader Impacts**

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

797 Answer: [Yes]

798 Justification: We discuss societal impacts in Section A of the appendix.

799 Guidelines:

- 800 • The answer NA means that there is no societal impact of the work performed.
- 801 • If the authors answer NA or No, they should explain why their work has no societal
802 impact or why the paper does not address societal impact.
- 803 • Examples of negative societal impacts include potential malicious or unintended uses
804 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
805 (e.g., deployment of technologies that could make decisions that unfairly impact specific
806 groups), privacy considerations, and security considerations.
- 807 • The conference expects that many papers will be foundational research and not tied
808 to particular applications, let alone deployments. However, if there is a direct path to
809 any negative applications, the authors should point it out. For example, it is legitimate
810 to point out that an improvement in the quality of generative models could be used to
811 generate deepfakes for disinformation. On the other hand, it is not needed to point out
812 that a generic algorithm for optimizing neural networks could enable people to train
813 models that generate Deepfakes faster.
- 814 • The authors should consider possible harms that could arise when the technology is
815 being used as intended and functioning correctly, harms that could arise when the
816 technology is being used as intended but gives incorrect results, and harms following
817 from (intentional or unintentional) misuse of the technology.
- 818 • If there are negative societal impacts, the authors could also discuss possible mitigation
819 strategies (e.g., gated release of models, providing defenses in addition to attacks,
820 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
821 feedback over time, improving the efficiency and accessibility of ML).

822 **11. Safeguards**

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

826 Answer: [NA]

827 Justification: We use public models that are fine-tuned for alignment on open-source datasets.
828 Our models do not contribute any additional risk over the base models as we are explicitly
829 training for alignment.

830 Guidelines:

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

841 **12. Licenses for existing assets**

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

845 Answer: [Yes]

846 Justification: The pretrained models in this work come from the Pythia family all
847 of which are classified under Apache License [https://huggingface.co/EleutherAI/pythia-](https://huggingface.co/EleutherAI/pythia-2.8b/tree/main)
848 2.8b/tree/main. The TL;DR comparison dataset used in this work uses a modified MIT
849 License <https://github.com/openai/summarize-from-feedback/blob/master/LICENSE>.

850 Guidelines:

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

866 13. New Assets

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

869 Answer: [NA]

870 Justification: We use open-source pretrained models and provide details to reproduce our
871 fine-tuning experiments.

872 Guidelines:

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

881 14. Crowdsourcing and Research with Human Subjects

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

885 Answer: [NA]

886 Justification: We do not use crowdsourcing or research with human subjects in this work

887 Guidelines:

- 888 • The answer NA means that the paper does not involve crowdsourcing nor research with
889 human subjects.
- 890 • Including this information in the supplemental material is fine, but if the main contribu-
891 tion of the paper involves human subjects, then as much detail as possible should be
892 included in the main paper.

893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914

- 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: We do not use crowdsourcing or research with human subjects in this work.

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.