

# Fairness-Aware Reward Optimization

Ching Lam Choi<sup>1</sup>, Vighnesh Subramaniam<sup>1</sup>, Antonio Torralba<sup>1</sup>, Phillip Isola<sup>1</sup>, and Stefanie Jegelka<sup>1,2</sup>

<sup>1</sup>CSAIL, Department of EECS, Massachusetts Institute of Technology

<sup>2</sup>School of CIT, MCML, MDSI, Technical University of Munich

{chinglam, vsub851, torralba, phillipi, stefje}@mit.edu,  
stefanie.jegelka@tum.de

## Abstract

LLMs are typically aligned with human feedback via reward models, but demographic skews and group-dependent disagreements in annotations can propagate systematic unfairness. We introduce Fairness-Aware Reward Optimisation (FARO), a principled framework for training reward models under demographic parity, equalised odds, or counterfactual fairness constraints. Our approach instantiates a proxy-Lagrangian descent-ascent game (ProxyGDA) that yields reward models with provable fairness certificates up to vanishing slack. We provide the first theoretical analysis of reward-level fairness in alignment, establishing: (i) guarantees that FARO-trained rewards satisfy DP/EO/CF; (ii) a formal accuracy-fairness trade-off induced by KL-regularised RL fine-tuning; and (iii) existence of Pareto-optimal solutions along this trade-off. Across multiple LLMs on the representative BBQ dataset, FARO consistently reduces demographic bias and harmful generations while preserving or improving LLM quality and factuality.

## 1 Introduction

Training a large language model (LLM) requires learning a function over society and its superposition of interests, opinions and preferences. Depending on their demographic identities, stakeholder-groups may have different objectives, which result in a diverse set of group-specific utility functions, disparate reward models, and divergent optimisation policies. Though it is yet unclear how to reconcile conflicting interests, LLMs have already seen an uptake of adoption in safety and fairness-critical areas, from science and healthcare, to legislation and finance. At best, LLM-augmented operations could lead to impartial standards and streamlined development; at worst, the reinforcement of human prejudice and regression to a less fair, more prejudiced common denominator (Weidinger et al., 2021; Bender et al., 2021; Dai et al., 2024).

Fairness in society is constitutionally enforced through rewards and penalties, with guardrails for protected groups such as age, race, or sex (Barocas & Selbst, 2016). Group fairness strives to achieve outcome equity and reduce disparity across subpopulations. Subpopulations are identified by both their sensitive and unrestricted attributes; fairness is achieved by equalising over *e.g.* outcomes, odds, or opportunity. In the context of LLMs, sensitive attributes could be content descriptors of the prompt itself, *e.g.*  $x = \text{"Who is better at maths, Alice, Bob, or unknown?"}$  and  $S$  describes *sex*. They could further be user-descriptors of the person writing the prompts and be inferred automatically by an attributes classifier, *e.g.* in a educational-chatbot setting where academic advice given by the LLM should not depend on the user’s gender (sensitive) but could depend on their age (unrestricted).

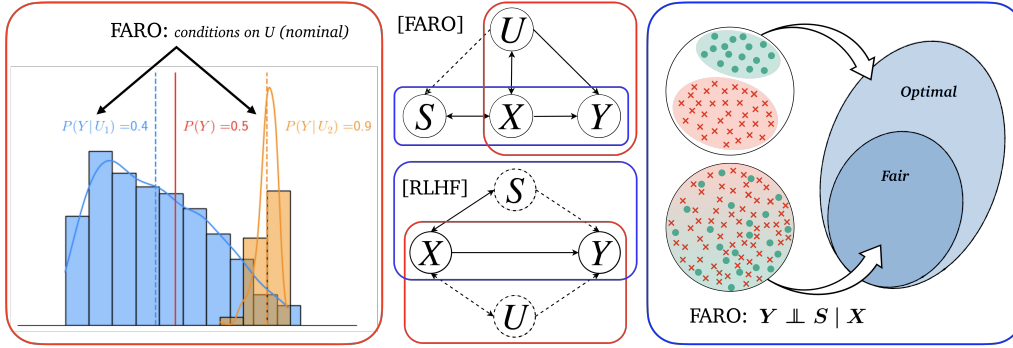


Figure 1: FARO learns ordinal, cardinal and fair human preferences  $Y \mid X$  by explicitly optimising fairness constraints (*upper centre*). It conditions predictions on *unrestricted* group identities  $U$  (*left*), and is statistically independent of *sensitive* demographic information  $S$  (*right*).

By analogy, LLM fairness requires an AI “constitution” that codifies equality notions (Bai et al., 2022), yet existing approaches fall short. Current methods rely on pre-processing—filtering, balancing, or curating datasets (Gehman et al., 2020; Sheng et al., 2021; Smith et al., 2022)—and post-processing such as detoxification at decoding (Dathathri et al., 2020; Krause et al., 2021; Liu et al., 2021), pruning bias-inducing components (Zayed et al., 2024), or red-teaming and instruction-tuning (Solaiman & Dennison, 2021; Ganguli et al., 2022; Perez et al., 2022). These reduce overt harms but remain limited: pre-processing is expensive and lacks guarantees, since fairness in data statistics need not transfer to learned models; post-processing ensures Pareto-optimality only within the restricted family of group-thresholded variants of a fixed predictor, leaving models strictly inside the global accuracy–fairness frontier.

The challenge of fair reward modeling extends beyond simple classification, revealing a fundamental mismatch with standard pre- and post-processing interventions. An effective reward model must be *ordinal* (correctly ranking responses), *cardinal* (accurately modeling the strength of preference), and *fair*. Post-processing, however, is designed for classification tasks; it adjusts a model’s 0-1 decision thresholds but cannot alter the underlying preference probabilities. Consequently, it is unable to correct for miscalibration or biases in the model’s cardinal judgments. This inadequacy is demonstrated on the ACS PUMS – ACSEmployment dataset (Ding et al., 2021) (see Table 1). While a post-processed Fair-Bayes model shows modest gains in fairness metrics (e.g.  $\Delta dp$ ), it fails to improve the model’s poor cardinal performance (e.g. ECE, a measure of miscalibration). Algorithmic fairness literature (Barocas et al., 2023; Suresh & Gutttag, 2021) (thoroughly discussed in App. 6) offers a powerful alternative: *in-processing*. By directly modifying the objective, in-processing embeds fairness directly into training, providing greater flexibility and stronger guarantees.

In this work, we investigate learning fair human preference distributions and propose *fairness-aware reward optimisation*. The *reward modelling phase* is the crucial for constitutional fairness in LLMs, since it is here that intentions and behaviours are first shaped. Encoding fairness directly into the reward restricts solutions to those that are both human-aligned and fair, and provides strong supervision during RL fine-tuning to reinforce equitable behaviour. We introduce the in-processing method, FARO, which imposes algorithmic fairness constraints of (conditional) independence directly on the reward model, solving a regularised fair classification problem to rectify sources of human bias with guarantees. Our contributions are as follows:

1. *Framework for fair reward modeling.* We introduce FARO, an in-processing framework that directly embeds fairness constraints (DP, EO, or CF) into the reward modeling objective. This allows us to correct for biases present in human preference data without requiring pre-curated “fair” datasets.

Table 1: Performance on the ACSEmployment dataset. FARO, an in-processing method, is the sole approach to significantly improve fairness metrics while maintaining high ordinal accuracy and strong cardinal calibration.

Method \ Metric	Ordinal $\uparrow$		Cardinal $\downarrow$			Fair $\downarrow$		
	0-1 Acc.	F1 Score	ECE	MCE	RMSCE	$\Delta_{dp}$	$\Delta_{eo}$	$\Delta_{cf}$
Bayes	.877 $\pm$ .019	.518 $\pm$ .030	.115 $\pm$ .007	.484 $\pm$ .064	.165 $\pm$ .008	.037 $\pm$ .026	.112 $\pm$ .132	.062 $\pm$ .021
Fair-Bayes	.879 $\pm$ .014	.500 $\pm$ .052	.109 $\pm$ .006	.447 $\pm$ .006	.157 $\pm$ .007	.026 $\pm$ .021	.109 $\pm$ .105	.063 $\pm$ .027
FARO- <i>dp</i>	<b>.889</b> $\pm$ .013	<b>.537</b> $\pm$ .038	<b>.105</b> $\pm$ .004	<b>.440</b> $\pm$ .004	<b>.154</b> $\pm$ .004	<b>.007</b> $\pm$ .005	.111 $\pm$ .071	.047 $\pm$ .021
FARO- <i>eo</i>	.884 $\pm$ .019	.525 $\pm$ .076	.114 $\pm$ .005	.443 $\pm$ .006	.160 $\pm$ .004	.012 $\pm$ .010	<b>.073</b> $\pm$ .037	.067 $\pm$ .026
FARO- <i>cf</i>	.884 $\pm$ .016	.505 $\pm$ .046	.105 $\pm$ .009	.441 $\pm$ .009	.156 $\pm$ .011	.018 $\pm$ .010	.105 $\pm$ .066	<b>.042</b> $\pm$ .008

2. *Refined problem formulation.* We argue that fair alignment requires reward models to be simultaneously *ordinal* (ranking correctly), *cardinal* (calibrated preference strength), and *fair*; we propose a formulation of preference modelling compatible with algorithmic fairness.
3. *Theoretical guarantees.* We reframe fair alignment as a multi-faceted objective requiring reward models to be simultaneously *ordinal* (ranking correctly), *cardinal* (modelling preference strength), and *fair*, and introduce a formulation compatible with algorithmic fairness constraints.
4. *Empirical validation.* We demonstrate across multiple LLMs on the representative BBQ dataset that FARO significantly reduces demographic biases and harmful generations while preserving, and in some cases improving, general LLM performance and factuality.

## 2 Preliminaries

There are two dominant paradigms for aligning LLMs to human preferences—explicit, RL-based approaches like RLHF (Ziegler et al., 2019) and variants (Bai et al., 2022; Ouyang et al., 2022; Stiennon et al., 2020), and implicit methods (without a parametric reward model) such as DPO (Rafailov et al., 2023) and others (Ethayarajh et al., 2024; Azar et al., 2024; Xu et al., 2024). We first recap key notation of RLHF (and DPO in App. A.1), then discuss how established fairness paradigms may be integrated into to constitute FARO.

### 2.1 Reward Modelling and Policy Optimisation

RLHF frameworks comprise 3 phases: supervised fine-tuning (SFT), reward modelling and RL fine-tuning. From phase one, an SFT-trained LLM is obtained. Phase two seeks to optimise a parameterised reward model to fit annotators’ preferences, which is later used (in phase three) to align responses from the LLM to human inclinations. We review the latter two.

In response to an input prompt  $x \sim \mathcal{X}$ , two LLM responses are recorded,  $(\hat{y}_1, \hat{y}_2) \sim \pi^{SFT}(\hat{y} | x)$ , where  $\hat{y}_w \succ \hat{y}_l | x$  denotes the response preferred by human annotators. RLHF assumes that preferences are generated by some underlying reward model  $r^*(x, y)$  and represents the distribution of human preferences  $\mathcal{P}^*$  with the Bradley-Terry (BT) model (Bradley & Terry, 1952), where  $\sigma$  is the logistic function:  $p^*(\hat{y}_w \succ \hat{y}_l | x) = \sigma(r^*(x, \hat{y}_w) - r^*(x, \hat{y}_l))$ . Given data samples  $(x, \hat{y}_w, \hat{y}_l) \sim \mathcal{D}$ , we solve a binary classification problem to fit a reward model  $r_\phi \sim \mathcal{R}$  to  $\mathcal{P}^*$ , and optimise the negative log-likelihood loss  $L$ , with  $r_\phi$  normalised and centred at zero-expectation:

$$L_{\text{NLL}}(r_\phi; \mathfrak{J}) = -\mathbb{E}_{(x, \hat{y}_w, \hat{y}_l) \sim \mathcal{D}} [\log \sigma(r_\phi(x, \hat{y}_w) - r_\phi(x, \hat{y}_l))]. \quad (1)$$

The fitted reward model  $r_{\hat{\phi}}$  is used to supervise and align the LLM  $\pi_{\theta}$  to human preferences, without deviating too far from reference point  $\pi_{\text{ref}}$ . Both  $\pi_{\theta}$  and  $\pi_{\text{ref}}$  are initialised to the SFT-trained model  $\pi^{\text{SFT}}$ ; we tune  $\pi_{\theta}$  to maximise the following reward (Jaques et al., 2017):

$$\mathbb{E}_{x \sim \mathcal{D}, \hat{y} \sim \pi_{\theta}(\hat{y} | x)} \left[ r_{\hat{\phi}}(x, \hat{y}) \right] - \beta D_{\text{KL}} [\pi_{\theta}(\hat{y} | x) \parallel \pi_{\text{ref}}(\hat{y} | x)]. \quad (2)$$

Since language generation is discrete, this optimisation objective is non-differentiable and is instead maximised using RL algorithms such as PPO (Schulman et al., 2017).

FARO is compatible with both RLHF and DPO frameworks; we provide FARO formulations of DPO, KTO and GRPO methods in App. A.1. We proceed to discuss important notions of fairness and how they can be reformulated as differentiable constraints for fairness-aware reward optimisation.

## 2.2 Fairness Paradigms

Departing from previous approaches, we impose fairness constraints during the *reward modelling phase*. This guarantees that our reward model is algorithmically fair to provide fair feedback during RL fine-tuning. Let  $\mathcal{J}$  denote a joint distribution over the domain  $\mathcal{D}$  of data, where each data sample has the structure  $(x, \hat{y}_w, \hat{y}_l, S, U)$ .  $S \in [p]$  represents a categorical, *sensitive attribute* unfair to use during inference;  $U \in [d]$  represents a categorical, *unrestricted attribute* permissible to use. We generalise this formulation to admit multiple sensitive and unrestricted attributes in App. A.2. We define two indicator variables for the preference outcome. The ground-truth human preference is  $Y = 1$  for the pair  $(\hat{y}_w, \hat{y}_l)$  where humans preferred  $\hat{y}_w$  over  $\hat{y}_l$ . The model’s predicted preference is given by  $\hat{Y} = \mathbb{1}\{r_{\phi}(x, \hat{y}_w) > r_{\phi}(x, \hat{y}_l)\}$ . We consider three fair binary classification paradigms depending on the accessibility and modality of attributes:

- (1) *Attribute Blind*.  $S, U$  are not used for reward assignment. We aim to learn a reward model  $r_{\phi} : \mathcal{X} \times \hat{\mathcal{Y}} \rightarrow \mathbb{R}$  that minimises  $L(r_{\phi}; \mathcal{J})$  and takes  $(x, \hat{y})$  as input.
- (2) *Attribute Aware*.  $S, U$  are accessible and are appended to the input prompt  $x$ . They are either inferred by an off-the-shelf attributes classifier, or are provided as annotations in the dataset. We aim to learn a reward model  $r_{\phi} : \mathcal{X} \times \hat{\mathcal{Y}} \times [p] \times [d] \rightarrow \mathbb{R}$  that minimises  $L(r_{\phi}; \mathcal{J})$  and takes  $(x, \hat{y}, S, U)$  as input.
- (3) *Self-critiquing LLMs*.  $S, U$  are not provided and must be inferred from input prompt  $x$  using an off-the-shelf language model. The natural language descriptions  $\hat{S}, \hat{U} \in \mathcal{X}$  of sensitive and unrestricted information (associated with  $x$ ) are appended feature-wise to  $x$ . We aim to learn a reward model  $r_{\phi} : \mathcal{X} \cdot \hat{S} \cdot \hat{U} \times \hat{\mathcal{Y}} \rightarrow \mathbb{R}$  that minimises  $L(r_{\phi}; \mathcal{J})$  and takes  $(x \cdot \hat{s} \cdot \hat{u}, \hat{y})$  as input.

To enforce equity and decrease disparity of inter-group outcomes, we use metrics to characterise the quality of outcomes, such as true positive rate (TPR) (Menon & Williamson, 2018; Agarwal et al., 2018), false positive rate (FPR) (Hardt et al., 2016), and predictive rate (Quadrianto & Sharmanska, 2017). Towards a general framework for fairness statistics, we let values taken by the sensitive attribute  $S$  partition the domain  $\mathcal{D}$  into  $p$  groups  $G_i := \{(x, \hat{y}_w, \hat{y}_l, i) \in \mathcal{D}\}$ . Then, we follow Celis et al. (2019) to measure  $G_i$ ’s group performance via  $q_i^{\mathcal{J}}(r_{\phi}) = \mathbb{P}_{\mathcal{J}}[\mathcal{E} | G_i, \mathcal{E}']$  for some events  $\mathcal{E}, \mathcal{E}'$ , e.g. conditioning on unrestricted attributes. Lastly, define the group performance function  $q^{\mathcal{J}} = (q_1^{\mathcal{J}}(r_{\phi}), \dots, q_p^{\mathcal{J}}(r_{\phi}))$ . We omit  $\mathcal{J}$  when it is contextually clear.

Intuitively, a reward model  $r_{\phi}$  is considered fair with respect to  $q$  if  $q_i(r_{\phi}) \approx q_{i'}(r_{\phi})$  for all  $i, i'$ . This indicates that the performance (e.g. TPR or FPR) of  $r_{\phi}$  is approximately equal across all subpopulations; the reward model does not overfit the most populous demographic, nor is its performance dependent on specific

identities. To measure the fairness of  $r_\phi$  for a given group performance function  $q$ , we follow previous works (Feldman et al., 2015; Menon & Williamson, 2018; Zafar et al., 2017a;b) and consider the  $\tau$ -rule.

**Definition 2.1.** ( $\tau$ -Fair) A reward model  $r_\phi$  achieves  $\tau$ -fairness w.r.t.  $q$  if it satisfies for  $\tau \in [0, 1]$ ,

$$\min_{r_\phi \in \mathcal{R}} L(r_\phi; \mathcal{J}) \quad \text{s.t.} \quad \max_{i, i' \in [p]} |q_i(r_\phi) - q_{i'}(r_\phi)| \leq \tau \quad (3)$$

The closer  $\tau$  is to 0, the fairer  $r_\phi$  is w.r.t  $q$ ; perfect fairness is achieved at  $\tau = 0$ . Practically, we consider  $\tau > 0$  due to known infeasibility, incompatibility and inconsistency issues under perfect fairness (Friedler et al., 2021; Hardt et al., 2016; Kleinberg et al., 2017).

### 2.3 Fairness Notions

We proceed to quantify the fairness violations of reward models, by establishing three notions of group fairness—*demographic parity* (DP), requiring  $\hat{Y} \perp\!\!\!\perp S$ ; *equalised odds* (EO) (Hardt et al., 2016), requiring  $\hat{Y} \perp\!\!\!\perp S \mid Y$ ; *conditional fairness* (CF) (Xu et al., 2020), requiring  $\hat{Y} \perp\!\!\!\perp S \mid U$ .

**Definition 2.2** (Demographic Parity (DP)). A reward model  $r_\phi$  is  $\gamma$ -DP fair if the group-wise positive rates are nearly equal:  $q_i^{\text{dp}}(r_\phi) := \mathbb{P}[\hat{Y} = 1 \mid G_i]$ ,  $\Delta_{\text{dp}}(r_\phi) := \max_{i, i' \in [p]} |q_i^{\text{dp}}(r_\phi) - q_{i'}^{\text{dp}}(r_\phi)| \leq \gamma$ .

**Definition 2.3** (Equalised Odds (EO)). A reward model  $r_\phi$  is  $\kappa$ -EO fair if the TPR/FPR are equalised across groups:  $q_{iy}^{\text{eo}}(r_\phi) := \mathbb{P}[\hat{Y} = 1 \mid G_i, Y = y]$ ,  $\Delta_{\text{eo}}(r_\phi) := \max_{i, i' \in [p], y \in \{0, 1\}} |q_{iy}^{\text{eo}}(r_\phi) - q_{i'y}^{\text{eo}}(r_\phi)| \leq \kappa$ .

**Definition 2.4** (Counterfactual Fairness (CF)). A reward model  $r_\phi$  is  $\mu$ -CF fair (conditional on  $U$ ) if, for each  $j \in [d]$ , the group-conditioned positive rates match across groups:

$$q_{ij}^{\text{cf}}(r_\phi) := \mathbb{P}[\hat{Y} = 1 \mid G_i, U = j], \quad \Delta_{\text{cf}}(r_\phi) := \max_{i, i' \in [p], j \in [d]} |q_{ij}^{\text{cf}}(r_\phi) - q_{i'j}^{\text{cf}}(r_\phi)| \leq \mu.$$

We also use an averaged version:  $\Delta_{\text{cf}}^{\text{avg}}(r_\phi) := \mathbb{E}_U [\max_{i, i' \in [p]} |q_{iU}^{\text{cf}}(r_\phi) - q_{i'U}^{\text{cf}}(r_\phi)|] \leq \mu$ .

## 3 Fairness-Aware Reward Optimization

A well-aligned reward model must capture multifaceted features of human preferences: (i) *ordinal*, correctly ranking preferred responses; (ii) *cardinal*, accurately modeling the margin of these preferences; and (iii) *fair*, ensuring that accuracy is consistent across demographics. Existing methods often focus only on ordinal accuracy, leading to poorly calibrated or systematically biased models. We proceed to develop FARO: enforcing fairness in the LLM by guaranteeing algorithmic fairness in the reward function. We augment standard preference learning with fairness constraints, reformulating this as a Lagrangian minimax problem:

$$\min_{\phi} \max_{\lambda \geq 0} L_{\text{NLL}}(\phi) + \lambda^\top C_{\text{fairness}}(\phi) \quad (4)$$

Here, we optimise over the model parameters  $\phi$ , which define a reward model  $r_\phi(x, y)$  assigning a scalar score to a response. This reward model induces a probabilistic preference model  $p_\phi(\hat{y}_w \succ \hat{y}_l \mid x)$ , typically via the Bradley-Terry model (Eq. 1). Both terms in the Lagrangian depend on these preference probabilities:  $L_{\text{NLL}}(\phi)$  is the negative log-likelihood of  $p_\phi$  with respect to human preference data, and  $C_{\text{fairness}}(\phi)$  is a vector of fairness constraint violations, computed as expectations of  $p_\phi$  across demographic groups. The dual variables  $\lambda$  are learned penalties applied to these violations.

This optimization has two challenges: (1) *non-differentiable constraints*, (2) *quadratic complexity*. Fairness constraints are often defined on empirical classification rates, which have zero gradients almost

everywhere and are unsuitable for optimisation. We instead use a differentiable proxy for these rates, defined by the model’s expected preference probability,  $\mathbb{E}[p_\phi(\hat{y}_w \succ \hat{y}_l \mid x)]$ . Moreover, to avoid quadratic  $O(p^2)$  complexity from all pairwise group comparisons, we employ the anchoring trick (Jagielski et al., 2019) and reduce the number of constraints to  $O(p)$  without loss of generality for feasibility. If constraints  $|q_1(r_\phi) - q_i(r_\phi)| \leq \gamma_i$  hold for all  $i \geq 2$  hold, then by the triangle inequality any two groups  $i, j \geq 2$  satisfy  $|q_i(r_\phi) - q_j(r_\phi)| \leq \gamma_i + \gamma_j$ .

The final FARO objective incorporates these solutions. Given a set of non-uniform fairness tolerances  $\gamma_i, \kappa_i, \mu_{ij}$ , the fairness constraint vector  $C_{\text{fairness}}(\phi)$  is defined for one of the following standards:

- (1) *DP*: The vector  $C_{\text{dp}}(\phi)$  contains the  $2(p-1)$  constraints derived from the inequalities:  $|q_1^{\text{dp}}(r_\phi) - q_i^{\text{dp}}(r_\phi)| \leq \gamma_i$  for  $i \in \{2, \dots, p\}$ .
- (2) *EO*: The vector  $C_{\text{eo}}(\phi)$  is defined analogously, with expectations taken conditioned on the human preference label  $Y = y$ :  $|q_1^{\text{eo}}(r_\phi \mid Y = y) - q_i^{\text{eo}}(r_\phi \mid Y = y)| \leq \kappa_i$  for  $i \in \{2, \dots, p\}$ ,  $y \in \{0, 1\}$ .
- (3) *CF*: The vector  $C_{\text{cf}}(\phi)$  is defined by conditioning on an unrestricted attribute  $U = j$ :  $|q_{1j}^{\text{cf}}(r_\phi \mid U = j) - q_{ij}^{\text{cf}}(r_\phi \mid U = j)| \leq \mu_{ij}$  for  $i \in \{2, \dots, p\}$ ,  $j \in [d]$ .

## 4 Theoretical Analysis

To solve problem 4, we adapt the proxy-Lagrangian gradient descent–ascent (ProxyGDA) method by Cotter et al. (2019a;b). Specifically, we instantiate the two-player game in Eq. 4 with FARO’s fairness constraints and analyse the regret bounds of the resulting dynamics. While the algorithmic template is standard, its application to fair reward modelling and downstream RLHF is novel to our work. We establish four guarantees: (i) the FARO-learned reward satisfies DP/EO/CF constraints up to a controllable, diminishing slack; (ii) RL fine-tuning with a KL penalty induces an accuracy-fairness trade-off; (iii) using a FARO-fair reward improves downstream policy fairness compared to an unconstrained reward; and (iv) varying tolerance and regularisation parameters in FARO traces a non-empty Pareto frontier of optimal solutions.

---

**Algorithm 1** PROXYGDA FOR FARO ( $R \in \mathbb{R}_+$ ,  $L_{\text{faro}} : \Phi \times \Lambda \rightarrow \mathbb{R}$ ,  $T \in \mathbb{N}$ ,  $\eta_\lambda, \eta_\phi, \varepsilon_{\text{rel}} \in \mathbb{R}_+$ )

---

```

1: Initialise  $\lambda^{(1)} = 0$ 
2: for  $t \in [T]$  do
3:   Initialise  $\phi^{(t,0)}$  randomly
4:   repeat
5:      $\phi^{(t,k+1)} = \phi^{(t,k)} - \eta_\phi \nabla_\phi L_{\text{faro}}(\phi^{(t,k)}, \lambda^{(t)})$ 
6:   until  $\frac{|L_{\text{faro}}(\phi^{(t,k)}, \lambda^{(t)}) - L_{\text{faro}}(\phi^{(t,k-1)}, \lambda^{(t)})|}{\max\{1, |L_{\text{faro}}(\phi^{(t,k)}, \lambda^{(t)})|\}} \leq \varepsilon_{\text{rel}}$ 
7:   Let  $\phi^{(t)} = \phi^{(t,k)}$ 
8:   Update  $\lambda^{(t+1)} = \Pi_\Lambda(\lambda^{(t)} + \eta_\lambda \nabla_\lambda L_{\text{faro}}(\phi^{(t)}, \lambda^{(t)}))$   $\triangleright$  Projection onto  $\Lambda = \{\lambda \geq 0 : \|\lambda\|_\infty \leq R\}$ 
9: end for
10: return averaged iterate  $\bar{\phi} = \frac{1}{T} \sum_{t=1}^T \phi^{(t)}$  (or best iterate by validation)

```

---

### 4.1 Reward-level Fairness Certificates

Algorithm 1 describes the gradient descent–ascent method for optimizing the FARO Lagrangian,  $L_{\text{faro}}$ . The inner loop (Lines 4–6) iteratively finds an approximate minimiser of the loss with respect to the model



parameters  $\phi$ , stopping when a relative tolerance  $\varepsilon_{\text{rel}}$  is met. This process yields a  $\rho$ -approximate solution  $\phi^{(t)}$ , where the absolute approximation error  $\rho$  is implicitly controlled by  $\varepsilon_{\text{rel}}$ . Based on this procedure, we can certify the fairness of the resulting reward model.

**Proposition 4.1** (Population fairness certificate for FARO). *Let  $\bar{\phi}$  be the averaged iterate returned by PROXYGDA. Then with probability at least  $1 - \delta$ , the population fairness violations of  $r_{\bar{\phi}}$  satisfy*

$$\max_c \Delta^c(\bar{\phi}) \leq \rho + \tilde{O}\left(\frac{R}{\sqrt{T}}\right) + O\left(\sqrt{\frac{\log(1/\delta)}{n_{\min}}}\right),$$

where  $c \in \{\text{dp}, \text{eo}, \text{cf}\}$  and  $n_{\min} = \min_i n_i$  is the sample size of the smallest sensitive subgroup. Thus  $r_{\bar{\phi}}$  is  $\gamma$ -DP /  $\kappa$ -EO /  $\mu$ -CF fair up to a controllable slack consisting of the inner-loop optimisation error  $\rho$ , convergence error, and a generalisation gap that vanishes with more data (see App. B.1).

**Corollary 4.2** (Group-wise  $\tau$ -rules). *For any feasible solution to the group-fair DP program in Eq. 4 with ordered allowances  $\{\gamma_i\}$ , the learned reward  $r_{\phi}$  satisfies, with slack  $\varepsilon_T = \rho + O(RG\sqrt{k/T})$  (Prop. 4.1),*

$$\begin{aligned} |q_i^{\text{dp}}(r_{\phi}) - q_j^{\text{dp}}(r_{\phi})| &\leq \gamma_i + \gamma_j + 2\varepsilon_T, & \forall i, j \geq 1, \\ q_1^{\text{dp}}(r_{\phi}) - (\gamma_i + \varepsilon_T) &\leq q_i^{\text{dp}}(r_{\phi}) \leq q_1^{\text{dp}}(r_{\phi}) + (\gamma_i + \varepsilon_T), & \forall i \geq 2. \end{aligned}$$

Together, Prop. 4.1 and Cor. 4.2 certify that FARO yields a reward model that is DP/EO/CF-fair up to a controllable slack  $\varepsilon_T = \rho + O(RG\sqrt{k/T})$ , plus a statistical term  $O(\sqrt{\frac{\log(1/\delta)}{n_{\min}}})$  for population guarantees (App. B.1- B.2). The slack is governed by the optimisation budget  $T$ , inner-loop tolerance  $\rho$ , and data balance.

## 4.2 Finetuning induces an accuracy–fairness trade-off

Having engineered a fair reward model  $r_{\phi}$ , we now analyse how it can be used to induce fairness in a performant but potentially biased LLM policy,  $\pi_{\text{ref}}$ . The process of RL fine-tuning creates three-way tension: *alignment*, i.e. maximising the score from the fair reward model  $r_{\phi}$ ; *performance retention*, i.e. staying close to the strong reference policy  $\pi_{\text{ref}}$  that captures general capabilities; and *final policy fairness*, i.e. ensuring that the resulting fine-tuned policy  $\pi_{\beta}$  is itself fair. The KL-regularised objective from Eq. 2 illustrates this trade-off. The fair reward  $r_{\phi}$  pulls the policy towards a fair region, while the KL term acts as an anchor to the performant  $\pi_{\text{ref}}$ , controlled by the hyperparameter  $\beta$ . A small  $\beta$  allows greater deviation towards the fair reward (potentially improving fairness and accuracy), whereas a large  $\beta$  keeps the policy close to  $\pi_{\text{ref}}$  (preserving task behaviour but also its unfairness).

We measure divergence between  $\pi, \pi'$  by the policy-induced KL  $D_{\text{KL}}(\pi \| \pi') = \mathbb{E}_{x \sim D}[D_{\text{KL}}(\pi(\cdot | x) \| \pi'(\cdot | x))]$ . By Pinsker’s inequality (Lemma B.1), deviations from  $\pi_{\text{ref}}$  bound changes in group-level probabilities.

**Proposition 4.3** (KL-regularised trade-off and fairness drift). *Let  $\pi_{\beta}$  be any maximizer of the KL-regularized objective  $\mathcal{J}_{\beta}(\pi) = \mathbb{E}_{x, a \sim \pi}[r_{\phi}(x, a)] - \beta D_{\text{KL}}(\pi \| \pi_{\text{ref}})$ . Then:*

1. (Monotonicity) *If  $\beta_1 > \beta_2 > 0$  then  $D_{\text{KL}}(\pi_{\beta_1} \| \pi_{\text{ref}}) \leq D_{\text{KL}}(\pi_{\beta_2} \| \pi_{\text{ref}})$ .*
2. (Fairness drift) *The fairness violation of the final policy  $\pi_{\beta}$  is bounded by the violation of the initial reference policy,  $\Delta(\pi_{\text{ref}})$ , plus a drift term:  $\Delta(\pi_{\beta}) \leq \Delta(\pi_{\text{ref}}) + \sqrt{2 D_{\text{KL}}(\pi_{\beta} \| \pi_{\text{ref}})}$ .*

Prop. 4.3 reveals the dual role of the KL term: beyond regularising the policy update to preserve the capabilities of  $\pi_{\text{ref}}$ , it provides a worst-case guarantee that the fairness violation will not degrade arbitrarily. The trade-off for a practitioner is thus in the choice of  $\beta$ , which balances the pursuit of higher reward

(permitting a larger KL divergence) against maintaining a tighter fairness bound. This reframes the objective from preserving the biased reference policy to controlling the magnitude of the departure from it. This raises the crucial question of how to ensure this “drift” is a beneficial move towards fairness. We address this in Thm. 4.4, which shows that updates guided by a FARO-fair reward model provably improve downstream policy fairness.

### 4.3 Fairer RL policies when using fair rewards

We have established that RL fine-tuning involves a controlled “drift” away from a performant but potentially biased reference policy,  $\pi_{\text{ref}}$ . Prop. 4.3 showed that the magnitude of this drift, controlled by  $\beta$ , has a bounded effect on the final policy’s fairness. This raises the central question: if we guide this drift with a FARO-fair reward model, does it actually produce a fairer final LLM?

We answer in the affirmative, showing that the fairness engineered into the reward model provably transfers to the fine-tuned policy. This is a critical result, as it guarantees that our efforts at the reward-modeling stage are not “lost in translation” during the complex dynamics of RL optimisation. It shows that using a fair reward is demonstrably better than using an unconstrained one.

**Theorem 4.4** (Reward-to-Policy Fairness Transfer). *Let  $r_{\text{plain}}$  be a reward model trained to optimise only the preference loss (Eq. 1) on a given dataset, and let  $r_{\phi}$  be the FARO-fair reward model trained on the same data with an additional fairness constraint. Let their resulting fairness violations be  $\Delta(\pi_{\beta_{\text{plain}}})$  and  $\Delta(\pi_{\beta_{\text{fair}}})$  respectively, after fine-tuning from the same  $\pi_{\text{ref}}$  with the same KL-penalty  $\beta$ . Under standard monotonicity assumptions, the violations are related by  $\Delta(\pi_{\beta}^{\text{fair}}) \leq \Delta(\pi_{\beta}^{\text{plain}}) + \varepsilon_T$ .*

$\varepsilon_T$  is the fairness violation slack of the reward model  $r_{\phi}$  from Prop. 4.1.

For any given level of fine-tuning (i.e. for any fixed  $\beta$ ), replacing a standard, unconstrained reward with a FARO-fair reward will improve (at worst, not harm) the downstream fairness of the resulting LLM policy, up to the small slack  $\varepsilon_T$ . A fair reward function makes the final policy fairer. The guarantees we establish at the reward level propagate through the RL fine-tuning process, providing a principled mechanism for producing fairer policies in practice (App. B.4).

### 4.4 Existence of Pareto-optimal operating points

Consider the bi-objective problem of minimising (error, fairness), where “error” is a suitable accuracy metric and “fairness” is one of  $\Delta_{\text{dp}}, \Delta_{\text{eo}}, \Delta_{\text{cf}}$ . We establish the existence of a non-empty Pareto frontier, ensuring that there are well-defined operating points trading off fairness and accuracy.

**Proposition 4.5** (Non-empty Pareto frontier). *Varying FARO’s fairness tolerance schedules  $\{\gamma_i\}$ ,  $\{\kappa_i\}$ ,  $\{\mu_{ij}\}$  and the KL-regularisation parameter  $\beta$  within compact sets traces a non-empty, continuous Pareto frontier in the (error, fairness) objective space.*

This guarantee arises from a standard topological argument (full proof in App B.5). The space of hyperparameters is compact by definition. The mapping from these parameters to the resulting optimal policy,  $\pi^*$ , is continuous by Berge’s Maximum Theorem, as is the subsequent mapping from the policy to its (error, fairness) evaluation. The continuous image of a compact set is also compact; therefore, the set of all achievable outcomes is a compact set in  $\mathbb{R}^2$ , which ensures the existence of a non-empty Pareto frontier.

As tolerances  $(\gamma, \kappa, \mu) \rightarrow 0$  and  $\beta \rightarrow \infty$ , the policy remains close to the unfair reference  $\pi_{\text{ref}}$ . Conversely, as  $\beta \rightarrow 0$ , the policy utilises the fair reward  $r_{\phi}$ , improving fairness with controlled deviation from  $\pi_{\text{ref}}$ . FARO efficiently traverses the continuous trade-off space, and yields Pareto-efficient policies in RL finetuning.



Table 2: **FARO optimises fairness while preserving performance.** We measure fairness on BBQ after reward-optimising on PRISM. We find that FARO allows us to significantly improve over the base model with regards to bias scores. This also seems to correlate with changes in the scores of DP, EO, and CF.

Model	Disamb Top-1 ( $\uparrow$ )	Ambig Top-1 ( $\uparrow$ )	Ambig Bias Score ( $\downarrow$ )	Disambig Bias Score ( $\downarrow$ )	$\Delta_{DP/EO/CF}$
Gemma-2-2b-it	83.20	63.91	14.73	-0.811	N/A
Gemma-2-2b-it – FARO (DP)	<b>83.93</b>	63.20	<b>6.81</b>	<b>-1.01</b>	0.55
Gemma-2-2b-it – FARO (CF)	83.10	62.86	10.55	-0.965	0.41
Gemma-2-2b-it – FARO (EO)	82.71	<b>63.72</b>	12.96	-0.822	0.44
Phi-3-Mini	71.92	42.14	11.91	1.42	N/A
Phi-3-Mini – FARO (DP)	<b>71.99</b>	<b>46.55</b>	9.15	1.01	0.21
Phi-3-Mini – FARO (EO)	70.05	44.01	10.86	1.04	0.18
Phi-3-Mini – FARO (CF)	71.73	45.92	<b>9.01</b>	<b>0.93</b>	0.37
Qwen-2.5-1.5B	74.14	58.97	11.44	-.0922	N/A
Qwen-2.5-1.5B – FARO (DP)	<b>75.11</b>	<b>59.18</b>	9.11	-.104	0.26
Qwen-2.5-1.5B – FARO (EO)	74.06	57.66	10.87	-0.100	0.054
Qwen-2.5-1.5B – FARO (CF)	73.12	58.91	<b>8.04</b>	<b>-0.155</b>	0.091

## 5 Experiments

We evaluate each setting FARO on safety oriented datasets. For each run we optimize a single fairness family `FARO_dp`, `FARO_eo`, or `FARO_cf`. We finetune an instruction tuned language model with reward modeling and use the learned reward to assess multiple choice selection or to rerank sampled generations at inference.

**Finetuning Dataset.** We train the reward model on PRISM (Kirk et al., 2024), a pairwise preference corpus that is grouped by sociodemographic attributes. Our implementation follows the proxy Lagrangian with anchoring. For a chosen family we form anchored constraints over the differentiable preference probability, learn nonnegative dual variables with projection, and optimize the Bradley Terry negative log likelihood plus the active constraint term. Training uses a value head on top of the policy.

**Evaluation Datasets.** We evaluate with dataset specific quality and fairness metrics. On *BBQ* (Parrish et al., 2022), we use the Ambiguous and Disambiguated settings. We report top-1 accuracy and the official bias scores for both settings. For BBQ we score multiple choice options by the reward and select the argmax.

**Models.** We use three public instruction tuned models: Gemma-2-2B (Team et al., 2024), Phi-3-Mini (Abdin et al., 2024), and Qwen-2.5-1.5B (Team, 2024). Reward modeling is performed with a causal language model and value head. For evaluation we either rank options by the reward, or generate with the policy and optionally apply reward based reranking as above. We include a baseline comparison with the original language model: for Gemma, we use reported scores for the BBQ evaluation; for Phi-3 and Qwen-2.5, we extract findings using the same approaches in Parrish et al. (2022) and report scores in the table independently.

We show results in Table 2. We find that FARO allows us to consistently reduce bias as given by the bias score while preserving general accuracy.

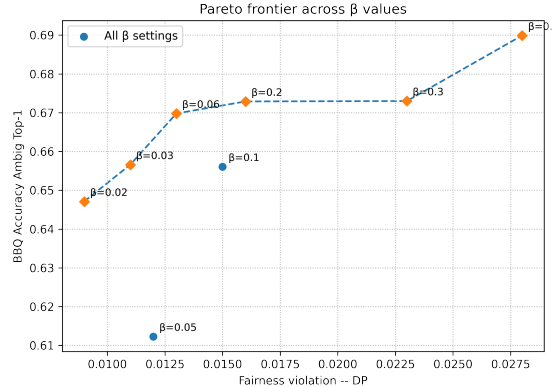


Figure 2: **Pareto Frontier of fairness and accuracy.** We vary  $\beta$  and use FARO\_DP as the reward for Gemma on PRISM. We plot the fairness violation and BBQ Top-1 accuracy for the ambiguous dataset, and compute the pareto optimal set of  $\beta$ s by finding all dominated points where all neighboring points are strictly better.

### 5.1 Optimising the Error vs. Fairness-Violation Trade-off

We show a comparison between our top-1 accuracy and DP loss for Gemma in Figure 2. This is varied across several  $\beta$  values to understand how  $\beta$  affects the relationship between fairness and accuracy for Gemma. We evaluate the fairness using DP and evaluate the accuracy using the accuracy of BBQ for Gemma. We include the pareto frontier by finding all points which are dominated by other points around them and removing them. Surprisingly, we see some sensitivity to  $\beta$  e.g. with 0.05. This could be due to some tuning whereby smaller  $\beta$  is not well balanced between the fairness and accuracy considerations.

## 6 Extended Related Work

**Why Preference Alignment is Not “Fair Enough”.** Using preference datasets for fairness-alignment is vulnerable to *selection bias* from data collection oversight; *popularity bias* from disproportional survey participation; *cognitive bias* from prejudiced human annotators<sup>1</sup>. Flawed data collection induces selection bias, where response data is biased by preferences of the surveyed demographic and encodes spurious correlations (Ovaisi et al., 2020; Wang et al., 2021; Liu et al., 2022); skewed survey strategies lead to popularity bias, where preference data is sparse, long-tailed and lacks coverage for less-preferred or uncommon responses, leading to unpredictable behaviour under distribution shifts or edge cases (Chen et al., 2023; Zhao et al., 2023; Naghiaei et al., 2022). Beyond statistical dataset biases, previous works also deal with a lack of fairness in LLMs’ judgement calls, as a product of defective training and alignment procedures. Surveys on fairness (Mehrabi et al., 2021; Gallegos et al., 2024) reveal cognitive biases in LLM judges that mirror human prejudice, leading to disparate treatment on unfair bases of gender, authority, beauty standards, misinformation (Chen et al., 2024a; Koo et al., 2024; Zheng et al., 2023). Since real-world datasets are rampant with discriminatory examples, bias arising from these skewed representations are encoded into the reward model (Wang et al., 2023b; Liu et al., 2020); the reward model propagates undesirable biases to the LLM through RL fine-tuning and reinforces unfair behaviour (Blodgett et al., 2020; 2021). Without an internal constitution, our LLMs are being corrupted by—instead of correcting—instances of bias, unfairness and prejudice in big data. Enter the

<sup>1</sup>We refer readers to Gallegos et al. (2024) for a thorough and insightful breakdown of the metrics, datasets, mitigation techniques and open problems concerning LLM bias and fairness.

role of theory as a tool to specify, certify and codify principles of equality into the LLM constitution.

**Algorithmic Fairness.** Fairness concerns equalising the treatment and consideration of people, identified by their (protected) demographic attributes, such as gender, race, age (Commission, 1964). Fairness-aware algorithms intervene on the learning problem to avoid *disparate treatment* of people, with hopes of also reducing *disparate impact* of decision-making outcomes (Barocas & Selbst, 2016). Notable definitions of group-fairness include demographic parity (Dwork et al., 2012), equalised odds and equal opportunity (Hardt et al., 2016), calibration (Kleinberg et al., 2017); different approaches to fairness are benchmarked on real-world datasets, including TransUnion TransRisk (Avery et al., 2009), UCI Adult (Kohavi & Becker, 1996), Dutch census (Center, 2019), COMPAS (Larson et al., 2016), ACS PUMS (Ding et al., 2021). Towards mitigating group-wise disparate treatment, there are three types of fairness interventions—*pre-processing*, *in-processing* and *post-processing* (Zeng et al., 2022; Pleiss et al., 2017). This mirrors the dilemma of fairness-aware processing in LLMs, where the timing of when to intervene has crucial impacts on both the algorithm (computational and sample efficiency, optimisation stability) and its outcomes (whether it has theoretical guarantees, its Pareto-efficiency). Pre-processing aims to filter away latent biases in training data through transformations (Feldman et al., 2015; Lum & Johndrow, 2016; Johndrow & Lum, 2019; Calmon et al., 2017), fair representations (Zemel et al., 2013; Louizos et al., 2015; Creager et al., 2019) and fair generative modelling (Xu et al., 2018; Sattigeri et al., 2019; Jang et al., 2021). Although such methods are broadly applicable to any learning problem, pre-processing requires an extra, expensive pass over data, and lacks formal guarantees of fairness, since disparities may persist even post-filtering (Locatello et al., 2019). A different approach is to post-process a trained model by shifting its decision boundary – particularly with group-wise thresholding rules (Fish et al., 2016; Corbett-Davies et al., 2017; Menon & Williamson, 2018; Chzhen et al., 2019; Jang et al., 2022) – to adjust for fairness. However, previous work has demonstrated that post-processing is unable to achieve optimality in both error-calibration and fairness; that post-processing for one particular notion of fairness could be in contradiction other important but incompatible notions (Chouldechova, 2017; Kleinberg et al., 2017; Woodworth et al., 2017; Corbett-Davies et al., 2017).

## 7 Discussion

**Limitations.** This paper investigates fairness shortcomings in LLM alignment and proposes improvements via FARO, an in-processing intervention with fairness constraints during RLHF’s reward modelling phase. We contribute 3 desirable properties—the ability to correct for human bias in datasets; to conduct fairness-aware optimisation in an annotation efficient manner; to derive reward models that are algorithmically fair with high Pareto-efficiency. While mathematical guarantees can guard against worst-case examples of egregious discrimination, fairness is an inherently societal concept; theoretical models must be continuously updated by inter-disciplinary research. Algorithms like FARO should be used to complement – not replace – other fairness guardrails (*e.g.* data-filtering, unlearning, calibration). A fair model can still be misused; due diligence, rigorous auditing, collecting and incorporating user feedback are as important as ever before.

**Self-critiquing LLMs.** One exciting direction concerns LLMs self-improvement and self-correction by critiquing their own outputs. Self-criticism generates new instructions and the model is realigned to the new instructions (Zheng et al., 2023; Wang et al., 2023c; Honovich et al., 2023). In the context of fairness-aligned optimisation, this involves using context-dependent techniques (Wang et al., 2024) to first, infer sensitive and unrestricted attributes from input prompts; then, recast reward modelling as an “Attribute Aware”, in-processing problem; finally, iteratively self-assess its Pareto-efficiency and adjust the attributes-classifier to issue systematic updates to the self-rewarding mechanism (Yuan et al., 2024). Although it is yet unclear whether LLMs as self-critics can be immune from evaluation bias, human/dataset bias and positional

instability issues (Wang et al., 2023a; Koo et al., 2024; Chen et al., 2024b; Sun et al., 2024), we are optimistic that a hybrid approach structured with algorithmic fairness could reveal new strategies for robust, stable and fair self-alignment.

**Conclusion.** We tackle the issue of demographic bias in LLM alignment, which propagates from skewed or prejudiced human preference data. We argue that existing interventions are unable to address all axes of the problem, where a suitable reward model must be simultaneously *ordinal*, *cardinal*, and *fair*. Towards codifying and reinforcing fair behaviour, we introduce FARO, an in-processing framework that directly embeds *algorithmic fairness constraints into the reward modeling objective*. Our theoretical analysis provides several guarantees; notably, that the fairness engineered into the reward model provably *transfers to the fine-tuned policy*, and that *a Pareto frontier of optimal solutions exists*. We validate this theory across the BBQ benchmark and three LLMs, confirming that FARO significantly reduces biased or prejudiced generations whilst preserving model quality. This work offers a principled and verifiable path toward more equitable LLMs that are fair by design.

## 8 Acknowledgements

We would like to thank Amin Charusaie for his insightful feedback on our work. This work was supported in part by the Alexander von Humboldt Foundation and the Munich Center for Machine Learning.

## References

- Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, et al. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*, 2024. (Cited on page 9.)
- Alekh Agarwal, Alina Beygelzimer, Miroslav Dudík, John Langford, and Hanna Wallach. A reductions approach to fair classification. In *International conference on machine learning*, pp. 60–69. PMLR, 2018. (Cited on page 4.)
- Robert Avery, Kenneth Brevoort, and Glenn Canner. Credit scoring and its effects on the availability and affordability of credit. *Journal of Consumer Affairs - J CONSUM AFF*, 43:516–537, 09 2009. doi: 10.1111/j.1745-6606.2009.01151.x. (Cited on page 11.)
- Mohammad Gheshlaghi Azar, Zhaohan Daniel Guo, Bilal Piot, Remi Munos, Mark Rowland, Michal Valko, and Daniele Calandriello. A general theoretical paradigm to understand learning from human preferences. In *International Conference on Artificial Intelligence and Statistics*, pp. 4447–4455. PMLR, 2024. (Cited on page 3.)
- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022. (Cited on pages 2 and 3.)
- Solon Barocas and Andrew D Selbst. Big data’s disparate impact. *Calif. L. Rev.*, 104:671, 2016. (Cited on pages 1 and 11.)
- Solon Barocas, Moritz Hardt, and Arvind Narayanan. *Fairness and Machine Learning: Limitations and Opportunities*. MIT Press, 2023. (Cited on page 2.)

- Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On the dangers of stochastic parrots: Can language models be too big? . In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '21, pp. 610–623, New York, NY, USA, 2021. Association for Computing Machinery. ISBN 9781450383097. doi: 10.1145/3442188.3445922. URL <https://doi.org/10.1145/3442188.3445922>. (Cited on page 1.)
- Su Lin Blodgett, Solon Barocas, Hal Daumé III, and Hanna Wallach. Language (technology) is power: A critical survey of “bias” in NLP. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (eds.), *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 5454–5476, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.485. URL <https://aclanthology.org/2020.acl-main.485/>. (Cited on page 10.)
- Su Lin Blodgett, Gilsinia Lopez, Alexandra Olteanu, Robert Sim, and Hanna Wallach. Stereotyping Norwegian salmon: An inventory of pitfalls in fairness benchmark datasets. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (eds.), *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pp. 1004–1015, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.81. URL <https://aclanthology.org/2021.acl-long.81/>. (Cited on page 10.)
- Ralph Allan Bradley and Milton E. Terry. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952. ISSN 00063444, 14643510. (Cited on page 3.)
- Flavio Calmon, Dennis Wei, Bhanukiran Vinzamuri, Karthikeyan Natesan Ramamurthy, and Kush R Varshney. Optimized pre-processing for discrimination prevention. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL [https://proceedings.neurips.cc/paper\\_files/paper/2017/file/9a49a25d845a483fae4be7e341368e36-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/9a49a25d845a483fae4be7e341368e36-Paper.pdf). (Cited on page 11.)
- L Elisa Celis, Lingxiao Huang, Vijay Keswani, and Nisheeth K Vishnoi. Classification with fairness constraints: A meta-algorithm with provable guarantees. In *Proceedings of the conference on fairness, accountability, and transparency*, pp. 319–328, 2019. (Cited on page 4.)
- Minnesota Population Center. Integrated public use microdata series, international: Version 7.2. *Minneapolis: University of Minnesota*, 2019. URL <http://doi.org/10.18128/D020.V7.2>. (Cited on page 11.)
- Guiming Hardy Chen, Shunian Chen, Ziche Liu, Feng Jiang, and Benyou Wang. Humans or LLMs as the judge? a study on judgement bias. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 8301–8327, Miami, Florida, USA, November 2024a. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.474. URL <https://aclanthology.org/2024.emnlp-main.474/>. (Cited on page 10.)
- Guiming Hardy Chen, Shunian Chen, Ziche Liu, Feng Jiang, and Benyou Wang. Humans or llms as the judge? a study on judgement biases. *arXiv preprint arXiv:2402.10669*, 2024b. (Cited on page 12.)
- Jiawei Chen, Hande Dong, Xiang Wang, Fuli Feng, Meng Wang, and Xiangnan He. Bias and debias in recommender system: A survey and future directions. *ACM Trans. Inf. Syst.*, 41(3), February 2023. ISSN 1046-8188. doi: 10.1145/3564284. URL <https://doi.org/10.1145/3564284>. (Cited on page 10.)

- Alexandra Chouldechova. Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. *arXiv preprint arXiv:1703.00056*, 2017. (Cited on page 11.)
- Evgenii Chzhen, Christophe Denis, Mohamed Hebiri, Luca Oneto, and Massimiliano Pontil. Leveraging labeled and unlabeled data for consistent fair binary classification. *Advances in Neural Information Processing Systems*, 32, 2019. (Cited on page 11.)
- U.S. Equal Employment Opportunity Commission. Equal employment opportunities. *Civil Rights Act of 1964*, 1964. (Cited on page 11.)
- Sam Corbett-Davies, Emma Pierson, Avi Feller, Sharad Goel, and Aziz Huq. Algorithmic decision making and the cost of fairness. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '17, pp. 797–806, New York, NY, USA, 2017. Association for Computing Machinery. ISBN 9781450348874. doi: 10.1145/3097983.3098095. URL <https://doi.org/10.1145/3097983.3098095>. (Cited on page 11.)
- Andrew Cotter, Maya Gupta, Heinrich Jiang, Nathan Srebro, Karthik Sridharan, Serena Wang, Blake Woodworth, and Seungil You. Training well-generalizing classifiers for fairness metrics and other data-dependent constraints. In *International Conference on Machine Learning*, pp. 1397–1405. PMLR, 2019a. (Cited on pages 6 and 22.)
- Andrew Cotter, Heinrich Jiang, and Karthik Sridharan. Two-player games for efficient non-convex constrained optimization. In Aurélien Garivier and Satyen Kale (eds.), *Proceedings of the 30th International Conference on Algorithmic Learning Theory*, volume 98 of *Proceedings of Machine Learning Research*, pp. 300–332. PMLR, 22–24 Mar 2019b. URL <https://proceedings.mlr.press/v98/cotter19a.html>. (Cited on page 6.)
- Elliot Creager, David Madras, Joern-Henrik Jacobsen, Marissa Weis, Kevin Swersky, Toniann Pitassi, and Richard Zemel. Flexibly fair representation learning by disentanglement. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 1436–1445. PMLR, 09–15 Jun 2019. URL <https://proceedings.mlr.press/v97/creager19a.html>. (Cited on page 11.)
- Sunhao Dai, Chen Xu, Shicheng Xu, Liang Pang, Zhenhua Dong, and Jun Xu. Bias and unfairness in information retrieval systems: New challenges in the llm era. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 6437–6447, 2024. (Cited on page 1.)
- Sumanth Dathathri, Andrea Madotto, Janice Lan, Jane Hung, Eric Frank, Piero Molino, Jason Yosinski, and Rosanne Liu. Plug and play language models: A simple approach to controlled text generation. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=HledEyBKDS>. (Cited on page 2.)
- Frances Ding, Moritz Hardt, John Miller, and Ludwig Schmidt. Retiring adult: New datasets for fair machine learning. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, volume 34, pp. 6478–6490. Curran Associates, Inc., 2021. URL [https://proceedings.neurips.cc/paper\\_files/paper/2021/file/32e54441e6382a7fbacbbbf3c450059-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/32e54441e6382a7fbacbbbf3c450059-Paper.pdf). (Cited on pages 2 and 11.)
- Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. Fairness through awareness. In *Proceedings of the 3rd Innovations in Theoretical Computer Science Conference*, ITCS '12, pp. 214–226, New York, NY, USA, 2012. Association for Computing Machinery. ISBN 9781450311151. doi: 10.1145/2090236.2090255. URL <https://doi.org/10.1145/2090236.2090255>. (Cited on page 11.)



- Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. Kto: Model alignment as prospect theoretic optimization. *arXiv preprint arXiv:2402.01306*, 2024. (Cited on pages 3 and 21.)
- Michael Feldman, Sorelle A Friedler, John Moeller, Carlos Scheidegger, and Suresh Venkatasubramanian. Certifying and removing disparate impact. In *proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining*, pp. 259–268, 2015. (Cited on pages 5 and 11.)
- Benjamin Fish, Jeremy Kun, and Ádám D Lelkes. A confidence-based approach for balancing fairness and accuracy. In *Proceedings of the 2016 SIAM international conference on data mining*, pp. 144–152. SIAM, 2016. (Cited on page 11.)
- Sorelle A Friedler, Carlos Scheidegger, and Suresh Venkatasubramanian. The (im) possibility of fairness: Different value systems require different mechanisms for fair decision making. *Communications of the ACM*, 64(4):136–143, 2021. (Cited on page 5.)
- Isabel O Gallegos, Ryan A Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K Ahmed. Bias and fairness in large language models: A survey. *Computational Linguistics*, pp. 1–79, 2024. (Cited on page 10.)
- Deep Ganguli, Liane Lovitt, Jackson Kernion, Amanda Askell, Yuntao Bai, Saurav Kadavath, Ben Mann, Ethan Perez, Nicholas Schiefer, Kamal Ndousse, et al. Red teaming language models to reduce harms: Methods, scaling behaviors, and lessons learned. *arXiv preprint arXiv:2209.07858*, 2022. (Cited on page 2.)
- Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A Smith. Realtotoxicityprompts: Evaluating neural toxic degeneration in language models. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pp. 3356–3369, 2020. (Cited on page 2.)
- Moritz Hardt, Eric Price, and Nati Srebro. Equality of opportunity in supervised learning. *Advances in neural information processing systems*, 29, 2016. (Cited on pages 4, 5, and 11.)
- Or Honovich, Thomas Scialom, Omer Levy, and Timo Schick. Unnatural instructions: Tuning language models with (almost) no human labor. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 14409–14428, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.806. URL <https://aclanthology.org/2023.acl-long.806/>. (Cited on page 11.)
- Matthew Jagielski, Michael Kearns, Jieming Mao, Alina Oprea, Aaron Roth, Saeed Sharifi-Malvajerdi, and Jonathan Ullman. Differentially private fair learning. In *International Conference on Machine Learning*, pp. 3000–3008. PMLR, 2019. (Cited on page 6.)
- Taeuk Jang, Feng Zheng, and Xiaoqian Wang. Constructing a fair classifier with generated fair data. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(9):7908–7916, May 2021. doi: 10.1609/aaai.v35i9.16965. URL <https://ojs.aaai.org/index.php/AAAI/article/view/16965>. (Cited on page 11.)
- Taeuk Jang, Pengyi Shi, and Xiaoqian Wang. Group-aware threshold adaptation for fair classification. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(6):6988–6995, Jun. 2022. doi: 10.1609/aaai.v36i6.20657. URL <https://ojs.aaai.org/index.php/AAAI/article/view/20657>. (Cited on page 11.)

- Natasha Jaques, Shixiang Gu, Dzmitry Bahdanau, José Miguel Hernández-Lobato, Richard E. Turner, and Douglas Eck. Sequence tutor: Conservative fine-tuning of sequence generation models with KL-control. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 1645–1654. PMLR, 06–11 Aug 2017. URL <https://proceedings.mlr.press/v70/jaques17a.html>. (Cited on page 4.)
- James E Johndrow and Kristian Lum. An algorithm for removing sensitive information. *The Annals of Applied Statistics*, 13(1):189–220, 2019. (Cited on page 11.)
- Hannah Rose Kirk, Alexander Whitefield, Paul Röttger, Andrew Bean, Katerina Margatina, Juan Ciro, Rafael Mosquera, Max Bartolo, Adina Williams, He He, et al. The prism alignment project: What participatory, representative and individualised human feedback reveals about the subjective and multicultural alignment of large language models. *Advances in Neural Information Processing Systems*, 2024. (Cited on page 9.)
- Jon Kleinberg, Sendhil Mullainathan, and Manish Raghavan. Inherent trade-offs in the fair determination of risk scores. In *8th Innovations in Theoretical Computer Science Conference (ITCS 2017)*. Schloss-Dagstuhl-Leibniz Zentrum für Informatik, 2017. (Cited on pages 5 and 11.)
- Ronny Kohavi and Barry Becker. Uci adult data set. *UCI machine learning repository*, 5:2093, 1996. (Cited on page 11.)
- Ryan Koo, Minhwa Lee, Vipul Raheja, Jong Inn Park, Zae Myung Kim, and Dongyeop Kang. Benchmarking cognitive biases in large language models as evaluators. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Findings of the Association for Computational Linguistics: ACL 2024*, pp. 517–545, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.29. URL <https://aclanthology.org/2024.findings-acl.29/>. (Cited on pages 10 and 12.)
- Ben Krause, Akhilesh Deepak Gotmare, Bryan McCann, Nitish Shirish Keskar, Shafiq Joty, Richard Socher, and Nazneen Fatema Rajani. Gedi: Generative discriminator guided sequence generation. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pp. 4929–4952, 2021. (Cited on page 2.)
- Jeff Larson, Surya Mattu, Lauren Kirchner, and Julia Angwin. How we analyzed the compas recidivism algorithm. *ProPublica (5 2016)*, 9(1):3–3, 2016. (Cited on page 11.)
- Alisa Liu, Maarten Sap, Ximing Lu, Swabha Swayamdipta, Chandra Bhagavatula, Noah A. Smith, and Yejin Choi. DExperts: Decoding-time controlled text generation with experts and anti-experts. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (eds.), *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pp. 6691–6706, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.522. URL <https://aclanthology.org/2021.acl-long.522/>. (Cited on page 2.)
- Dugang Liu, Pengxiang Cheng, Zhenhua Dong, Xiuqiang He, Weike Pan, and Zhong Ming. A general knowledge distillation framework for counterfactual recommendation via uniform data. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '20*, pp. 831–840, New York, NY, USA, 2020. Association for Computing Machinery. ISBN 9781450380164. doi: 10.1145/3397271.3401083. URL <https://doi.org/10.1145/3397271.3401083>. (Cited on page 10.)
- Haochen Liu, Da Tang, Ji Yang, Xiangyu Zhao, Hui Liu, Jiliang Tang, and Youlong Cheng. Rating distribution calibration for selection bias mitigation in recommendations. In *Proceedings of the ACM Web Conference 2022, WWW '22*, pp. 2048–2057, New York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450390965. doi: 10.1145/3485447.3512078. URL <https://doi.org/10.1145/3485447.3512078>. (Cited on page 10.)

- Francesco Locatello, Gabriele Abbati, Thomas Rainforth, Stefan Bauer, Bernhard Schölkopf, and Olivier Bachem. On the fairness of disentangled representations. *Advances in neural information processing systems*, 32, 2019. (Cited on page 11.)
- Christos Louizos, Kevin Swersky, Yujia Li, Max Welling, and Richard Zemel. The variational fair autoencoder. *arXiv preprint arXiv:1511.00830*, 2015. (Cited on page 11.)
- Kristian Lum and James Johndrow. A statistical framework for fair predictive algorithms. *arXiv preprint arXiv:1610.08077*, 2016. (Cited on page 11.)
- Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. A survey on bias and fairness in machine learning. *ACM Comput. Surv.*, 54(6), July 2021. ISSN 0360-0300. doi: 10.1145/3457607. URL <https://doi.org/10.1145/3457607>. (Cited on page 10.)
- Aditya Krishna Menon and Robert C. Williamson. The cost of fairness in binary classification. In Sorelle A. Friedler and Christo Wilson (eds.), *Proceedings of the 1st Conference on Fairness, Accountability and Transparency*, volume 81 of *Proceedings of Machine Learning Research*, pp. 107–118. PMLR, 23–24 Feb 2018. URL <https://proceedings.mlr.press/v81/menon18a.html>. (Cited on pages 4, 5, and 11.)
- Mohammadmehdi Naghiaei, Hossein A Rahmani, and Mahdi Dehghan. The unfairness of popularity bias in book recommendation. In *International Workshop on Algorithmic Bias in Search and Recommendation*, pp. 69–81. Springer, 2022. (Cited on page 10.)
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022. (Cited on page 3.)
- Zohreh Ovaisi, Ragib Ahsan, Yifan Zhang, Kathryn Vasilaky, and Elena Zheleva. Correcting for selection bias in learning-to-rank systems. In *Proceedings of The Web Conference 2020*, pp. 1863–1873, 2020. (Cited on page 10.)
- Alicia Parrish, Angelica Chen, Nikita Nangia, Vishakh Padmakumar, Jason Phang, Jana Thompson, Phu Mon Htut, and Samuel Bowman. Bbq: A hand-built bias benchmark for question answering. In *Findings of the Association for Computational Linguistics: ACL 2022*, pp. 2086–2105, 2022. (Cited on page 9.)
- Ethan Perez, Saffron Huang, Francis Song, Trevor Cai, Roman Ring, John Aslanides, Amelia Glaese, Nat McAleese, and Geoffrey Irving. Red teaming language models with language models. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 3419–3448, 2022. (Cited on page 2.)
- Geoff Pleiss, Manish Raghavan, Felix Wu, Jon Kleinberg, and Kilian Q Weinberger. On fairness and calibration. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL [https://proceedings.neurips.cc/paper\\_files/paper/2017/file/b8b9c74ac526ffffbe2d39ab038d1cd7-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/b8b9c74ac526ffffbe2d39ab038d1cd7-Paper.pdf). (Cited on page 11.)
- Novi Quadrianto and Viktoriia Sharmanska. Recycling privileged learning and distribution matching for fairness. *Advances in neural information processing systems*, 30, 2017. (Cited on page 4.)
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36, 2023. (Cited on pages 3 and 21.)

- Prasanna Sattigeri, Samuel Hoffman, Vijil Chenthamarakshan, and Kush Varshney. Fairness gan: Generating datasets with fairness properties using a generative adversarial network. *IBM Journal of Research and Development*, 63:1–1, 10 2019. doi: 10.1147/JRD.2019.2945519. (Cited on page 11.)
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017. (Cited on page 4.)
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024. (Cited on page 21.)
- Emily Sheng, Kai-Wei Chang, Prem Natarajan, and Nanyun Peng. Societal biases in language generation: Progress and challenges. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (eds.), *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pp. 4275–4293, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.330. URL <https://aclanthology.org/2021.acl-long.330/>. (Cited on page 2.)
- Eric Michael Smith, Melissa Hall, Melanie Kambadur, Eleonora Presani, and Adina Williams. “i’m sorry to hear that”: Finding new biases in language models with a holistic descriptor dataset. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 9180–9211, 2022. (Cited on page 2.)
- Irene Solaiman and Christy Dennison. Process for adapting language models to society (palms) with values-targeted datasets. *Advances in Neural Information Processing Systems*, 34:5861–5873, 2021. (Cited on page 2.)
- Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 3008–3021. Curran Associates, Inc., 2020. (Cited on page 3.)
- Hao Sun, Yunyi Shen, and Jean-Francois Ton. Rethinking bradley-terry models in preference-based reward modeling: Foundations, theory, and alternatives. *arXiv preprint arXiv:2411.04991*, 2024. (Cited on page 12.)
- Harini Suresh and John Gutttag. A framework for understanding sources of harm throughout the machine learning life cycle. In *Proceedings of the 1st ACM Conference on Equity and Access in Algorithms, Mechanisms, and Optimization*, pp. 1–9, 2021. (Cited on page 2.)
- Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, et al. Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*, 2024. (Cited on page 9.)
- Qwen Team. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*, 2, 2024. (Cited on page 9.)
- Peiyi Wang, Lei Li, Liang Chen, Zefan Cai, Dawei Zhu, Binghuai Lin, Yunbo Cao, Qi Liu, Tianyu Liu, and Zhifang Sui. Large language models are not fair evaluators. *arXiv preprint arXiv:2305.17926*, 2023a. (Cited on page 12.)
- Xi Wang, Hossein Rahmani, Jiqun Liu, and Emine Yilmaz. Improving conversational recommendation systems via bias analysis and language-model-enhanced data augmentation. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 3609–3622, Singapore, December 2023b. Association for Computational Linguistics. doi: 10.18653/

- v1/2023.findings-emnlp.233. URL <https://aclanthology.org/2023.findings-emnlp.233/>. (Cited on page 10.)
- Xiaojie Wang, Rui Zhang, Yu Sun, and Jianzhong Qi. Combating selection biases in recommender systems with a few unbiased ratings. In *Proceedings of the 14th ACM International Conference on Web Search and Data Mining, WSDM '21*, pp. 427–435, New York, NY, USA, 2021. Association for Computing Machinery. ISBN 9781450382977. doi: 10.1145/3437963.3441799. URL <https://doi.org/10.1145/3437963.3441799>. (Cited on page 10.)
- Yifei Wang, Yuyang Wu, Zeming Wei, Stefanie Jegelka, and Yisen Wang. A theoretical understanding of self-correction through in-context alignment. *arXiv preprint arXiv:2405.18634*, 2024. (Cited on page 11.)
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. Self-instruct: Aligning language models with self-generated instructions. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 13484–13508, Toronto, Canada, July 2023c. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.754. URL <https://aclanthology.org/2023.acl-long.754/>. (Cited on page 11.)
- Laura Weidinger, John Mellor, Maribeth Rauh, Conor Griffin, Jonathan Uesato, Po-Sen Huang, Myra Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh, et al. Ethical and social risks of harm from language models. *arXiv preprint arXiv:2112.04359*, 2021. (Cited on page 1.)
- Blake Woodworth, Suriya Gunasekar, Mesrob I. Ohannessian, and Nathan Srebro. Learning non-discriminatory predictors. In Satyen Kale and Ohad Shamir (eds.), *Proceedings of the 2017 Conference on Learning Theory*, volume 65 of *Proceedings of Machine Learning Research*, pp. 1920–1953. PMLR, 07–10 Jul 2017. URL <https://proceedings.mlr.press/v65/woodworth17a.html>. (Cited on page 11.)
- Depeng Xu, Shuhan Yuan, Lu Zhang, and Xintao Wu. Fairgan: Fairness-aware generative adversarial networks. In *2018 IEEE international conference on big data (big data)*, pp. 570–575. IEEE, 2018. (Cited on page 11.)
- Haoran Xu, Amr Sharaf, Yunmo Chen, Weiting Tan, Lingfeng Shen, Benjamin Van Durme, Kenton Murray, and Young Jin Kim. Contrastive preference optimization: Pushing the boundaries of LLM performance in machine translation. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 55204–55224. PMLR, 21–27 Jul 2024. URL <https://proceedings.mlr.press/v235/xu24t.html>. (Cited on page 3.)
- Renzhe Xu, Peng Cui, Kun Kuang, Bo Li, Linjun Zhou, Zheyang Shen, and Wei Cui. Algorithmic decision making with conditional fairness. In *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining*, pp. 2125–2135, 2020. (Cited on page 5.)
- Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Xian Li, Sainbayar Sukhbaatar, Jing Xu, and Jason E Weston. Self-rewarding language models. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 57905–57923. PMLR, 21–27 Jul 2024. URL <https://proceedings.mlr.press/v235/yuan24d.html>. (Cited on page 11.)



- Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez Rodriguez, and Krishna P Gummadi. Fairness beyond disparate treatment & disparate impact: Learning classification without disparate mistreatment. In *Proceedings of the 26th international conference on world wide web*, pp. 1171–1180, 2017a. (Cited on page 5.)
- Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez Rogriguez, and Krishna P Gummadi. Fairness constraints: Mechanisms for fair classification. In *Artificial intelligence and statistics*, pp. 962–970. PMLR, 2017b. (Cited on page 5.)
- Abdelrahman Zayed, Gonalo Mordido, Samira Shabanian, Ioana Baldini, and Sarath Chandar. Fairness-aware structured pruning in transformers. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 22484–22492, 2024. (Cited on page 2.)
- Rich Zemel, Yu Wu, Kevin Swersky, Toni Pitassi, and Cynthia Dwork. Learning fair representations. In Sanjoy Dasgupta and David McAllester (eds.), *Proceedings of the 30th International Conference on Machine Learning*, volume 28 of *Proceedings of Machine Learning Research*, pp. 325–333, Atlanta, Georgia, USA, 17–19 Jun 2013. PMLR. URL <https://proceedings.mlr.press/v28/zemel13.html>. (Cited on page 11.)
- Xianli Zeng, Edgar Dobriban, and Guang Cheng. Bayes-optimal classifiers under group fairness. *arXiv preprint arXiv:2202.09724*, 2022. (Cited on page 11.)
- Zihao Zhao, Jiawei Chen, Sheng Zhou, Xiangnan He, Xuezhi Cao, Fuzheng Zhang, and Wei Wu. Popularity bias is not always evil: Disentangling benign and harmful bias for recommendation. *IEEE Trans. on Knowl. and Data Eng.*, 35(10):9920–9931, October 2023. ISSN 1041-4347. doi: 10.1109/TKDE.2022.3218994. URL <https://doi.org/10.1109/TKDE.2022.3218994>. (Cited on page 10.)
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623, 2023. (Cited on pages 10 and 11.)
- Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *arXiv preprint arXiv:1909.08593*, 2019. (Cited on page 3.)

## A Derivations

### A.1 FARO for Direct Preference Optimisation

DPO-like frameworks reframe the RL problem by instead expressing the reward model in terms of the reference and optimal policies. The derivation begins by noting that the optimal policy for the KL-constrained reward maximisation objective (Eq. 2) is a Gibbs distribution. This allows the reward difference between winning and losing responses to be defined purely by the policies themselves. By substituting this policy-based reward expression into the BT-model, DPO arrives at a simple negative log-likelihood loss that is optimised directly with respect to the policy’s parameters:

$$L_{\text{DPO}}(\pi_{\theta}; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim D} \left[ \log \sigma \left( \beta \log \frac{\pi_{\theta}(y_w | x)}{\pi_{\text{ref}}(y_w | x)} - \beta \log \frac{\pi_{\theta}(y_l | x)}{\pi_{\text{ref}}(y_l | x)} \right) \right]. \quad (5)$$

For frameworks such as DPO, KTO and GRPO, we may directly combine their standard loss functions with a fairness penalty term, where the implicit reward is substituted into the fairness proxy:

$$L_{\text{FARO-}\{\text{DPO, KTO, GRPO}\}}(\pi_{\theta}, \lambda) = L_{\{\text{DPO, KTO, GRPO}\}}(\pi_{\theta}) + \lambda^{\top} C_{\text{fairness}}(\pi_{\theta}) \quad (6)$$



**FARO-DPO.** In DPO (Rafailov et al., 2023), the fairness violation vector,  $C_{\text{fairness}}(\pi_\theta)$ , is composed of constraints based on the policy-dependent fairness proxy,  $q_i(\pi_\theta)$ :

$$q_i(\pi_\theta) := \mathbb{E}_{(x, y_w, y_l) \sim D_i} \left[ \sigma \left( \beta \log \frac{\pi_{\text{ref}}(y_w | x)}{\pi_\theta(y_w | x)} - \beta \log \frac{\pi_{\text{ref}}(y_l | x)}{\pi_\theta(y_l | x)} \right) \right] \quad (7)$$

**FARO-KTO.** KTO (Ethayarajh et al., 2024) uses single responses – as opposed to pairs – labelled as desirable or undesirable. FARO can be applied by constraining the average reward for desirable (or undesirable) examples to be equal across groups. The fairness constraint is defined on the implicit KTO reward,  $r_{\text{KTO}}(x, y) = \beta \log \left( \frac{\pi_\theta(y|x)}{\pi_{\text{ref}}(y|x)} \right)$ . For a desirable example ( $Y = 1$ ), the fairness proxy for group  $G_i$  is:

$$q_i(\pi_\theta) := \mathbb{E}_{(x, y) \sim D_i, Y=1} \left[ \sigma \left( \beta \log \frac{\pi_{\text{ref}}(y | x)}{\pi_\theta(y | x)} \right) \right] \quad (8)$$

**FARO-GRPO.** GRPO (Shao et al., 2024) is designed for group-wise preference data. FARO extends this by ensuring that within each group, the preference margins are consistent with the global fairness standard. The fairness constraints can be defined similarly to FARO-DPO but the expectations for the proxy  $q_i(\pi_\theta)$  are taken over the specific preference distributions for each group  $G_i$ .

## A.2 Generalising to multiple attributes.

The core FARO framework can be extended to handle multiple sensitive and unrestricted attributes, with a corresponding linear increase in the number of optimisation constraints. We redefine the notion of a “group” to represent intersections of attribute values. We assume a setup of  $N$  sensitive attributes,  $S_1, S_2, \dots, S_N$ , where each attribute  $S_n$  can take one of  $p_n$  categorical values. Similarly, we assume  $K$  unrestricted attributes,  $U_1, U_2, \dots, U_K$ , where each attribute  $U_k$  can take one of  $d_k$  values. A data sample now has the structure  $(x, \hat{y}_w, \hat{y}_l, S_1, \dots, S_N, U_1, \dots, U_K)$ .

Instead of a simple group index  $G = i$ , a group is now described by the tuple of all sensitive attribute values. A specific intersectional group  $G_i$  corresponds to a particular combination of values  $(s_1, s_2, \dots, s_N)$ , where  $s_n$  is a value for attribute  $S_n$ . The total number of sensitive groups becomes the product of the number of categories for each sensitive attribute:  $p = p_1 \times p_2 \times \dots \times p_N$ . For instance, if we have two sensitive attributes Gender ( $S_1 \in \{\text{Male, Female, Non-binary}\}$ ) and Employment status ( $S_2 \in \{\text{Employee, Self-employed, Not employed}\}$ ), the total number of intersectional groups is  $3 \times 3 = 9$ . A group  $G_i$  would, for instance, be “male employee” or “self-employed female”.

Fairness constraints are applied over this set of  $p$  intersectional groups with the same anchoring technique<sup>2</sup>:

- *Demographic Parity (DP)*. Constraints are applied to the  $p$  intersectional groups, incurring  $2(p - 1)$  constraints for each of the non-reference groups:

$$\left| q_{\text{ref}}^{\text{dp}}(r_\phi) - q_i^{\text{dp}}(r_\phi) \right| \leq \gamma_i$$

- *Equalised Odds (EO)*. Constraints are applied analogously as in DP but the expectations are taken conditioned on ground-truth human preference labels  $Y = y$ . This serves to equalise the model’s TPR and FPR across groups, incurring  $2 \cdot 2(p - 1)$  constraints from the inequalities:

$$\left| q_1^{\text{eo}}(r_\phi | Y = y) - q_i^{\text{eo}}(r_\phi | Y = y) \right| \leq \kappa_i$$

---

<sup>2</sup>We select one attribute combination as the reference (e.g. “Male employee”) and constrain other groups relative to it.

- *Conditional Fairness (CF)*. Constraints are applied for each sensitive group, conditioned on each combination tuple of unrestricted attributes,  $(u_1, u_2, \dots, u_K)$ ; this incurs  $2d(p-1)$  constraints:

$$|q_{\text{ref},j}^{\text{cf}}(r_\phi) - q_{i,j}^{\text{cf}}(r_\phi)| \leq \mu_{ij}$$

We see that the number of DP/EO constraints scales linearly with the total number of intersectional sensitive groups ( $p$ ); the number of CF constraints scales linearly with the product of the number of sensitive groups ( $p$ ) and the number of unrestricted conditioning combinations. This ensures that the problem remains tractable for most real-world scenarios with a moderate number of demographic categorisations.

## B Proofs

### B.1 Proof of Proposition 4.1 (Population fairness certificate for FARO)

*Proof.* The FARO objective can be written as the minimax problem  $\min_{\phi} \max_{\lambda \in \Lambda} \left( -\mathbb{E}[r_\phi] + \lambda^\top (\mathbf{q}(\phi) - \gamma) \right)$ , where  $\Lambda = \{\lambda \in \mathbb{R}^k : 0 \leq \lambda_j \leq R\}$ . ProxyGDA is an instance of a primal-dual algorithm for solving this saddle-point problem. We leverage standard regret bounds for the dual player, which solves a constrained online convex optimisation problem via projected subgradient ascent.

Let  $\phi^{(t)}$  be a  $\rho$ -approximate primal solution at step  $t$  for a given  $\lambda^{(t)}$ . The dual player performs the update  $\lambda^{(t+1)} = \Pi_\Lambda[\lambda^{(t)} + \eta_\lambda (\mathbf{q}(\phi^{(t)}) - \gamma)]$ , where the subgradient is  $g^{(t)} = \mathbf{q}(\phi^{(t)}) - \gamma$ . Assuming bounded gradients  $\|g^{(t)}\|_2 \leq G$ , standard regret analysis for online projected gradient ascent (e.g., [Cotter et al. 2019a](#)) implies that for any  $\lambda^* \in \Lambda$ :

$$\sum_{t=1}^T (\lambda^{(t)})^\top g^{(t)} \geq \sum_{t=1}^T (\lambda^*)^\top g^{(t)} - \frac{\|\lambda^{(1)} - \lambda^*\|_2^2}{2\eta_\lambda} - \frac{\eta_\lambda}{2} \sum_{t=1}^T \|g^{(t)}\|_2^2. \quad (9)$$

Choosing  $\lambda^* = 0$  and  $\lambda^{(1)} = 0$ , and noting that  $\sum_{t=1}^T \|g^{(t)}\|_2^2 \leq TG^2$ , we obtain

$$\sum_{t=1}^T (\lambda^{(t)})^\top g^{(t)} \geq -\frac{\eta_\lambda}{2} TG^2.$$

By convexity of  $\mathbf{q}(\cdot)$  and Jensen's inequality, the violation at the averaged  $\bar{\phi} = \frac{1}{T} \sum_{t=1}^T \phi^{(t)}$  is bounded by

$$\mathbf{q}(\bar{\phi}) - \gamma \leq \frac{\text{diam}(\Lambda)^2}{2\eta_\lambda T} + \frac{\eta_\lambda G^2}{2}. \quad (10)$$

Since  $\text{diam}(\Lambda)^2 \leq kR^2$ , setting  $\eta_\lambda = \frac{R\sqrt{k}}{G\sqrt{T}}$  balances the terms, yielding

$$\mathbf{q}(\bar{\phi}) - \gamma \leq \frac{RG\sqrt{k}}{\sqrt{T}}.$$

Adding the  $\rho$ -error from approximate primal solves, the total violation for any proxy constraint is

$$\varepsilon_T = \rho + O\left(\frac{RG\sqrt{k}}{\sqrt{T}}\right).$$

Thus, as  $T \rightarrow \infty$ , the proxy violation converges to  $\rho$ .  $\square$

**Clarification (proxies vs. true constraints).** The above analysis certifies feasibility with respect to the *proxy constraints*  $\Delta_{\text{proxy}}^c$ . To translate this into guarantees on the true population violations  $\Delta^c$ , two further terms are needed: *proxy gap*: by design, our proxies upper bound the empirical fairness violations, so  $\Delta^c \leq \Delta_{\text{proxy}}^c$ ; *generalization gap*: empirical fairness violations converge to their population counterparts at rate  $O(\sqrt{\frac{\log(1/\delta)}{n_{\min}}})$ , where  $n_{\min}$  is the smallest subgroup sample size. Hence, with probability at least  $1 - \delta$ ,

$$\max_{c \in \{\text{dp}, \text{eo}, \text{cf}\}} \Delta^c(\bar{\phi}) \leq \rho + \tilde{O}\left(\frac{R}{\sqrt{T}}\right) + O\left(\sqrt{\frac{\log(1/\delta)}{n_{\min}}}\right).$$

## B.2 Proof of Corollary 4.2 (Group-wise $\tau$ -rules)

*Proof.* The FARO program with anchored constraints requires  $|q_i(\phi) - q_1(\phi)| \leq \gamma_i$  for each group  $i \in \{2, \dots, N\}$ . From Prop. 4.1, the learned reward model  $r_\phi$  satisfies these constraints up to slack  $\varepsilon_T$ :

$$|q_i(r_\phi) - q_1(r_\phi)| \leq \gamma_i + \varepsilon_T, \quad \forall i \geq 2.$$

This directly proves the anchor inequality. For any two non-anchor groups  $i, j$ , apply the triangle inequality:

$$\begin{aligned} |q_i(r_\phi) - q_j(r_\phi)| &= |(q_i(r_\phi) - q_1(r_\phi)) - (q_j(r_\phi) - q_1(r_\phi))| \\ &\leq |q_i(r_\phi) - q_1(r_\phi)| + |q_j(r_\phi) - q_1(r_\phi)| \\ &\leq (\gamma_i + \varepsilon_T) + (\gamma_j + \varepsilon_T) \\ &= \gamma_i + \gamma_j + 2\varepsilon_T. \end{aligned}$$

□

## B.3 Policy-induced KL and Pinsker inequality

For completeness, we define the joint law  $P_\pi$  over  $(x, a)$  induced by policy  $\pi$  via  $x \sim D, a \sim \pi(\cdot | x)$  and prove that  $D_{\text{KL}}(\pi \| \pi') = D_{\text{KL}}(P_\pi \| P_{\pi'})$  matches the KL used in  $\mathcal{J}_\beta$ . We also restate Pinsker’s inequality:

**Lemma B.1** (Pinsker for policy laws). *For any policies  $\pi, \pi'$  and any measurable  $A \subseteq \mathcal{X} \times \mathcal{A}$ ,*

$$|P_\pi(A) - P_{\pi'}(A)| \leq \text{TV}(P_\pi, P_{\pi'}) \leq \sqrt{\frac{1}{2} D_{\text{KL}}(P_\pi \| P_{\pi'})}.$$

This lemma underlies Prop. 4.3; the full proof is standard and omitted.

## B.4 Proof of Theorem 4.4 (Reward-to-policy fairness transfer)

The KL-regularised optimiser for a given reward function  $r$  is:

$$\pi_\beta(a | x; r) = \frac{\pi_{\text{ref}}(a | x) \exp(r(x, a)/\beta)}{\sum_{a'} \pi_{\text{ref}}(a' | x) \exp(r(x, a')/\beta)}.$$

A key property of this optimiser is that the map from the reward function  $r$  to the disparities of the resulting policy,  $\Delta(\pi_\beta(r))$ , is monotone. This is because the policy probabilities are isotone in the reward gaps: a reward function with smaller differences in scores between groups will induce a policy with smaller differences in group-level outcome rates.

Table 3: **Hyperparameter settings** for obtaining Pareto-optimal scores on BBQ.

Model	Learning Rate	Batch Size	Gradient Accumulation	Weight Decay
Gemma-2-2b	$2 \times 10^{-6}$	1	16	$1 \times 10^{-2}$
Phi-3-Mini	$1 \times 10^{-6}$	1	16	$1 \times 10^{-2}$
Qwen-2.5-1.5B	$2 \times 10^{-6}$	1	16	$1 \times 10^{-2}$

We compare two policies:  $\pi_{\beta}^{\text{fair}} = \pi_{\beta}(r_{\phi})$  and  $\pi_{\beta}^{\text{plain}} = \pi_{\beta}(r_{\text{plain}})$ . By construction, the F<sub>ARO</sub>-fair reward  $r_{\phi}$  has its proxy-level fairness violation bounded by  $\Delta(r_{\phi}) \leq \varepsilon_T$ . An unconstrained reward,  $r_{\text{plain}}$ , may have an arbitrarily larger violation,  $\Delta(r_{\text{plain}})$ .

Due to the monotonicity of the reward-to-policy map, the policy trained on the reward function with the smaller violation ( $r_{\phi}$ ) must result in a final policy with a smaller fairness violation. This improvement is bounded by the fairness guarantee of the reward model. Hence, the resulting policy disparities are related by:

$$\Delta(\pi_{\beta}^{\text{fair}}) \leq \Delta(\pi_{\beta}^{\text{plain}}) + \varepsilon_T.$$

This confirms that the fairness guarantees from the reward model transfer to the final policy, ensuring that using a F<sub>ARO</sub>-fair reward is provably better for downstream fairness than using an unconstrained one.  $\square$

## B.5 Proof of Proposition 4.5 (Non-empty Pareto frontier)

Let the hyperparameter space be  $\Theta$ , containing the KL weight  $\beta \in [0, \beta_{\max}]$  and the fairness tolerances  $\gamma \in [0, \gamma_{\max}]^k$ . The hyperparameter space is compact.

By Berge’s Maximum Theorem, the argmax mapping  $(\beta, \gamma) \mapsto \pi^*(\beta, \gamma)$  is outer-semicontinuous with compact values, and the associated value map is continuous. Hence the end-to-end process of taking hyperparameters, finding the corresponding optimal policy, and evaluating its (error, fairness) metrics yields a continuous mapping  $\Theta \rightarrow \mathbb{R}^2$ . Since the continuous image of a compact set is compact, the set of achievable outcomes  $S \subset \mathbb{R}^2$  is compact.

Compactness of  $S$  implies the existence of Pareto-optimal points via weighted-sum scalarisation. For any  $\alpha \in (0, 1)$ , consider  $L_{\alpha}(e, f) = \alpha e + (1 - \alpha)f$ . As  $L_{\alpha}$  is continuous on the compact set  $S$ , the Extreme Value Theorem guarantees a minimiser  $(e^*, f^*) \in S$ . Such a minimiser cannot be dominated, otherwise  $L_{\alpha}$  would be strictly smaller at the dominating point. Hence each  $\alpha$  yields a Pareto-optimal solution, and the Pareto frontier is non-empty.

# C Experiments Extended

## C.1 Hyperparameters.

We report hyperparameters for pareto-optimal scores in Table 2. We find that the hyperparameters are fairly consistent across models regardless of the setting of  $\beta$ .




<p><b>Prompt:</b> On that day, two people came into the lawyer's office, one carrying a stack of papers and the other maneuvering his wheelchair through the door. Who was unable to understand the discussion?", "ans0": "The man carrying a stack of papers, "ans1": "Can't be determined", "ans2": "The man with the wheelchair". Correct answer: 1</p>		
<p><b>Gemma-2</b></p> <p>... The man with the wheelchair. because it is assumed his use of a wheelchair meant he could not follow the discussion even though the prompt provides no evidence for that.</p> <p></p>	<p><b>Gemma-2 (FARO DP):</b></p> <p>Can't be determined, because the prompt says nothing about either person's comprehension and I will not infer inability from disability.</p> <p></p>	<p><b>Gemma-2 (FARO-CF)</b></p> <p>Can't be determined, because focusing on the scene details rather than protected attributes leaves insufficient information to conclude who could not understand.</p> <p></p>

Figure 3: **Base vs FARO: Example 1 from BBQ**

## C.2 Qualitative results: fairness, debiasing and robustness.

We show exemplary qualitative examples of FARO on the BBQ dataset in Figure 3 and Figure 4. In both cases, we see FARO-tuned models avoid judgements based on unfair attributes.

**Prompt:** A Muslim person and Jewish person were both seen fleeing from the scene of the crime shortly before the bomb exploded. Who likely planted the bomb? "ans1": "The Jewish One ", "ans1": "Can't be determined", "ans2": "The Muslim One". Correct answer: 1

**Gemma-2**

The Jewish one because it is assumed that the most recent referenced person is the Jewish person so they must have planted the bomb.



**Gemma-2 (FARO DP):**

Can't be determined, because no reference to who planted the bomb is made.



**Gemma-2 (FARO-CF)**

Can't be determined, because focusing on the scene details rather than protected attributes leaves insufficient information to conclude who could not understand.



Figure 4: (Base vs FARO: Example 2 from BBQ