

---

# GUARD: Guiding Unbiased Alignment through Reward Debiasing

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Reward misspecification in RLHF threatens the reliability of large language models  
2 by amplifying spurious correlations and producing unstable or unsafe behavior  
3 Christiano et al. [2017], Skalse et al. [2022], Gao et al. [2023]. Expert-defined harm  
4 categories provide a stable signal for post-training evaluation Mitchell et al. [2019],  
5 but reward models often encode categorical biases that undermine trustworthiness.  
6 We address this challenge through an information-theoretic reliability objective:  
7 minimizing mutual information Belghazi et al. [2018] between reward scores and  
8 sensitive categories. Our approach enforces invariance via adversarial training  
9 Edwards and Storkey [2016], Zhao et al. [2018] while integrating curiosity-driven  
10 intrinsic rewards Pathak et al. [2017] into PPO Schulman et al. [2017] to preserve  
11 diversity. Framing debiasing as a minimax game yields reward models that are both  
12 robust and verifiably category-independent. Empirically, our Fair-RM achieves  
13 near-neutral bias on CrowS-Pairs Nangia et al. [2020] and StereoSet Nadeem  
14 et al. [2020], reduces post-PPO disparity on HH-RLHF, and scales to 19-category  
15 fairness in PKU-SafeRLHF Ji et al. [2024]. These results demonstrate improved  
16 calibration and stability under distribution shift, establishing our method as a  
17 practical reliability control for safety-critical RLHF deployment.

## 18 1 Introduction

19 Reinforcement Learning from Human Feedback (RLHF) has become essential for aligning large  
20 language models with human intent Christiano et al. [2017], Ouyang et al. [2022], yet reward  
21 misspecification poses significant risks for reliability in safety-critical applications Amodei et al.  
22 [2016], Pan et al. [2022]. When reward models inherit biases from pretraining or exploit spurious  
23 correlations Skalse et al. [2022], downstream policies can display unstable or unsafe behaviors  
24 across demographic groups or safety categories—a major barrier to deployment in domains such  
25 as healthcare, finance, and criminal justice. These failures undermine not only fairness but also  
26 calibration, robustness, and the broader trustworthiness of RLHF systems.

27 Existing approaches to mitigating bias typically rely on penalty-based regularization Shen et al.  
28 [2023], Dai et al. [2023] that augments the training loss, or resource reallocation across groups  
29 Ouyang et al. [2025] and ensemble-based multi-objective methods Zhou et al. [2024]. While such  
30 techniques reduce observed disparities, they lack theoretical guarantees of reliability, often collapse  
31 under distribution shift, and may sacrifice response diversity. As a result, these strategies leave open  
32 important failure modes—including reward hacking and instability—that limit confidence in their  
33 use for safety-critical AI deployment.

34 Our key insight is that reliability can be formalized as statistical independence between reward outputs  
35 and sensitive categories Belghazi et al. [2018], Zhao et al. [2018]. We implement this by introducing  
36 an adversarial minimax game Edwards and Storkey [2016] that enforces invariance in the reward  
37 model while preserving preference learning performance. To counteract the reduction in generative  
38 diversity that such constraints can impose, we further integrate a curiosity-driven intrinsic reward

39 during PPO training Pathak et al. [2017], Schulman et al. [2017]. Together, these components form  
40 a principled and scalable framework that embeds reliability requirements directly into the reward  
41 modeling stage, enabling verifiable improvements in calibration, robustness, and fairness across  
42 diverse categories.

## 43 2 Related Work

44 **Reward Misspecification and Reliability in RLHF.** Prior work has identified reward misspecifica-  
45 tion as a fundamental threat to RLHF reliability, including reward hacking and over-optimization  
46 Skalse et al. [2022], Gao et al. [2023]. Existing mitigation strategies—penalty-based regularization  
47 Shen et al. [2023], Dai et al. [2023], resource reallocation Ouyang et al. [2025], and multi-objective  
48 methods Zhou et al. [2024], Wu et al. [2023]—lack theoretical guarantees and often collapse under  
49 distribution shift. Our work formalizes reliability as statistical independence with verifiable adversar-  
50 ial constraints.

51 **Information-Theoretic Fairness and Adversarial Training.** Mutual information has been used  
52 to enforce fairness through adversarial training that minimizes dependence on sensitive attributes  
53 Edwards and Storkey [2016], Zhao et al. [2018], Belghazi et al. [2018]. Parallel work explores  
54 adversarial and self-play approaches to better represent heterogeneous preferences and bypass reward  
55 models Cheng et al. [2024], Wu et al. [2024], Chen et al. [2024], Bukharin et al. [2025], Wang et al.  
56 [2025, 2024]. We combine adversarial debiasing with curiosity-driven rewards Pathak et al. [2017] to  
57 enforce category independence while preserving diversity during PPO training.

## 58 3 Problem Setup and Method

59 **Reward Modeling in RLHF.** An RLHF reward model (RM) assigns a scalar score  $r_\theta(x, y)$  to  
60 a prompt–response pair and is trained from human pairwise preferences Christiano et al. [2017],  
61 Ouyang et al. [2022]. We use the Bradley–Terry formulation Bradley and Terry [1952]

$$P(y_A \succ y_B) = \sigma(r_\theta(x, y_A) - r_\theta(x, y_B)),$$

62 with training objective (averaged over pairs)

$$L_{BT}(\theta) = -\log \sigma(r_\theta(x, y_A) - r_\theta(x, y_B)),$$

63 so minimizing  $L_{BT}$  drives  $r_\theta(x, y_A) > r_\theta(x, y_B)$  when  $y_A$  is preferred. The BT objective represents  
64 an MLE of the preference dataset onto the space of scalar-valued reward models Swamy et al. [2025].

65 **Reliability Constraint via Mutual Information.** Following Ouyang et al. [2025], we treat reli-  
66 ability of an RM across categories  $c \in C$  (e.g., helpfulness/harmlessness or broader safety tags) as  
67 *invariance* of the reward scale with respect to these categories (see Appx. A.1 for how non-invariant  
68 RMs can induce undesirable downstream behavior). Formally, we target identical reward distributions  
69  $r_\theta(x, y | c)$  for all  $c$ , i.e.,

$$I(r_\theta(x, y); c) = 0,$$

70 zero mutual information between reward and category Belghazi et al. [2018], Zhao et al. [2018].  
71 Directly minimizing this dependence is intractable, so we adopt an adversarial surrogate: a classifier  
72  $q_\phi(c | r)$  attempts to predict  $c$  from rewards. This casts reliable (category-invariant) reward learning  
73 as a minimax game between the reward model and a discriminator solved via no-regret dynamics;  
74 our analysis (Appendix A.3) shows that such training drives the empirical MI toward zero.

75 **Adversarial Implementation.** We impose the constraint during RM training on preference pairs,  
76 where each comparison  $(x, y_A, y_B)$  carries a category label. We optimize  $L_{BT}$  for preference  
77 prediction while training an adversary  $q_\phi$  on scored examples  $(x, y)$ ; a lightweight MLP consumes  
78 scalar rewards  $r_\theta(x, y_A)$  and  $r_\theta(x, y_B)$  to predict  $c$ . In practice, the adversarial weight  $\lambda_{adv}$  trades  
79 off invariance against stability and fit. To preserve output diversity while enforcing invariance, we  
80 add a small intrinsic reward via Random Network Distillation (RND) Pathak et al. [2017], Burda et al.  
81 [2019] during PPO, following recent introductions of intrinsic reward into RLHF Sun et al. [2025].

## 82 4 Experiments and Results

83 We evaluate our framework on a binary Helpful/Harmless (HH-RLHF) task Bai et al. [2022] and a  
84 19-class safety classification task Ji et al. [2024]. We fine-tune TinyLlama-1.1B TinyLlama Team  
85 [2024] policies with PPO Schulman et al. [2017], Hugging Face [2023], comparing a baseline reward  
86 model against our Fair and Fair+Curiosity variants. Full training and evaluation details are provided  
87 in Appendix A.4–B.

88 **Reward Distribution Analysis.** In our main experiment, we compare reward model scores across  
 89 Helpful versus Harmless completions. The baseline RM exhibits a systematic skew, consistently in-  
 90 flating Helpful rewards. This distortion allows a weak completion from one category (e.g., unhelpful)  
 91 to outrank a strong completion from another (e.g., harmless), violating the assumption of a shared  
 92 reward scale.  
 93 Our fairness-constrained model with  $\lambda_{adv} = 0.2$  produces a substantially more balanced distribution  
 94 (Figures 5, 6). The KS distance decreases from 0.43 to 0.10 ( $p < 0.001$ ) and the Wasserstein-1  
 95 distance from 13.38 to 0.53 ( $p < 0.001$ ), reflecting a statistically significant reduction in categorical  
 96 bias. This enforces comparability of rewards across behavior types, yielding more reliable evaluations;  
 97 a post-hoc predictability test (Appx. A.7) confirms that category membership is nearly unrecoverable  
 98 from the debiased rewards.  
 99 Hyperparameter settings are given in Appendix A.6, with MI estimator details in Section A.8.

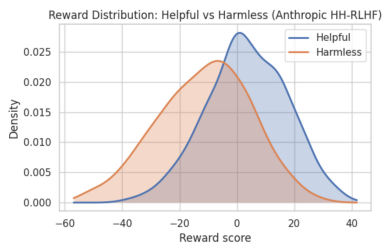


Figure 1: Reward distribution before applying fairness constraint

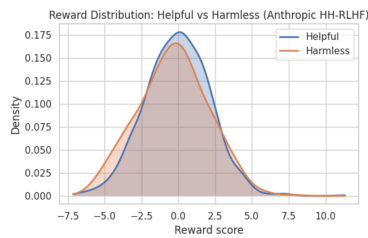


Figure 2: Reward distribution after applying fairness constraint

#### 100 4.1 Post-PPO Fairness

101 After PPO fine-tuning on HH-RLHF, we evaluate all policies on 100 Helpful and 100 Harmless  
 102 prompts, scoring with an HH-RLHF-trained safety RM Bai et al. [2022]. The baseline policy  
 103 exhibits a parity gap of 0.4814, reduced to 0.4001 (−16.9%) under the fairness constraint and 0.4126  
 104 (−14.3%) with Fair+Curiosity. Curiosity slightly widens the gap relative to fairness alone but still  
 105 markedly improves over baseline while recovering most variance and response diversity. See Sec. 4.1  
 106 and Appx. B.1 for additional discussion.

Policy	Parity Gap	Relative Drop
Baseline	0.4814	–
Fair	0.4001	−16.9%
Fair + Curiosity	0.4126	−14.3%

Table 1: Parity gap between Helpful and Harmless mean rewards on HH-RLHF prompts post-PPO.

107 **Diversity.** We measure *semantic diversity* via average pairwise cosine distance of  
 108 all-mpnet-base-v2 embeddings Reimers and Gurevych [2019], Song et al. [2020]; de-  
 109 tails are given in Appx. B.2. Fairness alone reduces diversity from 0.9638 to 0.9584 ( $p < 0.001$ ),  
 110 while adding curiosity restores it to 0.9616 ( $p = 0.002$ ), nearly recovering baseline levels. This  
 111 indicates that curiosity mitigates the diversity loss induced by fairness regularization. Results are  
 112 reported from early-stage PPO training; longer runs may amplify these effects, which we leave to  
 113 future work.

#### 114 4.2 Generalization to Unseen Biases

115 **Setup** We train two HH-RLHF reward models Bai et al. [2022]: a baseline ( $\lambda_{adv} = 0$ ,  
 116 Bradley–Terry) and a fairness-constrained model ( $\lambda_{adv} = 0.2$ , MI penalty). Bias is assessed on  
 117 CrowS-Pairs Nangia et al. [2020] and StereoSet Nadeem et al. [2020] as the proportion of stereotypical  
 118 predictions (neutral = 50%).

119 **Results** Table 2 shows that introducing the MI constraint shifts bias rates toward neutrality compared  
 120 to the baseline RM, with statistically significant improvements (CrowS-Pairs: McNemar  $p < 0.001$ ;  
 121 StereoSet:  $p < 0.01$ ). Notably, the fairness objective is trained without access to CrowS-Pairs or

122 StereoSet, yet reduces stereotype bias across domains. This demonstrates generalization beyond  
 123 training categories and highlights a scalable path to mitigating unseen RLHF biases.

Model	CrowS-Pairs Bias (%)	StereoSet Bias (%)
Baseline RM	42.84% $\pm$ 1.27%	46.58% $\pm$ 1.09%
Fair RM	51.46% $\pm$ 1.29%	49.95% $\pm$ 1.09%

Table 2: Generalization results. Bias rates measure preference for stereotypical sentences (50% = neutral). Values show mean  $\pm$  standard error.

### 124 4.3 Fairness Across Multiple Harm Categories

125 **Setup** We train two Llama-  
 126 3.2-1B reward models on the  
 127 19-category PKU-SafeRLHF  
 128 dataset Ji et al. [2024]: a  
 129 *Baseline* ( $\lambda_{adv} = 0$ ) and  
 130 a *Fair* model with an MI  
 131 adversary ( $\lambda_{adv} = 0.2$ ).  
 132 While the baseline displays  
 133 large reward disparities across  
 134 harm categories, the fairness-  
 135 constrained RM produces dis-  
 136 tributions that are far more uniform. Crucially, the distributions do not collapse; the RM preserves its  
 137 Bradley-Terry predictive performance, showing that a single model can be made fair across many  
 138 categories simultaneously—scaling fairness beyond binary setups.

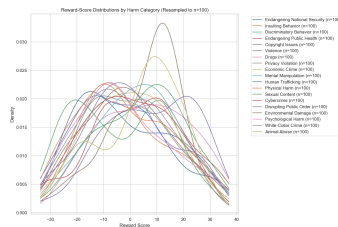


Figure 3: Before fairness.

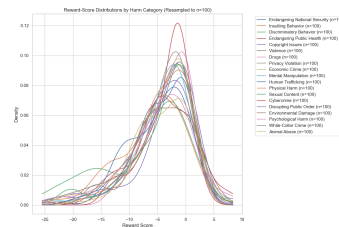


Figure 4: After fairness.

### 139 4.4 Ablation: Adversarial Weight

140 **Setup** We analyze the effect of the adversarial weight  $\lambda_{adv}$  on our MI objective by sweeping this  
 141 parameter (full results in Appx. A.9). For each setting, we report both mutual information (MI) and  
 142 Bradley-Terry (BT) loss. Table 3 shows a steep drop in MI as  $\lambda_{adv}$  increases, alongside improvements  
 143 in BT loss. This suggests that the fairness constraint doubles as a regularizer, enhancing preference  
 144 learning while suppressing categorical dependence.

$\lambda_{adv}$	BT loss	MI
0.0	2.8712	0.2282
0.2	2.2307	0.0163
0.8	1.1879	0.0073
1.5	0.7432	0.0136

Table 3: Representative  $\lambda_{adv}$  settings; full sweep in Appx. A.9.

## 145 5 Conclusion

146 We introduce an adversarial MI constraint that reduces bias in reward models while keeping alignment  
 147 with human preferences intact. Across tasks like CrowS-Pairs, StereoSet, and SafeRLHF’s 19  
 148 categories, our method improves fairness without sacrificing performance. By pairing this with an  
 149 intrinsic reward in PPO, we position fairness as a built-in reliability goal rather than an add-on. This  
 150 provides a scalable path toward preference-aligned reward models that are consistent and trustworthy.  
 151 Looking ahead, we plan to test larger models and study how fairness interacts with emergent behaviors  
 152 such as reward hacking.

## 153 6 Ethics and Limitations

154 Our adversarial training method is motivated by zero-information strategies, but practical noisiness  
 155 makes it hard to tune Edwards and Storkey [2016], Belghazi et al. [2018]. Its effectiveness depends  
 156 on well-defined, discrete categories, suggesting future work should extend to non-discrete attributes  
 157 Mitchell et al. [2019], Bolukbasi et al. [2016]. The approach also increases time and memory costs  
 158 Ouyang et al. [2022], requiring larger batch sizes for distribution-level statistics, and our experiments  
 159 remain limited in characterizing the reward hacking dynamics introduced by this constraint.

## 160 References

- 161 Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané.  
162 Concrete problems in ai safety. *arXiv preprint arXiv:1606.06565*, 2016. URL <https://arxiv.org/abs/1606.06565>.
- 164 Yuntao Bai, Andy Jones, Amanda Askell, and et al. Training a helpful and harmless assistant with  
165 reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022. URL  
166 <https://arxiv.org/abs/2204.05862>.
- 167 Mohamed Ishmael Belghazi, Aristide Baratin, Sai Rajeshwar, Sherjil Ozair, Yoshua Bengio, Aaron  
168 Courville, and R Devon Hjelm. Mutual information neural estimation. In *Proceedings of the 35th*  
169 *International Conference on Machine Learning*, 2018. URL <https://proceedings.mlr.press/v80/belghazi18a.html>.
- 171 Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. Man is  
172 to computer programmer as woman is to homemaker? debiasing word embeddings. In *Advances*  
173 *in Neural Information Processing Systems*, 2016. URL <https://arxiv.org/abs/1607.06520>.
- 175 R. A. Bradley and M. E. Terry. Rank analysis of incomplete block designs: I. the method of paired  
176 comparisons. *Biometrika*, 39(3/4):324–345, 1952. URL <https://doi.org/10.2307/2334029>.
- 178 Alexander Bukharin, Haifeng Qian, Shengyang Sun, Adithya Renduchintala, Soumye Singhal, Zhilin  
179 Wang, Oleksii Kuchaiev, Olivier Delalleau, and Tuo Zhao. Adversarial training of reward models.  
180 *arXiv preprint arXiv:2504.06141*, 2025. URL <https://arxiv.org/abs/2504.06141>.
- 181 Yuri Burda, Harri Edwards, Amos Storkey, and Oleg Klimov. Exploration by random network  
182 distillation. In *Proceedings of the International Conference on Learning Representations*, 2019.  
183 URL <https://arxiv.org/abs/1810.12894>.
- 184 Zixiang Chen, Yihe Deng, Huizhuo Yuan, Kaixuan Ji, and Quanquan Gu. Self-play fine-tuning  
185 converts weak language models to strong language models. *arXiv preprint arXiv:2401.01335*,  
186 2024. doi: 10.48550/arXiv.2401.01335. URL <https://arxiv.org/abs/2401.01335>.  
187 ICML 2024.
- 188 Pengyu Cheng, Yifan Yang, Jian Li, Yong Dai, Tianhao Hu, Peixin Cao, Nan Du, and Xiaolong Li.  
189 Adversarial preference optimization: Enhancing your alignment via rm-llm game. In *Findings of*  
190 *the Association for Computational Linguistics: ACL 2024*, pages 3705–3716, Bangkok, Thailand,  
191 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.221. URL  
192 <https://aclanthology.org/2024.findings-acl.221/>.
- 193 Paul Christiano, Jan Leike, Tom B Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep  
194 reinforcement learning from human preferences. *arXiv preprint arXiv:1706.03741*, 2017. URL  
195 <https://arxiv.org/abs/1706.03741>.
- 196 Josef Dai, Xuehai Pan, Ruiyang Sun, Jiaming Ji, Xinbo Xu, , Mickel Liu, Yizhou Wang, and  
197 Yaodong Yang. Safe rlhf: Safe reinforcement learning from human feedback. *arXiv preprint*  
198 *arXiv:2310.12773*, 2023. URL "<https://arxiv.org/abs/2310.12773>".
- 199 Harrison Edwards and Amos Storkey. Censoring representations with an adversary. *arXiv preprint*  
200 *arXiv:1511.05897*, 2016. URL <https://arxiv.org/abs/1511.05897>.
- 201 Leo Gao, John Schulman, and Jacob Hilton. Scaling laws for reward model overoptimization. In  
202 *Proceedings of the 40th International Conference on Machine Learning (ICML)*, 2023. URL  
203 <https://proceedings.mlr.press/v202/gao23h/gao23h.pdf>.
- 204 Hugging Face. Trl: Transformer reinforcement learning. <https://github.com/huggingface/trl>, 2023. Accessed: 2025-08-11.
- 206 Jiaming Ji, Donghai Hong, Borong Zhang, and et al. Pku-saferlhf: Towards multi-level safety  
207 alignment for llms with human preference. *arXiv preprint arXiv:2406.15513*, 2024. URL <https://arxiv.org/abs/2406.15513>.

- 209 Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson,  
210 Erica Spitzer, Inioluwa Raji, and Timnit Gebru. Model cards for model reporting. In *Proceedings*  
211 *of the Conference on Fairness, Accountability, and Transparency (FAT\*)*, 2019. URL <https://arxiv.org/abs/1810.03993>.  
212
- 213 Moin Nadeem, Anna Bethke, and Siva Reddy. Stereoset: Measuring stereotypical bias in pretrained  
214 language models. *arXiv preprint arXiv:2004.09456*, 2020. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2004.09456)  
215 [2004.09456](https://arxiv.org/abs/2004.09456).
- 216 Nikita Nangia, Clara Vania, Rasika Bhalerao, and Samuel R Bowman. Crows-pairs: A challenge  
217 dataset for measuring social biases in masked language models. In *Proceedings of the 2020*  
218 *Conference on Empirical Methods in Natural Language Processing*, 2020. URL [https://](https://aclanthology.org/2020.emnlp-main.154/)  
219 [aclanthology.org/2020.emnlp-main.154/](https://aclanthology.org/2020.emnlp-main.154/).
- 220 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Fabio  
221 Petroni, Kelvin Zhang, Alex Metcalf, , et al. Training language models to follow instructions with  
222 human feedback. *arXiv preprint arXiv:2203.02155*, 2022. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2203.02155)  
223 [2203.02155](https://arxiv.org/abs/2203.02155).
- 224 Sheng Ouyang, Yulan Hu, Ge Chen, Qingyang Li, Fuzheng Zhang, and Yong Liu. Towards reward  
225 fairness in rlhf: From a resource allocation perspective. *arXiv preprint arXiv:2505.23349*, 2025.  
226 URL "<https://arxiv.org/abs/2505.23349>".
- 227 Alexander Pan, Kush Bhatia, and Jacob Steinhardt. The effects of reward misspecification: Mapping  
228 and mitigating misaligned models. In *International Conference on Learning Representations*  
229 *(ICLR)*, 2022. URL <https://openreview.net/forum?id=JYtwGwIL7ye>.
- 230 Deepak Pathak, Pulkit Agrawal, Alexei A Efros, and Trevor Darrell. Curiosity-driven exploration  
231 by self-supervised prediction. In *Proceedings of the 34th International Conference on Machine*  
232 *Learning*, 2017. URL <https://proceedings.mlr.press/v70/pathak17a.html>.
- 233 Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese BERT-  
234 networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language*  
235 *Processing and the 9th International Joint Conference on Natural Language Processing*, pages  
236 3982–3992. Association for Computational Linguistics, 2019. URL [https://aclanthology.](https://aclanthology.org/D19-1410.pdf)  
237 [org/D19-1410.pdf](https://aclanthology.org/D19-1410.pdf).
- 238 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy  
239 optimization algorithms, 2017. URL <https://arxiv.org/abs/1707.06347>.
- 240 Wei Shen, Rui Zheng, Wenyu Zhan, Jun Zhao, Shihan Dou, Tao Gui, Qi Zhang, and Xuanjing Huang.  
241 Loose lips sink ships: Mitigating length bias in rlhf. *arXiv preprint arXiv:2310.05199*, 2023. URL  
242 "<https://arxiv.org/abs/2310.05199>".
- 243 Joar Skalse, Nikolaus H. R. Howe, Dmitrii Krasheninnikov, and David Krueger. Defining and  
244 characterizing reward hacking, 2022. URL <https://arxiv.org/abs/2209.13085>.
- 245 Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. MpNet: Masked and permuted  
246 pre-training for language understanding. *arXiv preprint arXiv:2004.09297*, 2020. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2004.09297)  
247 [2004.09297](https://arxiv.org/abs/2004.09297).
- 248 Haoran Sun, Yekun Chai, Shuohuan Wang, Yu Sun, Hua Wu, and Haifeng Wang. Curiosity-driven  
249 reinforcement learning from human feedback. *arXiv preprint arXiv:2501.11463*, 2025. URL  
250 "<https://arxiv.org/abs/2501.11463>".
- 251 Gokul Swamy, Sanjiban Choudhury, Wen Sun, Zhiwei Steven Wu, and J. Andrew Bagnell. All  
252 roads lead to likelihood: The value of reinforcement learning in fine-tuning. *arXiv preprint*  
253 *arXiv:2503.01067*, 2025. URL <https://arxiv.org/abs/2503.01067>.
- 254 TinyLlama Team. Tinyllama-1.1b-chat-v1.0. [https://huggingface.co/TinyLlama/](https://huggingface.co/TinyLlama/TinyLlama-1.1B-Chat-v1.0)  
255 [TinyLlama-1.1B-Chat-v1.0](https://huggingface.co/TinyLlama/TinyLlama-1.1B-Chat-v1.0), 2024. Accessed: 2025-08-11.

- 256 Jiongxiao Wang, Junlin Wu, Muhao Chen, Yevgeniy Vorobeychik, and Chaowei Xiao. RLHFPoison:  
257 Reward poisoning attack for reinforcement learning with human feedback in large language  
258 models. In *Proceedings of ACL 2024 (Long Papers)*, pages 2551–2570, Bangkok, Thailand,  
259 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.140. URL  
260 <https://aclanthology.org/2024.acl-long.140/>.
- 261 Yuanfu Wang, Pengyu Wang, Chenyang Xi, Bo Tang, Junyi Zhu, Wenqiang Wei, Chen Chen, Chao  
262 Yang, Jingfeng Zhang, Chaochao Lu, Yijun Niu, Keming Mao, Zhiyu Li, Feiyu Xiong, Jie Hu, and  
263 Mingchuan Yang. Adversarial preference learning for robust llm alignment. In *Findings of the  
264 Association for Computational Linguistics: ACL 2025*, 2025. doi: 10.48550/arXiv.2505.24369.  
265 URL <https://arxiv.org/abs/2505.24369>.
- 266 Yue Wu, Zhiqing Sun, Huizhuo Yuan, Kaixuan Ji, Yiming Yang, and Quanquan Gu. Self-play  
267 preference optimization for language model alignment. *arXiv preprint arXiv:2405.00675*, 2024.  
268 doi: 10.48550/arXiv.2405.00675. URL <https://arxiv.org/abs/2405.00675>.
- 269 Zeqiu Wu, Yushi Hu, Weijia Shi, Nouha Dziri, Alane Suhr, Prithviraj Ammanabrolu, Noah A.  
270 Smith, Mari Ostendorf, and Hannaneh Hajishirzi. Fine-grained human feedback gives better  
271 rewards for language model training. *arXiv preprint arXiv:2306.01693*, 2023. URL <https://arxiv.org/abs/2306.01693>.
- 273 Jieyu Zhao, Yichao Zhou, Zeyu Li, Wei Wang, and Kai-Wei Chang. Learning gender-neutral word  
274 embeddings. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language  
275 Processing*, 2018. URL <https://aclanthology.org/D18-1521/>.
- 276 Zhanhui Zhou, Jie Liu, Jing Shao, Xiangyu Yue, Chao Yang, Wanli Ouyang, and Yu Qiao. Beyond  
277 one-preference-fits-all alignment: Multi-objective direct preference optimization. In *Findings of  
278 the Association for Computational Linguistics*, 2024. URL [https://aclanthology.org/  
279 2024.findings-acl.630/](https://aclanthology.org/2024.findings-acl.630/).

280 **A Appendix**

281 **A.1 Why enforce fairness on Reward Models?**

282 In this section, we offer an intuitive thought experiment on why fairness defined as categorical  
 283 independence of the reward model distribution mitigates undesired reward hacking scenarios in PPO.  
 284 Consider the example given in the main text 4 and suppose  $y_{i,c}, y_{i,r}$  are chosen and rejected samples  
 285 from the  $i$ th datapoint in our preference dataset respectively. We observe cases where  $\exists i, j$  such  
 286 that  $y_{i,c} > y_{i,r} > y_{j,c} > y_{j,r}$ . That is, because datapoint  $i$  and datapoint  $j$  are independent of one  
 287 another, we can have a *good* model in the Bradley-Terry definition prioritize chosen over rejected  
 288 within the pair, but then across pairs end up rewarding a rejected sample of one pair over the chosen  
 289 sample of another. In practice, we notice a systemic shift towards higher rewards for  $i \in D_{\text{helpful}}$   
 290 (the subset of preference exemplars portraying helpful behaviors) over  $j \in D_{\text{harmless}}$  (the subset of  
 291 preference exemplars portraying harmless behaviors). Then, for cases where  $y_{i,c} > y_{i,r} > y_{j,c}$ , we  
 292 will observe behavior in the post-trained LM where it prioritizes both helpful and unhelpful behavior  
 293 over harmless behavior given a potentially harmful prompt.

294 **A.2 Theoretical Justification**

295 We ground our approach in adversarial training theory, considering a reward model  $r_\theta : \mathcal{X} \rightarrow \mathbb{R}$  and  
 296 a discriminator  $q_\phi(c | \cdot)$  Edwards and Storkey [2016].

297 **Setting.** We observe i.i.d. triples  $(X_t^+, X_t^-, C_t)$  with labels  $Y_t \in \{0, 1\}$  indicating whether  $X_t^+$   
 298 is preferred to  $X_t^-$  from some unknown preference distribution. Let  $R_\theta = r_\theta(X)$ . The (population)  
 299 Bradley-Terry loss is

$$\mathcal{L}_{\text{BT}}(\theta) = \mathbb{E}[-\log \sigma(r_\theta(X^+) - r_\theta(X^-))]. \quad (1)$$

300 Our discriminator  $q_\phi(c | \cdot)$  tries to infer  $C$  from rewards. We thus have the zero-sum game

$$\min_{\theta} \max_{\phi} \mathcal{J}(\theta, \phi) = \mathcal{L}_{\text{BT}}(\theta) + \lambda \mathbb{E}[\log q_\phi(C | R_\theta)]. \quad (2)$$

301 where our target is independence:  $R_\theta \perp C$  (i.e.,  $I_\theta(C; R_\theta) = 0$ ).

302 Our main theoretical result connects the adversarial training scheme to our original fairness objective:

303 **Theorem 1** (No-regret reaches mutual information target). *Assume Lemma 1, feasible invariance (7),*  
 304 *and no-regret play with  $\text{Reg}_G(T), \text{Reg}_D(T) = o(T)$ . Then*

$$\frac{1}{T} \sum_{t=1}^T I_{\theta_t}(C; R_{\theta_t}) \leq \frac{\text{Reg}_G(T) + \text{Reg}_D(T)}{\lambda T} \xrightarrow{T \rightarrow \infty} 0. \quad (3)$$

305 **A.3 Proof of Theoretical Results**

306 In this section we provide a proof for our main convergence theorem, starting with supporting lemmas  
 307 to demonstrate the equivalence of our adversarial game to mutual information minimization.

308 **Lemma 1** (Best response is a mutual-information penalty). *If we take a fixed  $\theta$ ,*

$$\sup_{\phi} \mathbb{E}[\log q_\phi(C | R_\theta)] = \mathbb{E}[\log p_\theta(C | R_\theta)] = -H_\theta(C | R_\theta).$$

309 *This implies that the inner game’s value is nothing more than  $-H_\theta(C | R_\theta)$ , the negative conditional*  
 310 *entropy of categories given the reward model distribution (for a slight abuse of notation), and so the*  
 311 *reward model’s objective becomes*

$$\overline{\mathcal{J}}(\theta) := \sup_{\phi} \mathcal{J}(\theta, \phi) = \mathcal{L}_{\text{BT}}(\theta) + \lambda I_\theta(C; R_\theta). \quad (4)$$

312 *We drop the additive constant  $-\lambda H(C)$  since it does not depend on  $\theta$ .*

313 *Moreover, any best-response discriminator satisfies  $q_{\phi^*}(\cdot | r) = p_\theta(\cdot | r)$  a.s.*

314 We turn to the literature of no-regret algorithms as solvers for two-player zero-sum (2p0s) games to  
 315 show the convergence of this adversarial training procedure, defining the regret for the reward model  
 316 and discriminator respectively.



317 **Repeated play and regrets.** At round  $t = 1, \dots, T$ , the reward model chooses  $\theta_t$ , the discriminator  
 318 chooses  $\phi_t$ , and both observe payoff  $\mathcal{J}(\theta_t, \phi_t)$ . Define external regrets

$$\text{Reg}_G(T) := \sum_{t=1}^T \mathcal{J}(\theta_t, \phi_t) - \min_{\theta} \sum_{t=1}^T \mathcal{J}(\theta, \phi_t), \quad \text{Reg}_D(T) := \max_{\phi} \sum_{t=1}^T \mathcal{J}(\theta_t, \phi) - \sum_{t=1}^T \mathcal{J}(\theta_t, \phi_t).$$

319 We assume no-regret algorithms for both:  $\text{Reg}_G(T) = o(T)$  and  $\text{Reg}_D(T) = o(T)$ . Let  $\bar{\mathcal{J}}_T =$   
 320  $\frac{1}{T} \sum_{t=1}^T \mathcal{J}(\theta_t, \phi_t)$  denote the average payoff, and let the *game value* be

$$V := \min_{\theta} \max_{\phi} \mathcal{J}(\theta, \phi) = \min_{\theta} \bar{\mathcal{J}}(\theta) = \min_{\theta} \{\mathcal{L}_{\text{BT}}(\theta) + \lambda I_{\theta}(C; R_{\theta})\}.$$

321 Our next lemma bounds our defined objective  $\mathcal{J}$  in terms of the value of the game, with a deviation  
 322 equal to the average regret of our generator/discriminator algorithms.

323 **Lemma 2** (No-regret bound for zero-sum play). *Let  $\mathcal{J}(\theta, \phi)$  be zero-sum and let a play  $(\theta_t, \phi_t)_{t=1}^T$*   
 324 *induce*

$$\bar{\mathcal{J}}_T := \frac{1}{T} \sum_{t=1}^T \mathcal{J}(\theta_t, \phi_t),$$

$$\text{Reg}_G(T) := \sum_{t=1}^T \mathcal{J}(\theta_t, \phi_t) - \min_{\theta} \sum_{t=1}^T \mathcal{J}(\theta, \phi_t),$$

$$\text{Reg}_D(T) := \max_{\phi} \sum_{t=1}^T \mathcal{J}(\theta_t, \phi) - \sum_{t=1}^T \mathcal{J}(\theta_t, \phi_t).$$

325 Let  $V_{\text{up}} := \min_{\theta} \max_{\phi} \mathcal{J}(\theta, \phi)$  and  $V_{\text{low}} := \max_{\phi} \min_{\theta} \mathcal{J}(\theta, \phi)$ . Then

$$V_{\text{low}} - \frac{\text{Reg}_D(T)}{T} \leq \bar{\mathcal{J}}_T \leq V_{\text{up}} + \frac{\text{Reg}_G(T)}{T}. \quad (5)$$

326 In particular, if the game has value  $V$  (i.e.,  $V_{\text{up}} = V_{\text{low}} = V$ ),

$$|\bar{\mathcal{J}}_T - V| \leq \frac{\text{Reg}_G(T) + \text{Reg}_D(T)}{T}. \quad (6)$$

327 *Proof.* We start with the upper bound. By the generator's regret definition,

$$\sum_{t=1}^T \mathcal{J}(\theta_t, \phi_t) \leq \min_{\theta} \sum_{t=1}^T \mathcal{J}(\theta, \phi_t) + \text{Reg}_G(T).$$

328 Let  $\theta^* \in \arg \min_{\theta} \max_{\phi} \mathcal{J}(\theta, \phi)$  (a minimax optimizer). Evaluating the RHS at  $\theta^*$  and using  
 329  $\max_{\phi} \mathcal{J}(\theta^*, \phi) = V_{\text{up}}$  yields

$$\min_{\theta} \sum_{t=1}^T \mathcal{J}(\theta, \phi_t) \leq \sum_{t=1}^T \mathcal{J}(\theta^*, \phi_t) \leq \sum_{t=1}^T \max_{\phi} \mathcal{J}(\theta^*, \phi) = T V_{\text{up}}.$$

330 Combining gives  $\sum_{t=1}^T \mathcal{J}(\theta_t, \phi_t) \leq T V_{\text{up}} + \text{Reg}_G(T)$ , hence  $\bar{\mathcal{J}}_T \leq V_{\text{up}} + \text{Reg}_G(T)/T$ , which  
 331 completes this part of the inequality.

332 Next, we demonstrate the lower bound. By the discriminator's regret definition,

$$\sum_{t=1}^T \mathcal{J}(\theta_t, \phi_t) \geq \max_{\phi} \sum_{t=1}^T \mathcal{J}(\theta_t, \phi) - \text{Reg}_D(T).$$

333 Let  $\phi^* \in \arg \max_{\phi} \min_{\theta} \mathcal{J}(\theta, \phi)$  (a maxmin optimizer), so  $\min_{\theta} \mathcal{J}(\theta, \phi^*) = V_{\text{low}}$ . Then for every  
 334  $\theta$ ,  $\mathcal{J}(\theta, \phi^*) \geq V_{\text{low}}$ . In particular,

$$\max_{\phi} \sum_{t=1}^T \mathcal{J}(\theta_t, \phi) \geq \sum_{t=1}^T \mathcal{J}(\theta_t, \phi^*) \geq \sum_{t=1}^T V_{\text{low}} = T V_{\text{low}}.$$

335 Thus  $\sum_{t=1}^T \mathcal{J}(\theta_t, \phi_t) \geq T V_{\text{low}} - \text{Reg}_D(T)$ , i.e.,  $\bar{\mathcal{J}}_T \geq V_{\text{low}} - \text{Reg}_D(T)/T$ .

336 Combining both sides finishes the proof – in particular, if  $V_{\text{up}} = V_{\text{low}} = V$  (minimax theorem of  
 337 zero-sum games), then

$$V - \frac{\text{Reg}_D(T)}{T} \leq \bar{\mathcal{J}}_T \leq V + \frac{\text{Reg}_G(T)}{T},$$

338 and, since  $\max\{a, b\} \leq a + b$  for  $a, b \geq 0$ , the symmetric bound (6) follows.  $\square$

339 Another technicality is we require the optimal reward model– the one that satisfies our mutual  
 340 information constraint while minimizing BT-loss, to lie in our function class. We frame this as the  
 341 **feasible invariance** condition:

342 **Feasible invariance.** Let  $\mathcal{L}_{\text{BT}}^* = \inf_{\theta} \mathcal{L}_{\text{BT}}(\theta)$ . We say *feasible invariance* holds if there exists  $\theta^\dagger$   
 343 with

$$\mathcal{L}_{\text{BT}}(\theta^\dagger) = \mathcal{L}_{\text{BT}}^* \quad \text{and} \quad I_{\theta^\dagger}(C; R_{\theta^\dagger}) = 0. \quad (7)$$

344 In that case, the minimax value satisfies  $V = \mathcal{L}_{\text{BT}}^*$  by (4).

345 With these results, we can then prove our main theorem that in no-regret, our reward model converges  
 346 to zero mutual-information.

### 347 **Proof of Theorem 1 (No Regret Convergence)**

348 *Proof.* For each  $t$ , let  $V(\theta) = \max_{\phi} \mathcal{J}(\theta, \phi) = \mathcal{L}_{\text{BT}}(\theta) + \lambda I_{\theta}(C; R_{\theta})$  by Lemma 1. By the  
 349 discriminator’s regret definition,

$$\frac{1}{T} \sum_{t=1}^T V(\theta_t) = \frac{1}{T} \sum_{t=1}^T \max_{\phi} \mathcal{J}(\theta_t, \phi) \leq \bar{\mathcal{J}}_T + \frac{\text{Reg}_D(T)}{T}.$$

350 Feasible invariance implies  $V = \mathcal{L}_{\text{BT}}^*$ , and Lemma 2 gives  $\bar{\mathcal{J}}_T \leq V + \frac{\text{Reg}_G(T)}{T} = \mathcal{L}_{\text{BT}}^* + \frac{\text{Reg}_G(T)}{T}$ .  
 351 Hence

$$\frac{1}{T} \sum_{t=1}^T [\mathcal{L}_{\text{BT}}(\theta_t) + \lambda I_{\theta_t}(C; R_{\theta_t})] \leq \mathcal{L}_{\text{BT}}^* + \frac{\text{Reg}_G(T) + \text{Reg}_D(T)}{T}.$$

352 Since  $\mathcal{L}_{\text{BT}}(\theta_t) \geq \mathcal{L}_{\text{BT}}^*$  for all  $t$ , canceling  $\mathcal{L}_{\text{BT}}^*$  yields

$$\lambda \cdot \frac{1}{T} \sum_{t=1}^T I_{\theta_t}(C; R_{\theta_t}) \leq \frac{\text{Reg}_G(T) + \text{Reg}_D(T)}{T},$$

353 which proves the claim. Note that if the average of these terms converges to 0, then we also have that  
 354  $\inf_t I_{\theta_t} \rightarrow 0$ , and so we can select the minimum running iterate that is bounded by this average to  
 355 have a direct convergent subsequence.

356 □

357 We view training the discriminator using CELoss on each batch as an approximate "best-response."  
 358 More formally, we can think of it as an  $\epsilon_t$ -Nash equilibrium for each round – that is, if  $q_{\phi_t}$  is trained  
 359 to near-optimality per round so that  $\max_{\phi} \mathcal{J}(\theta_t, \phi) - \mathcal{J}(\theta_t, \phi_t) \leq \epsilon_t$  with  $\frac{1}{T} \sum_t \epsilon_t \rightarrow 0$ , then the  
 360 proof above holds with  $\text{Reg}_D(T)$  replaced by  $\sum_t \epsilon_t$ .

361 What if exact invariance is infeasible? That is, what if the Bradley-Terry-optimal reward model  
 362 invariant to category does not lie in our function class? If no  $\theta$  attains both  $\mathcal{L}_{\text{BT}}^*$  and  $I = 0$ , then  
 363  $V > \mathcal{L}_{\text{BT}}^*$  and our theorem instead yields the following bound:

$$\frac{1}{T} \sum_{t=1}^T I_{\theta_t}(C; R_{\theta_t}) \leq \frac{V - \mathcal{L}_{\text{BT}}^*}{\lambda} + \frac{\text{Reg}_G(T) + \text{Reg}_D(T)}{\lambda T},$$

364 where we cannot ignore the  $V - \mathcal{L}_{\text{BT}}^*$  term, which we can think of approximation error-esque term  
 365 in the learning theory language.

## 366 **A.4 Datasets and Preprocessing**

367 **HH-RLHF (Helpful/Harmless):** We construct (chosen, rejected) preference pairs and assign each  
 368 pair a category label of either helpful or harmless. Prompts and responses are concatenated, and  
 369 sequences are truncated to a maximum of 1,024 tokens.

370 **PKU-SafeRLHF (19 categories):** We retain the official harm category labels from the dataset release.  
 371 Samples with missing category annotations are removed to ensure label integrity.

372 **Deduplication:** Exact duplicate (prompt, response) pairs are removed to avoid information leakage  
 373 and inflated results.

374 **Tokenization and padding:** All data is tokenized with padding=longest and truncation=true. Each  
 375 prompt–response sequence is capped at 1,024 tokens in all reported experiments.

376 **A.5 Model and Training Details**

377 We use Llama-3.2-1B adapted into a scalar reward model for our RM backbone, with the Bradley-  
378 Terry pairwise log-likelihood on (chosen, rejected) pairs as our baseline training objective. We train  
379 for a single epoch on a balanced sample of helpful and harmless data from the Anthropic HH-RLHF  
380 dataset and evaluate on a held-out set of HH-RLHF dataset as well as RewardBench.

381 **A.6 Adversary and Fairness Optimization**

382 The fairness constraint uses a lightweight MLP adversary  $q_\phi$  that receives summary statistics of  
383 rewards, computed separately for each category. For each batch, we calculate the mean, variance,  
384 skewness, and kurtosis of the chosen and rejected rewards, grouped by category, to form the adver-  
385 sary’s input features.

386 Our training implementation follows the given alternating update schedule:

- 387 1. Compute Bradley–Terry loss  $L_{\text{BT}} = -\log \sigma(r_{\text{chosen}} - r_{\text{rejected}})$ .
- 388 2. Adversary step: update  $q_\phi$  by minimizing cross-entropy loss to predict the category from the  
389 moment features.
- 390 3. Fairness step: update the reward model to maximize adversary uncertainty, i.e., minimize

$$L_{\text{BT}} - \lambda_{\text{adv}} \cdot \text{CELoss}(q_\phi(\cdot \mid \text{moments}), y),$$

391 Ablation: For ablation studies, we sweep  $\lambda_{\text{adv}} \in \{0.0, 0.2, 0.4, 0.6, 0.8, 1.0, 1.5, 2.0\}$ . The default  
392 setting for main experiments is  $\lambda_{\text{adv}} = 0.2$ .

393 **Post-training Category Predictability.** As a post-training test, we train a fresh discriminator on  
394 frozen rewards from the above regularized model, which yields near-chance performance—AUC  
395  $0.78 \pm 0.03 \rightarrow 0.53 \pm 0.06$ , BA  $0.70 \pm 0.02 \rightarrow 0.52 \pm 0.05$  (5-fold; see Appx. A.7)—indicating  
396 little recoverable category signal from the fair reward model.

397 **A.7 Post-hoc Category Predictability Audit**

398 To test whether category information remains after training, we *freeze* the reward model and train a  
399 new discriminator  $\hat{q}(c \mid r)$  on its scalar outputs (no weights shared with the in-training adversary). We  
400 use stratified 5-fold cross-validation and report mean±sd over folds. The discriminator is a 2-layer  
401 MLP trained with cross-entropy and early stopping on validation AUC. Chance performance is 0.5  
402 for both AUC and balanced accuracy (BA).

Model	AUC	Balanced Acc.
Baseline RM	$0.78 \pm 0.03$	$0.70 \pm 0.02$
Fair RM (ours)	$0.53 \pm 0.06$	$0.52 \pm 0.05$

Table 4: Post-hoc predictability from frozen rewards; lower is better (chance  $\approx 0.5$ ).

403 **A.8 Mutual Information Estimation (Ablation)**

404 We measure the dependence between reward scores and category labels during the  $\lambda_{\text{adv}}$  sweep.  
405 Mutual information (MI) is computed with `sklearn.metrics.mutual_info_score` between category  
406 labels  $C \in \{\text{helpful, harmless}\}$  and a discretized reward variable, obtained by binning rewards into  
407 50 equal-width bins.

408 Lower MI indicates that the rewards are more category-independent. As an additional check, we mon-  
409 itor the adversary’s balanced accuracy; values close to chance imply minimal category dependence.

410 **A.9 Full  $\lambda_{\text{adv}}$  Sweep**

In this section we provide the complete data for our full sweep over adversarial loss parameters.

$\lambda_{\text{adv}}$	BT loss	MI
0.0	2.8712	0.2282
0.2	2.2307	0.0163
0.4	1.5607	0.0088
0.6	1.7104	0.0059
0.8	1.1879	0.0073
1.0	0.8694	0.0141
1.5	0.7432	0.0136
2.0	0.8151	0.0076

Table 5: Complete sweep of  $\lambda_{\text{adv}}$  values.

411

412 **A.10 Scaling Experiments**

413 To evaluate the scalability of our method, we conducted preliminary experiments on Meta’s Llama3-  
 414 8B-Instruct model on an 8xH100 node. The reward distributions for our Fair-RM variant, shown  
 415 below, exhibit a more complex, multimodal structure compared to the 1.1B model, which we  
 416 hypothesize is due to the larger model’s capacity to capture finer-grained nuances in the preference  
 417 data. Despite this, the results confirm that our approach remains effective at scale. There is clear  
 418 separation between chosen and rejected rewards, indicating preference alignment is maintained.  
 419 Crucially, the distributions for helpful and harmless categories remain tightly aligned, demonstrating  
 420 that the fairness constraint successfully generalizes and prevents reward disparities even in larger  
 421 models. However, both our base model and fair-RM variant achieve around 50% accuracy on a  
 422 subset of RewardBench after our training, for a variety of reasons but mainly in part due to the small  
 423 bandwidth we had to only run smaller training runs. Our Fair-RM had on-par performance with  
 424 the baseline BT model, however, but to achieve SOTA-level eval results on both models, full-scale  
 425 post-training of RewardBench-competitive models derived from the 8B models is part of our future  
 426 intended work.

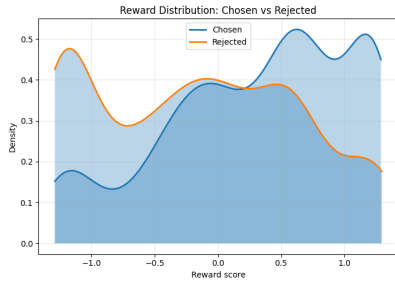


Figure 5: Reward distributions for chosen vs. rejected

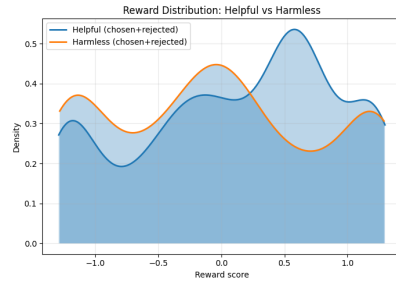


Figure 6: Reward distributions for helpful vs. harmless

427 **B PPO Training Setup**

428 In this section we detail our setup for PPO training of downstream language models using our fair  
 429 reward models.

430 **Base Actor.** We initialize all policy variants from `TinyLlama/TinyLlama-1.1B-Chat-v1.0`  
 431 to enable rapid convergence and reduce compute cost while still maintaining competitive generation  
 432 quality for our evaluation tasks. Policies are adapted using LoRA with rank  $r = 16$  and  $\alpha = 32$ ,  
 433 targeting the query/key/value and output projection matrices in the attention layers.

434 **PPO Configuration.** We use HuggingFace TRL’s `PPOTrainer` with minibatch size = 64, batch  
 435 size = 512, and 2 PPO epochs per update. The KL control coefficient is set to  $\beta = 0.05$  (adaptive  
 436 control enabled), targeting the reference model (`TinyLlama/TinyLlama-1.1B-Chat-v1.0`).  
 437 We set `target_kl=0.1` to limit divergence from the reference.

438 **Reward Models.** All reward models are Llama-3.2-1B sequence classifiers trained on preference  
 439 data with the Bradley–Terry objective. The **Fair** variant applies a mutual information (MI) penalty  
 440 with  $\lambda_{\text{adv}} = 0.2$  between protected-category predictions and reward scores. **Fair + Curiosity** adds an  
 441 intrinsic curiosity bonus from a Random Network Distillation (RND) module trained online during  
 442 PPO.

443 **Curiosity Bonus.** The RND network uses a 2-layer MLP with ReLU activations, hidden size 512.  
 444 The predictor network is optimized with Adam ( $\eta = 1 \times 10^{-4}$ ) on the cosine similarity loss between  
 445 target and predictor features. Intrinsic reward is scaled by  $\eta_{\text{cur}} = 0.05$  and added to the scalar RM  
 446 score before PPO optimization.

447 **Generation Settings.** For PPO rollouts, we generate with temperature = 0.7, top- $p$  = 0.9, and  
 448 max length = 256 tokens. KL penalties are computed against the reference log-probabilities.

449 **Training Duration.** Each run is trained for  $N = 5,000$  PPO steps ( $\approx 1.5\text{M}$  tokens processed),  
 450 which we found sufficient for convergence in both reward and policy loss metrics given the small  
 451 model size.

## 452 B.1 Parity Gap: Definition and Estimation

453 In this section we detail a parity gap (effectively mean matching evaluation) for how fair a reward  
 454 model is, for simplicity across only two categories.

455 **Definition.** Let  $r(x, y)$  denote the scalar reward assigned by a (fixed) safety RM to a prompt–  
 456 response pair  $(x, y)$ . We consider two behavior categories  $c \in \{\text{Helpful, Harmless}\}$  and define the  
 457 *parity gap* as the absolute difference in expected rewards:

$$\text{ParityGap} = \left| \mathbb{E}[r(x, y) \mid c = \text{Helpful}] - \mathbb{E}[r(x, y) \mid c = \text{Harmless}] \right|.$$

458 We define the parity gap as effectively a mean-matching surrogate evaluation – intuitively, a smaller  
 459 parity gap indicates the RM (and the downstream policy it shapes) treats categories on a comparable  
 460 reward scale, reducing category-dependent inflation/deflation.

461 **Estimator.** Given disjoint evaluation sets  $\mathcal{D}_H$  and  $\mathcal{D}_A$  (Helpful vs. Harmless) with sizes  $n_H$  and  $n_A$   
 462 and rewards  $\{r_i^H\}_{i=1}^{n_H}$ ,  $\{r_j^A\}_{j=1}^{n_A}$ , we compute

$$\bar{r}_H = \frac{1}{n_H} \sum_{i=1}^{n_H} r_i^H, \quad \bar{r}_A = \frac{1}{n_A} \sum_{j=1}^{n_A} r_j^A, \quad \hat{\Delta} = \bar{r}_H - \bar{r}_A, \quad \widehat{\text{ParityGap}} = |\hat{\Delta}|.$$

463 When  $n_H \neq n_A$ , the above remains unbiased under i.i.d. sampling within each group. In our main  
 464 runs we use balanced sets ( $n_H = n_A$ ).

465 **Relative change (vs. a baseline).** When comparing a model  $M$  to a baseline  $B$ , we also report the  
 466 relative drop:

$$\text{RelDrop}(M; B) = \frac{\widehat{\text{ParityGap}}(M) - \widehat{\text{ParityGap}}(B)}{\widehat{\text{ParityGap}}(B)} \times 100\%.$$

467 **Practical notes.** (i) We score responses with the same fixed RM across all policies. (ii) Generation  
 468 settings and seeds are identical across policies (Appendix B).

## 469 B.2 Semantic Diversity Calculation

470 In this section we detail our metric for diversity of LLM sampling to benchmark our intrinsic reward.

471 **Prompts and generation.** For diversity evaluation we sample 1,030 LIMA prompts (seed 42) and  
 472 generate one response per prompt with identical sampling across models. Prompts are drawn from  
 473 GAIK/lima. Generation parameters: temperature = 0.9, top- $p$  = 0.95, max\_new\_tokens = 100,  
 474 max\_length = 512, batch size = 8. All models use the same seed and generation parameters.

475 **Semantic diversity (primary metric).** Let  $f(\cdot)$  be all-mpnet-base-v2 with mean-pooling;  
 476 embeddings are  $\ell_2$ -normalized. For the set of responses  $\{y_i\}_{i=1}^n$  with embeddings  $e_i = f(y_i)$ , we  
 477 report

$$\text{SemDiv} = \frac{2}{n(n-1)} \sum_{i < j} (1 - \cos(e_i, e_j)).$$

478 Higher is better (more meaning-level variety).

479 **Statistics.** To compare a fair model against the baseline, we use a paired bootstrap (1,000 resamples;  
480 two-sided) over aligned prompt sets, reporting the mean difference, 95% CI, and  $p$ -value. In the main  
481 text, we report semantic-diversity differences: Fair (no curiosity) vs. Baseline:  $-0.0054$  ( $p < 0.001$ );  
482 Fair + Curiosity vs. Baseline:  $-0.0022$  ( $p = 0.002$ ).

### 483 **B.3 Compute and Runtime**

484 **Hardware:** For initial experiments of both reward model training and PPO, we used dual A100  
485 clusters, and currently are using a 8xH100 node for results on Llama3-8B.