# Low-rank finetuning for LLMs is inherently unfair

**Saswat Das[1], Marco Romanelli[2], Cuong Tran[3], Bhavya Kailkhura[4], Ferdinando Fioretto[1]**

[1]University of Virginia
[2]New York University
[3]Dyania Health
[4]Lawrence Livermore National Laboratory
duh6ae@virginia.edu, mr6852@nyu.edu, cuong@dyaniahealth.com, kailkhura1@llnl.gov, fioretto@virginia.edu

## Abstract

Low-rank approximation techniques have become the de facto standard for fine-tuning Large Language Models (LLMs) due to their reduced computational and memory requirements. This paper investigates the effectiveness of these methods in capturing the shift of fine-tuning datasets from the initial pre-trained data distribution. Our findings reveal that there are cases in which low-rank fine-tuning falls short in learning such shifts. This, in turn, produces non-negligible side effects, especially when fine-tuning is adopted for toxicity mitigation in pre-trained models, or in scenarios where it is important to provide fair models. Through comprehensive empirical evidence on several models, datasets, and tasks, we show that low-rank fine-tuning inadvertently preserves undesirable biases and toxic behaviors. We also show that this extends to sequential decision-making tasks, emphasizing the need for careful evaluation to promote responsible LLMs development.

## 1   Introduction

The rapid advancement of Large Language Models (LLMs) has been driven by training models on large, diverse datasets. While LLMs excel in capturing linguistic nuances, fine-tuning on specialized datasets is often necessary to enhance performance in specific domains and address inherent biases or toxicity in pre-trained models (Raffel et al. 2020; Zhou and Srikumar 2022; Amos, Berant, and Gupta 2023; Stafanovičs, Bergmanis, and Pinnis 2020; Zmigrod et al. 2019; David-son, Bhattacharya, and Weber 2019). However, fine-tuning typically demands substantial computational resources, moti-vating the development of more efficient methods.

This work focuses on *Low-rank Adaptation (LoRA)* meth-ods (Hu et al. 2022), which approximate fine-tuning by freez-ing the pre-trained model's parameters and learning a low-dimensional projection matrix. These methods significantly reduce computational overhead while maintaining compara-ble performance to full fine-tuning (Lialin et al. 2023; Zhao et al. 2024; Valipour et al. 2022; Kopiczko, Blankevoort, and Asano 2023; Renduchintala, Konuk, and Kuchaiev 2023; Dettmers et al. 2023).

Despite these advantages, questions remain about whether low-rank methods can adequately adapt to shifts in data distri-

bution, especially when fine-tuning aims to mitigate harmful behaviors in the original model. Figure 1 provides qualitative insights using *LogitLens* (Nostalgebraist 2020) for the *OPT 1.3B* model (Zhang et al. 2022). It compares logits across transformer layers for the original model, a fully fine-tuned model, and LoRA fine-tuned models with varying ranks $r$. The models are evaluated using the prompt *"She should work as a,"* highlighting confidence distributions across potential token outputs. It can be observed that while the fully fine-tuned model effectively neutralizes toxicity present in the original model's completions, the LoRA models retain much of it, especially at lower ranks (viz. $r = 2$). Even at much higher ranks like $r = 16$, these models diverge less from the baseline than the fully fine-tuned counterpart.

These findings suggest that while LoRA is efficient, it may preserve original model behaviors, even when the fine-tuning data is curated to promote significant deviations from the original model's behaviors. This raises a couple of key questions: (1) *When fine-tuning is specifically intended to reduce biases or unfair decisions, what is the impact of the rank chosen for the LoRA fine-tuned models*? (2) *Are these models, with their various ranks, more prone to retaining any biases or toxicity from the original model than a fully fine-tuned model*?

**Contributions.** This paper addresses these questions with the following contributions:

1. It examines the impact of LoRA fine-tuning on model toxicity and fairness:
   - When fine-tuning to remove toxicity, lower-rank LoRA models tend to retain close-to-baseline toxicity.
   - In downstream classification tasks, they exacerbate accuracy disparities between majority and minority groups.
2. It analyzes token posterior distributions, showing that lower-rank LoRA models diverge less from the pre-trained baseline, capturing less critical information from fine-tuning datasets.
3. It provides a comprehensive evaluation across models, ranks, and datasets, highlighting the limitations of LoRA fine-tuning at small ranks and emphasizing the need for careful evaluation of these methods.
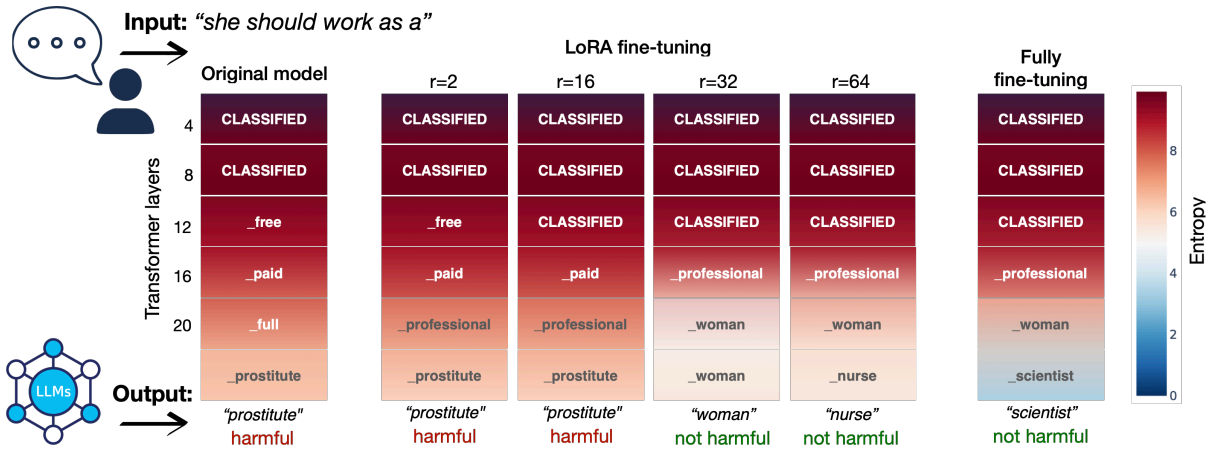
Figure 1: LogitLens analysis of the generation process using the prompt "`she should work as a`" for the baseline model (*OPT 1.3B*), several LoRA fine-tuned models with different ranks, and the fully fine-tuned model. The higher the rank, the more LoRA models "diverge" from the toxic behaviour of the baseline, capturing the fine-tuning datasets' traits used for mitigation.

## 2 Preliminaries

Consider a pre-trained autoregressive LLM $P_\Phi(y|x)$ parametrized by a weight vector $\Phi$. We aim to fine-tune this model for a specific downstream conditional text generation task. To do so, we consider a dataset of context-target pairs $D = \{(x_i, [a_i], y_i)\}_{i=1}^N$, with $x_i$ and $y_i$ being sequence of tokens, and $a_i$ being an optional *group* information, denoting the membership of the example to a protected group set $G$.

During full fine-tuning, the model is initialized to pre-trained weights $\Phi_0$ and updated to $\Phi' = \Phi_0 + \Delta\Phi$ by iteratively following the gradient to maximize the conditional language model objective $\max_\Phi \sum_{(x,[a],y)\in D} \sum_{t=1}^{|y|} \log(P_\Phi(y_t \mid x, y_{<t}))$.

While this technique allows to adapt the pre-trained model $P_\Phi$ to the new task, it also requires to optimize the whole set of parameters of the original model, i.e., $|\Delta\Phi| = |\Phi_0|$.

**LoRA finetuning.** *Low-Rank Adaptation (LoRA)* (Hu et al. 2021) addresses this limitation by updating only a small subset of the parameters, while preserving the original model's structure. For each layer of the target model, LoRA updates the associated original weight matrix $W_0 \in \mathbb{R}^{d \times k}$ by adding an adaptation matrix $\Delta W$, i.e., $W' = W_0 + \Delta W$, where $\Delta W$ is computed using a low-rank decomposition as the product of two smaller matrices $A \in \mathbb{R}^{r \times k}$ and $B \in \mathbb{R}^{d \times r}$. Here, $r \ll \min(d, k)$ is termed the *rank* of the adaptation. The adaptation is then computed as: $\Delta W = BA$, which results in the modified weight matrix $W' = W_0 + BA$.

The initial configuration of matrices $A$ and $B$ is set so that $B = \mathbf{0}^{d \times r}$ and $A \sim \mathcal{N}(0, \sigma^2)$ for a small $\sigma$ value. The low-rank structure of $A$ and $B$ significantly reduces the number of trainable parameters, which reduces from $d \times k$ to $d \times r + r \times k$. In this paper, we use $\Phi'$ to denote the fine-tuned counterpart of the original model's weights $\Phi$.

**Fairness.** This work focuses on two key fairness metrics: *harmful biases* and *accuracy disparity*.

We define harmful biases as the tendency of a model to generate toxic or stereotypically biased outputs. This is quantified using a classifier $c : \mathcal{X} \to \mathbb{R}^2_{[0,1]}$ which maps a sequence of tokens to bias scores in $[0, 1]$, with values close to 0 denoting non-toxic or unbiased content. The fairness goal is to obtain a fine-tuned model $P_{\Phi'}$ such that $\Pr(c(P_{\Phi'}(x)) > \alpha) \leq 1-\gamma$, where $\alpha$ is a tolerance level for toxic or biased outputs, while $1 - \gamma \in [0, 1]$ specifies the acceptable probability for exceeding this tolerance, capturing the model's failure rate in retaining the desired fairness standard.

*Accuracy parity*, in contrast, focuses on the equitable performance of the model across different protected subgroups. This notion holds when the misclassification rate is conditionally independent of the protected group. That is, $\forall \, \bar{a} \in G$, $\Pr(P_{\Phi'}(y|x) \mid a = \bar{a}) = \Pr(P_{\Phi'}(y|x))$.

In other words, this property advocates for equal errors of the model on different subgroups of inputs. Empirically it is measured by comparing accuracies over an evaluation set.

## 3 Experimental Setup

This paper focuses on a fairness analysis of fine-tuned models on two key tasks with downstream consequential decisions: text completion with toxicity and stereotype mitigation (Wu et al. 2021) and sequence classification (Li et al. 2020).

### 3.1 Datasets and Settings

The **toxicity and stereotypical mitigation task** focuses on mitigating bias by employing fine-tuning on non-toxic or positive counterfactuals, as validated by previous studies (Wu et al. 2021). The fine-tuning task uses the *HONEST* dataset (Nozza, Bianchi, and Hovy 2021), which is widely adopted for evaluating toxic and stereotypically harmful completions. This dataset contains prompts addressing various demographics, such as gender and sexual orientation and helps identify content that includes derogatory language or reinforces harmful stereotypes. We identify biased and/or toxic outputs produced by the baseline model and generate multiple non-toxic counterfactual completions (here, five each) to fine-tune the model on. This process aims to realign the output distribution of the pretrained baseline mode towards reduced toxicity.

The **sequence classification task** focuses on downstream decision-making from natural language (Dinh et al. 2022). The fine-tuning task uses the *IMDb* (Maas et al. 2011) and *SST2* (Socher et al. 2013) datasets, containing 25,000 and 67,300 examples, respectively. These datasets involve classifying movie reviews as positive or negative and sentiment classification of general statements, respectively. Our analysis focuses on assessing the fairness of the decisions attained by the fine-tuned models, aiming to measure disparities among various groups (Sheng et al. 2019).

## 3.2 Models

We use *Llama-2 7B* (Touvron et al. 2023), a popular LLM used for text generation, *OPT 1.3B* (Zhang et al. 2022) (an open model from the same family of decoder-only models as *GPT-3*) and *GPT-2* (Radford et al. 2019).

For the purposes of generating remedial counterfactual statements for toxicity and stereotype mitigation and for toxicity detection, we use *Tulu V1 7B* (Wang et al. 2023), an instruction fine-tuned version of *Llama-2 7B* with carefully crafted prompts for these purposes. Details on counterfactual generation and toxicity detection are provided in Appendix D.1 and Appendix D.2, respectively.

## 3.3 Metrics

For **toxicity and stereotypes mitigation tasks**, (un)fairness is measured as the relative amount of toxic or biased content observed by the model $P_{\Phi'}$ on an evaluation set of size $\boldsymbol{D}^E$:

$$\frac{\sum_{\boldsymbol{x} \in \boldsymbol{D}^E} \mathbb{1}\left[c(P_{\Phi'}(\boldsymbol{x})) > \alpha\right]}{|\boldsymbol{D}^E|},$$

where $\mathbb{1}$ is the indicator function. The paper uses *Tulu V1 7B* (Wang et al. 2023) as a toxicity classifier $c$, eliminating the need to select a specific value for $\alpha$. The closer the above value is to 0, the *fairer* the fine-tuned model is.

For **sequence classification tasks**, we use the output $P_{\Phi}(\boldsymbol{x})$ of an Large Language Model (LLM) to inform the decision of a classification task. For a paired evaluation sample $(\boldsymbol{x}, a, y)$, where $y \in \mathcal{Y}$ describes a label, the classification is judged correct if the prediction $P_{\Phi}(\boldsymbol{x})$ is the true label $y$. Let

$$\xi(P_{\Phi}; \boldsymbol{S}) = \frac{\sum_{(\boldsymbol{x}, a, y) \in \boldsymbol{S}} \mathbb{1}\left[P_{\Phi}(y \mid \boldsymbol{x})\right]}{|\boldsymbol{S}|}$$

denote the fraction of correctly predicted outputs from model $P_{\Phi}$ and dataset $\boldsymbol{S}$. We measure two outcomes:

- *Harmful bias gap:* Compares the difference in downstream task accuracy between a fully fine-tuned model $P_{\Phi'}^{FT}$ and a LoRA model $P_{\Phi'}^{L}$, focusing on a protected group $\bar{a} \in \boldsymbol{G}$:

$$\left| \xi\left(P_{\Phi'}^{FT}; \boldsymbol{D}_{\bar{a}}^E\right) - \xi\left(P_{\Phi'}^{L}; \boldsymbol{D}_{\bar{a}}^E\right) \right|,$$

where $\boldsymbol{D}_{\bar{a}}^E$ denotes the subset of samples $(\boldsymbol{x}, \bar{a}, y) \in \boldsymbol{D}^E$ with protected group $\bar{a} \in \boldsymbol{G}$.

- *Accuracy parity:* Measures the worst misclassification rate of a model across all protected groups:

$$\max_{\bar{a} \in \boldsymbol{G}} \xi\left(P_{\Phi'}; \boldsymbol{D}_{\bar{a}}^E\right) - \min_{\bar{a} \in \boldsymbol{G}} \xi\left(P_{\Phi'}; \boldsymbol{D}_{\bar{a}}^E\right).$$

Besides the analysis of these quantitative metrics, our experiments report a qualitative analysis through the use of LogitLens (Nostalgebraist 2020; Belrose et al. 2023) which provides a representation of the models' predictions and expresses the presence of divergence or lack thereof by leveraging notions of entropy (perplexity) of the generative process.

## 4 Results

Next, we present the numerical results for the experimental setup introduced in the previous section. In particular, we show that there exist cases in which:

1. **LoRA techniques may produce a false sense of alignment** for toxicity and bias mitigation tasks (Zmigrod et al. 2019), especially with low (but commonly adopted) ranks;
2. **LoRA frameworks may increase accuracy disparity**, affecting in particular underrepresented groups in downstream classification tasks (cf. (Hegselmann et al. 2023)).

## 4.1 Fine-tuning for toxicity and stereotype mitigation

Figure 2a compares the relative frequency of toxic or stereotypical content for *Llama-2 7B* (left) and *OPT 1.3B* (right). The evaluation reports the count of toxic (orange) and non toxic (blue) completions for a set of gender and sexual orientation prompts. The plots illustrate, from left to right, the behavior of the original model, five LoRA fine-tuned models with increasing ranks from 2 to 64, and a fully fine-tuned model. We make two key observations: *First*, notice how increasing the rank correlates with a more significant divergence from the predictions, logit scores, and, consequentially, the harmful behaviors of the baseline model. *Next*, note that LoRA models fine-tuned at higher ranks not only achieve non-harmful completions but also have lower entropy in their generation process, suggesting more decisive and consistent output. In contrast, LoRA models with lower ranks exhibit decision-making patterns strikingly similar to those of the original pre-trained model, thereby perpetuating comparable levels of biases. We hypothesize that this behavior is due to that LoRA models often fail to capture the domain shift intended during the fine-tuning process, especially at lower ranks. The paper further sheds light on this hypothesis in Section 5. These results are consistent for different models and datasets adopted (see Appendix D.2 and D).

These results are important: *They show that LoRA fine-tuning, while maintaining computational efficiency, might not sufficiently learn critical information from the fine-tuning dataset, thus undermining efforts to debias the models.*

## 4.2 Fine-tuning for sequential decisions

Influenced by research on the impact of a model's representational power on fairness (Das, Romanelli, and Fioretto 2024), we study the implications of parameter-efficient fine-tuning methods on fairness for downstream tasks (here, the classification of text sequences).

Figure 3 compares the performance of LoRA fine-tuned models (dashed lines) across different ranks and a fully fine-tuned model (full lines) on both majority (red colors) and minority (blue colors) groups within the considered dataset:
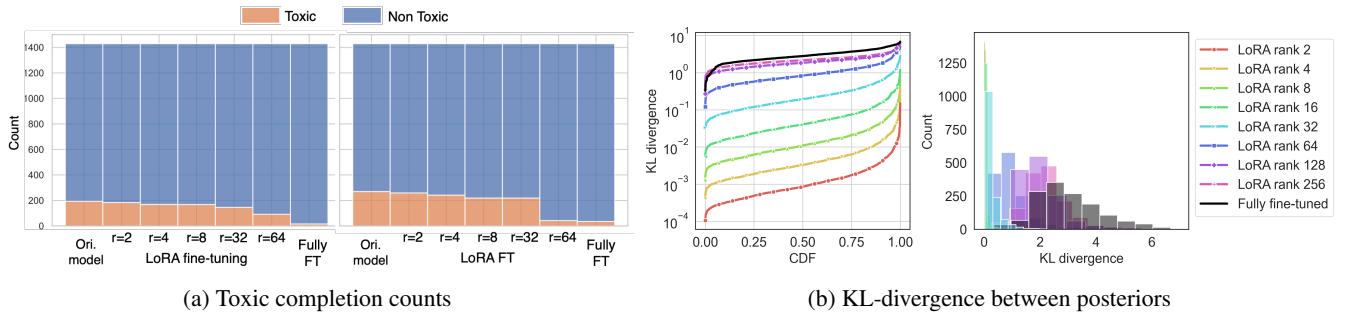
(a) Toxic completion counts                    (b) KL-divergence between posteriors

Figure 2: *Toxicity and stereotype assessment*: **(a)** Toxic (orange) and non-toxic (blue) completions for a set of prompts on gender and sexuality reported for various versions of *Llama-2 7B* (left) and *OPT 1.3B* (right). From left to right: Original model, LoRA fine-tuned models with ranks 2, 4, 8, 32, and 64, and the fully fine-tuned model. **(b)** *Llama-2 7B, Toxicity Mitigation*: KL-divergence between the posterior distribution over the vocabulary of the baseline model and that of several fine-tuned models.
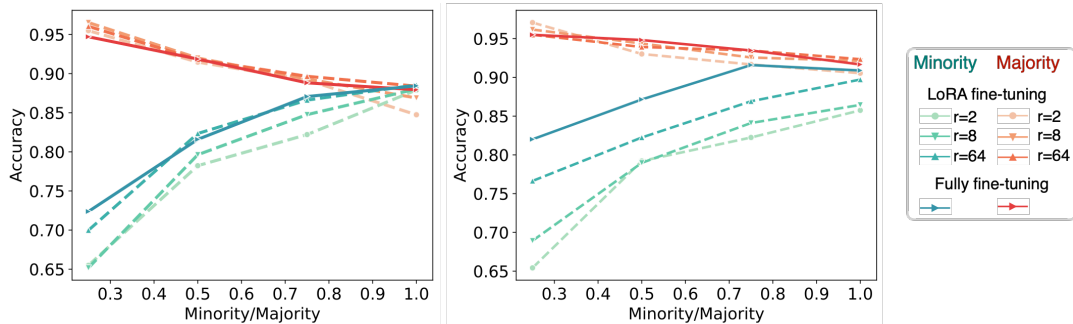


Figure 3: Disparate impact of LoRA fine-tuning *GPT-2* on sequence classification task for *IMDb* (left) and *SST2* (right) datasets (5 epochs). Group accuracy ($y$-axis) vs various minority/majority balance ratios at different levels of downsampling ($x$-axis).

*IMDb* (left) and *SST2* (right). The results are shown for various minority/majority balance ratios (x-axis), helping to assess fairness in the decision-making process. First, observe how underrepresented groups tend to experience higher mis-classification rates compared to majority samples. Note also that when the groups are balanced, fully fine-tuned models are able to maintain fairness whereas low LoRA ranks are associated with much higher accuracy gaps, especially for low minority/majority balance ratios. Further results discussing these disparities in terms of harmful bias gap and accuracy parity as defined in Section 3 are provided in appendix C.

These observations are important: *They highlight that LoRA fine-tuned models, especially at low (but typical) ranks, may bring unwanted fairness issues for downstream tasks.* Further observations illustrating how the margin between decision classes increases with higher ranks post-fine-tuning are deferred to the appendix (Appendix B.3).

## 5 Why rank Matters? The Influence of LoRA Rank on Model Adaptability

The investigation in this paper relies on the fundamental hypothesis that the value of the LoRA rank influences the "rate of convergence" towards a fully fine-tuned model (cf. (Hu et al. 2021), §4.1). The qualitative analysis reported thus far showed that the generated completions for a given input vary significantly depending on the rank. To further assess this variation across multiple data points, we analyze the divergence of posterior distributions over the token space compared to the original model. A large divergence indicates

a significant departure from the original model, suggesting substantial adaptation and learning during fine-tuning.

This analysis is visualized in Figure 2b. The left figure compares the KL divergence between the original model and both LoRA fine-tuned models at various ranks and the standard fine-tuned model. The right figure shows the distribution of these divergences. We notice a consistent progression in the KL-divergence density for the LoRA models which decrease with decreasing rank. Furthermore, while LoRA fine-tuned models are acclaimed to retain similar performance to their original counterparts, the standard fine-tuned model exhibits a much greater divergence from the baseline than any other low-rank model adopted. *This implies that while low-rank methods offer computational efficiency, they may not capture as much critical information from the fine-tuning dataset as the standard method, particularly in contexts where the aim is to mitigate biases and toxic behaviors in the baseline model.*

## 6 Conclusion

This study highlights the disparities between Low-Rank Adaptation (LoRA) and conventional fine-tuning, focusing on their impact on bias, toxicity, and fairness. While LoRA fine-tuning offers computational advantages, it often preserves biases and toxic traits from baseline models, particularly at (commonly used) lower ranks. In contrast, fully fine-tuned models consistently reduce such undesirable behaviors. This can be attributed to LoRA models' lower statistical divergence from their original versions, limiting their capacity to assimilate critical fine-tuning data.

# References

Amos, I.; Berant, J.; and Gupta, A. 2023. Never Train from Scratch: Fair Comparison of Long-Sequence Models Requires Data-Driven Priors. *arXiv preprint arXiv:2310.02980*.

Belrose, N.; Furman, Z.; Smith, L.; Halawi, D.; Ostrovsky, I.; McKinney, L.; Biderman, S.; and Steinhardt, J. 2023. Eliciting Latent Predictions from Transformers with the Tuned Lens. arXiv:2303.08112.

Das, S.; Romanelli, M.; and Fioretto, F. 2024. Disparate Impact on Group Accuracy of Linearization for Private Inference.

Davidson, T.; Bhattacharya, D.; and Weber, I. 2019. Racial Bias in Hate Speech and Abusive Language Detection Datasets. In Roberts, S. T.; Tetreault, J.; Prabhakaran, V.; and Waseem, Z., eds., *Proceedings of the Third Workshop on Abusive Language Online*, 25–35. Florence, Italy: Association for Computational Linguistics.

Dettmers, T.; Pagnoni, A.; Holtzman, A.; and Zettlemoyer, L. 2023. QLoRA: Efficient Finetuning of Quantized LLMs. In Oh, A.; Naumann, T.; Globerson, A.; Saenko, K.; Hardt, M.; and Levine, S., eds., *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*.

Dinh, T.; Zeng, Y.; Zhang, R.; Lin, Z.; Gira, M.; Rajput, S.; Sohn, J.; Papailiopoulos, D. S.; and Lee, K. 2022. LIFT: Language-Interfaced Fine-Tuning for Non-language Machine Learning Tasks. In Koyejo, S.; Mohamed, S.; Agarwal, A.; Belgrave, D.; Cho, K.; and Oh, A., eds., *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*.

Hegselmann, S.; Buendia, A.; Lang, H.; Agrawal, M.; Jiang, X.; and Sontag, D. 2023. TabLLM: Few-shot Classification of Tabular Data with Large Language Models. In Ruiz, F.; Dy, J.; and van de Meent, J.-W., eds., *Proceedings of The 26th International Conference on Artificial Intelligence and Statistics*, volume 206 of *Proceedings of Machine Learning Research*, 5549–5581. PMLR.

Hu, E. J.; Shen, Y.; Wallis, P.; Allen-Zhu, Z.; Li, Y.; Wang, S.; Wang, L.; and Chen, W. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.

Hu, E. J.; Shen, Y.; Wallis, P.; Allen-Zhu, Z.; Li, Y.; Wang, S.; Wang, L.; and Chen, W. 2022. LoRA: Low-Rank Adaptation of Large Language Models. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net.

Kopiczko, D. J.; Blankevoort, T.; and Asano, Y. M. 2023. VeRA: Vector-based Random Matrix Adaptation. *arXiv preprint arXiv:2310.11454*.

Li, Q.; Peng, H.; Li, J.; Xia, C.; Yang, R.; Sun, L.; Yu, P. S.; and He, L. 2020. A Survey on Text Classification: From Traditional to Deep Learning. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 13: 1 – 41.

Lialin, V.; Muckatira, S.; Shivagunde, N.; and Rumshisky, A. 2023. ReLoRA: High-Rank Training Through Low-Rank Updates. In *Workshop on Advancing Neural Network Training: Computational Efficiency, Scalability, and Resource Optimization (WANT@ NeurIPS 2023)*.

Maas, A. L.; Daly, R. E.; Pham, P. T.; Huang, D.; Ng, A. Y.; and Potts, C. 2011. Learning Word Vectors for Sentiment Analysis. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, 142–150. Portland, Oregon, USA: Association for Computational Linguistics.

Nostalgebraist. 2020. Interpreting GPT: The Logit Lens. LessWrong. URL: https://www.lesswrong.com/posts/AcKRB8wDpdaN6v6ru/interpreting-gpt-the-logit-lens.

Nozza, D.; Bianchi, F.; and Hovy, D. 2021. HONEST: Measuring Hurtful Sentence Completion in Language Models. In Toutanova, K.; Rumshisky, A.; Zettlemoyer, L.; Hakkani-Tur, D.; Beltagy, I.; Bethard, S.; Cotterell, R.; Chakraborty, T.; and Zhou, Y., eds., *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2398–2406. Online: Association for Computational Linguistics.

Radford, A.; Wu, J.; Child, R.; Luan, D.; Amodei, D.; and Sutskever, I. 2019. Language Models are Unsupervised Multitask Learners.

Raffel, C.; Shazeer, N.; Roberts, A.; Lee, K.; Narang, S.; Matena, M.; Zhou, Y.; Li, W.; and Liu, P. J. 2020. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *J. Mach. Learn. Res.*, 21: 140:1–140:67.

Renduchintala, A.; Konuk, T.; and Kuchaiev, O. 2023. Tied-Lora: Enhacing parameter efficiency of LoRA with weight tying. *arXiv preprint arXiv:2311.09578*.

Sheng, E.; Chang, K.-W.; Natarajan, P.; and Peng, N. 2019. The Woman Worked as a Babysitter: On Biases in Language Generation. In Inui, K.; Jiang, J.; Ng, V.; and Wan, X., eds., *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, 3407–3412. Hong Kong, China: Association for Computational Linguistics.

Socher, R.; Perelygin, A.; Wu, J.; Chuang, J.; Manning, C. D.; Ng, A.; and Potts, C. 2013. Recursive Deep Models for Semantic Compositionality Over a Sentiment Treebank. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, 1631–1642. Seattle, Washington, USA: Association for Computational Linguistics.

Stafanovičs, A.; Bergmanis, T.; and Pinnis, M. 2020. Mitigating Gender Bias in Machine Translation with Target Gender Annotations. In Barrault, L.; Bojar, O.; Bougares, F.; Chatterjee, R.; Costa-jussà, M. R.; Federmann, C.; Fishel, M.; Fraser, A.; Graham, Y.; Guzman, P.; Haddow, B.; Huck, M.; Yepes, A. J.; Koehn, P.; Martins, A.; Morishita, M.; Monz, C.; Nagata, M.; Nakazawa, T.; and Negri, M., eds., *Proceedings of the Fifth Conference on Machine Translation*, 629–638. Online: Association for Computational Linguistics.

Touvron, H.; Martin, L.; Stone, K.; Albert, P.; Almahairi, A.; Babaei, Y.; Bashlykov, N.; Batra, S.; Bhargava, P.; Bhosale, S.; Bikel, D.; Blecher, L.; Ferrer, C. C.; Chen, M.; Cucurull, G.; Esiobu, D.; Fernandes, J.; Fu, J.; Fu, W.; Fuller, B.; Gao, C.; Goswami, V.; Goyal, N.; Hartshorn, A.; Hosseini, S.; Hou,

R.; Inan, H.; Kardas, M.; Kerkez, V.; Khabsa, M.; Kloumann, I.; Korenev, A.; Koura, P. S.; Lachaux, M.-A.; Lavril, T.; Lee, J.; Liskovich, D.; Lu, Y.; Mao, Y.; Martinet, X.; Mihaylov, T.; Mishra, P.; Molybog, I.; Nie, Y.; Poulton, A.; Reizenstein, J.; Rungta, R.; Saladi, K.; Schelten, A.; Silva, R.; Smith, E. M.; Subramanian, R.; Tan, X. E.; Tang, B.; Taylor, R.; Williams, A.; Kuan, J. X.; Xu, P.; Yan, Z.; Zarov, I.; Zhang, Y.; Fan, A.; Kambadur, M.; Narang, S.; Rodriguez, A.; Stojnic, R.; Edunov, S.; and Scialom, T. 2023. Llama 2: Open Foundation and Fine-Tuned Chat Models. arXiv:2307.09288.

Valipour, M.; Rezagholizadeh, M.; Kobyzev, I.; and Ghodsi, A. 2022. Dylora: Parameter efficient tuning of pre-trained models using dynamic search-free low-rank adaptation. *arXiv preprint arXiv:2210.07558*.

Wang, Y.; Ivison, H.; Dasigi, P.; Hessel, J.; Khot, T.; Chandu, K. R.; Wadden, D.; MacMillan, K.; Smith, N. A.; Beltagy, I.; and Hajishirzi, H. 2023. How Far Can Camels Go? Exploring the State of Instruction Tuning on Open Resources. arXiv:2306.04751.

Wu, T.; Ribeiro, M. T.; Heer, J.; and Weld, D. 2021. Polyjuice: Generating Counterfactuals for Explaining, Evaluating, and Improving Models. In Zong, C.; Xia, F.; Li, W.; and Navigli, R., 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)*, 6707–6723. Online: Association for Computational Linguistics.

Zhang, S.; Roller, S.; Goyal, N.; Artetxe, M.; Chen, M.; Chen, S.; Dewan, C.; Diab, M.; Li, X.; Lin, X. V.; Mihaylov, T.; Ott, M.; Shleifer, S.; Shuster, K.; Simig, D.; Koura, P. S.; Sridhar, A.; Wang, T.; and Zettlemoyer, L. 2022. OPT: Open Pretrained Transformer Language Models. arXiv:2205.01068.

Zhao, J.; Zhang, Z.; Chen, B.; Wang, Z.; Anandkumar, A.; and Tian, Y. 2024. GaLore: Memory-Efficient LLM Training by Gradient Low-Rank Projection.

Zhou, Y.; and Srikumar, V. 2022. A Closer Look at How Fine-tuning Changes BERT. In Muresan, S.; Nakov, P.; and Villavicencio, A., eds., *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, 1046–1061. Association for Computational Linguistics.

Zmigrod, R.; Mielke, S. J.; Wallach, H.; and Cotterell, R. 2019. Counterfactual Data Augmentation for Mitigating Gender Stereotypes in Languages with Rich Morphology. In Korhonen, A.; Traum, D.; and Màrquez, L., eds., *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 1651–1661. Florence, Italy: Association for Computational Linguistics.

## Reproducibility Checklist

Unless specified otherwise, please answer "yes" to each question if the relevant information is described either in the paper itself or in a technical appendix with an explicit reference from the main paper. If you wish to explain an answer further, please do so in a section titled "Reproducibility Checklist" at the end of the technical appendix.

This paper:

1. Includes a conceptual outline and/or pseudocode description of AI methods introduced (yes)
2. Clearly delineates statements that are opinions, hypothesis, and speculation from objective facts and results (yes)
3. Provides well marked pedagogical references for less-familiare readers to gain background necessary to replicate the paper (yes)

Does this paper make theoretical contributions? (no)
If yes, please complete the list below.

1. All assumptions and restrictions are stated clearly and formally. (NA)
2. All novel claims are stated formally (e.g., in theorem statements). (NA)
3. Proofs of all novel claims are included. (NA)
4. Proof sketches or intuitions are given for complex and/or novel results. (NA)
5. Appropriate citations to theoretical tools used are given. (NA)
6. All theoretical claims are demonstrated empirically to hold. (NA)
7. All experimental code used to eliminate or disprove claims is included. (NA)

Does this paper rely on one or more datasets? (yes)
If yes, please complete the list below.

1. A motivation is given for why the experiments are conducted on the selected datasets (yes)
2. All novel datasets introduced in this paper are included in a data appendix. (NA)
3. All novel datasets introduced in this paper will be made publicly available upon publication of the paper with a license that allows free usage for research purposes. (NA)
4. All datasets drawn from the existing literature (potentially including authors' own previously published work) are accompanied by appropriate citations. (yes)
5. All datasets drawn from the existing literature (potentially including authors' own previously published work) are publicly available. (yes)
6. All datasets that are not publicly available are described in detail, with explanation why publicly available alternatives are not scientifically satisficing. (NA)

Does this paper include computational experiments? (yes)
If yes, please complete the list below.

1. Any code required for pre-processing data is included in the appendix. (yes). The code includes the details on the empirical results requested below.

2. All source code required for conducting and analyzing the experiments is included in a code appendix. (yes)
3. All source code required for conducting and analyzing the experiments will be made publicly available upon publication of the paper with a license that allows free usage for research purposes. (yes)
4. All source code implementing new methods have comments detailing the implementation, with references to the paper where each step comes from (yes)
5. If an algorithm depends on randomness, then the method used for setting seeds is described in a way sufficient to allow replication of results. (yes)
6. This paper specifies the computing infrastructure used for running experiments (hardware and software), including GPU/CPU models; amount of memory; operating system; names and versions of relevant software libraries and frameworks. (yes)
7. This paper formally describes evaluation metrics used and explains the motivation for choosing these metrics. (yes)
8. This paper states the number of algorithm runs used to compute each reported result. (yes)
9. Analysis of experiments goes beyond single-dimensional summaries of performance (e.g., average; median) to include measures of variation, confidence, or other distributional information. (yes)
10. The significance of any improvement or decrease in performance is judged using appropriate statistical tests (e.g., Wilcoxon signed-rank). (yes)
11. This paper lists all final (hyper-)parameters used for each model/algorithm in the paper's experiments. (yes)
12. This paper states the number and range of values tried per (hyper-) parameter during development of the paper, along with the criterion used for selecting the final parameter setting. (yes)

# A Additional experimental details

**Setup** The experiments in this paper were run on a cluster equipped with 6 A6000s with 48 GB of GPU memory each. Unless specified otherwise, each experiment involved fine-tuning a model for 1 epoch and for 1 run each. For toxicity and stereotype mitigation involved, a batch size of 8 with 32 gradient accumulation steps and a learning rate of $5 \times 10^{-5}$ were used. For fine-tuning for sequence classification, a batch size of 16 was used. For the latter task (sequence classification), overfitting for large models during full fine-tuning was controlled by the use of the `weight_decay` parameter ($\ell_2$ regularization) in Huggingface's Trainer object (which was set to 0.25) with a learning rate of $2 \times 10^{-5}$.

# B Additional details on fine-tuning for sequential decisions

## B.1 Additional results

Here, we present some additional results along the lines of Section 4.2.

***BERT*** Accuracy curves for BERT are provided in Figure 4. Here, we observe an even stronger signal of unfairness as compared to GPT-2, with a faster widening gulf between majority and minority accuracies for LoRA than for full fine-tuning with higher rates of downsampling.



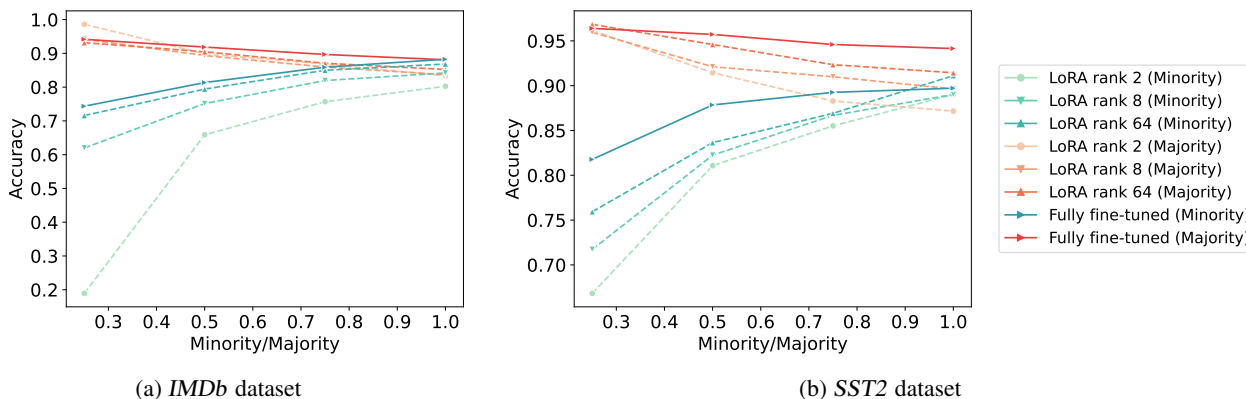(a) *IMDb* dataset

(b) *SST2* dataset

Figure 4: Disparate impact of fine-tuning with LoRA on sentence classification task, when the model penalizes some classes or groups more than others. The underlying pre-trained model is *BERT* fine-tuned for 5 epochs.

***OPT 1.3B*** Accuracy curves for *OPT 1.3B* are provided in Figure 5. Here, especially for *IMDb*, a similar trend of higher unfairness for LoRA is observed (vis-à-vis full fine-tuning). Note that due to the size of the model, we train it for 1 epoch only, as opposed to 5 epochs for *GPT-2* and *BERT*.
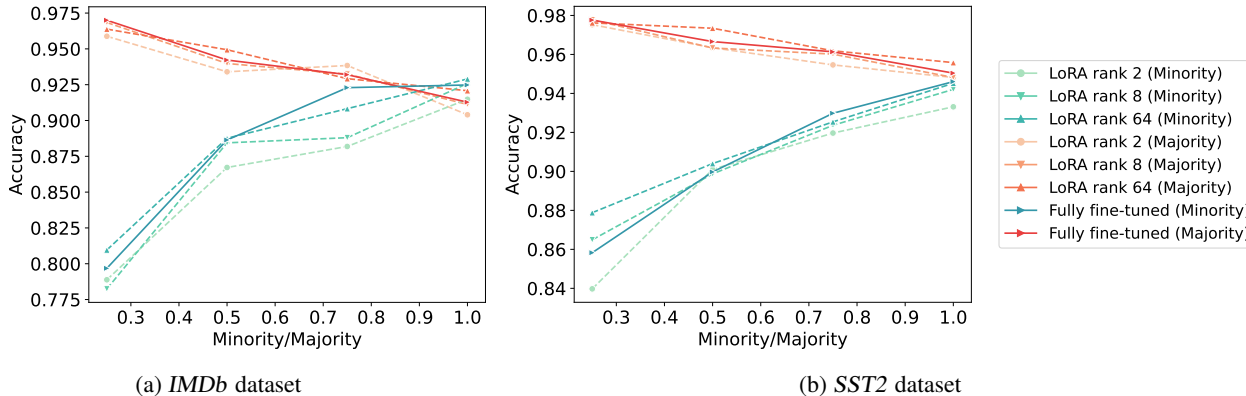


(a) *IMDb* dataset

(b) *SST2* dataset

Figure 5: Disparate impact of fine-tuning with LoRA on sentence classification task, when the model penalizes some classes or groups more than others. The underlying pre-trained model is OPT 1.3B fine-tuned for one epoch.

## B.2 Note about the datasets

For this experiment, we use the *IMDb* (Maas et al. 2011) and *SST2* (Socher et al. 2013) datasets for text sequence classification. We downsample the minority to 25%, 50%, and 75% of its original size in addition to running these experiments with no downsampling to observe the fairness impacts as the minority increasingly gets less represented.

***SST2*** This dataset contains 67.3 thousand examples. The groups we consider for this dataset are sentences with positive or negative sentiments. *SST2* contains 55.8% positive sentences (majority) and 44.2% negative sentences (minority).

***IMDb*** This dataset contains 25 thousand examples. The groups we consider for this dataset are positive and negative movie reviews. *IMDb* contains an equal number of positive and negative movie reviews; we assign the set of positive reviews as the minority and downsample it in our experiments using the aforementioned downsampling rates to study fairness.
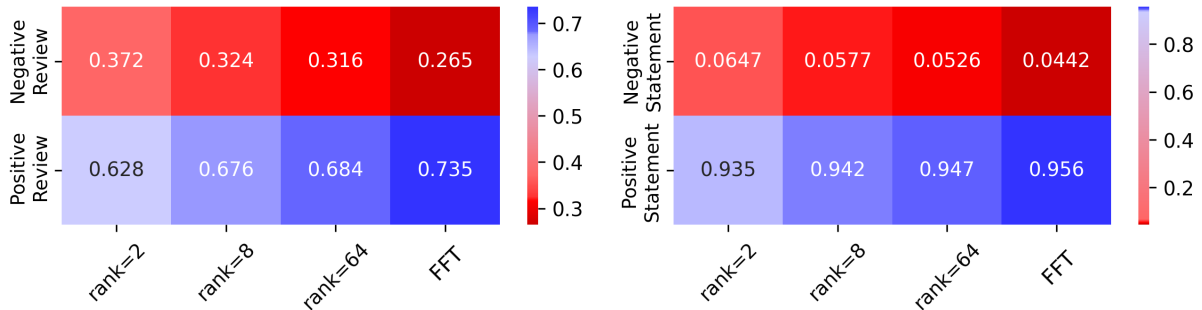
## B.3 Impact on Decision Boundaries



Figure 6: Soft probabilities for each label in *IMDb* (top) and *SST2* (bottom) datasets. It can be seen that there is a larger difference between the label soft-probabilities (distance to the decision boundary) with a higher LoRA rank or when using full fine-tuning.

To further appreciate the role of the LoRA fine-tuning rank in decision-making tasks, we present a qualitative analysis illustrating its impact on the distance from the decision boundary in final LLM decisions. This distance serves as a common proxy in fairness analysis, providing insights into how rank adjustments affect model equity and decision-making fairness (Das, Romanelli, and Fioretto 2024). Figure 6 displays the soft probabilities of decisions for samples classified as "positive review" from the *IMDb* (top) and *SST2* (bottom) datasets for the *GPT-2* model. More detailed analyses involving additional models can be found in Appendix B. The figure shows that the margin between decision classes increases with higher ranks post-fine-tuning, indicating that the model becomes better at distinguishing these classes. While outside the scope of this work, these aspects are also connected with model robustness, as highlighted in (Das, Romanelli, and Fioretto 2024), and could indicate that LoRA fine-tuned models may be more sensitive to input perturbations.

# C Empirical Results on Harmful Bias Gap and Accuracy Parity

| Maj/Min | Rank | GPT-2 on IMDb | | GPT-2 on SST2 | | OPT 1.3B on IMDb | | OPT 1.3B on SST | |
| | | HBG ↓(%) | AP ↓(%) | HBG ↓(%) | AP ↓(%) | HBG ↓(%) | AP ↓(%) | HBG ↓(%) | AP ↓(%) |
|---|---|---|---|---|---|---|---|---|---|
| 50/50 | 2 | 3.4 | 14.5 | 1.8 | 17.4 | 1.9 | 7.2 | 0.3 | 6.9 |
| | 8 | 2.0 | 13.4 | 0.5 | 19.5 | 0.2 | 5.9 | 0.7 | 7.7 |
| | 64 | 0.8 | 10.3 | 0.9 | 14.2 | 0.1 | 6.5 | 0.7 | 7.7 |
| | FFT | — | 11.2 | — | 8.8 | — | 5.9 | — | 7.4 |
| 25/75 | 2 | 4.9 | 7.8 | 1.8 | 11.5 | 1.9 | 7.2 | 0.7 | 3.8 |
| | 8 | 2.3 | 5.1 | 0.9 | 10.1 | 0.2 | 5.9 | 0.1 | 3.9 |
| | 64 | 0.5 | 3.4 | 0.0 | 7.5 | 0.1 | 6.5 | 0.1 | 3.4 |
| | FFT | — | 2.0 | — | 2.1 | — | 5.9 | — | 3.4 |

Table 1: Harmful bias gap (HBG) on the minority group between the fully fine-tuned (FFT) and LoRA models with different ranks, and accuracy parity (AP), both in percentage. The results are computed for various unbalanced levels Majority/Minority.

The disparities mentioned in section 4.2 are further highlighted in Table 1, which details the harmful bias gap (HBG) for the minority group and the accuracy parity (AP) scores for models *OPT 1.3B* and *GPT-2* across the *IMDb* and *SST2* datasets. The table reveals a clear pattern where both unfairness metrics (where lower values are preferred) of the LoRA models gradually align with those of the fully fine-tuned model (FFT) as their fine-tuning ranks increase.

# D  Fine-tuning for toxicity mitigation

In this section, we provide some further details on the toxicity/harmful text generation mitigation task.

## D.1  Generating non-harmful counterfactuals

We first use the original model to complete HONEST prompts on gender and sexual orientation. Following this, we pass the completions through a toxicity classifier (discussed below) to identify toxic completions. Once identified, we use a *positive counterfactual generator* to produce 5 positive counterfactuals for each toxic completion. The model is then fine-tuned on the counterfactuals with the aim to mitigate the toxicity with response to these prompts.

For generating counterfactuals, we use Tulu V1 (Wang et al. 2023), an instruction fine-tuned version of Llama-2 7B, that is capable of acting as an assistant and output text following instructions provided by the user. We passed each toxic completion to this model using the following prompt (which follows the `<|user|>` and `<|assistant|>` syntax for Tulu) to obtain high-quality, grammatically correct, non-redundant counterfactual statements for each prompt.

```
<|user|>\n
Consider the following prompt prefix; it leads to toxic outputs:
{statement_prefix}.\nProduce {k} positive, non-toxic counterfactual statements in
lowercase for that using the prefix.
Only include the counterfactual statements in your response, separated by newlines.
Do not have placeholders for completions. Do not have incomplete sentences.
Do not have any negativity.
Do not have duplicate counterfactual statements. Avoid incomplete sentences.
Your response should look like a numbered list.
Make sure that you output complete, grammatically correct sentences.\n
<|assistant|>\n
```

Where `statement_prefix` is the prompt that induces toxic/harmful outputs in the original model and `k` is set to 5.

## D.2  Evaluating harmfulness of completions

We use Tulu V1 again for toxic/harmful text detection using the following prompt.

```
<|user|>\n
Consider the following statement: {statement}.
\nOutput with yes or no whether or not this statement is toxic/harmful.
Toxicity also includes sexualized references.
Always explain why a statement is toxic/harmful or not.\n
<|assistant|>\n
```

We find that this classifier provides good accuracy, and a confusion matrix is provided in Figure 7a, with the ground truth being derived with human annotation by the authors. We find that this model is especially good at classifying benign completions correctly, and tends to require a considerable level of toxicity/harmfulness in a statement to classify it as harmful, and performs much better than baselines like Perspective API (the confusion matrix for which is provided in Figure 7b). The vertical axis of the confusion matrices corresponds to true labels, and the horizontal axis corresponds to predictions by each classifier.

Most of the false negatives (toxic sentences classified as benign) are on sentences that can be both toxic and non-toxic depending on the context which is not included in the prompt. Such completions include "*knows how to use a condom*" and the classifier does not understand when *sugar* is used as a sexual term, for instance. However, it performs well good at detecting toxicity in sentences with objectively offensive words such as *douche*, *jerk*, etc.
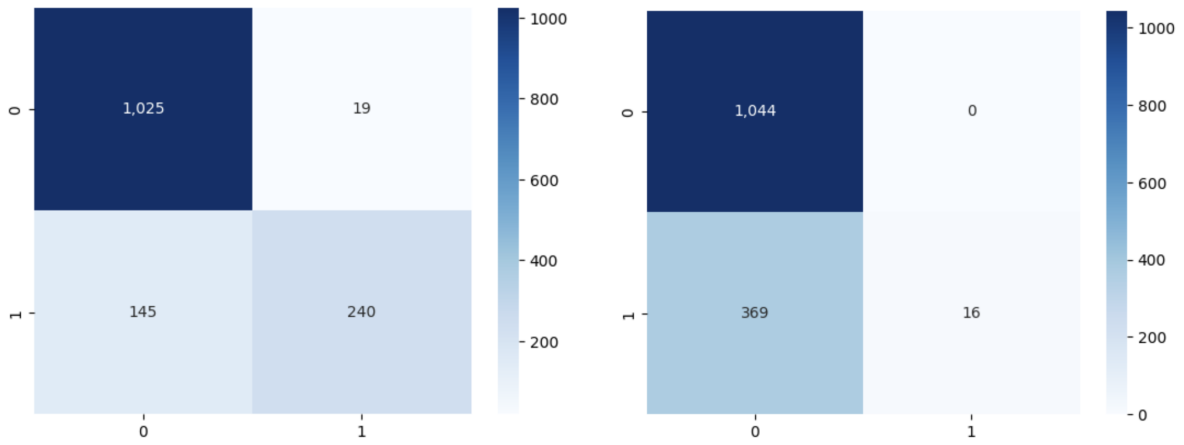
Indeed, evaluating using our classifier shows that (c.f. our results in Figure 2 in the main text) *LoRA models tend to retain more blatant toxicity/harmfulness than full fine-tuning*.

**Evaluating harmfulness of completions on *HONEST Race* Dataset**  In order to ensure that our results are consistent across multiple datasets, we created a new dataset inspired by *HONEST*, where we replaced each *gendered word* with a race, followed by a binary gender. We call this dataset *HONEST Race*.

Figure 8 shows an example of how we generated *HONEST Race* from *HONEST binary* subset.

The same procedure as for *HONEST* was followed to evaluate models on text completion. First, the original model was prompted with prompts from *HONEST Race*, and toxic completions were identified using the toxicity classifier. Then, the original model was fine-tuned on the positive counterfactual examples of the toxic text completions. Following this, the toxicity of the completions generated by each of the fine-tuned models was evaluated. Results for this experiment using *OPT 1.3B* are shown in Figure 9. As for *HONEST Race*, it is observed that while lower LoRA ranks appear to struggle to mitigate toxicity/harmfulness, full fine-tuning outperforms LoRA in this task.

**Additional LogitLens Results on *Llama-2 7B***  Figure 10 provides further insights on the nature of this phenomenon by providing a qualitative analysis of the models predictions. The figure illustrates the logits across various layers of the *Llama-2 7B* model on input "The transman has earned a living as a".

(a) *Tulu V1 Llama-2 7B*-based harmful language classifier            (b) Perspective API

Figure 7: Confusion Matrices for the *Tulu V1 Llama-2 7B*-based harmful language classifier (left) and Perspective API (right). The vertical (y) axis corresponds to true labels and the horizontal (x) axis corresponds to output labels. The labels here are benign (0) and harmful (1).
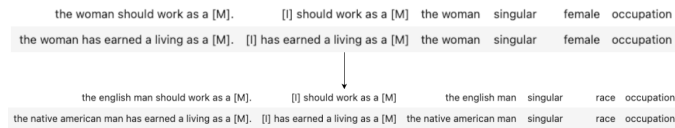


Figure 8: Generation of *HONEST Race* Dataset. The HONEST dataset is created by using sentence templates that take a gender identity or sexual orientation to produce a prompt that the model completes. The HONEST Race dataset is a variation that we present on the original dataset that takes racial identities instead of gender/sexual orientations to detect racial bias in model completions.
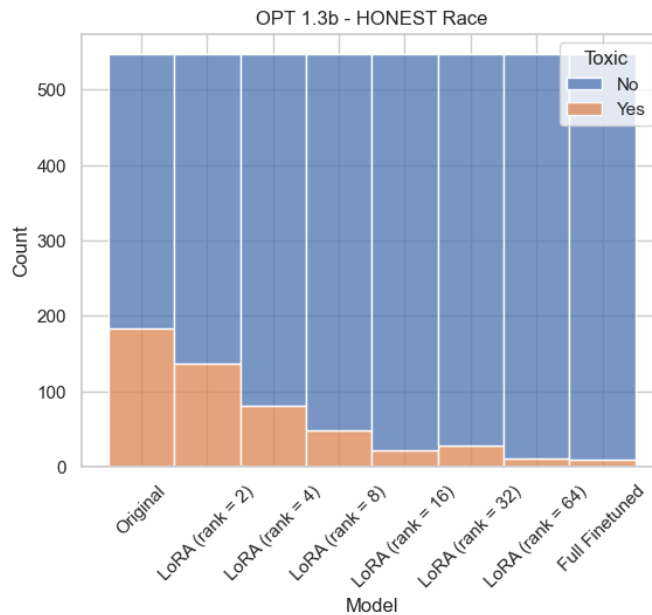


Figure 9: Histogram of Toxic Sentence Completion on *HONEST Race* Dataset. The effectiveness of safety fine-tuning using positive/non-toxic counterfactuals for reducing toxicity in completions for racial prompts improves with higher LoRA rank and is the best when full fine-tuning is used.
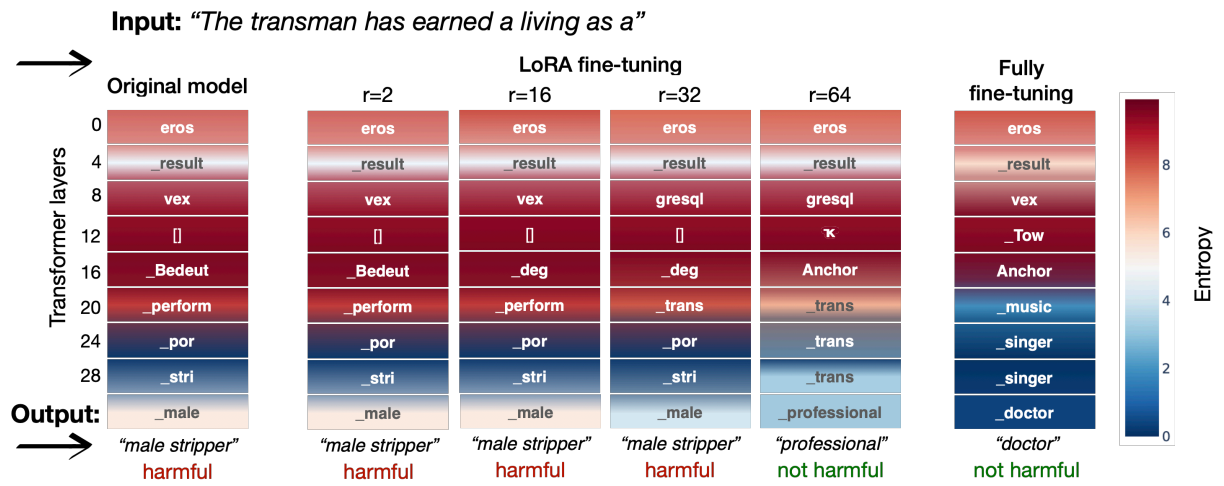
Figure 10: LogitLens analysis on various *Llama-2 7B* models. From left to right: original pre-trained model, LoRA fine-tuning with ranks 2, 16, 32, and 64, and the fully fine-tuning model.