## Direct Alignment with Heterogeneous Preferences

Ali Shirali\*<sup>†1</sup>, Arash Nasr-Esfahany\*<sup>2</sup>, Abdullah Alomar<sup>2,3</sup>, Parsa Mirtaheri<sup>4</sup>, Rediet Abebe<sup>‡5</sup>, and Ariel Procaccia<sup>‡5</sup>

<sup>1</sup>University of California Berkeley <sup>2</sup>Massachusetts Institute of Technology <sup>3</sup>Ikigai Labs <sup>4</sup>University of California San Diego <sup>5</sup>Harvard University

February 25, 2025

#### Abstract

Alignment with human preferences is commonly framed using a universal reward function, even though human preferences are inherently heterogeneous. We formalize this heterogeneity by introducing user types and examine the limits of the homogeneity assumption. We show that aligning to heterogeneous preferences with a single policy is best achieved using the average reward across user types. However, this requires additional information about annotators. We examine improvements under different information settings, focusing on direct alignment methods. We find that minimal information can yield first-order improvements, while full feedback from each user type leads to consistent learning of the optimal policy. Surprisingly, however, no sample-efficient consistent direct loss exists in this latter setting. These results reveal a fundamental tension between consistency and sample efficiency in direct policy alignment.

#### 1 Introduction

Human rewards and preferences are heterogeneous [1, 2, 3, 4, 5]. Despite this, learning from preference data often bypasses this insight, relying on what we dub the *preference homogeneity assumption*. This tension in assumptions is readily apparent in standard human-AI alignment methods—such as reinforcement learning from human feedback (RLHF) [6, 7, 8] and direct preference optimization (DPO) [9]—which assume a single reward function captures the interests of the entire population.

We examine the limits of the preference homogeneity assumption when individuals belong to user types, each characterized by a specific reward function. Recent work has shown that in this setting, the homogeneity assumption can lead to unexpected behavior [10, 11, 12]. One challenge is that, under this assumption, learning from human preferences becomes unrealizable, as a single reward function cannot capture the complexity of population preferences with multiple reward functions [13, 14]. Both RLHF and DPO rely on maximum likelihood estimation (MLE) to optimize the reward or policy. Unrealizability implies their likelihood functions cannot fully represent the underlying preference data distribution, resulting in a nontrivial optimal MLE solution. From another perspective, learning a universal reward or policy from a heterogeneous population inherently involves an aggregation of diverse interests, and this aggregation is nontrivial.

In the quest for a single policy that accommodates a heterogeneous population with multiple user types, we show that the only universal reward yielding a well-defined alignment problem is an affine

<sup>\*</sup>Equal contribution

<sup>&</sup>lt;sup>†</sup>Work done while visiting Harvard

<sup>&</sup>lt;sup>‡</sup>Equal advising

aggregation of the reward functions across user types, with the average reward as a natural choice. However, standard methods like DPO do not maximize this user-weighted average reward. Building on insights by Siththaranjan, Laidlaw, and Hadfield-Menell [15], we show that DPO implicitly maximizes Borda count, which comes with unexpected drawbacks, e.g., the optimal solution depends on how alternative responses are sampled, even for infinite data.

We observe that learning the average reward over user types—or equivalently, a policy that maximizes it—from anonymous data is impossible. Focusing on *direct alignment methods*, which avoid explicit reward modeling, we study the benefits of using annotator data for a range of information settings. We show that improving DPO with a first-order correction to its objective is possible with minimal annotator information. Specifically, we design an approximate direct alignment method when each preference data point is paired with another one labeled by the same user.

On the other hand, we find that there are limits to what is possible even with significant annotator information. In particular, we propose a consistent loss function for direct alignment when we have feedback on each data point from each user type. But this loss is sample-inefficient, using only data where all annotator types agree. Surprisingly, we prove that no consistent loss uses the rest of the data.

In summary, the homogeneity assumption leads to undesirable outcomes when aligning a single AI agent to diverse preferences. Our analysis shows that there is a limited class of reward aggregation that results in a valid objective for alignment, with average reward over user types emerging as the natural candidate. This requires some annotator data, though small amounts of data can yield significant improvements. Our findings, however, uncover a fundamental tension between consistency and sample efficiency in direct alignment. To achieve both sample efficiency and consistency, we must forgo the benefits of direct optimization and instead train individualized reward models, which inevitably incurs significant training and storage costs.

#### 2 Preliminaries

In the alignment problem, we consider a setting where a reward function  $r^*$  evaluates responses to queries. Formally,  $r^*([x, y])$  is the reward value of responding y to a query x.

The alignment problem involves designing a *policy*, which chooses high-reward responses. Let  $\pi$  denote a policy, defining a probability distribution over possible responses: i.e., given a query  $\boldsymbol{x}$ ,  $\pi$  returns  $\boldsymbol{y}$  with probability  $\pi(\boldsymbol{y} \mid \boldsymbol{x})$ . Commonly, we start with a reference policy  $\pi_{\text{ref}}$ , which serves as a prior over  $\boldsymbol{y}$  [16]. The goal is then to find a new policy  $\pi$  that, for every  $\boldsymbol{x}$ , maximizes

$$\mathbb{E}_{\boldsymbol{y} \sim \pi(\cdot | \boldsymbol{x})} [r^*([\boldsymbol{x}, \boldsymbol{y}])] - \beta D_{\mathrm{KL}} (\pi(\cdot | \boldsymbol{x}); \pi_{\mathrm{ref}}(\cdot | \boldsymbol{x})). \tag{1}$$

We denote the optimal policy by  $\pi^*$ . In practice,  $\pi_{\text{ref}}$  is often a pretrained language model, and the regularization parameter  $\beta$  controls deviation from it. Eq. (1) often includes  $\mathbb{E}_{\boldsymbol{x}}$ , which is important in practice but does not affect  $\pi^*$  in theory.

When  $r^*$  is explicitly known, we can directly apply RL to maximize Eq. (1). In many real-world settings, however, we do not know  $r^*$  and must estimate it. In such cases, we can collect human feedback to infer the reward function, after which we can use RL to optimize the policy, commonly known as RLHF. While RLHF is widely used, tuning this approach can be challenging due to the inherent complexities of RL. Recently, direct alignment with preferences has gained popularity as an alternative approach [17, 9, 18]. Unlike RLHF, direct alignment methods bypass explicit reward modeling to instead train a policy directly from human feedback.

**Preference Model.** Both direct alignment and RLHF rely on a model of human preference to relate reward values with observed preference data. Consider the case where responses  $y_1, y_2$  are generated for a given query x. We express the probability that  $y_2$  is preferred to  $y_1$  as

$$\Pr(y_2 \succ y_1 \mid x; r^*) = \sigma(r^*([x, y_2]) - r^*([x, y_1])),$$
 (2)

where  $\sigma$  is a non-decreasing function in [0,1] [19, 20, 21]. A widely-used choice for  $\sigma$  is the sigmoid function corresponding to the well-known Bradley-Terry (BT) model [22].

**Direct Preference Optimization.** Among direct alignment methods, DPO has emerged as the most widely used approach. It leverages a closed-form solution to Eq. (1), which allows it to link any reward directly to its optimal policy. Thereby, rather than explicitly estimating the reward, DPO optimizes a policy whose *induced reward* best explains the observed preferences. We derive this connection below:

First, maximizing Eq. (1) has a well-known solution [23]. The optimal policy  $\pi^*$  takes the form:

$$\pi^*(\boldsymbol{y} \mid \boldsymbol{x}) = \frac{1}{Z(\boldsymbol{x})} \pi_{\text{ref}}(\boldsymbol{y} \mid \boldsymbol{x}) \cdot \exp\left(\frac{1}{\beta} r^*([\boldsymbol{x}, \boldsymbol{y}])\right). \tag{3}$$

Here,  $Z(\boldsymbol{x}) = \sum_{\boldsymbol{y}'} \pi_{\text{ref}}(\boldsymbol{y}' \mid \boldsymbol{x}) \cdot \exp(\frac{1}{\beta}r^*([\boldsymbol{x}, \boldsymbol{y}']))$  is the partition function. Eq. (3) establishes a direct relationship between policy ratios and reward differences:

$$r^*(\boldsymbol{y}_2) - r^*(\boldsymbol{y}_1) = \beta \log \frac{\pi^*(\boldsymbol{y}_2)}{\pi_{\text{ref}}(\boldsymbol{y}_2)} - \beta \log \frac{\pi^*(\boldsymbol{y}_1)}{\pi_{\text{ref}}(\boldsymbol{y}_1)}.$$
 (4)

Henceforth, we omit  $\boldsymbol{x}$  when we can do so without ambiguity. This equation shows that the difference in rewards between two responses is fully captured by the difference in their policy ratios. Using this formulation, we can define the *induced reward* of a policy  $\pi$  by  $\beta \log \frac{\pi(\boldsymbol{y})}{\pi_{\text{ref}}(\boldsymbol{y})}$ . The induced reward of  $\pi$  is the reward for which  $\pi$  is the optimal policy.

The difference in rewards of  $y_1$  and  $y_2$  is sufficient to express the likelihood of  $y_2 > y_1$  in Eq. (2). Using Eq. (4), we can therefore write the likelihood as a function of  $\pi^*$ :

$$\Pr(\mathbf{y}_2 \succ \mathbf{y}_1 \mid \pi^*) = \sigma \left(\beta \log \frac{\pi^*(\mathbf{y}_2)}{\pi_{\text{ref}}(\mathbf{y}_2)} - \beta \log \frac{\pi^*(\mathbf{y}_1)}{\pi_{\text{ref}}(\mathbf{y}_1)}\right). \tag{5}$$

For any policy  $\pi$ , we can define  $\Pr(y_2 \succ y_1 \mid \pi)$  similarly. We can then estimate  $\pi^*$  using MLE: Given a dataset  $\mathcal{D}$  with query and response pairs  $(x, y_l, y_w)$ , where  $y_w \succ y_l$ , DPO finds  $\pi^*$  by maximizing the log-likelihood

$$\sum_{(\boldsymbol{x}, \boldsymbol{y}_l, \boldsymbol{y}_w) \in \mathcal{D}} \log \Pr(\boldsymbol{y}_w \succ \boldsymbol{y}_l \mid \boldsymbol{x}; \pi), \qquad (6)$$

or equivalently, minimizing the cross-entropy loss. Under mild assumptions, MLE is a consistent estimator of  $\pi^*$ .

#### 3 Problem Formulation

The alignment problem is traditionally framed under a preference homogeneity assumption, where a single reward is presumed to capture all individual interests. In practice, people's preferences can differ significantly. To better capture real-world settings, we formalize preference heterogeneity by allowing reward functions to vary across user types.

Heterogeneous Preferences. The influential study of "individual choice behavior" by Luce [20] and other foundational works on human decision-making in mathematical psychology such as Shepard [24] focus on *individual* preference models. Luce [20] uses an axiomatic approach to establish the existence of a value function for each individual that, once normalized, explains the individual's choice probabilities. The widely used BT is one such example.

In practice, we often cannot observe individuals' identities. Therefore, standard approaches in preference modeling use a single reward function across the entire population. This homogeneity

assumption makes preference learning *unrealizable*: Even if a specific model family can explain preferences of every individual, we cannot ensure that a model from the same family explains population-level choices. For example, we cannot represent a mixture of BT models with a single BT (we prove this in Proposition C.1 for completeness).

To account for heterogeneity, we need to define individual rewards. However, learning at scale with this level of granularity is impractical, especially when working with finite data. Hence, we group individuals into multiple *user types*, denoted by  $\mathcal{U}$ . Individuals with the same type have similar rewards, but this need not hold across types.

For a given user of type  $u \in \mathcal{U}$ , we denote the corresponding reward function by  $r^*(\cdot; u)$ . This function assigns a scalar reward  $r^*([\boldsymbol{x}, \boldsymbol{y}]; u)$  for every response  $\boldsymbol{y}$  to a given query  $\boldsymbol{x}$ . We model the preferences of an individual of type u with

$$\Pr(\boldsymbol{y}_2 \succ \boldsymbol{y}_1 \mid r^*, u) = \sigma(r^*(\boldsymbol{y}_2; u) - r^*(\boldsymbol{y}_1; u)), \qquad (7)$$

The population-level preferences are therefore given by:

$$\Pr(\boldsymbol{y}_2 \succ \boldsymbol{y}_1 \mid r^*) = \mathbb{E}_u \left[ \sigma \left( r^*(\boldsymbol{y}_2; u) - r^*(\boldsymbol{y}_1; u) \right) \right]. \tag{8}$$

The Extended Alignment Problem. Our goal is to find a single policy that can effectively accommodate a heterogeneous population. This is essential when user types are not observable during inference. Furthermore, a universal policy may be preferable when personalization comes with significant drawbacks: e.g., cases where prioritizing a broadly accepted notion of truth or safety is more important than catering to individual preferences [25, 26].

Deriving a universal policy requires aggregation of diverse rewards. As we show next, an affine combination is the only form of aggregation that guarantees a well-defined problem, i.e., a problem that yields the same optimal policy for every reward that is consistent with the preference data.

**Proposition 3.1.** Consider an aggregation  $f: \mathbb{R}^{\mathcal{U}} \to \mathbb{R}$ . If  $f(\{r(\boldsymbol{y};u)\}_{u \in \mathcal{U}})$  induces the same ordering over  $\boldsymbol{y}$  for every reward r consistent with the preferences distribution, then under weak regulatory assumptions, f must be affine.

See proof on page 20. This result rules out many commonly-used aggregations, such as Max-Min [27] or Nash social welfare [28]. The expected reward across user types emerges as a natural choice here. Any other affine combination would weigh people unequally, which requires strong justifications and is rare in practice.

To summarize, our objective is to maximize

$$\mathbb{E}_{\boldsymbol{y} \sim \pi(\cdot | \boldsymbol{x})} \left[ \mathbb{E}_{u} \left[ r^{*}([\boldsymbol{x}, \boldsymbol{y}]; u) \right] \right] - \beta D_{\mathrm{KL}} \left( \pi(\cdot | \boldsymbol{x}); \pi_{\mathrm{ref}}(\cdot | \boldsymbol{x}) \right)$$
(9)

for every prompt x. With this extended framework in mind, we next discuss why standard approaches like RLHF or DPO do not necessarily yield the optimal policy.

## 4 Implications of Homogeneity Assumption

With heterogeneous preferences, standard RLHF or DPO cannot yield the optimal policy  $\pi^*$  that maximizes Eq. (9). If they did, it would also be possible to learn the user-weighted average reward as the induced reward of  $\pi^*$ . However, as we show in Proposition 5.1 and was previously observed by Siththaranjan, Laidlaw, and Hadfield-Menell [15] and Procaccia, Schiffer, and Zhang [29], learning the expected reward from anonymous preferences is impossible.

To explain DPO's failure in finding  $\pi^*$ , we extend its derivation to the heterogeneous setting in Section 4.1. This analysis lays the foundations to account for heterogeneity in DPO later on. In Section 4.2, we show that DPO's policy aligns with Borda count and, in Section 4.3, highlight its limitations. While our analysis focuses on DPO, similar insights extend to RLHF by substituting the policy with its induced reward.

#### 4.1 Objective is Not the Expected Reward

We follow DPO's derivation from Section 2 but under heterogeneity. We show the closed-form connection between  $\pi^*$  and  $r^*$  is no longer sufficient to express the likelihood function. Beginning with Eq. (9), the optimal policy is

$$\pi^*(\boldsymbol{y}) = \frac{1}{Z(\boldsymbol{x})} \, \pi_{\text{ref}}(\boldsymbol{y}) \cdot \exp\left(\frac{1}{\beta} \, \mathbb{E}_u \big[ r^*(\boldsymbol{y}; u) \big] \right). \tag{10}$$

Define  $\Delta r^*(\boldsymbol{y}_1, \boldsymbol{y}_2; u) := r^*(\boldsymbol{y}_2; u) - r^*(\boldsymbol{y}_1; u)$ . The policy ratios of  $\pi^*$  are related to the expected difference in rewards:

$$\mathbb{E}_{u}\left[\Delta r^{*}(\boldsymbol{y}_{1}, \boldsymbol{y}_{2}; u)\right] = \beta \log \frac{\pi^{*}(\boldsymbol{y}_{2})}{\pi_{\text{ref}}(\boldsymbol{y}_{2})} - \beta \log \frac{\pi^{*}(\boldsymbol{y}_{1})}{\pi_{\text{ref}}(\boldsymbol{y}_{1})}.$$
(11)

In the homogeneous case,  $\Delta r^*$  was sufficient to describe the likelihood of  $y_2 > y_1$ . However, with heterogeneous preferences,  $\mathbb{E}_u[\Delta r^*]$  alone does not suffice to write the likelihood function in Eq. (8). It is only under the approximation

$$\mathbb{E}_{u}\left[\sigma\left(\Delta r^{*}(\boldsymbol{y}_{1},\boldsymbol{y}_{2};u)\right)\right] \approx \sigma\left(\mathbb{E}_{u}\left[\Delta r^{*}(\boldsymbol{y}_{1},\boldsymbol{y}_{2};u)\right]\right)$$
(12)

that we can write  $\Pr(y_2 \succ y_1 \mid r^*)$  in terms of policy ratios as in Eq. (5), and minimize DPO's loss to find  $\pi^*$ .

#### 4.2 Ordinal Consistency with Borda Count

If DPO were the answer, what would the question be? We partially answer this question by an adaptation of a result from Siththaranjan, Laidlaw, and Hadfield-Menell [15] which we restate for completeness. First, define Borda count as follows.

**Definition 4.1** (Normalized Borda count). For a prompt x, let  $\mathcal{D}(\cdot \mid x)$  denote the distribution of alternative responses sampled for x. The Normalized Borda Count (NBC) of y at x is the probability that an annotator with a random type prefers y over an alternative response  $y' \sim \mathcal{D}(\cdot \mid x)$ :

$$NBC(\boldsymbol{y} \mid \boldsymbol{x}) := \mathbb{E}_{\boldsymbol{y}' \sim \mathcal{D}(\cdot \mid \boldsymbol{x})} \left[ \Pr(\boldsymbol{y} \succ \boldsymbol{y}' \mid \boldsymbol{x}; r^*) \right]. \tag{13}$$

We next show that DPO's policy ratios are ordinally consistent with the normalized Borda count.

**Proposition 4.2.** Suppose responses to  $\mathbf{x}$  in the preference dataset are drawn from  $\mathcal{D}(\cdot \mid \mathbf{x})$ . In the limit of many data points, DPO's induced reward, or equivalently  $\frac{\pi_{\mathrm{DPO}}(\cdot \mid \mathbf{x})}{\pi_{\mathrm{ref}}(\cdot \mid \mathbf{x})}$ , has the same ordering over responses as  $\mathrm{NBC}(\cdot \mid \mathbf{x})$ .

See proof on page 20. Proposition 4.2 also applies to the homogeneous setting. In this case, however, NBC aligns with  $r^*$ . It is worth mentioning that DPO is not the only method consistent with NBC; identity preference optimization (IPO) [30] uses NBC as its objective. We next highlight key differences between NBC and the user-weighted expected reward along with DPO's drawbacks in practice.

<sup>&</sup>lt;sup>1</sup>We can view NBC( $y \mid x$ ) as an aggregation of rewards at y. One can verify that NBC meets the order consistency condition of Proposition 3.1. However, it uses the reward value at  $y' \neq y$  to define the aggregated reward at y and thus does not fall under Proposition 3.1. In fact, this interdependency causes the issues we discuss Section 4.3.

#### 4.3 Practical Drawbacks

In case of heterogeneous preferences, Borda count can significantly diverge from the user-weighted expected reward. This is studied under *distortion* in social choice problems [31]. Notably, NBC in Eq. (13) depends on  $\mathcal{D}$ . Therefore, although data collection is irrelevant in defining the optimal policy, it does affect  $\pi_{\text{DPO}}$ .

Next, we illustrate two key differences between  $\pi^*$  and  $\pi_{DPO}$  using examples. Unless otherwise stated, we assume  $\mathcal{D}(\cdot \mid \boldsymbol{x})$  and  $\pi_{ref}(\cdot \mid \boldsymbol{x})$  are uniform, and annotators follow BT. Refer to Appendix A for further drawbacks of DPO (minority suppression and IIA violation).

Sensitivity to Preference Dataset Distribution. Suppose  $\mathcal{U} = \{A, B\}$  and types are equally represented. Given three possible responses, type A prefers  $y_1$  but type B prefers  $y_2$ :

$$r^*(\mathbf{y}_1; A) = 6, r^*(\mathbf{y}_2; A) = 1, r^*(\mathbf{y}_3; A) = 4,$$
  
 $r^*(\mathbf{y}_1; B) = 3, r^*(\mathbf{y}_2; B) = 9, r^*(\mathbf{y}_3; B) = 4.$ 

One can verify when  $\mathcal{D}(y_1) = \mathcal{D}(y_2)$ , increasing  $\mathcal{D}(y_3)$  from 0.02 to 0.04 changes  $\pi_{DPO}$ 's preference from  $y_2$  to  $y_1$ .

DPO's policy is also sensitive to the preference model. Consider a variation of BT with a temperature of 2:  $\sigma_2(z) := (1 + \exp(-z/2))^{-1}$ . For the same users and uniform sampling of alternatives, increasing the temperature from 1 to 2 flips  $\pi_{\text{DPO}}$ 's ranking over  $y_1$  and  $y_2$  while the preference model has no effect on  $\pi^*$ .

We have to emphasize that the dependence of NBC, and consequently  $\pi_{DPO}$ , on the dataset sampling distribution  $\mathcal{D}$  is not due to finite-sample limitations or insufficient offline dataset support. This issue persists even with complete data coverage and in the limit of infinite data.

**Mediocrity Promotion.** Consider the task of summarization. Suppose  $\mathcal{U} = \{A, B, C\}$  and types are equally represented. Type A(B) strongly favors longer (shorter) summaries while type C slightly prefers medium-length ones:

$$r^*(\text{short}; A) = 0, r^*(\text{med}; A) = 1, r^*(\log; A) = 4,$$
  
 $r^*(\text{short}; B) = 4, r^*(\text{med}; B) = 1, r^*(\log; B) = 0,$   
 $r^*(\text{short}; C) = 0, r^*(\text{med}; C) = 1, r^*(\log; C) = 0.$ 

In this case,  $\pi^*(\text{short}) = \pi^*(\text{long}) > \pi^*(\text{med})$ , however, NBC(short) = NBC(long) < NBC(med). DPO prefers medium-length summaries not strongly favored by any type.

Real-World Examples. The examples above are not contrived; in real-world cases, NBC can produce rankings different from  $\pi^*$  and is sensitive to dataset distribution as extensively studied under distortion of social choice rules. To show this with a real example, we use Pew Research Center surveys and analyze a question to 5101 participants: "The next time you purchase a vehicle, how likely are you to consider purchasing an electric vehicle?" (options from A: very likely to D: not at all likely). We discuss how we select this question in Appendix E. Responses come from groups of different political leanings: Republican (45%), Democratic (48%), and Neither/refused (7%).

Assuming the Luce-Shepard model [24] (see Eq. (27)), we estimate the reward for each group to calculate NBC and a user-weighted average reward. To find NBC, we use two distributions for alternatives: a uniform distribution  $\mathcal{D}_U$  and a slightly altered distribution  $\mathcal{D}_a$  with 0.2 total variation distance (TV) from  $\mathcal{D}_U$ . As shown in Fig. 1, NBC (with  $\mathcal{D}_U$ ) ranks option C first despite its mediocrity: it is the second or third preference for the three groups (see Fig. 8). In contrast, the user-weighted average reward favors D: the first and second preference for Republicans and the no-lean groups, respectively. Notably, altering  $\mathcal{D}_U$  to  $\mathcal{D}_a$  flips NBC's top ranking, highlighting NBC's sensitivity to dataset distribution. Similar discrepancies appear in other Pew surveys (see Appendix E).

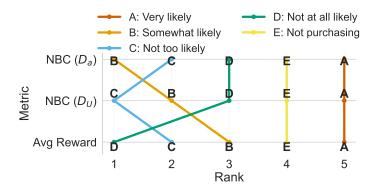


Figure 1: "The next time you purchase a vehicle, how likely are you to seriously consider purchasing an EV?" NBC ranking differs from the user-weighted average reward and is sensitive to the dataset distribution.

## 5 Approximate Direct Alignment with Minimal Annotator Information

The failure of standard alignment methods to find the optimal policy  $\pi^*$  raises the question of whether it is even possible to design a method that identifies  $\pi^*$ . The answer is no without annotator information. To show this, it suffices to prove that the ranking based on the user-weighted average reward is not learnable. This implies that  $\pi^*$  is also not learnable since its induced reward corresponds to this average reward by definition. We defer the formal definition of learnability to Definition C.2 in the appendix. Based on this definition, we prove the following impossibility:

**Proposition 5.1.** If there are at least two alternatives and two user types with a continuous preference model, the ranking based on the user-weighted expected reward is not learnable without annotator information.

See proof on page 21. Siththaranjan, Laidlaw, and Hadfield-Menell [15] (Theorem 3.4) presented a version of Proposition 5.1 in case of two alternatives and two types with  $\sigma(\Delta r) = \mathbb{I}\{\Delta r > 0\}$ . Procaccia, Schiffer, and Zhang [29] (Theorem 2.2) generalized this to BT. Proposition 5.1 presents a fresh perspective by generalizing the impossibility to any continuous preference model and presenting multiple proof strategies, including one that draws on a robust version of Arrow's theorem [32].

To circumvent the impossibility in Proposition 5.1, we must either relax the requirement of exactly identifying  $\pi^*$  or collect some information from the annotators. This section focuses on the former, and the latter is the subject of Section 6. Next, we introduce an approximate alignment objective, along with the required information and algorithms to solve it.

#### 5.1 First-Order Approximation

The approximation in Eq. (12) is equivalent to using a zeroth-order Taylor expansion of  $\sigma(\cdot)$  around the average reward to calculate the likelihood function. To improve it, we extend DPO by incorporating an additional non-zero term from the expansion, which we call *first-order corrected DPO*. The derivation is as follows. Expanding  $\sigma(\Delta r^*(y_1, y_2; u))$  around  $\Delta \bar{r}^*(y_1, y_2) := \mathbb{E}_u[\Delta r^*(y_1, y_2; u)]$  up to the second order gives the below approximation for the likelihood:

$$\mathbb{E}_{u}\left[\sigma\left(\Delta r^{*}(\boldsymbol{y}_{1},\boldsymbol{y}_{2};u)\right)\right] \approx \sigma\left(\Delta \bar{r}^{*}(\boldsymbol{y}_{1},\boldsymbol{y}_{2})\right) + \frac{1}{2}\sigma''\left(\Delta \bar{r}^{*}(\boldsymbol{y}_{1},\boldsymbol{y}_{2})\right) \cdot \operatorname{Var}_{u}\left[\Delta r^{*}(\boldsymbol{y}_{1},\boldsymbol{y}_{2};u)\right].$$
(14)

Note that Eq. (12) is loose when  $\sigma$  is nonlinear and preferences have high variance. We improve likelihood approximation by incorporating the variance term from Eq. (14). To calculate Eq. (14), we can substitute  $\Delta \bar{r}^*$  in by the difference in log policy ratios (Eq. (11)). We then need to estimate the variance term. Section 5.3 offers a variance estimator.

Once the variance is estimated by a function  $V(y_1, y_2)$ , first-order corrected DPO estimates  $\Pr(y_2 \succ y_1 \mid \pi)$  using

$$\sigma(h(\boldsymbol{y}_1, \boldsymbol{y}_2; \pi)) + \frac{\alpha}{2} \sigma''(h(\boldsymbol{y}_1, \boldsymbol{y}_2; \pi)) \cdot V(\boldsymbol{y}_1, \boldsymbol{y}_2).$$

Here,  $\alpha > 0$  determines the strength of correction, and

$$h(\boldsymbol{y}_1, \boldsymbol{y}_2; \pi) := \beta \log \frac{\pi(\boldsymbol{y}_2)}{\pi_{\text{ref}}(\boldsymbol{y}_2)} - \beta \log \frac{\pi(\boldsymbol{y}_1)}{\pi_{\text{ref}}(\boldsymbol{y}_1)}$$
(15)

denotes the difference of  $\pi$ 's induced rewards. Given  $\Pr(\mathbf{y}_2 \succ \mathbf{y}_1 \mid \pi)$  and a preference dataset  $\mathcal{D}$ , we can maximize the log-likelihood similar to Eq. (6). For numerical stability, we use a stable logarithm  $\log(z) \coloneqq \log(\max\{z,\epsilon\})$  in computations. Note that our theory suggests  $\alpha = 1$  while DPO uses  $\alpha = 0$ . Our empirical findings in Section 7.2 show that larger values of  $\alpha$  improve the effectiveness of the correction. We next discuss the estimation of  $V(\mathbf{y}_1, \mathbf{y}_2)$ .

#### 5.2 Impossibility without Annotator Information

If we limit our algorithms to M-estimators, which encompass most practical learning methods, consistent estimation of the variance term is impossible with anonymous data:

**Proposition 5.2.** There is no M-estimator that can estimate  $V(x, y_1, y_2) := \text{Var}_u \left[ \Delta r^*(x, y_1, y_2; u) \right]$  consistently without annotator information.

See proof on page 23. While Proposition 5.1 already implies that we cannot learn  $\pi^*$  without annotation information, Proposition 5.2 goes further, showing that even improving DPO with a first-order approximation is practically impossible. Next, we show how minimal annotation information can overcome this impossibility.

#### 5.3 Using Paired Preferences

We can get around Proposition 5.2 by collecting additional information on annotators. Specifically, it suffices to have a dataset  $\mathcal{D}$  of pairs of preferences in the form  $\{(\boldsymbol{x},\boldsymbol{y}_1,\boldsymbol{y}_2,o),(\boldsymbol{x}',\boldsymbol{y}_1',\boldsymbol{y}_2',o')\}$  where  $o = \mathbb{1}\{\boldsymbol{y}_2 \succ \boldsymbol{y}_1\}$ ,  $o' = \mathbb{1}\{\boldsymbol{y}_2' \succ \boldsymbol{y}_1'\}$  are labeled by the *same person*. Using  $\mathcal{D}$ , we can train a *joint likelihood model*  $J(\boldsymbol{x},\boldsymbol{y}_1,\boldsymbol{y}_2,\boldsymbol{x}',\boldsymbol{y}_1',\boldsymbol{y}_2')$  by minimizing cross-entropy between J and  $(o \cdot o')$  as the label. The joint likelihood model consistently estimates

$$\mathbb{E}_u \Big[ \sigma \big( \Delta r^*(\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2; u) \big) \cdot \sigma \big( \Delta r^*(\boldsymbol{x}', \boldsymbol{y}_1', \boldsymbol{y}_2'; u) \big) \Big] \; .$$

This is in fact sufficient to estimate the variance term:

**Lemma 5.3.** Using  $J_1$  and  $J_2$  as shorthands for  $J(x, y_1, y_2, x, y_1, y_2)$  and  $J(x, y_1, y_2, x, y_2, y_1)$ , we can use the following to estimate the variance term:

$$V(y_1, y_2) = \frac{J_1 - (J_1 + J_2)^2}{\sigma'(\Delta \bar{r}^*(y_1, y_2))^2}.$$
 (16)

See proof on page 24. Note that we can substitute  $\Delta \bar{r}^*(y_1, y_2)$  in terms of log policy ratios from Eq. (11). Thus, we have all the elements to calculate V in Eq. (16). This completes our derivation of first-order corrected DPO.

### 6 Direct Alignment with Maximum Annotator Information

Recall that learning the optimal policy from anonymous data is impossible, and an approximate improvement to DPO requires only minimal information about the annotations. But what if we collect richer data? Can we design a direct alignment method that consistently learns the optimal policy  $\pi^*$ ? To explore this, we consider a dataset where every sample is labeled by representatives of all user types. We show that consistent direct alignment is possible using this dataset.

Suppose user types are in a finite set  $\mathcal{U}$  and equally represented. This assumption makes our negative results stronger. Consider a rich data collection: for every context and candidate pairs  $(\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2)$ , we collect one preference data point from each user type. Let  $\boldsymbol{o} \in \{0, 1\}^{|\mathcal{U}|}$  be the vector that indicates preferences where  $o_u = 1$  if user of type  $u \in \mathcal{U}$  has preferred  $\boldsymbol{y}_2$ , and 0 otherwise. Given such a dataset  $\mathcal{D}$  with context, candidates, and preferences represented as  $(\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2, \boldsymbol{o})$ , our goal is to design a loss function

$$\mathcal{L}(\mathcal{D}; \pi) = \frac{1}{|\mathcal{D}|} \sum_{(\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2, \boldsymbol{o}) \in \mathcal{D}} l(\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2, \boldsymbol{o}; \pi)$$
(17)

such that  $\arg\min_{\pi} \mathcal{L}(\mathcal{D}; \pi)$  is a consistent estimator of  $\pi^*$ . Designing such a loss is, in fact, possible. For instance, suppose we only look into the agreement cases in  $\mathcal{D}$  where o is either all one or zero. Conditioned on agreement, we will show that the probability of  $y_2 \succ y_1$  is proportional to  $\exp\left(\sum_u \Delta r^*(y_1, y_2; u)\right)$ . We can write this likelihood in terms  $\pi^*$  directly as we have a correspondence between  $\pi^*$  and the difference in user-weighted average rewards (see Eq. (11)). We formally show this possibility through a temperature-adjusted DPO:

**Proposition 6.1.** Defining l in Eq. (17) as follows results in a consistent estimation of the optimal policy when preferences follow the BT model:

$$l(\boldsymbol{y}_1, \boldsymbol{y}_2, \boldsymbol{o}; \pi) = \begin{cases} -\log \sigma(|\mathcal{U}| \cdot h(\boldsymbol{y}_1, \boldsymbol{y}_2; \pi)), & \boldsymbol{o} = \vec{\boldsymbol{1}}, \\ -\log \sigma(|\mathcal{U}| \cdot h(\boldsymbol{y}_2, \boldsymbol{y}_1; \pi)), & \boldsymbol{o} = \vec{\boldsymbol{0}}, \\ 0 & o.w. \end{cases}$$

Here, h is the difference of  $\pi$ 's induced rewards (Eq. (15)).

See proof on page 24. Consistent loss function is not unique. We give another example in Proposition C.3, and a systematic way to find such losses in Lemma D.2. In both examples, loss functions reduce to the standard DPO loss when  $|\mathcal{U}| = 1$ .

While the loss function in Proposition 6.1 benefits from consistency, it only uses samples where all user types have agreed. In other words, it discards a sample with any disagreement. A natural question arises: Can we design a loss function that uses all data, including those with disagreement, while maintaining consistency? Surprisingly, the answer is no:

**Theorem 6.2.** Suppose l in Eq. (17) only depends on  $(x, y_1, y_2)$  through  $\pi$  and  $\pi_{ref}$ . If there are more than three types of user and the preferences follow BT, any loss  $\mathcal{L}$  that allows a consistent estimation of the optimal policy discards samples with disagreement, i.e., those with  $o \notin \{0, 1\}$ .

See proof on page 25. This theorem highlights a tension: To improve efficiency, one must compromise either consistency or direct optimization. The approximate direct alignment method proposed in Section 5 exemplifies forgoing consistency. Next, we discuss an alternative that favors consistency.

An Indirect Practical Solution: Averaging Personalized Rewards. The tradeoff between sample efficiency and consistency arises from the requirement for direct optimization. To regain sample efficiency, we may relax the requirement for direct alignment by training reward models while still avoiding RL. Specifically, we can learn personalized reward models  $r(\cdot; u)$  for different user types  $u \in \mathcal{U}$ , calculate a user-weighted expected reward, and use it to relabel a preference dataset. A dataset labeled

with this average reward makes any direct alignment method applicable and avoids RL. It is both consistent and sample-efficient when personalized reward learning is feasible, but comes at the cost of additional training for each user type and memory to store multiple models. Relabeling followed by DPO is found effective in practice [33].

#### 7 Experiments

We provide empirical evidence for our claims throughout the paper. Section 7.1 extends our sensitivity example in Section 4.3 to a real-world preference dataset. In Section 7.2, we simulate DPO and our proposed improvements in a synthetic, small-scale environment where we can visualize and compare the resulting policies. Finally, we scale this experiment in Section 7.3 by fine-tuning large language models, illustrating the extent of improvement over DPO. Our code for reproducing the experimental results is publicly available at https://github.com/arashne/dahp.

#### 7.1 NBC Sensitivity to Sampling Distribution

Recall from Section 4.2 that the common practice of alignment assuming homogeneity results in ordinal consistency with NBC. Here, we analyze the sensitivity of NBC to the distribution of pairwise preference datasets in real-world cases by using Pew surveys [34], the same dataset used in Section 4.3. Specifically, we address two questions: (i) Across the questions in the Pew surveys, how often would NBC rankings change when the sampling distribution of alternatives varies (while retaining support over all alternatives)? (ii) How much must the sampling distribution deviate from uniform to alter NBC rankings?

To answer these questions, we first estimate the reward of each option in each question (see Section 4.3 and Appendix E for further details). Given the rewards, we can calculate NBC under any sampling distribution using Eq. (13).

For question (i), among 1519 questions from 19 Pew surveys, NBC rankings change due to changing the sampling distribution from uniform in 20% of cases (306 questions), with the preferred choice changing in 136 cases. For question (ii), we find that a modest change in the sampling distribution suffices; in half the cases, a total variation (TV) distance of less than 0.23 from uniform alters the rankings. The cumulative distribution function (CDF) of the minimum TV distances required to change NBC rankings is shown in Fig. 2. Further experimental details are in Appendix E.

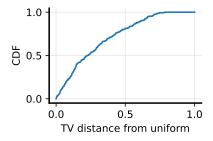


Figure 2: CDF of the minimum TV distances (from uniform) required to change the NBC order in the Pew surveys. A change of 0.23 is sufficient to change the order in half the questions.

#### 7.2 Synthetic Experiments

We generalize the discrete environment of Xu et al. [35] with a heterogeneous population. This environment enables us to visualize the differences between DPO's policy and the optimal policy, as well as to evaluate the effectiveness of applying a first-order correction (Section 5) and using a consistent loss function (Section 6).

**Environment.** A prompt x can take a value from 1 to n. There are also only n possible responses to each x. The reward for responding y to x for a type u is  $r^*([x,y];u) = R_u(\operatorname{dist}(x+u,y))$ , where dist is a circular distance, and  $R_u$  is a linearly decreasing function floored at zero. In this experiment, we set n = 40 and consider three equally represented types:  $\mathcal{U} = \{-10, 0, 10\}$ . We assume BT annotators.

Since the reward (and thus the policies) depends only on y-x, we can reduce everything to a 1D representation by setting  $\delta \coloneqq y-x$  and averaging over x. For example, for a policy  $\pi$ , define a 1D policy  $\pi(\delta) \coloneqq \frac{1}{n} \sum_{x \in [n]} \pi(x+\delta \mid x)$ , and similarly a 1D reward  $r^*(\delta; u)$ . We also compute standard errors of these 1D representations across x. Fig. 3 shows our choice of rewards as well as the expected reward across user types.

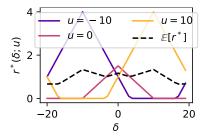


Figure 3: Rewards in the synthetic experiments

**Policies.** For a uniform  $\pi_{\text{ref}}$  and  $\beta = 1$ , Eq. (10) implies  $\pi^*(y \mid x) \propto \exp\left(\frac{1}{3}\sum_{u \in \mathcal{U}} r^*([x,y];u)\right)$ . We generate a large dataset of preferences under uniform context and alternative distributions and use the Adam optimizer to minimize the loss for different methods. For the first-order correction of DPO, we additionally train a joint likelihood model J to estimate the variance term V from Eq. (16). We use the loss from Proposition 6.1 as our choice for the consistent loss. This loss effectively utilizes less data compared to other methods since it throws away data points with disagreement. Refer to the accompanying code for more details.

**Results.** Fig. 4(a) presents  $\pi_{DPO}$  along with  $\pi^*$  and NBC. Unlike the optimal policy which prefers  $\delta$  around -10 and 10, DPO prefers  $\delta \approx 0$ . To a large extent, DPO's policy is ordinally consistent with NBC.

Fig. 4(b) shows that increasing correction strength  $\alpha$  brings the corrected DPO policy closer to  $\pi^*$ . In particular, at  $\alpha=1$ , the corrected DPO already favors alternatives with  $\delta\in\{-10,10\}$ , consistent with  $\pi^*$ . Furthermore, increasing  $\alpha$  makes these alternatives even more favorable. In the full-information setting, Fig. 4(c) shows that minimizing the consistent loss largely leads to  $\pi^*$ . Note that minor deviations from theoretical derivations are likely due to limited data and imperfect optimization in these experiments.

#### 7.3 Semi-Synthetic Experiments

To demonstrate our findings in a realistic setup, we use LoRA [36] to fine-tune Llama-3-8B [37] for both reward learning and direct alignment on two relabeled variations of the HH-RLHF dataset [38], which contains user prompts with pairs of chatbot responses. To simulate heterogeneous preferences, we define three user types with distinct length-based rewards: the first type prefers long, the second type prefers short, and the third type prefers mid-length prompt response combinations (see Appendix F for details). Recall from Section 3 that we argued for the average reward across user types as the proper objective for alignment with heterogeneous preferences. Therefore, we use agreement with the ground-truth average reward on the test set as the success metric.

First, we use vanilla reward learning and DPO with the homogeneity assumption on an anonymous preference dataset where every sample is labeled with a random user type. As Fig. 5 shows in blue,

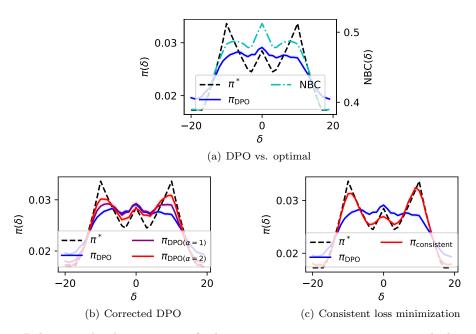


Figure 4: Policies explicitly accounting for heterogeneity are more consistent with the average reward across types in a synthetic setup.

their induced order agrees with the average reward in 89.6% and 67.4% of the test cases, respectively. Next, we use the loss function in Proposition 6.1 on a dataset with maximum annotator information (Section 6). To create this dataset, for every response pair to a prompt, we sample the preferences of the three user types until a consensus is reached, using the agreed-upon preference as the label. As Fig. 5 shows in red, the learned reward and the induced reward by the learned policy agree with the average reward in 93.9% and 71.7% of the test cases. In summary, explicitly accounting for heterogeneity increases the agreement with the average reward across user types by 4.3%—an additional 368 test cases—for both reward learning and direct alignment.

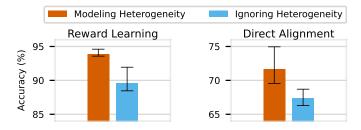


Figure 5: In the presence of preference labels from every user type, our proposed loss function produces reward models (left) and aligned policies (right) that are more consistent with the average reward across user types, compared to typical approaches that overlook heterogeneity. Bars show the mean, and whiskers denote the second and third quartiles across five random seeds.

#### 8 Related Work

Aligning models to "serve pluralistic human values" [12] can involve personalization to the user's specific reward [39, 40] or aggregation of diverse rewards. The latter, which is the subject of our

study, can use insights from social choice theory [11, 10].

Closest to our work, EM-DPO [41] simultaneously learns the distribution of user types and their corresponding policies. However, EM introduces significant complexities and lacks guarantees. Moreover, the identifiability of types requires additional assumptions. MODPO [42] applies DPO for each user type while utilizing estimated rewards from other user types to maximize a linear combination of rewards. Neither method obtains a policy by directly minimizing a loss over preference data. For a more extensive related work, refer to Appendix B.

#### 9 Discussion and Conclusion

Aligning a single policy to the average reward across user types in a heterogeneous population requires collecting annotator information. This can range from minimal information such as linking two instances labeled by the same annotator, to richer information like using questionnaires to infer annotator types. We improved DPO using the former and introduced a consistent loss when annotators from all user types label every data point. With additional assumptions, unsupervised methods might be able to identify annotator types from anonymous datasets [43]. Further research should explore the additional structures that, when used during data collection, can help with identifiability.

Our results revealed a tension between consistency and sample efficiency in direct alignment. Thus, an alternative approach—individual reward training and aggregation—may be more practical for addressing heterogeneity when individual rewards are identifiable.

We use the average reward across user types, as the natural choice among aggregations that define a well-defined alignment problem from pairwise preferences. However, the choice of aggregation is inherently a social and policy question rather than purely a technical one. In certain contexts, the policymaker might prefer to give higher weight to disadvantaged people to address issues such as inequality. Additionally, in some cases, the very existence of a reward function may be questionable, requiring objectives to be defined in terms of choice probabilities rather than rewards.

We believe that trained policies should not be used to elicit or represent aggregate preferences, even when reward aggregation is appropriate and estimation is consistent. While such policies may capture certain patterns in user behavior, they do not necessarily reflect the underlying interests or values of the population. In other words, we view the resulting policy as a functional tool for decision-making rather than a true representation of users' collective interests.

In summary, while preference heterogeneity is well recognized in mathematical psychology, standard methods often bypass this complexity. As we showed, accounting for heterogeneity, even when the goal remains the same as in the homogeneous setting—to derive a single policy—can render common techniques inefficient or inapplicable. Understanding these limitations calls for new approaches that explicitly incorporate heterogeneity while effectively balancing efficiency, consistency, and practicality.

## Acknowledgment

We thank Mohammad Alizadeh, Pouya Hamadanian, Moritz Hardt, Erfan Jahanparast, and Itai Shapira for discussions and feedback on earlier versions of this paper. Abebe was partially supported by the Andrew Carnegie Fellowship Program. Procaccia was partially supported by the National Science Foundation under grants IIS-2147187 and IIS-2229881; by the Office of Naval Research under grants N00014-24-1-2704 and N00014-25-1-2153; and by a grant from the Cooperative AI Foundation.

### References

[1] Lora Aroyo and Chris Welty. "Truth is a lie: Crowd truth and the seven myths of human annotation". In: AI Magazine 36.1 (2015), pp. 15–24.

- [2] Ellie Pavlick and Tom Kwiatkowski. "Inherent disagreements in human textual inferences". In: Transactions of the Association for Computational Linguistics 7 (2019), pp. 677–694.
- [3] Remi Denton et al. "Whose ground truth? Accounting for individual and collective identities underlying dataset annotation". In: arXiv preprint arXiv:2112.04554 (2021).
- [4] Marta Sandri et al. "Why Don't You Do It Right? Analysing Annotators' Disagreement in Subjective Tasks". In: Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics. Ed. by Andreas Vlachos and Isabelle Augenstein. Dubrovnik, Croatia: Association for Computational Linguistics, May 2023, pp. 2428-2441. DOI: 10.18653/v1/2023.eacl-main.178. URL: https://aclanthology.org/2023.eacl-main.178/.
- [5] Michael JQ Zhang et al. "Diverging Preferences: When do Annotators Disagree and do Models Know?" In: arXiv preprint arXiv:2410.14632 (2024).
- [6] Daniel M Ziegler et al. "Fine-tuning language models from human preferences". In: arXiv preprint arXiv:1909.08593 (2019).
- [7] Nisan Stiennon et al. "Learning to summarize with human feedback". In: Advances in Neural Information Processing Systems 33 (2020), pp. 3008–3021.
- [8] Long Ouyang et al. "Training language models to follow instructions with human feedback". In: Advances in Neural Information Processing Systems 35 (2022), pp. 27730–27744.
- [9] Rafael Rafailov et al. "Direct preference optimization: Your language model is secretly a reward model". In: Advances in Neural Information Processing Systems 36 (2024).
- [10] Vincent Conitzer et al. "Social Choice Should Guide AI Alignment in Dealing with Diverse Human Feedback". In: arXiv preprint arXiv:2404.10271 (2024).
- [11] Luise Ge et al. "Axioms for AI Alignment from Human Feedback". In: Advances in Neural Processing Systems 36 (2024).
- [12] Taylor Sorensen et al. "Position: A Roadmap to Pluralistic Alignment". In: Forty-first International Conference on Machine Learning. 2024. URL: https://openreview.net/forum?id=gQpBnRHwxM.
- [13] Vincent Dumoulin et al. "A density estimation perspective on learning from pairwise human preferences". In: arXiv preprint arXiv:2311.14115 (2023).
- [14] Chanwoo Park et al. "RLHF from heterogeneous feedback via personalization and preference aggregation". In: ICML 2024 Workshop: Aligning Reinforcement Learning Experimentalists and Theorists. 2024.
- [15] Anand Siththaranjan, Cassidy Laidlaw, and Dylan Hadfield-Menell. "Distributional preference learning: Understanding and accounting for hidden context in RLHF". In: arXiv preprint arXiv:2312.08358 (2023).
- [16] Tomasz Korbak, Ethan Perez, and Christopher Buckley. "RL with KL penalties is better viewed as Bayesian inference". In: Findings of the Association for Computational Linguistics: EMNLP 2022. 2022, pp. 1083–1091.
- [17] Yao Zhao et al. "SLiC-HF: Sequence likelihood calibration with human feedback". In: arXiv preprint arXiv:2305.10425 (2023).
- [18] Shangmin Guo et al. "Direct language model alignment from online AI feedback". In: arXiv preprint arXiv:2402.04792 (2024).
- [19] LL Thurstone. "A law of comparative judgment." In: Psychological Review 34.4 (1927), p. 273.
- [20] R Duncan Luce. Individual choice behavior. Vol. 4. Wiley New York, 1959.
- [21] John I Yellott Jr. "The relationship between Luce's choice axiom, Thurstone's theory of comparative judgment, and the double exponential distribution". In: *Journal of Mathematical Psychology* 15.2 (1977), pp. 109–144.

- [22] Ralph Allan Bradley and Milton E Terry. "Rank analysis of incomplete block designs: I. The method of paired comparisons". In: *Biometrika* 39.3/4 (1952), pp. 324–345.
- [23] Brian D Ziebart. Modeling purposeful adaptive behavior with the principle of maximum causal entropy. Carnegie Mellon University, 2010.
- [24] Roger N. Shepard. "Stimulus and response generalization: A stochastic model relating generalization to distance in psychological space". In: *Psychometrika* 22.4 (Dec. 1957), pp. 325–345. ISSN: 1860-0980. DOI: 10.1007/BF02288967. URL: https://doi.org/10.1007/BF02288967.
- [25] Lucas Monteiro Paes et al. "On the epistemic limits of personalized prediction". In: Advances in Neural Information Processing Systems 35 (2022), pp. 1979–1991.
- [26] Hannah Rose Kirk et al. "The benefits, risks and bounds of personalizing the alignment of large language models to individuals". In: *Nature Machine Intelligence* (2024), pp. 1–10.
- [27] Souradip Chakraborty et al. "MaxMin-RLHF: Towards equitable alignment of large language models with diverse human preferences". In: arXiv preprint arXiv:2402.08925 (2024).
- [28] Mamoru Kaneko and Kenjiro Nakamura. "The Nash social welfare function". In: *Econometrica: Journal of the Econometric Society* (1979), pp. 423–435.
- [29] Ariel D. Procaccia, Benjamin Schiffer, and Shirley Zhang. "Clone-Robust AI Alignment". In: arXiv preprint arXiv:2501.09254 (2025).
- [30] Mohammad Gheshlaghi Azar et al. "A general theoretical paradigm to understand learning from human preferences". In: *International Conference on Artificial Intelligence and Statistics*. PMLR. 2024, pp. 4447–4455.
- [31] Elliot Anshelevich et al. "Distortion in social choice problems: The first 15 years and beyond". In: arXiv preprint arXiv:2103.00911 (2021).
- [32] Ehud Friedgut, Gil Kalai, and Assaf Naor. "Boolean functions whose Fourier transform is concentrated on the first two levels". In: Advances in Applied Mathematics 29.3 (2002), pp. 427–437.
- [33] Evan Frick et al. "How to Evaluate Reward Models for RLHF". In: arXiv preprint arXiv:2410.14872 (2024).
- [34] Pew Research Center. About Pew Research Center. Accessed: 2025-01-20. 2025. URL: https://www.pewresearch.org/about/.
- [35] Shusheng Xu et al. "Is DPO Superior to PPO for LLM Alignment? A Comprehensive Study". In: Forty-first International Conference on Machine Learning. 2024.
- [36] Edward J Hu et al. "LoRA: Low-Rank Adaptation of Large Language Models". In: International Conference on Learning Representations. 2022. URL: https://openreview.net/forum?id=nZeVKeeFYf9.
- [37] AI@Meta. "Llama 3 Model Card". In: (2024). URL: https://github.com/meta-llama/llama3/blob/main/MODEL\_CARD.md.
- [38] Yuntao Bai et al. "Training a helpful and harmless assistant with reinforcement learning from human feedback". In: arXiv preprint arXiv:2204.05862 (2022).
- [39] Sriyash Poddar et al. "Personalizing Reinforcement Learning from Human Feedback with Variational Preference Learning". In: *The Thirty-eighth Annual Conference on Neural Information Processing Systems*. 2024. URL: https://openreview.net/forum?id=gRG6SzbW9p.
- [40] Daiwei Chen et al. "PAL: Pluralistic Alignment Framework for Learning from Heterogeneous Preferences". In: arXiv preprint arXiv:2406.08469 (2024).
- [41] Keertana Chidambaram, Karthik Vinay Seetharaman, and Vasilis Syrgkanis. Direct Preference Optimization With Unobserved Preference Heterogeneity. 2024. URL: https://openreview.net/forum?id=NQZNNUsutn.

- [42] Zhanhui Zhou et al. "Beyond one-preference-fits-all alignment: Multi-objective direct preference optimization". In: Findings of the Association for Computational Linguistics ACL 2024. 2024, pp. 10586–10613.
- [43] Xiaomin Zhang et al. "On the identifiability of mixtures of ranking models". In: arXiv preprint arXiv:2201.13132 (2022).
- [44] Usman Anwar et al. "Foundational Challenges in Assuring Alignment and Safety of Large Language Models". In: *Transactions on Machine Learning Research* (2024). Survey Certification, Expert Certification. ISSN: 2835-8856. URL: https://openreview.net/forum?id=oVTkOs8Pka.
- [45] Stephen Casper et al. "Open Problems and Fundamental Limitations of Reinforcement Learning from Human Feedback". In: *Transactions on Machine Learning Research* (2023). Survey Certification, Featured Certification. ISSN: 2835-8856. URL: https://openreview.net/forum?id=bx24KpJ4Eb.
- [46] Corby Rosset et al. "Direct Nash optimization: Teaching language models to self-improve with general preferences". In: arXiv preprint arXiv:2404.03715 (2024).
- [47] Zhaolin Gao et al. "Rebel: Reinforcement learning via regressing relative rewards". In: arXiv preprint arXiv:2404.16767 (2024).
- [48] Haoxian Chen et al. "Mallows-DPO: Fine-Tune Your LLM with Preference Dispersions". In: arXiv preprint arXiv:2405.14953 (2024).
- [49] Yunhao Tang et al. "Generalized Preference Optimization: A Unified Approach to Offline Alignment". In: Forty-first International Conference on Machine Learning. 2024. URL: https://openreview.net/forum?id=gu3nacA9AH.
- [50] Yu Meng, Mengzhou Xia, and Danqi Chen. "SimPO: Simple Preference Optimization with a Reference-Free Reward". In: *The Thirty-eighth Annual Conference on Neural Information Processing Systems*. 2024. URL: https://openreview.net/forum?id=3Tzcot1LKb.
- [51] Xinyu Li, Zachary C Lipton, and Liu Leqi. "Personalized language modeling from personalized human feedback". In: arXiv preprint arXiv:2402.05133 (2024).
- [52] Nishant Balepur et al. Whose Boat Does it Float? Improving Personalization in Preference Tuning via Inferred User Personas. 2025. arXiv: 2501.11549 [cs.CL]. URL: https://arxiv. org/abs/2501.11549.
- [53] Seongyun Lee et al. "Aligning to Thousands of Preferences via System Message Generalization". In: The Thirty-eighth Annual Conference on Neural Information Processing Systems. 2024. URL: https://openreview.net/forum?id=recsheQ7e8.
- [54] Meihua Dang et al. "Personalized Preference Fine-tuning of Diffusion Models". In:  $arXiv\ preprint\ arXiv:2501.06655\ (2025)$ .
- [55] Joel Jang et al. "Personalized soups: Personalized large language model alignment via post-hoc parameter merging". In: arXiv preprint arXiv:2310.11564 (2023).
- [56] Allison Lau et al. "Personalized Adaptation via In-Context Preference Learning". In: arXiv preprint arXiv:2410.14001 (2024).
- [57] Tianyi Qiu. "Representative Social Choice: From Learning Theory to AI Alignment". In: arXiv preprint arXiv:2410.23953 (2024).
- [58] Parand A Alamdari, Soroush Ebadian, and Ariel D Procaccia. "Policy aggregation". In: Advanced in Neural Information Processing Systems 36 (2024).
- [59] Jessica Dai and Eve Fleisig. "Mapping social choice theory to RLHF". In: arXiv preprint arXiv:2404.13038 (2024).
- [60] Gokul Swamy et al. "A Minimaximalist Approach to Reinforcement Learning from Human Feedback". In: Forty-first International Conference on Machine Learning. 2024. URL: https://openreview.net/forum?id=5kVgd2MwMY.

- [61] Hadassah Harland et al. "Adaptive Alignment: Dynamic Preference Adjustments via Multi-Objective Reinforcement Learning for Pluralistic AI". In: arXiv preprint arXiv:2410.23630 (2024).
- [62] Haoxiang Wang et al. "Arithmetic control of LLMs for diverse user preferences: Directional preference alignment with multi-objective rewards". In: arXiv preprint arXiv:2402.18571 (2024).
- [63] Huiying Zhong et al. "Provable multi-party reinforcement learning with diverse human feedback". In: arXiv preprint arXiv:2403.05006 (2024).
- [64] Dexun Li et al. "Aligning crowd feedback via distributional preference reward modeling". In: arXiv preprint arXiv:2402.09764 (2024).
- [65] Ryan Boldi et al. "Pareto-optimal learning from preferences with hidden context". In: arXiv preprint arXiv:2406.15599 (2024).
- [66] Alexandre Rame et al. "Rewarded soups: Towards Pareto-optimal alignment by interpolating weights fine-tuned on diverse rewards". In: Advances in Neural Information Processing Systems 36 (2024).
- [67] Angelica Chen et al. "Preference Learning Algorithms Do Not Learn Preference Rankings". In: The Thirty-eighth Annual Conference on Neural Information Processing Systems. 2024. URL: https://openreview.net/forum?id=YkJ5BuEXdD.
- [68] Dun Zeng et al. On Diversified Preferences of Large Language Model Alignment. 2024. arXiv: 2312.07401 [cs.AI]. URL: https://arxiv.org/abs/2312.07401.
- [69] Hritik Bansal, John Dang, and Aditya Grover. "Peering Through Preferences: Unraveling Feedback Acquisition for Aligning Large Language Models". In: *The Twelfth International Conference on Learning Representations*. 2024. URL: https://openreview.net/forum?id=dK161MwbCy.
- [70] Shibani Santurkar et al. "Whose opinions do language models reflect?" In: *International Conference on Machine Learning*. PMLR. 2023, pp. 29971–30004.
- [71] Michiel Bakker et al. "Fine-tuning language models to find agreement among humans with diverse preferences". In: Advances in Neural Information Processing Systems 35 (2022), pp. 38176–38189.
- [72] Liwei Jiang et al. "Can Language Models Reason about Individualistic Human Values and Preferences?" In: arXiv preprint arXiv:2410.03868 (2024).
- [73] Thomas P. Zollo et al. PersonalLLM: Tailoring LLMs to Individual Preferences. 2024. arXiv: 2409.20296 [cs.LG]. URL: https://arxiv.org/abs/2409.20296.
- [74] Ilya Loshchilov and Frank Hutter. "Decoupled Weight Decay Regularization". In: International Conference on Learning Representations. 2019. URL: https://openreview.net/forum?id=Bkg6RiCqY7.

## A Additional Drawbacks of DPO under Heterogeneity

Violating Independence of Irrelevant Alternatives (IIA). Suppose  $\mathcal{U} = \{A, B\}$  and types are equally represented. Given two possible responses  $y_1$  and  $y_2$ , type A prefers  $y_1$  but type B prefers  $y_2$ :

$$r^*(\mathbf{y}_1; A) = 6, r^*(\mathbf{y}_2; A) = 1,$$
  
 $r^*(\mathbf{y}_1; B) = 3, r^*(\mathbf{y}_2; B) = 9.$ 

A direct calculation shows  $\mathbb{E}_u[r^*(\boldsymbol{y}_2)] > \mathbb{E}_u[r^*(\boldsymbol{y}_1)]$  and NBC( $\boldsymbol{y}_2$ ) > NBC( $\boldsymbol{y}_1$ ). So, both  $\pi^*$  and  $\pi_{DPO}$  prefer  $\boldsymbol{y}_2$ . Let's consider another possible response  $\boldsymbol{y}_3$  which is not the preferred response for any user type:

$$r^*(\mathbf{y}_3; A) = r^*(\mathbf{y}_3; B) = 2.$$

While  $\pi^*$  still prefers  $y_2$  to  $y_1$ , now NBC( $y_1$ )  $\approx 0.62 > \text{NBC}(y_2) \approx 0.55$ , so, introducing an irrelevant alternative can alter DPO's ranking over existing alternatives.

**Tyranny of Majority.** Suppose  $\mathcal{U} = \{A, B\}$  with type A shaping 90% of the population. Given two responses  $y_1, y_2$ , type A slightly favors  $y_2$  but type B finds  $y_2$  offensive:

$$r^*(\mathbf{y}_1; A) = 0.5, r^*(\mathbf{y}_2; A) = 1,$$
  
 $r^*(\mathbf{y}_1; B) = 0.5, r^*(\mathbf{y}_2; B) = -10.$ 

In this case,  $\pi^*$  prefers  $y_1$  even though type B is a minority. In contrast, we have  $NBC(y_1) \approx 0.47$ ,  $NBC(y_2) \approx 0.53$ , which implies that the majority dominates in DPO.

#### **B** Additional Related Work

The challenge of handling heterogeneous preferences in alignment has been recognized as a significant problem in alignment research [44, 45, 11, 12]. This problem has attracted considerable attention from researchers in the field. Here, we highlight a few representative works that address key directions in tackling this challenge.

Analysis of DPO. Our study of how standard preference learning methods, such as DPO, behave in the presence of heterogeneous preferences was inspired by Siththaranjan, Laidlaw, and Hadfield-Menell [15]'s result, which shows that RLHF aggregates preferences according to a well-known voting rule called Borda count. Chakraborty et al. [27] highlights the impossibility of aligning with a singular reward model in RLHF by providing a lower bound on the gap between the optimal policy and a subpopulation's optimal policy. Dumoulin et al. [13] adopts a density estimation perspective on learning from human feedback to illustrate the challenges of preference learning from a population of annotators with diverse viewpoints. Rosset et al. [46] and Gao et al. [47] point out the limitations of point-wise reward models in expressing complex, intransitive preferences that may arise due to the aggregation of diverse preferences. Additionally, frameworks that generalize DPO and unify different alignment methods have been proposed to analyze current approaches and explore possible alternatives [48, 49, 30, 50].

**Policy Personalization.** Many works in the literature have proposed personalization as a solution to the problem of pluralistic alignment. Poddar et al. [39] propose a latent variable formulation of the problem and learn rewards and policies conditioned on it. Chen et al. [40] use an ideal point model for preferences and learn latent spaces representing different preferences. Mapping user information to user representations, Li, Lipton, and Leqi [51] perform personalized DPO to jointly learn a user model and a personalized language model. Balepur et al. [52] use abductive reasoning to infer user personas and train models to tailor responses accordingly. Lee et al. [53] explore the possibility of

steering a language model to align with a user's intentions through system messages. Dang et al. [54] extend personalized alignment to text-to-image diffusion models. Jang et al. [55] perform personalized alignment by decomposing preferences into multiple dimensions. Lau et al. [56] dynamically adapt the model to individual preferences using in-context learning.

Preference Aggregation. Closely aligned with our goal of serving the entire population with a single policy, several works have explored ways to aggregate diverse preferences. The rich literature on social choice theory has proven to be a valuable source of inspiration for studying existing preference learning approaches and proposing new ones [10, 57, 58, 11, 59]. Drawing insights from social choice theory, robustness to approximate clones has been proposed as a desirable property of RLHF algorithms, which current methods lack [29]. The Minimax Winner, a concept in preference aggregation, has inspired the use of the proportion of wins as the reward for a particular trajectory to align a model through self-play [60]. The impact of heterogeneity on strategic behavior in feedback and its effects on aggregation are also explored in Park et al. [14], which further examines the use of different social welfare functions for preference aggregation.

Methods. Solutions proposed to address different formulations of the problem span a wide range of methods. Siththaranjan, Laidlaw, and Hadfield-Menell [15] estimate a distribution of scores for alternatives to account for heterogeneity as hidden context. Chidambaram, Seetharaman, and Syrgkanis [41] propose an Expectation-Maximization (EM) version of DPO to minimize a notion of worst-case regret. Multi-objective reinforcement learning [61, 55] and its direct optimization variant [42] have also been proposed to align with diverse preferences. Wang et al. [62] train a multi-objective reward model to capture diverse preferences. Zhong et al. [63] use meta-learning to learn diverse preferences and aggregate them using different social welfare functions. Li et al. [64] design an optimal-transport-based loss to calibrate their model with the categorical distribution of preferences. Producing a Pareto front of models has also been explored as a solution. Boldi et al. [65] employ an iterative process to select solutions, while Rame et al. [66] interpolate the weights of independent networks linearly to achieve a Pareto-optimal generalization across preferences.

Empirical Observations. Empirical studies of alignment methods have had a significant impact on the study of preference learning. Zhang et al. [5] demonstrate the Bradley-Terry model's failure to distinguish between unanimous agreement among annotators and the majority opinion in cases of diverging user preferences. Chen et al. [67] show that RLHF and DPO struggle to improve ranking accuracy. Zeng et al. [68] study the role of model size and data size in the impact of diversified human preferences. Bansal, Dang, and Grover [69] demonstrate the significant influence of feedback protocol choice on alignment evaluation. Santurkar et al. [70] explore the opinions reflected by a language model, while Bakker et al. [71] investigate a language model's ability to generate consensus statements by training it to predict individual preferences. Jiang et al. [72] propose individualistic alignment to predict an individual's values, and Zollo et al. [73] introduce the PersonalLLM benchmark to measure a model's adaptation to a particular user's preferences.

#### C Additional Statements

**Proposition C.1.** There exists a mixture of BTs that a single BT cannot represent.

See proof on page 26.

**Definition C.2** (Learnability). Denote by  $\mathcal{D}_{r,\sigma}$  an i.i.d. sampled pairwise preference dataset labeled by random users with reward r and preference model  $\sigma$ . Let  $\bar{r}(\boldsymbol{y}) := \mathbb{E}_u[r(\boldsymbol{y};u)]$ . We say that the ranking based on  $\bar{r}$  is (weakly) learnable if, for some  $\epsilon > 0$ , there exists an algorithm with a bounded sample complexity m, such that for every reward r, when given a dataset  $\mathcal{D}_{r,\sigma}$  of size  $|\mathcal{D}_{r,\sigma}| \ge m(\epsilon,\bar{r})$ , it outputs a ranking consistent with  $\bar{r}$  with a probability at least  $\epsilon$  above the chance level.

**Proposition C.3.** Defining l in Eq. (17) as follows results in a consistent estimation of the optimal policy when preferences follow the BT model:

$$l(\mathbf{y}_1, \mathbf{y}_2, \mathbf{o}; \pi) = \begin{cases} -\sigma(h(\mathbf{y}_1, \mathbf{y}_2; \pi)) - I(\sigma(h(\mathbf{y}_1, \mathbf{y}_2; \pi))), & \mathbf{o} = \mathbf{1}, \\ -\sigma(h(\mathbf{y}_2, \mathbf{y}_1; \pi)) - I(\sigma(h(\mathbf{y}_2, \mathbf{y}_1; \pi))), & \mathbf{o} = \mathbf{0}, \\ 0 & o.w. \end{cases}$$

Here, we define  $I(\theta) := \int_1^{\theta} \left(\frac{1}{\theta'} - 1\right)^{|\mathcal{U}|} d\theta'$ , and h is the difference of  $\pi$ 's induced rewards (Eq. (15)). See proof on page 26.

### D Missing Proofs

**Proposition 3.1.** Consider an aggregation  $f: \mathbb{R}^{\mathcal{U}} \to \mathbb{R}$ . If  $f(\{r(\boldsymbol{y};u)\}_{u \in \mathcal{U}})$  induces the same ordering over  $\boldsymbol{y}$  for every reward r consistent with the preferences distribution, then under weak regulatory assumptions, f must be affine.

Proof of Proposition 3.1. First of all, if a reward function  $r^*(y; u)$  can explain the preferences of a user type u, any other reward function  $r(y; u) := r^*(y; u) + c(u)$  induces the same preference distribution:

$$\Pr(\mathbf{y}_2 \succ \mathbf{y}_1 \mid r; u) = \sigma(r(\mathbf{y}_2; u) - r(\mathbf{y}_1; u)) = \sigma(r^*(\mathbf{y}_2; u) - r^*(\mathbf{y}_1; u)) = \Pr(\mathbf{y}_2 \succ \mathbf{y}_1 \mid r^*; u).$$

Therefore,  $r^*$  is identifiable up to a bias term that can depend on the context and user type.

Consider a reward aggregation function  $f: \mathbb{R}^{\mathcal{U}} \to \mathbb{R}$ . Denoting all the rewards from different user types by a vector  $\mathbf{r}(\mathbf{y}) \in \mathbb{R}^{\mathcal{U}}$ , the aggregation  $f(\mathbf{r}(\mathbf{y}))$  should induce the same ranking for every  $\mathbf{r}$  consistent with (possibly infinite) preference data. Our above argument then implies that  $f(\mathbf{r}^*(\mathbf{y}) + \mathbf{c})$  should induce a consistent ranking for every  $\mathbf{c} \in \mathbb{R}^{\mathcal{U}}$ . For a sufficiently large space of alternatives, where  $\mathbf{r}^*(\mathbf{y})$  can take any value within a closed interval of  $\mathbb{R}$ , this is possible only if there exists a function  $\psi: \mathbb{R}^{\mathcal{U}} \to \mathbb{R}$  such that

$$f(\mathbf{r}_2 + \mathbf{c}) - f(\mathbf{r}_1 + \mathbf{c}) = \psi(\mathbf{r}_2 - \mathbf{r}_1),$$

for every c,  $r_1$ , and  $r_2$  in  $\mathbb{R}^{\mathcal{U}}$ . Choosing  $c = -r_1$ , this implies  $f(r) = f(0) + \psi(r)$  for every r. Therefore, we have the following Cauchy functional equation for  $\psi$ :

$$\psi(\mathbf{r} + \mathbf{\Delta}) = \psi(\mathbf{r}) + \psi(\mathbf{\Delta}).$$

Under weak regularity conditions such as the monotonicity or continuity of f, it is well-known that  $\psi$  has to be a linear function. This implies that f has to be an affine function which completes the proof.

**Proposition 4.2.** Suppose responses to  $\mathbf{x}$  in the preference dataset are drawn from  $\mathcal{D}(\cdot \mid \mathbf{x})$ . In the limit of many data points, DPO's induced reward, or equivalently  $\frac{\pi_{\mathrm{DPO}}(\cdot \mid \mathbf{x})}{\pi_{\mathrm{ref}}(\cdot \mid \mathbf{x})}$ , has the same ordering over responses as  $\mathrm{NBC}(\cdot \mid \mathbf{x})$ .

Proof of Proposition 4.2. We start from DPO's objective in Eq. (6). For notational simplicity, we assume  $\pi(\boldsymbol{y} \mid \boldsymbol{x})$  already contains a normalization by  $\pi_{\text{ref}}(\boldsymbol{y} \mid \boldsymbol{x})$ . In the limit of many data points, we can rewrite DPO's objective as the minimization of a cross-entropy loss

$$\mathcal{L}_{\text{DPO}}(\pi) := -\mathbb{E}_{\boldsymbol{x},\boldsymbol{y},\boldsymbol{y}'} \left[ \bar{\sigma} \left( \Delta r^*(\boldsymbol{x},\boldsymbol{y}',\boldsymbol{y}) \right) \cdot \log \sigma \left( \beta \log \frac{\pi(\boldsymbol{y} \mid \boldsymbol{x})}{\pi(\boldsymbol{y}' \mid \boldsymbol{x})} \right) + \left( 1 - \bar{\sigma} \left( \Delta r^*(\boldsymbol{x},\boldsymbol{y}',\boldsymbol{y}) \right) \right) \cdot \log \left( 1 - \sigma \left( \beta \log \frac{\pi(\boldsymbol{y} \mid \boldsymbol{x})}{\pi(\boldsymbol{y}' \mid \boldsymbol{x})} \right) \right) \right],$$

<sup>&</sup>lt;sup>2</sup>We can view NBC( $y \mid x$ ) as an aggregation of rewards at y. One can verify that NBC meets the order consistency condition of Proposition 3.1. However, it uses the reward value at  $y' \neq y$  to define the aggregated reward at y and thus does not fall under Proposition 3.1. In fact, this interdependency causes the issues we discuss Section 4.3.

where  $\bar{\sigma}(\Delta r^*(\boldsymbol{x}, \boldsymbol{y}', \boldsymbol{y}))$  is shorthand for  $\Pr(\boldsymbol{y} \succ \boldsymbol{y}' \mid \boldsymbol{x}; r^*) = \mathbb{E}_u[\sigma(r^*([\boldsymbol{x}, \boldsymbol{y}]; u) - r^*([\boldsymbol{x}, \boldsymbol{y}']; u))]$ . The minimizer of  $\mathcal{L}_{\text{DPO}}$  should meet the first-order condition:  $\frac{\partial \mathcal{L}_{\text{DPO}}}{\partial \pi(\boldsymbol{y} \mid \boldsymbol{x})} = 0$ , for every  $\boldsymbol{x}$  and  $\boldsymbol{y}$ . Then, a direct calculation shows that the optimal policy  $\pi^*$  meets

$$\mathbb{E}_{\boldsymbol{y}' \sim \mathcal{D}(\cdot|\boldsymbol{x})} \left[ \sigma \left( \beta \log \frac{\pi^*(\boldsymbol{y} \mid \boldsymbol{x})}{\pi^*(\boldsymbol{y}' \mid \boldsymbol{x})} \right) \right] - \mathbb{E}_{\boldsymbol{y}' \sim \mathcal{D}(\cdot|\boldsymbol{x})} \left[ \bar{\sigma} \left( \Delta r^*(\boldsymbol{x}, \boldsymbol{y}', \boldsymbol{y}) \right) \right] = 0.$$
 (18)

Recognize that the second term is  $NBC(y \mid x)$ :

$$NBC(\boldsymbol{y} \mid \boldsymbol{x}) = \mathbb{E}_{\boldsymbol{y}' \sim \mathcal{D}(\cdot \mid \boldsymbol{x})} \Big[ \mathbb{E}_u \Big[ \sigma \big( r^*([\boldsymbol{x}, \boldsymbol{y}]; u) - r^*([\boldsymbol{x}, \boldsymbol{y}']; u) \big) \Big] \Big].$$

In the absence of heterogeneity, we have  $\bar{\sigma} = \sigma$ , so setting  $\beta \log \pi^*(\boldsymbol{y} \mid \boldsymbol{x}) = r^*([\boldsymbol{x}, \boldsymbol{y}]) + C(\boldsymbol{x})$  for a normalizing C would solve Eq. (18). In general, we are not aware of any closed-form solution. However, we can still infer the ordering the optimal policy induces from Eq. (18): Since the first term is increasing in  $\pi^*(\boldsymbol{y} \mid \boldsymbol{x})$ , the optimal policy will be monotone in NBC( $\boldsymbol{y} \mid \boldsymbol{x}$ ). This completes the proof.

**Proposition 5.1.** If there are at least two alternatives and two user types with a continuous preference model, the ranking based on the user-weighted expected reward is not learnable without annotator information.

*Proof of Proposition 5.1.* We give three related proof strategies. The first strategy works for every preference model. The second strategy draws on a connection to a robust version of Arrow's impossibility theorem [32]. The third strategy is inspired by Procaccia, Schiffer, and Zhang [29]. We start with the notation and definitions specific to this proof.

Notation and Definitions. Consider a fixed prompt  $\boldsymbol{x}$  with a set of possible responses  $\mathcal{Y}$ . Let R denote a *complete ranking* over  $\mathcal{Y}$ , where  $\boldsymbol{y}_2 R \boldsymbol{y}_1$  indicates whether  $\boldsymbol{y}_2 \succ \boldsymbol{y}_1$  or vice versa. A *profile* refers to a set of complete rankings. For a heterogeneous reward function  $r(\boldsymbol{y};u)$  and a prior  $\mathcal{P}$  over user types  $\mathcal{U}$ , let  $R_{\bar{r}}$  be the ranking according to  $\bar{r}(\boldsymbol{y}) := \mathbb{E}_{u \sim \mathcal{P}}[r(\boldsymbol{y};u)]$ .

A pairwise preference dataset  $\mathcal{D}$  consists of tuples  $(y_1, y_2, o)$ , where  $o := \mathbb{1}\{y_2 \succ y_1\}$ . We assume that  $y_1$  and  $y_2$  are i.i.d. draws. When a random user with a reward function r labels each instance in the dataset, we denote the resulting dataset by  $\mathcal{D}_r$ . A pairwise learning algorithm  $\mathcal{A}$  produces a complete ranking over  $\mathcal{Y}$  based on the pairwise preference dataset  $\mathcal{D}$ .

**Proof Strategy 1.** Suppose there exists an algorithm  $\mathcal{A}$  such that for some  $\mathcal{Y}$  with  $|\mathcal{Y}| \geq 2$ , for any reward function r and any preference dataset  $\mathcal{D}_r$  with  $|\mathcal{D}_r| \geq n_{\bar{r}}$ , it outputs  $R_{\bar{r}}$  on  $\mathcal{Y}$  with a probability of at least  $\frac{1}{|\mathcal{Y}|!} + \epsilon$ .

Suppose r is a heterogeneous reward function that its expectation induces a complete ranking  $R_{\bar{r}}$  with no tie. Define a new heterogeneous reward function  $r_{\gamma}$  as follows. Consider a new user type  $0 \notin \mathcal{U}$ . For some  $\gamma > 1$ , let  $r_{\gamma}(\boldsymbol{y}; u) = \gamma r(\boldsymbol{y}; u)$  when  $u \neq 0$ , and  $r_{\gamma}(\boldsymbol{y}; 0) = 0$ . Define a new user distribution  $\mathcal{P}_{\gamma}(u) \coloneqq (1 - \frac{1}{\gamma}) \mathbb{I}\{u = 0\} + \frac{1}{\gamma} \mathcal{P}(u)$ . It is straightforward to verify  $\bar{r}_{\gamma} \coloneqq \mathbb{E}_{u \sim \mathcal{P}_{\gamma}}[r_{\gamma}(\boldsymbol{y}; u)] = \bar{r}$ . Therefore, with high probability,  $\mathcal{A}$  outputs  $R_{\bar{r}_{\gamma}} = R_{\bar{r}}$  from  $\mathcal{D}_{r_{\gamma}}$  for every  $\gamma > 1$ :

$$\Pr\left(\mathcal{A}(\mathcal{D}_{r_{\gamma}}) = R_{\bar{r}}\right) \ge \frac{1}{|\mathcal{Y}|!} + \epsilon.$$

As we increase  $\gamma$ , for any continuous preference model  $\sigma$ , the pairwise preference dataset  $\mathcal{D}_{r_{\gamma}}$  approaches a uniform preference dataset  $\mathcal{D}_{\text{unif}}$  labeled mostly by an indifferent annotator of type u = 0. So, we have

$$\Pr\left(\mathcal{A}(\mathcal{D}_{\text{unif}}) = R_{\bar{r}}\right) \ge \frac{1}{|\mathcal{Y}|!} + \epsilon. \tag{19}$$

This is true for every r. For different choices of r, agreements with  $R_{\bar{r}}$  are disjoint events. Since there are  $|\mathcal{Y}|!$  different rankings overall, the pigeonholed principle implies  $\epsilon = 0$ .

**Proof Strategy 2.** The proof is by contradiction. Suppose there exists an algorithm  $\mathcal{A}$  that for any reward function r and any preference dataset  $\mathcal{D}_r$  with  $|\mathcal{D}_r| \geq n_{\bar{r}}$ , it outputs  $R_{\bar{r}}$  with a probability of at least  $1 - \epsilon$ . We follow Friedgut, Kalai, and Naor [32] and define a social choice function as a function that yields an asymmetric relation on the alternatives given a profile. A social choice is rational if it is an order relation on the alternatives, and is neutral if it is invariant under permutations of alternatives.

Let  $\mathcal{Y}_3 = \{y_1, y_2, y_3\}$  be an arbitrary subset of  $\mathcal{Y}$  with size 3. Next, we construct a neutral social choice function f acting on  $\mathcal{Y}_3$  that is independent of irrelevant alternatives (IIA). Here is how we design f: Let  $P_r$  be a profile of size  $n \geq n_{\bar{r}}$  at the input, where every  $R \in P_r$  is an i.i.d. draw from a Plackett–Luce (PL) ranking model with  $\exp(r)$  as the weight of the alternatives. Note that the marginal distribution induced by PL on any two alternatives follows BT. Create three pairwise preference datasets  $\mathcal{D}_{r,12}$ ,  $\mathcal{D}_{r,23}$ , and  $\mathcal{D}_{r,13}$  where  $\mathcal{D}_{r,ij} = \{y_i R y_j \mid R \in P_r\}$ . The social function f applies  $\mathcal{A}$  to every dataset to obtain a relation over  $\mathcal{Y}_3$ . By construction, f is neutral and IIA.

By assumption, for every internal dataset  $\mathcal{D}_{r,ij}$  we have  $\Pr\left(\mathcal{A}(\mathcal{D}_{r,ij}) = R_{r,ij}\right) \geq 1 - \epsilon$ , where  $R_{r,ij}$  is the projection of  $R_{\bar{r}}$  to only two alternatives  $y_i$  and  $y_j$ . Using union bound, we have

$$\Pr\left(f(P_r) = R_{\bar{r}}\right) \ge 1 - 3\epsilon$$
.

Similar to strategy 1, we can define a new reward function  $r_{\gamma}$  with  $\bar{r}_{\gamma} = \bar{r}$  such that the above holds for every  $r_{\gamma}$  with  $\gamma > 1$ . Then, by increasing  $\gamma$ , the profile  $P_{r_{\gamma}}$  approaches a uniformly distributed profile  $P_{\text{unif}}$ . In this case, since  $R_{\bar{r}}$  is an order relation, we have  $\Pr\left(f \text{ is rational}\right) \geq 1 - 3\epsilon$ . Then Theorem 1.3 of Friedgut, Kalai, and Naor [32] implies that for some global constant K,

$$\Pr(f \text{ is dictatorship}) \ge 1 - 3K\epsilon$$
.

On the other hand, we know that the order that  $R_{\bar{r}}$  induces for different r is not a dictatorship, so, we have

$$\Pr(f \text{ is dictatorship}) < 3\epsilon.$$

Putting these together, we obtain a lower bound on  $\epsilon$ :

$$\epsilon \ge \frac{1}{3(K+1)} > 0.$$

The rest of the proof is similar for the second and third strategies. We can use a boosting argument to show that from any weak pairwise learner, we can obtain an arbitrarily strong weak learner corresponding at the cost of collecting a larger dataset. We show this when for the weak learner  $\epsilon < \frac{1}{2}$  but a weaker condition of choosing  $R_{\bar{r}}$  better than chance level is sufficient for our argument. Consider a pairwise preference dataset  $\mathcal{D}_r$  of size  $m \, n_{\bar{r}}$ . Partition  $\mathcal{D}_r$  into m equal-size datasets. Let  $R_i$  be the output of  $\mathcal{A}$  on the  $i^{\text{th}}$  dataset. Since samples in  $\mathcal{D}_r$  are independently generated,  $R_i$ s err independently. Construct a meta-algorithm  $\mathcal{A}_{\text{maj}}$  that outputs the majority winner of  $R_i$ s. A standard Hoeffding bound implies

$$\Pr\left(\mathcal{A}_{\mathrm{maj}}(\mathcal{D}_r) \neq R_{\bar{r}}\right) \leq \exp\left(-2(1/2 - \epsilon)^2 m\right).$$

A simple calculation then shows that for any arbitrarily small  $\epsilon' > 0$ , by choosing  $m = O(\log(\frac{1}{\epsilon'})(\frac{1}{2} - \epsilon)^{-2})$  the majority-winner algorithm agrees with  $R_{\bar{r}}$  with probability of at least  $1 - \epsilon'$ . This contradicts the lower bound we established earlier and completes the proof.

**Proof Strategy 3.** The proof is by contradiction and is inspired by Procaccia, Schiffer, and Zhang [29]. This proof requires at least four different user types. Suppose there are two equally represented user types  $\mathcal{U} = \{A, B\}$  who follow BT. For some arbitrary response  $y_0 \in \mathcal{Y}$  and  $0 < \tau < \frac{1}{3}$ , consider

the following reward function:

$$r_{\tau}(\boldsymbol{y}; u) = \begin{cases} 0, & \boldsymbol{y} \neq \boldsymbol{y}_{0}, \\ \sigma^{-1}(\frac{2}{3} + \tau) = \log \frac{\frac{2}{3} - \tau}{\frac{1}{3} + \tau}, & \boldsymbol{y} = \boldsymbol{y}_{0}, u = A, \\ \sigma^{-1}(\frac{2}{3} - \tau) = \log \frac{\frac{2}{3} + \tau}{\frac{1}{3} - \tau}, & \boldsymbol{y} = \boldsymbol{y}_{0}, u = B. \end{cases}$$
(20)

One can see  $\Pr(\boldsymbol{y}_0 \succ \boldsymbol{y} \mid r_{\tau}) = \frac{2}{3}$  for every  $\boldsymbol{y} \neq \boldsymbol{y}_0$ . Therefore,  $r_{\tau}$  induces the same pairwise preference distribution for every  $\tau$ . On the other hand,  $\bar{r}_{\tau}(\boldsymbol{y}_0) := \mathbb{E}_u[r_{\tau}(\boldsymbol{y}_0; u)] = \log \frac{\frac{4}{9} - \tau^2}{\frac{1}{9} - \tau^2} > 0$  is an increasing function of  $\tau$ .

Consider two arbitrary  $\tau_1$  and  $\tau_2$  such that  $0 < \tau_1 < \tau_2 < \frac{1}{3}$ . Sample a pairwise preference dataset at follows. Draw a random u and a random permutation  $\rho$  over  $\mathcal{Y}$ . If  $\rho$  is not identity, ask an annotator of type u with reward  $r_{\tau_1}(\cdot;u)$  to label this sample and permute the ranking with  $\rho$ . If  $\rho$  is identity, ask an annotator of type u with reward  $r_{\tau_2}(\cdot;u)$  to label. This sampling is equivalent to sampling from  $2 * |\mathcal{Y}|!$  different user types. By symmetry, this pairwise preference dataset is distributionally equivalent to a dataset  $\mathcal{D}_{\text{unif}}$  with indifferent preferences. However, our construction implies that  $y_0$  has the highest expected reward and other alternatives have similar rewards:

$$ar{r}(oldsymbol{y}) = rac{1}{|\mathcal{Y}|} egin{cases} ar{r}_{ au_1}(oldsymbol{y})\,, & oldsymbol{y} 
eq oldsymbol{y}_0\,, \ ar{r}_{ au_2}(oldsymbol{y})\,, & oldsymbol{y} = oldsymbol{y}_0\,. \end{cases}$$

Denote the ranking based on  $\bar{r}$  above by  $R_0$ .

Similar to the second strategy, suppose there exists an algorithm  $\mathcal{A}$  that for any reward function r and any preference dataset  $\mathcal{D}_r$  with  $|\mathcal{D}_r| \geq n_{\bar{r}}$ , it outputs  $R_{\bar{r}}$  with a probability of at least  $1 - \epsilon$ . Collect a preference dataset as explained above with at least  $n_{\bar{r}}$  samples. Then, by assumption,

$$\Pr\left(\mathcal{A}(\mathcal{D}_{\text{unif}}) = R_0\right) \ge 1 - \epsilon$$
.

Note that our choice of  $y_0$  could be any of the alternatives. Therefore, the above should be true when any of the alternatives has the highest expected reward. This implies a lower bound on  $\epsilon$ :

$$\epsilon \ge 1 - \frac{1}{|\mathcal{Y}|} > 0.$$

The rest of the proof is similar to the second strategy.

**Proposition 5.2.** There is no M-estimator that can estimate  $V(\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2) := \operatorname{Var}_u \left[ \Delta r^*(\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2; u) \right]$  consistently without annotator information.

Proof of Proposition 5.2. Consider a dataset  $\mathcal{D}$  of context, candidate pairs, and preference represented as  $(\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2, o)$ , where  $o = \mathbb{1}\{\boldsymbol{y}_2 \succ \boldsymbol{y}_1\}$ , and  $\boldsymbol{y}_1, \boldsymbol{y}_2$  are independently drawn. Then consider an Mestimator

$$\underset{V}{\arg\min} \sum_{(\boldsymbol{x},\boldsymbol{y}_{l},\boldsymbol{y}_{w}) \in \mathcal{D}} \rho(\boldsymbol{x},\boldsymbol{y}_{l},\boldsymbol{y}_{w};V) = \sum_{(\boldsymbol{x},\boldsymbol{y}_{1},\boldsymbol{y}_{2},o) \in \mathcal{D}} o \cdot \rho(\boldsymbol{x},\boldsymbol{y}_{1},\boldsymbol{y}_{2};V) + (1-o) \cdot \rho(\boldsymbol{x},\boldsymbol{y}_{2},\boldsymbol{y}_{1};V) \,.$$

Under preference model of Eq. (2), for a reward r, we have  $\mathbb{E}[o \mid \boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2; u] = \sigma(\Delta r(\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2; u))$ . In the limit of a large dataset, the M-estimator solves

$$\underset{V}{\operatorname{arg \,min}} \ \mathbb{E}_{\boldsymbol{x},\boldsymbol{y}_{1},\boldsymbol{y}_{2},o,u} \Big[ o \cdot \rho(\boldsymbol{x},\boldsymbol{y}_{1},\boldsymbol{y}_{2};V) + (1-o) \cdot \rho(\boldsymbol{x},\boldsymbol{y}_{2},\boldsymbol{y}_{1};V) \Big] \\
= \mathbb{E}_{\boldsymbol{x},\boldsymbol{y}_{1},\boldsymbol{y}_{2},u} \Big[ \sigma \big( \Delta r(\boldsymbol{x},\boldsymbol{y}_{1},\boldsymbol{y}_{2};u) \big) \cdot \rho(\boldsymbol{x},\boldsymbol{y}_{1},\boldsymbol{y}_{2};V) + \sigma \big( \Delta r(\boldsymbol{x},\boldsymbol{y}_{2},\boldsymbol{y}_{1};u) \big) \cdot \rho(\boldsymbol{x},\boldsymbol{y}_{2},\boldsymbol{y}_{1};V) \Big] \qquad (V \text{ has no } o) \\
= \mathbb{E}_{\boldsymbol{x},\boldsymbol{y}_{1},\boldsymbol{y}_{2}} \Big[ \mathbb{E}_{u} \Big[ \sigma \big( \Delta r(\boldsymbol{x},\boldsymbol{y}_{1},\boldsymbol{y}_{2};u) \big) \Big] \cdot \rho(\boldsymbol{x},\boldsymbol{y}_{1},\boldsymbol{y}_{2};V) + \mathbb{E}_{u} \Big[ \sigma \big( \Delta r(\boldsymbol{x},\boldsymbol{y}_{2},\boldsymbol{y}_{1};u) \big) \Big] \cdot \rho(\boldsymbol{x},\boldsymbol{y}_{2},\boldsymbol{y}_{1};V) \Big] \qquad (V \text{ has no } u) \\
= 2 \mathbb{E}_{\boldsymbol{x},\boldsymbol{y}_{1},\boldsymbol{y}_{2}} \Big[ \mathbb{E}_{u} \Big[ \sigma \big( \Delta r(\boldsymbol{x},\boldsymbol{y}_{1},\boldsymbol{y}_{2};u) \big) \Big] \cdot \rho(\boldsymbol{x},\boldsymbol{y}_{1},\boldsymbol{y}_{2};V) \Big] . \qquad (\boldsymbol{y}_{1},\boldsymbol{y}_{2} \text{ are i.i.d.})$$

Here, we used the fact that V does not depend on u. We also relied on the assumption that  $y_1$  and  $y_2$  are identically and independently distributed. The above equation suggests that regardless of how  $\rho$  is designed,  $V(x, y_1, y_2)$  can only depend on u's distribution through  $\mathbb{E}_u\left[\sigma\left(\Delta r(x, y_1, y_2; u)\right)\right]$ . Therefore, no consistent M-estimator can generally estimate  $\operatorname{Var}_u\left[\Delta r(x, y_1, y_2; u)\right]$  even with the availability of infinite preference data.

**Lemma 5.3.** Using  $J_1$  and  $J_2$  as shorthands for  $J(x, y_1, y_2, x, y_1, y_2)$  and  $J(x, y_1, y_2, x, y_2, y_1)$ , we can use the following to estimate the variance term:

$$V(\mathbf{y}_1, \mathbf{y}_2) = \frac{J_1 - (J_1 + J_2)^2}{\sigma'(\Delta \bar{r}^*(\mathbf{y}_1, \mathbf{y}_2))^2}.$$
 (16)

Proof of Lemma 5.3. First of all, J can give us the likelihood itself:

$$J(\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2, \boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2) + J(\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2, \boldsymbol{x}, \boldsymbol{y}_2, \boldsymbol{y}_1) = \mathbb{E}_u \left[ \sigma \left( \Delta r^*(\boldsymbol{y}_1, \boldsymbol{y}_2; u) \right)^2 + \sigma \left( \Delta r^*(\boldsymbol{y}_1, \boldsymbol{y}_2; u) \right) \cdot \sigma \left( \Delta r^*(\boldsymbol{y}_2, \boldsymbol{y}_1; u) \right) \right]$$
$$= \mathbb{E}_u \left[ \sigma \left( \Delta r^*(\boldsymbol{y}_1, \boldsymbol{y}_2; u) \right) \right].$$

Here, we used the property  $\sigma(\Delta r^*(\boldsymbol{y}_2, \boldsymbol{y}_1; u)) = 1 - \sigma(\Delta r^*(\boldsymbol{y}_1, \boldsymbol{y}_2; u))$ . We also dropped  $\boldsymbol{x}$  from the notation for simplicity. Since J can give us both the first and second moments, we can use it to find  $\operatorname{Var}_u[\sigma(\Delta r^*(\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2; u))]$  as follows:

$$\operatorname{Var}_{u}\left[\sigma\left(\Delta r^{*}(\boldsymbol{x},\boldsymbol{y}_{1},\boldsymbol{y}_{2};u)\right)\right] = \mathbb{E}_{u}\left[\sigma\left(\Delta r^{*}(\boldsymbol{x},\boldsymbol{y}_{1},\boldsymbol{y}_{2};u)\right)^{2}\right] - \mathbb{E}_{u}\left[\sigma\left(\Delta r^{*}(\boldsymbol{x},\boldsymbol{y}_{1},\boldsymbol{y}_{2};u)\right)\right]^{2}$$

$$= J(\boldsymbol{x},\boldsymbol{y}_{1},\boldsymbol{y}_{2},\boldsymbol{x},\boldsymbol{y}_{1},\boldsymbol{y}_{2}) - \left(J(\boldsymbol{x},\boldsymbol{y}_{1},\boldsymbol{y}_{2},\boldsymbol{x},\boldsymbol{y}_{1},\boldsymbol{y}_{2}) + J(\boldsymbol{x},\boldsymbol{y}_{1},\boldsymbol{y}_{2},\boldsymbol{x},\boldsymbol{y}_{2},\boldsymbol{y}_{1})\right)^{2}.$$

In the last piece of the proof, we connect  $\operatorname{Var}_u[\sigma(\Delta r^*)]$  with  $\operatorname{Var}_u[\Delta r^*]$ . The Taylor expansion of  $\sigma(\Delta r^*)$  around  $\Delta \bar{r}^* := \mathbb{E}_u[\Delta r^*]$  gives

$$\operatorname{Var}_{u}\left[\sigma(\Delta r^{*})\right] = \operatorname{Var}_{u}\left[\sigma(\Delta \bar{r}^{*}) + \sigma'(\Delta \bar{r}^{*}) \cdot (\Delta r^{*} - \Delta \bar{r}^{*}) + O((\Delta r^{*} - \Delta \bar{r}^{*})^{2})\right]$$
$$= \sigma'(\Delta \bar{r}^{*})^{2} \cdot \operatorname{Var}_{u}\left[\Delta r^{*}\right] + O\left(\mathbb{E}_{u}\left[(\Delta r^{*} - \Delta \bar{r}^{*})^{3}\right]\right).$$

We can neglect the third-order term in calculations as first-order correction uses up to  $O\left(\mathbb{E}_u\left[(\Delta r^* - \Delta \bar{r}^*)^2\right]\right)$  in its approximation.

**Proposition 6.1.** Defining l in Eq. (17) as follows results in a consistent estimation of the optimal policy when preferences follow the BT model:

$$l(\mathbf{y}_1, \mathbf{y}_2, \mathbf{o}; \pi) = \begin{cases} -\log \sigma(|\mathcal{U}| \cdot h(\mathbf{y}_1, \mathbf{y}_2; \pi)), & \mathbf{o} = \vec{\mathbf{1}}, \\ -\log \sigma(|\mathcal{U}| \cdot h(\mathbf{y}_2, \mathbf{y}_1; \pi)), & \mathbf{o} = \vec{\mathbf{0}}, \\ 0 & o.w. \end{cases}$$

Here, h is the difference of  $\pi$ 's induced rewards (Eq. (15)).

*Proof of Proposition 6.1.* The proof follows similar steps as the derivation of DPO. First of all, conditioned on agreement, the likelihood of observing  $y_2 > y_1$  under the BT model is

$$\Pr(\mathbf{y}_{2} \succ \mathbf{y}_{1} \mid r^{*}, \text{agreement}) = \frac{\Pi_{u}\sigma(\Delta r^{*}(\mathbf{y}_{1}, \mathbf{y}_{2}; u))}{\Pi_{u}\sigma(\Delta r^{*}(\mathbf{y}_{1}, \mathbf{y}_{2}; u)) + \Pi_{u}\sigma(\Delta r^{*}(\mathbf{y}_{2}, \mathbf{y}_{1}; u))}$$

$$= \frac{\exp(\sum_{u} r^{*}(\mathbf{y}_{2}; u))}{\exp(\sum_{u} r^{*}(\mathbf{y}_{1}; u)) + \exp(\sum_{u} r^{*}(\mathbf{y}_{2}; u))}$$

$$= \sigma(|\mathcal{U}| \cdot \mathbb{E}_{u}[\Delta r^{*}(\mathbf{y}_{1}, \mathbf{y}_{2})]).$$

On the other hand, Eq. (11) allows us to write  $\mathbb{E}_u[\Delta r^*(\boldsymbol{y}_1, \boldsymbol{y}_2)]$  with difference  $\pi$ 's induced rewards, i.e., h:

$$\Pr(\mathbf{y}_2 \succ \mathbf{y}_1 \mid \pi^*, \text{agreement}) = \sigma(h(\mathbf{y}_1, \mathbf{y}_2; \pi^*)).$$

We can define the likelihood in this way for every policy  $\pi$ . Then, the proposed loss function is equivalent to maximizing log-likelihood, which under mild conditions is a consistent estimator for  $\pi^*$ .

**Theorem 6.2.** Suppose l in Eq. (17) only depends on  $(\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2)$  through  $\pi$  and  $\pi_{\text{ref}}$ . If there are more than three types of user and the preferences follow BT, any loss  $\mathcal{L}$  that allows a consistent estimation of the optimal policy discards samples with disagreement, i.e., those with  $\mathbf{o} \notin \{\mathbf{0}, \mathbf{1}\}$ .

Proof of Theorem 6.2. The proof involves three steps: First, the next lemma shows that any loss function l in Eq. (17) with the desired consistency property can only depend on  $\pi$  through the ratio  $\frac{\pi(y_2|x)}{\pi(y_1|x)}$ .

**Lemma D.1.** Suppose l in Eq. (17) only depends on  $(\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2)$  through  $\pi$  and  $\pi_{ref}$ . Then, for any l that gives a consistent estimation of the optimal policy in Eq. (10), there exists an equivalent loss  $\tilde{l}$  such that

$$l(\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2, \boldsymbol{o}; \pi) = \tilde{l}(\boldsymbol{o}; h(\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2; \pi)),$$

where h is defined in Eq. (15).

See proof on page 27.

In the second step, we further limit the search space of  $\tilde{l}$  (as introduced by Lemma D.1) to those that meet certain first- and second-order conditions:

**Lemma D.2.** Any loss  $\tilde{l}$  as in Lemma D.1 that leads to a consistent estimation of the optimal policy meets

$$\begin{split} \sum_{\boldsymbol{o} \in \{0,1\}^{\mathcal{U}}} \frac{\partial \tilde{l}}{\partial \theta} \big( \boldsymbol{o}; \theta^*(\boldsymbol{z}) \big) \cdot \chi_{\boldsymbol{o}} \big( \boldsymbol{z} \big) &= 0 \,, \\ \sum_{\boldsymbol{o} \in \{0,1\}^{\mathcal{U}}} \frac{\partial^2 \tilde{l}}{\partial \theta^2} \big( \boldsymbol{o}; \theta^*(\boldsymbol{z}) \big) \cdot \chi_{\boldsymbol{o}} \big( \boldsymbol{z} \big) &\geq 0 \,, \end{split}$$

for every  $z \in [0,1]^{\mathcal{U}}$ . Here, we define

$$\chi_{\boldsymbol{o}}(\boldsymbol{z}) := \prod_{u \in \mathcal{U}} z_u^{o_u} (1 - z_u)^{1 - o_u} .$$

and

$$\theta^*(z) \coloneqq rac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \sigma^{-1}(z_u).$$

See proof on page 29.

Finally, we show that when preferences follow the BT model and there are more than three user types, all  $\tilde{l}(\boldsymbol{o};\theta)$  terms corresponding to  $\boldsymbol{o} \notin \mathbf{0}, \mathbf{1}$  do not depend on  $\theta$ . Therefore, these terms do not depend on  $\pi$  and can be removed from the loss function, thereby completing the proof.

**Lemma D.3.** If  $|\mathcal{U}| > 3$ , for any loss  $\tilde{l}$  that meets the first-order condition of Lemma D.2, we have  $\frac{\partial \tilde{l}}{\partial \theta}(\mathbf{o};\theta) = 0$  for every  $\mathbf{o} \notin \{\mathbf{0},\mathbf{1}\}$ .

See proof on page 29.

**Proposition C.1.** There exists a mixture of BTs that a single BT cannot represent.

Proof of Proposition C.1. Suppose the pairwise comparison distribution over a set of alternatives  $(y_1, y_2, y_3, ...)$  satisfies the Bradley-Terry (BT) model; i.e.  $\Pr(y_i \succ y_j) = \sigma(r^*(y_2) - r^*(y_1))$ . Then:

$$\Pr(\mathbf{y}_{1} \succ \mathbf{y}_{2}) \Pr(\mathbf{y}_{2} \succ \mathbf{y}_{3}) \Pr(\mathbf{y}_{3} \succ \mathbf{y}_{1}) = \frac{\prod_{i=1}^{3} \exp(r^{*}(\mathbf{y}_{i}))}{\prod_{i=1}^{3} \left( \exp(r^{*}(\mathbf{y}_{i})) + \exp(r^{*}(\mathbf{y}_{(i+1) \bmod 3+1})) \right)}$$
$$= \Pr(\mathbf{y}_{1} \succ \mathbf{y}_{3}) \Pr(\mathbf{y}_{3} \succ \mathbf{y}_{2}) \Pr(\mathbf{y}_{2} \succ \mathbf{y}_{1}).$$

Now, consider two BT models corresponding to  $u_1$  and  $u_2$ , with a uniform mixture over them. For the mixture:

$$\Pr(\mathbf{y}_i \succ \mathbf{y}_j) = \frac{\Pr(\mathbf{y}_i \succ \mathbf{y}_j \mid u_1) + \Pr(\mathbf{y}_i \succ \mathbf{y}_j \mid u_2)}{2}.$$

The probability of cyclic preferences in one direction is given by

$$\Pr(\boldsymbol{y}_1 \succ \boldsymbol{y}_2) \Pr(\boldsymbol{y}_2 \succ \boldsymbol{y}_3) \Pr(\boldsymbol{y}_3 \succ \boldsymbol{y}_1) = \frac{\sum_{s \in \{1,2\}^3} \prod_{i=1}^3 \Pr(\boldsymbol{y}_i \succ \boldsymbol{y}_{(i+1) \bmod 3+1} \mid u_{s_i})}{8},$$

which is not necessarily equal to the probability of the cyclic preferences in the reverse direction:

$$\Pr(\boldsymbol{y}_1 \succ \boldsymbol{y}_3) \Pr(\boldsymbol{y}_3 \succ \boldsymbol{y}_2) \Pr(\boldsymbol{y}_2 \succ \boldsymbol{y}_1) = \frac{\sum_{s \in \{1,2\}^3} \prod_{i=1}^3 \Pr(\boldsymbol{y}_{(i+1) \bmod 3+1} \succ \boldsymbol{y}_i \mid u_{s_i})}{8}.$$

To verify this, consider specific examples such as  $\Pr(\boldsymbol{y}_i \succ \boldsymbol{y}_j \mid u_k) = \frac{\exp(r_k^i)}{\exp(r_k^i) + \exp(r_k^j)}$  with  $r_1 = (1, 2, 3)$  and  $r_2 = (1, 2, 4)$ . More generally, the BT assumption implies that, for a fixed reward  $r^*$ , the likelihood of a set of pairwise comparisons  $\{(\boldsymbol{y}_{p,1} > \boldsymbol{y}_{p,2})\}_{p \in [P]}$  is proportional to  $\prod_i \exp(r^*(\boldsymbol{y}_i))^{|\{p \in [P] \mid \boldsymbol{y}_{p,1} = i\}|}$  and depends only on the number of times each option is preferred in the comparisons. However, as demonstrated above, this property does not hold for a mixture of BT models.

**Proposition C.3.** Defining l in Eq. (17) as follows results in a consistent estimation of the optimal policy when preferences follow the BT model:

$$l(\mathbf{y}_1, \mathbf{y}_2, \mathbf{o}; \pi) = \begin{cases} -\sigma(h(\mathbf{y}_1, \mathbf{y}_2; \pi)) - I(\sigma(h(\mathbf{y}_1, \mathbf{y}_2; \pi))), & \mathbf{o} = \mathbf{1}, \\ -\sigma(h(\mathbf{y}_2, \mathbf{y}_1; \pi)) - I(\sigma(h(\mathbf{y}_2, \mathbf{y}_1; \pi))), & \mathbf{o} = \mathbf{0}, \\ 0 & o.w. \end{cases}$$

Here, we define  $I(\theta) \coloneqq \int_1^\theta \left(\frac{1}{\theta'} - 1\right)^{|\mathcal{U}|} d\theta'$ , and h is the difference of  $\pi$ 's induced rewards (Eq. (15)).

Proof of Proposition C.3. Recall  $\Pr(o_u = 1 \mid \boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2) = \sigma(\Delta r^*(\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2; u))$ . We use  $z_u(\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2)$  as a shorthand for this quantity and will drop the dependence on  $(\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2)$  whenever it is clear from the context. We also use s as a shorthand for  $\sigma(h)$ . In the limit of a very large dataset, the proposed loss approaches

$$\mathcal{L}(s) = -\mathbb{E}_{\boldsymbol{x},\boldsymbol{y}_1,\boldsymbol{y}_2} \left[ \left( \prod_{u \in \mathcal{U}} z_u \right) \left( s + I(s) \right) + \left( \prod_{u \in \mathcal{U}} (1 - z_u) \right) \left( 1 - s + I(1 - s) \right) \right].$$

Note that we wrote  $\mathcal{L}$  as a function of s instead of  $\pi$  since s is the only place that  $\pi$  appears. We first show that  $\mathcal{L}(s)$  has a unique global minimizer. To show an s is a global minimizer of  $\mathcal{L}$ , it suffices to show that s minimizes the term inside expectation for every  $(x, y_1, y_2)$ . Such a minimizer meets the first-order condition:

$$\left(\prod_{u \in \mathcal{U}} z_u\right) \left(1 + \left(\frac{1}{s} - 1\right)^{|\mathcal{U}|}\right) + \left(\prod_{u \in \mathcal{U}} (1 - z_u)\right) \left(-1 - \left(\frac{1}{1 - s} - 1\right)^{|\mathcal{U}|}\right) = 0.$$

Here, we used  $\frac{dI}{d\theta} = (\frac{1}{\theta} - 1)^{|\mathcal{U}|}$ . Define  $w := (\frac{1-s}{s})^{|\mathcal{U}|}$ . Then, the above condition reduces to a quadratic equation in terms of w:

$$1 + w - \left(\prod_{u \in \mathcal{U}} \left(\frac{1}{z_u} - 1\right)\right) (1 + w^{-1}) = (1 + w^{-1}) \left[w - \prod_{u \in \mathcal{U}} \left(\frac{1}{z_u} - 1\right)\right] = 0.$$

Solving for w, we obtain

$$s^* = \frac{1}{1 + \left(\prod_{u \in \mathcal{U}} \left(\frac{1}{z_u} - 1\right)\right)^{\frac{1}{|\mathcal{U}|}}}$$
.

For the BT model, a direct calculation then shows

$$s^*(\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2) = \sigma\left(\frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \Delta r(\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2; u)\right). \tag{21}$$

In fact,  $s^*$  is the only global minimizer of  $\mathcal{L}(s)$ . This is because  $\mathcal{L}(s)$  is convex in s:

$$\frac{\mathrm{d}^{2}\mathcal{L}}{\mathrm{d}s^{2}} = \mathbb{E}_{\boldsymbol{x},\boldsymbol{y}_{1},\boldsymbol{y}_{2}} \left[ \left( \prod_{u \in \mathcal{U}} z_{u} \right) \cdot \frac{|\mathcal{U}|}{s^{2}} \left( \frac{1}{s} - 1 \right)^{|\mathcal{U}| - 1} + \left( \prod_{u \in \mathcal{U}} (1 - z_{u}) \right) \cdot \frac{|\mathcal{U}|}{(1 - s)^{2}} \left( \frac{1}{1 - s} - 1 \right)^{|\mathcal{U}| - 1} \right] \geq 0.$$

Finally, one can verify that the policy that results in  $s^*$  (Eq. (21)) is the optimal policy  $\pi^*$ . This completes the proof that the proposed loss is a consistent loss for  $\pi^*$ .

**Lemma D.1.** Suppose l in Eq. (17) only depends on  $(\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2)$  through  $\pi$  and  $\pi_{ref}$ . Then, for any l that gives a consistent estimation of the optimal policy in Eq. (10), there exists an equivalent loss l such that

$$l(\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2, \boldsymbol{o}; \pi) = \tilde{l}(\boldsymbol{o}; h(\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2; \pi)),$$

where h is defined in Eq. (15).

Proof of Lemma D.1. Since l only depends on  $(\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2)$  through  $\pi$  and  $\pi_{\text{ref}}$ , we overload the notation and use  $l(\boldsymbol{o}; \pi(\boldsymbol{y}_1 \mid \boldsymbol{x}), \pi(\boldsymbol{y}_2 \mid \boldsymbol{x}))$  to denote the loss from  $(\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2, \boldsymbol{o})$ . In the limit of many data points,  $\mathcal{L}(\mathcal{D}; \pi)$  converges to

$$\mathcal{L}(\pi) = \mathbb{E}_{oldsymbol{x}, oldsymbol{y}_1, oldsymbol{y}_2, oldsymbol{o}} ig[ lig( oldsymbol{o}; \pi(oldsymbol{y}_1 \mid oldsymbol{x}), \pi(oldsymbol{y}_2 \mid oldsymbol{x}) ig) ig]$$
 .

Using  $Pr(o_u = 1 \mid \boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2) = \sigma(\Delta r^*(\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2; u))$ , the tower rule implies

$$\mathcal{L}(\pi) = \mathbb{E}_{\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2} \Big[ \mathbb{E}_{\boldsymbol{o}} \Big[ l \big( \boldsymbol{o}; \pi(\boldsymbol{y}_1 \mid \boldsymbol{x}), \pi(\boldsymbol{y}_2 \mid \boldsymbol{x}) \big) \mid \boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2 \Big] \Big]$$

$$= \mathbb{E}_{\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2} \Big[ \sum_{\boldsymbol{o} \in \{0,1\}^{\mathcal{U}}} l \big( \boldsymbol{o}; \pi(\boldsymbol{y}_1 \mid \boldsymbol{x}), \pi(\boldsymbol{y}_2 \mid \boldsymbol{x}) \big) \cdot \prod_{u \in \mathcal{U}} \sigma \big( \Delta r^*(\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2; u) \big)^{o_u} \big( 1 - \sigma \big( \Delta r^*(\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2; u) \big) \big)^{1 - o_u} \Big].$$

For notational simplicity, let's define

$$z_u(\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2) := \sigma(\Delta r^*(\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2; u)),$$
$$\chi_{\boldsymbol{o}}(\boldsymbol{z}) := \prod_{u \in \mathcal{U}} z_u^{o_u} (1 - z_u)^{1 - o_u}.$$

Note that  $\chi_o$  implicitly depends on  $(x, y_1, y_2)$  through z, which we drop from the notation when it is clear from the context. Using this notation, we can rewrite the  $\mathcal{L}(\pi)$ 's expansion as follows:

$$\mathcal{L}(\pi) = \mathbb{E}_{\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2} \left[ \sum_{\boldsymbol{o} \in \{0, 1\}^{\mathcal{U}}} l(\boldsymbol{o}; \pi(\boldsymbol{y}_1 \mid \boldsymbol{x}), \pi(\boldsymbol{y}_2 \mid \boldsymbol{x})) \cdot \chi_{\boldsymbol{o}} (\boldsymbol{z}(\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2)) \right]. \tag{22}$$

Overloading notation, we can always equivalently represent  $(\pi(y_1 \mid x), \pi(y_2 \mid x))$  as  $(\pi(y_1, y_2 \mid x), \pi(y_2 \mid y_1, y_2, x))$ , that is, with the probability that either of the two responses is chosen and the probability that the second one is preferred. Therefore, there exists a loss l' such that

$$l(\boldsymbol{o}; \pi(\boldsymbol{y}_1 \mid \boldsymbol{x}), \pi(\boldsymbol{y}_2 \mid \boldsymbol{x})) = l'(\boldsymbol{o}; \pi(\{\boldsymbol{y}_1, \boldsymbol{y}_2\} \mid \boldsymbol{x}), \pi(\boldsymbol{y}_2 \mid \{\boldsymbol{y}_1, \boldsymbol{y}_2\}, \boldsymbol{x}))$$
.

If  $\pi$  is optimal,  $\pi(y_2 \mid y_1, y_2, x)$  should also be optimal. Since  $\pi(y_2 \mid y_1, y_2, x)$  appears in only one term of the expectation in Eq. (22), we can conclude that

$$\underset{\boldsymbol{\theta}'}{\arg\min} \sum_{\boldsymbol{o} \in \{0,1\}^{\mathcal{U}}} l'\big(\boldsymbol{o}; \pi^*(\{\boldsymbol{y}_1, \boldsymbol{y}_2\} \mid \boldsymbol{x}), \boldsymbol{\theta}'\big) \cdot \chi_{\boldsymbol{o}}\big(\boldsymbol{z}(\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2)\big)$$

is the optimal  $\pi(y_2 \mid \{y_1, y_2\}, x)$  for every optimal  $\pi(\{y_1, y_2\} \mid x)$ . On the other hand, a property of the optimal policy  $\pi^*$  is that

$$\pi^*(\boldsymbol{y}_2 \mid \{\boldsymbol{y}_1, \boldsymbol{y}_2\}, \boldsymbol{x}) = \frac{\pi^*(\boldsymbol{y}_2 \mid \boldsymbol{x})}{\pi^*(\boldsymbol{y}_1 \mid \boldsymbol{x})} = \frac{\pi_{\mathrm{ref}}(\boldsymbol{y}_2 \mid \boldsymbol{x})}{\pi_{\mathrm{ref}}(\boldsymbol{y}_1 \mid \boldsymbol{x})} \cdot \exp\left(\frac{1}{\beta} \mathbb{E}_u \left[\Delta r^*(\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2; u)\right]\right).$$

Therefore, for every optimal policy  $\pi^*(\{y_1, y_2\} \mid x)$ , we have

$$\frac{\pi_{\text{ref}}(\boldsymbol{y}_{2} \mid \boldsymbol{x})}{\pi_{\text{ref}}(\boldsymbol{y}_{1} \mid \boldsymbol{x})} \cdot \exp\left(\frac{1}{\beta} \mathbb{E}_{u} \left[\Delta r^{*}(\boldsymbol{x}, \boldsymbol{y}_{1}, \boldsymbol{y}_{2}; u)\right]\right) = 
\underset{\boldsymbol{\theta'}}{\arg \min} \sum_{\boldsymbol{o} \in \{0,1\}^{\mathcal{U}}} l'(\boldsymbol{o}; \pi^{*}(\{\boldsymbol{y}_{1}, \boldsymbol{y}_{2}\} \mid \boldsymbol{x}), \boldsymbol{\theta'}) \cdot \chi_{\boldsymbol{o}}(\boldsymbol{z}(\boldsymbol{x}, \boldsymbol{y}_{1}, \boldsymbol{y}_{2})).$$
(23)

Recall from the optimal policy (Eq. (10)) that we can modify the reward function for responses other than  $(x, y_1, y_2)$  while keeping  $\Delta r^*(x, y_1, y_2; u)$  constant. This allows arbitrary changes to  $\pi^*(\{y_1, y_2\} \mid x)$  without altering the rest of Eq. (23). So, we can argue that l' does not depend on  $\pi(\{y_1, y_2\} \mid x)$  and we drop it from l' notation. Define a new loss based on l':

$$ilde{l}(oldsymbol{o}; heta)\coloneqq l'\Big(oldsymbol{o};rac{\pi_{ ext{ref}}(oldsymbol{y}_2\midoldsymbol{x})}{\pi_{ ext{ref}}(oldsymbol{y}_1\midoldsymbol{x})}\cdot\expig(rac{1}{eta} hetaig)\Big)\,.$$

Note that  $\tilde{l}$  implicitly depends on  $(\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2)$  through  $\pi_{\text{ref}}$  which we dropped from notation. Using  $\tilde{l}$ , we can write the original loss l as

$$l(\boldsymbol{o}; \pi(\boldsymbol{y}_1 \mid \boldsymbol{x}), \pi(\boldsymbol{y}_2 \mid \boldsymbol{x})) = l'(\boldsymbol{o}; \pi(\boldsymbol{y}_2 \mid \{\boldsymbol{y}_1, \boldsymbol{y}_2\}, \boldsymbol{x}))$$

$$= l'(\boldsymbol{o}; \frac{\pi(\boldsymbol{y}_2 \mid \boldsymbol{x})}{\pi(\boldsymbol{y}_1 \mid \boldsymbol{x})})$$

$$= \tilde{l}(\boldsymbol{o}; \beta \log \frac{\pi(\boldsymbol{y}_2 \mid \boldsymbol{x})}{\pi(\boldsymbol{y}_1 \mid \boldsymbol{x})} - \beta \log \frac{\pi_{\text{ref}}(\boldsymbol{y}_2 \mid \boldsymbol{x})}{\pi_{\text{ref}}(\boldsymbol{y}_1 \mid \boldsymbol{x})})$$

$$= \tilde{l}(\boldsymbol{o}; h(\boldsymbol{x}, \boldsymbol{y}_1, \boldsymbol{y}_2; \pi)).$$

This completes the proof.

**Lemma D.2.** Any loss  $\tilde{l}$  as in Lemma D.1 that leads to a consistent estimation of the optimal policy meets

$$\sum_{\boldsymbol{o} \in \{0,1\}^{\mathcal{U}}} \frac{\partial \tilde{l}}{\partial \theta} (\boldsymbol{o}; \theta^*(\boldsymbol{z})) \cdot \chi_{\boldsymbol{o}}(\boldsymbol{z}) = 0,$$
$$\sum_{\boldsymbol{o} \in \{0,1\}^{\mathcal{U}}} \frac{\partial^2 \tilde{l}}{\partial \theta^2} (\boldsymbol{o}; \theta^*(\boldsymbol{z})) \cdot \chi_{\boldsymbol{o}}(\boldsymbol{z}) \ge 0,$$

for every  $z \in [0,1]^{\mathcal{U}}$ . Here, we define

$$\chi_{\boldsymbol{o}}(\boldsymbol{z}) \coloneqq \prod_{u \in \mathcal{U}} z_u^{o_u} (1 - z_u)^{1 - o_u}.$$

and

$$\theta^*(z) \coloneqq \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \sigma^{-1}(z_u).$$

*Proof of Lemma D.2.* We will refer to the proof of Lemma D.1 in this proof. Using  $\tilde{l}$  in place of l' in Eq. (23), since  $\exp(\cdot)$  is monotone increasing, we have

$$\mathbb{E}_{u}\left[\Delta r^{*}(\boldsymbol{x},\boldsymbol{y}_{1},\boldsymbol{y}_{2};u)\right] = \operatorname*{arg\,min}_{\theta} \sum_{\boldsymbol{o} \in \{0,1\}^{\mathcal{U}}} \tilde{l}(\boldsymbol{o};\theta) \cdot \chi_{\boldsymbol{o}}(\boldsymbol{z}).$$

On the other hand, using the fact that user types in  $\mathcal{U}$  are equiprobable, we can write

$$\mathbb{E}_{u}\left[\Delta r^{*}(\boldsymbol{x},\boldsymbol{y}_{1},\boldsymbol{y}_{2};u)\right] = \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \sigma^{-1}(z_{u}).$$

Putting these together, it is necessary to have

$$\frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \sigma^{-1}(z_u) = \arg\min_{\theta} \sum_{o \in \{0,1\}^{\mathcal{U}}} \tilde{l}(o; \theta) \cdot \chi_o(z)$$

for every  $z \in [0,1]^{\mathcal{U}}$ . The rest of the proof is straightforward.

**Lemma D.3.** If  $|\mathcal{U}| > 3$ , for any loss  $\tilde{l}$  that meets the first-order condition of Lemma D.2, we have  $\frac{\partial \tilde{l}}{\partial \theta}(\mathbf{o};\theta) = 0$  for every  $\mathbf{o} \notin \{0,1\}$ .

Proof of Lemma D.3. First of all, for the BT model, a direct calculation shows

$$\theta^*(\boldsymbol{z}) := \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \sigma^{-1}(z_u) = \frac{1}{|\mathcal{U}|} \log \prod_{u \in \mathcal{U}} \left(\frac{z_u}{1 - z_u}\right) = \frac{1}{|\mathcal{U}|} \log \left(\frac{\chi_1}{\chi_0}\right).$$

Since  $\theta^*(z)$  depends on z only through  $\frac{\chi_0}{\chi_1}$ , we denote  $\frac{\partial \tilde{l}}{\partial \theta}(o; \theta^*)$  by  $g(o; \frac{\chi_0}{\chi_1})$ . The proof has two steps: In the first step, we relate  $g(o; \frac{\chi_0}{\chi_1})$  to  $g(o^{\oplus u'}; \frac{\chi_0}{\chi_1})$ , where we define

$$oldsymbol{o}^{\oplus u'} \coloneqq egin{cases} o_u \,, & u 
eq u' \,, \ 1 - o_u \,, & u = u' \,. \end{cases}$$

Using this connection, in the second step, we will show that  $g(\mathbf{o}; \frac{\chi_{\mathbf{o}}}{\chi_{\mathbf{1}}}) = 0$  for any  $\mathbf{o} \in \{\mathbf{0}, \mathbf{1}\}$  when  $|\mathcal{U}| \geq 4$ .

**Step 1.** Consider any  $o \notin \{0,1\}$ . When  $|\mathcal{U}| \geq 3$ , there exists  $u' \in \mathcal{U}$  such that  $o^{\oplus u'} \notin \{0,1\}$ . For such o and u', we define two non-empty sets

$$S^{1} := \{ u \mid u \in \mathcal{U}, u \neq u', o_{u} = 1 \},$$
  
$$S^{0} := \{ u \mid u \in \mathcal{U}, u \neq u', o_{u} = 0 \}.$$

We set z such that

$$\prod_{u \in \mathcal{S}^1} z_u = \prod_{u \in \mathcal{S}^0} (1 - z_u),$$

$$z_u \to 1^-, \ \forall u \in \mathcal{S}^1,$$

$$z_u \to 0^+, \ \forall u \in \mathcal{S}^0.$$
(24)

For this choice of z, asymptotically, we have

$$\begin{split} \frac{\chi_{\mathbf{0}}}{\chi_{\mathbf{1}}} &= \frac{1 - z_{u'}}{z_{u'}} \,, \\ \chi_{\mathbf{o}} &= z_{u'}^{o_{u'}} (1 - z_{u'})^{1 - o_{u'}} \,, \\ \chi_{\mathbf{o}^{\oplus u'}} &= z_{u'}^{1 - o_{u'}} (1 - z_{u'})^{o_{u'}} \,, \\ \chi_{\mathbf{o}'} &= 0 \,, \, \forall \mathbf{o}' \notin \{\mathbf{o}, \mathbf{o}^{\oplus u'}\} \,. \end{split}$$

Using the above, we can simplify the first-order condition in Lemma D.2 as

$$g(\mathbf{o}; \frac{1-z_{u'}}{z_{u'}}) \cdot z_{u'}^{o_{u'}} (1-z_{u'})^{1-o_{u'}} + g(\mathbf{o}^{\oplus u'}; \frac{1-z_{u'}}{z_{u'}}) \cdot z_{u'}^{1-o_{u'}} (1-z_{u'})^{o_{u'}} = 0.$$

This condition should be held for every  $z_{u'} \in [0,1]$ . Therefore, we can conclude

$$g(\mathbf{o}^{\oplus u'}; \alpha) = -g(\mathbf{o}; \alpha) \cdot \alpha^{1 - 2o_{u'}}, \qquad (25)$$

for every  $\alpha \in \mathbb{R}$ . This completes the first part of the proof.

Step 2. When  $o \notin \{0,1\}$  and  $|\mathcal{U}| \geq 4$ , there exist distinct user types u' and u'' such that none of  $o^{\oplus u'}$ ,  $o^{\oplus u''}$ , and  $o^{\oplus (u',u'')}$  are in  $\{0,1\}$ . Here, we used  $o^{\oplus (u',u'')}$  as a shorthand for  $(o^{\oplus u'})^{\oplus u''}$ . For such o, u', and u'', we define two non-empty sets

$$S^{1} := \{ u \mid u \in \mathcal{U} \setminus \{u', u''\}, o_{u} = 1 \},$$
  
$$S^{0} := \{ u \mid u \in \mathcal{U} \setminus \{u', u''\}, o_{u} = 0 \}.$$

We set z according to Eq. (24). Then, asymptotically,

$$\chi_{\boldsymbol{o}'} = 0, \ \forall \boldsymbol{o}' \notin \{\boldsymbol{o}, \boldsymbol{o}^{\oplus u'}, \boldsymbol{o}^{\oplus u''}, \boldsymbol{o}^{\oplus (u', u'')}\}.$$

Using the above, we can simplify the first-order condition in Lemma D.2 as

$$g(\boldsymbol{o}; \frac{\chi_{\mathbf{0}}}{\chi_{\mathbf{1}}}) \cdot \chi_{\boldsymbol{o}} + g(\boldsymbol{o}^{\oplus u'}; \frac{\chi_{\mathbf{0}}}{\chi_{\mathbf{1}}}) \cdot \chi_{\boldsymbol{o}^{\oplus u'}} + g(\boldsymbol{o}^{\oplus u''}; \frac{\chi_{\mathbf{0}}}{\chi_{\mathbf{1}}}) \cdot \chi_{\boldsymbol{o}^{\oplus u''}} + g(\boldsymbol{o}^{\oplus (u', u'')}; \frac{\chi_{\mathbf{0}}}{\chi_{\mathbf{1}}}) \cdot \chi_{\boldsymbol{o}^{\oplus (u', u'')}} = 0. \quad (26)$$

Because of the symmetry of this equation, we can assume without loss of generality that  $o_{u'} = 0$  and  $o_{u''} = 1$ . Therefore, asymptotically, we have

$$\begin{split} \frac{\chi_{\mathbf{0}}}{\chi_{\mathbf{1}}} &= \big(\frac{1-z_{u'}}{z_{u'}}\big) \big(\frac{1-z_{u''}}{z_{u''}}\big)\,,\\ \chi_{\mathbf{o}} &= \big(1-z_{u'}\big) \cdot z_{u''}\,,\\ \chi_{\mathbf{o}^{\oplus u'}} &= z_{u'} \cdot z_{u''}\,,\\ \chi_{\mathbf{o}^{\oplus u''}} &= \big(1-z_{u'}\big) \cdot \big(1-z_{u''}\big)\,,\\ \chi_{\mathbf{o}^{\oplus (u',u'')}} &= z_{u'} \cdot \big(1-z_{u''}\big)\,. \end{split}$$

Eq. (25) also implies

$$g(\mathbf{o}^{\oplus u'}; \alpha) = -g(\mathbf{o}; \alpha) \cdot \alpha ,$$
  

$$g(\mathbf{o}^{\oplus u''}; \alpha) = -g(\mathbf{o}; \alpha) \cdot \alpha^{-1} ,$$
  

$$g(\mathbf{o}^{\oplus (u', u'')}; \alpha) = g(\mathbf{o}; \alpha) .$$

Plugging these into Eq. (26) and simplifying equations, we obtain

$$g\left(\mathbf{o}; \left(\frac{1-z_{u'}}{z_{u'}}\right) \left(\frac{1-z_{u''}}{z_{u''}}\right)\right) \cdot (2z_{u'}-1)(2z_{u''}-1) = 0.$$

This equation should be held for every  $z_{u'}$  and  $z_{u''}$ . By appropriately setting  $z_{u'}$  and  $z_{u''}$ , we can conclude that  $g(\mathbf{o}; \alpha)$  should be zero for every  $\alpha$ . This completes this proof.

# E PEW Surveys Experiments: Details and Additional Examples

In this section, we expand on our PEW surveys experiment where we used polling data on key political and social issues to show: (i) how NBC rankings can differ from those maximizing the average reward; (ii) How sensitive NBC is to the sampling distribution of the pairwise preference data.

**Data.** We use several Pew Research Center surveys, specifically the American Trends Panel surveys number 35, 52, 79, 83, 99, 109, 111, 112, 114, 119, 120, 121, 126, 127, 128, 129, 130, 131, and 132. The choices are a mix of recent surveys and those relevant to science, technology, data and AI. Each survey include questions asked to thousands of participants. We categorize participants by political party leanings to define types. When processing the questions, we discard responses that are empty, as well as discarding the option "Refused". We note that discarding the option "Refused" had no effect on the results as it is not frequently chosen.

**Reward Estimation.** Although we observe how often each group selects a particular option, we don't directly observe respondents' internal rewards. To estimate this, we apply the Luce-Shepherd model [24, 20]:

$$\Pr\left(\text{option } i \text{ is chosen from } \mathcal{S}\right) = \frac{\exp\left(r(i;u)\right)}{\sum_{j\in\mathcal{S}}\exp\left(r(j;u)\right)},\tag{27}$$

where S is the set of options, and  $r(\cdot;u)$  is the reward for type u. This allows us to estimate each option's reward (up to a constant additive term) for each type. From these estimates and observed probabilities, we compute both the expected reward and the NBC metric, where in the latter we assume the uniform probability for alternatives unless specified otherwise.

**Sensitivity Experiments.** In Section 7.1, we estimate the sensitivity of NBC rankings to the sampling distribution of pairwise preference data by determining the minimum Total Variation (TV) distance required, if possible, to alter NBC rankings from that attained under the uniform distribution. Recall that NBC is defined as:

$$\mathrm{NBC}(\boldsymbol{y}; \mathcal{D}) \coloneqq \mathbb{E}_{\boldsymbol{y}' \sim \mathcal{D}(\cdot)} \big[ \Pr(\boldsymbol{y} \succ \boldsymbol{y}' \mid \boldsymbol{x}; r) \big].$$

To compute this, we first estimate the reward function r, then evaluate  $\Pr(\mathbf{y} \succ \mathbf{y}'; r)$  for all  $(\mathbf{y}, \mathbf{y}') \in \mathcal{Y} \times \mathcal{Y}$ , where  $\mathcal{Y}$  is the set of alternatives. Next, consider the feasibility of swapping the ranking induced by the uniform distribution  $\mathcal{D}_U$  for alternatives  $\mathbf{y}_i$  and  $\mathbf{y}_j$  with a new distribution  $\mathcal{D}_a$ , assuming  $\operatorname{NBC}(\mathbf{y}_i; \mathcal{D}_U) > \operatorname{NBC}(\mathbf{y}_i; \mathcal{D}_U)$ . This is equivalent to solving the following linear program:

$$\begin{aligned} & \text{minimize} & & \frac{1}{2} \mathbf{1}^{\top} \mathbf{s}, \\ & \text{subject to:} & & \mathbf{q} > \epsilon \mathbf{1}, \\ & & & \mathbf{s} \geq \frac{1}{N} \mathbf{1} - \mathbf{q}, \\ & & & & & \mathbf{s} \geq \mathbf{q} - \frac{1}{N} \mathbf{1}, \\ & & & & & & \mathbf{P}_i \mathbf{q} < \mathbf{P}_j \mathbf{q} + \delta, \\ & & & & & & & \mathbf{1}^{\top} \mathbf{q} = 1, \\ & & & & & & & & \mathbf{q} > \mathbf{0}. \end{aligned}$$

Here,  $N = |\mathcal{Y}|$ , and labeling the alternatives as  $\mathbf{y}_1, \dots, \mathbf{y}_N$ , we define  $\mathbf{P}_{ij} = \Pr(\mathbf{y}_i \succ \mathbf{y}_j; r)$ ,  $q_i = \mathcal{D}_a(\mathbf{y}_i)$ , and  $\mathbf{P}_k$  as the k-th row of  $\mathbf{P}$ . The parameter  $\epsilon > 0$  ensures support for all alternatives, and

 $\delta > 0$  controls the required magnitude of change in NBC beyond what is required for the swap. We set  $\epsilon = \delta = 10^{-5}$ .

To compute the minimum TV distance, we solve the program for all pairs (i, j) where  $NBC(y_i; \mathcal{D}_U) > NBC(y_j; \mathcal{D}_U)$  and record the smallest objective value. In this analysis, we group respondents by political leaning (specifically, the column F\_PARTYSUM\_FINAL). We also note that in this analysis we exclude survey questions with fewer than three options.

Additional Examples. We highlight some of the examples of discrepancies that we find in some of the surveys listed above in Figures 6, 7, 8, 9, and 10. We note though that while we indeed find a few examples of the discrepancy, NBC rankings is actually aligned with average reward rankings in most cases. One possible explanation for this is that in many questions, the distributions of responses across types were very similar, suggesting that a homogeneous reward model would have been appropriate, causing NBC and average reward to align.

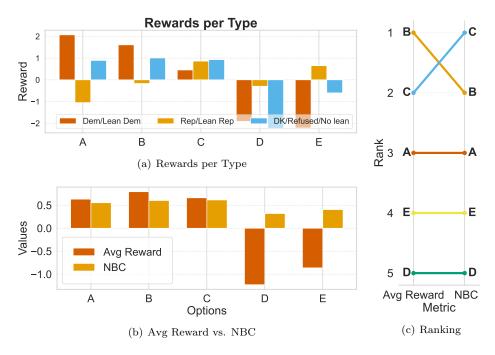


Figure 6: Do you think the plans and policies of the Biden administration will make the country's response to the coronavirus outbreak: A: A lot better; B: A little better; C: Not much different; D: A little worse; E: A lot worse

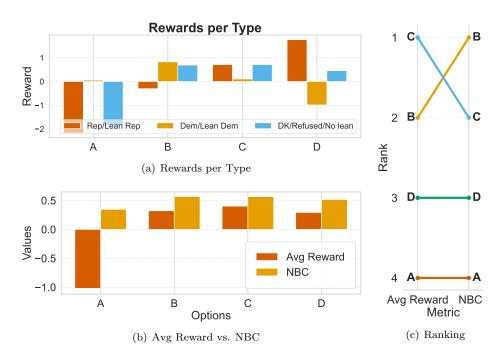


Figure 7: How would you rate the job Joe Biden is doing responding to the coronavirus outbreak? A: Excellent; B: Good; C: Only fair; D: Poor

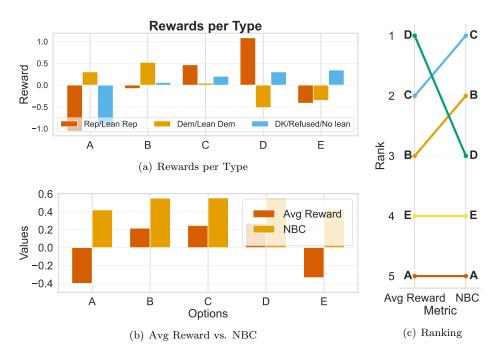


Figure 8: The next time you purchase a vehicle, how likely are you to seriously consider purchasing an electric vehicle? A: Very likely; B: Somewhat likely; C: Not too likely; D: Not at all likely; E: I do not expect to purchase a vehicle

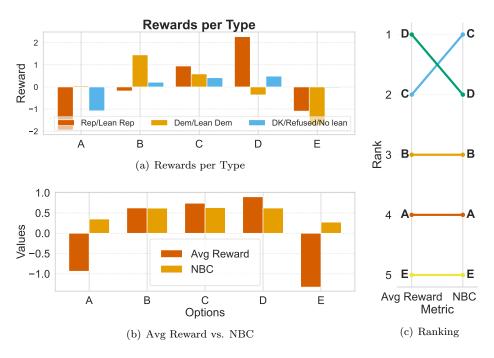


Figure 9: What is your overall opinion of Kamala Harris? A: Very favorable; B: Mostly favorable; C: Mostly unfavorable; D: Very unfavorable; E: Never heard of this person.

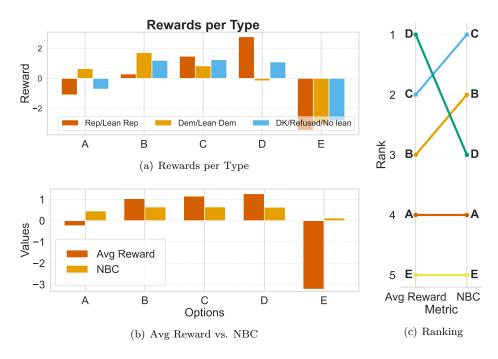


Figure 10: What is your overall opinion of Joe Biden? A: Very favorable; B: Mostly favorable; C: Mostly unfavorable; D: Very unfavorable; E: Never heard of this person.

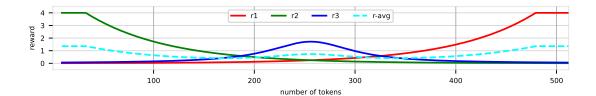


Figure 11: Reward definition for three user types in semi-synthetic experiments (Section 7.3) based on the length of prompt response combination. The first user type prefers long prompt response combinations, the second user type prefers short prompt response combinations, and the third user type prefers mid-length prompt response combinations. The dashed cyan line shows the average reward across the three user types.

## F Semi-Synthetic Experiment: Fine-Tuning Llama-3-8B on HH-RLHF

**Reward Models.** Fig. 11 shows the three distinct rewards we use for the three user types along with their average. In order to have a reliable ground-truth reward which we can rely on in evaluation, we define these rewards as functions of the number of tokens in prompt-response combinations.

Anonymous Dataset. We use prompts and response pairs from both helpfulness and harmlessness subsets of Anthropic's HH-RLHF dataset [38] and relabel the *chosen* and *rejected* responses manually. We filter for data points in which the sum of the number of tokens in the prompt and the number of tokens in the longer response do not exceed 512. This leaves us with 160, 800 training and 17, 104 test data points. For every data point (a prompt with a pair of responses), we sample one of the three user types uniformly at random. Given the type of user, we sample a preference based on BT [22] to label the two alternatives.

Dataset with Maximum Annotator Information. We use prompts and response pairs from both helpfulness and harmlessness subsets of Anthropic's HH-RLHF dataset [38] and relabel the *chosen* and *rejected* responses manually. We filter for data points in which the sum of the number of tokens in the prompt and the number of tokens in the longer response do not exceed 512. This leaves us with 160, 800 training and 17, 104 test data points. For every data point (a prompt with a pair of responses), we keep sampling BT [22] preferences from all user types until they agree with each other. Once the consensus is achieved, we stop sampling and use the agreed-upon preference as the label for this data point.

**Fine-Tuning Details.** We fine-tune Llama-3-8B [37] base model with LoRA [36]. We fine-tune for one epoch with a batch size of 2, and use a linear learning rate schedule that starts with  $3 \times 10^{-5}$  and decreases to zero. We use the Adam optimizer with a weight decay of 0.001 [74]. Regarding LoRA's hyper-parameters, we use the matrix rank of r = 8,  $\alpha = 32$ , and the dropout probability of 0.1.

For direct alignment experiments, we use a uniform reference policy. When ignoring heterogeneity, we do vanilla DPO over the anonymous dataset. When modeling heterogeneity, we use the loss function we propose in Proposition 6.1 over the dataset with maximum annotator information. We use the ordinal agreement between the ground-truth average reward and the reward induced by the aligned policy as the measure of accuracy.

For the reward learning experiments, we fine-tune the Llama-3-8B as a reward model. When ignoring heterogeneity, we assume BT and maximize the probability of the anonymous preference dataset under the learned reward model. When modeling heterogeneity, we use the loss function in

Table 1: Raw Accuracy (%) in Alignment Experiments

SEED	IGNORING HOMOGENEITY	Modeling Heterogeneity
0	65.61	66.63
1	67.55	69.55
2	68.77	75.21
3	68.70	72.04
4	66.28	74.95

Table 2: Raw Accuracy (%) in Reward Learning Experiments

SEED	IGNORING HOMOGENEITY	Modeling Heterogeneity
0	92.33	95.26
1	85.38	92.01
2	88.46	94.6
3	89.69	93.57
4	91.94	93.87

Proposition 6.1 over the dataset with maximum annotator information, but replace  $h(y_1, y_2; \pi)$  with the difference in rewards, i.e.,  $r(y_2) - r(y_1)$ . We use the ordinal agreement between the ground-truth average reward and the learned reward as the measure of accuracy.

**Detailed Results.** We conduct every experiment with five different random seeds. Fig. 5 shows the average,  $25^{\rm th}$  percentile, and  $75^{\rm th}$  percentile of accuracy across the five random seeds. Tables 1 and 2 show the raw accuracy numbers across the five random seeds.