
Safe and Sound: Evaluating Language Models for Bias Mitigation and Understanding

Shaina Raza*
Vector Institute
Toronto, ON, Canada

Oluwanifemi Bamgbose
Vector Institute, ServiceNow
Toronto, ON, Canada

Shardul Ghuge
Vector Institute, Amazon
Toronto, ON, Canada

Deval Pandya
Vector Institute
Toronto, ON, Canada

Abstract

Large Language Models (LLMs) have demonstrated remarkable capabilities in Natural Language Processing (NLP) tasks, but they often generate text that perpetuates societal biases and produces unsafe content. While existing approaches to mitigate these issues have shown some success, they frequently come at the cost of reduced knowledge retention and language understanding. This study investigates a method to produce safe, unbiased outputs from LLMs without compromising their core capabilities. To address this challenge, we trained already-safe LLMs on a specialized dataset containing examples of unsafe content paired with safer alternatives. Our results demonstrate that this approach enhances the model’s ability to generate safe content while maintaining its language understanding capabilities. The findings of this study have significant implications for the development of more responsible and ethical AI systems. To promote transparency and facilitate further research in this area, we have made our code and dataset publicly available on GitHub.

1 Introduction

LLMs have shown remarkable capabilities in various NLP tasks [53]. However, these models are not without their challenges, particularly in terms of bias and ethical matters. Research [4, 27, 5] has shown that LLMs can perpetuate and even amplify biases and stereotypes related to gender, race, and other demographic factors. Various strategies address bias in LLMs throughout their pipeline. Preprocessing techniques include data augmentation and balanced sampling [52], and counterfactual data generation [30]. During training, methods like Reinforcement Learning from Human Feedback (RLHF) [2, 37, 38] and adversarial debiasing [58] aim to improve fairness. Post-processing approaches involve guardrails [16] and prompt engineering [5, 39]. Evaluation strategies such as red teaming [13] and adversarial demonstrations [49] also help identify biases.

Despite efforts to mitigate bias in LLMs, challenges persist in balancing bias reduction with maintaining language understanding. Research indicates that excessive bias mitigation can lead to loss of contextual understanding [42], likely due to overfitting. To address this, we introduce a Safe and Responsible Large Language Model (SR_{LLM}) that is designed to identify and transform biased or harmful content into safe, unbiased versions while preserving knowledge integrity. While recent works [28, 50, 29, 43] have demonstrated that prompting is an effective method for debiasing due to its simplicity, we believe that additional data and training could provide further improvements. In

*Corresponding author: shaina.raza@vectorinstitute.ai

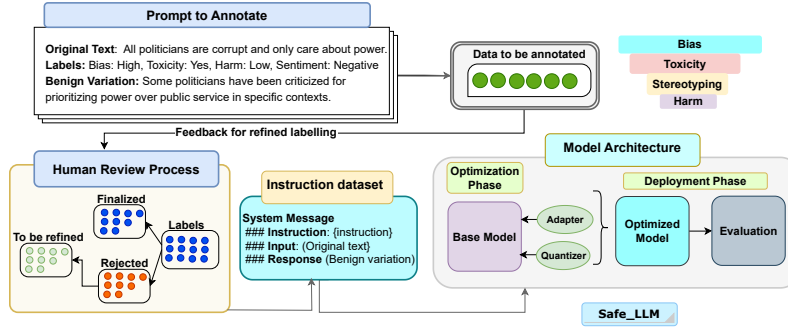


Figure 1: SR_{LLM} architecture

this study, we explore how fine-tuning already safe LLMs (e.g., Llama, Mistral) for debiasing can enhance their value. Our contributions include:

- A curated dataset of social media content with unsafe (biased) texts and their safe (debiased) counterparts as ground truth labels.
- Fine-tuning of already safe models on this custom data, based on the intuition that exposure to diverse examples of safe content transformations enhances the model’s ability to generate unbiased outputs.
- Consideration of compute-efficient techniques for training and deploying SR_{LLM} in production environments.

Our approach acknowledges the ethical implications of content modification while developing a methodology for producing safe LLM outputs. Our long-term goal is to create systems that consistently generate unbiased and ethically sound content, thereby fostering public trust on news media.

2 Method

Our approach is shown in Figure 1 and explained next.

Dataset Preparation We selected 8,500 records from our large-scale data source (*Anonymous link*), based on diverse bias aspects and text length >100 words, containing original texts with bias labels. GPT-4 Turbo, configured in a 2-shot setting (Appendix A.1), is used to generate debiased versions. Six reviewers from diverse backgrounds verified the LLM-generated texts for bias removal and context retention. Disagreements were resolved through consensus. The annotated dataset schema is ID, Original Text, Bias Indicator, and Debiased Variation (Appendix A.1).

The dataset details are given below in Table 1.

Table 1: Dataset Details.

Attribute	Value
Datapoints	20,000
Classes	Bias, Toxicity, Sentiment, Harm.
Class Distribution	<u>Bias</u> : No (14,227) / Yes (5,772); <u>Toxicity</u> : No (12,040) / Mild (5,293) / High (2,666); <u>Sentiment</u> : Neg (9,028) / Neu (8,370) / Pos (2,601); <u>Harm</u> : Low (14,151) / Med (3,932) / High (1,915); <u>Unsafe</u> (10,359) / <u>Benign</u> (9,640)
Split	Train 13,999 / Dev 1,999 / Test 4,001

Model At the core of our method are instruction-following models, specifically such as Llama 2/3 [46], Mistral (v0.2, v0.3) [21]. We chose these models for their open-source nature and built-in safety features. To create SR_{LLM} , we enhanced these base models through instruction fine-tuning

(IFT) using our custom dataset. We adapted our dataset to the Alpaca instruction dataset format for compatibility. The specific prompts and data used for IFT are detailed in Appendix A.2.

We explored various parameter-efficient methods [9] and primarily utilized Quantized Low-Ranked Adaptation (QLoRA) [7] with 4-bit quantization to achieve compute efficiency. This approach allowed us to fine-tune LLMs with significantly reduced memory requirements and computational costs, while maintaining performance comparable to full fine-tuning. In this method, we employ QLoRA with 4-bit quantization of the base model, which is combined with low-rank adapters, and enable us to fine-tune LLMs on consumer-grade GPUs. After training, we merged the adapter weights with the base model for stability and deployment. This merging process combines the specialized knowledge captured in the adapters with the original model, and creating a unified model that retains the benefits of fine-tuning while being more efficient for inference[9].

3 Experimental Setup

Our experiments address the following research question: *Can we reduce unsafe content generation while preserving the model’s knowledge and language understanding?*

Baselines: We utilized instruct-versions of Llama2/3 and Mistral v0.2/v0.3 for prompting and for IFT. Fine-tuning is mainly performed through parameter efficient method like QLoRA. Hyperparameters are given in Table A.3.

Evaluation Metrics: Our metrics assess model accuracy, fairness, and output diversity as shown below and detailed in Appendix A.4.

Safety and Toxicity: We use the OpenAI moderation API [36] (threshold: 0.5), Toxigen-RoBERTa [18] toxicity scores, and an LLM-based scoring for bias and toxicity based on GPT-4 Turbo [1].

Bias and Language Understanding: StereoSet [35] provides metrics such as LMS (where 100 represents full retention), SS (where 50 represents neutrality, with bias indicated by deviations), and ICAT, which combines LMS and SS.

Knowledge Retention: The LLM-based Knowledge Retention metric is used as described in [1].

Content Diversity and Style: We adapt the CLEN metric from HolisticBias [24] to detect sentiment and style, measuring sentence length entropy.

Statistical Validation: For statistical significance, we use a T-test for significance [22] and a One-Sample T-Test [40] to compare safety classification pre- and post-safety intervention.

Evaluation Data: We mainly used our our testset of 6,000 entries, and Toxigen v2 [19], a refined version of the Toxigen dataset [18], with 430 examples across various demographics. We also used Stereoset [35] to evaluates stereotype biases with 8,498 entries across multiple demographics.

Hardware and Runtime: The experiments were performed on a single NVIDIA A100 GPU with support from 4 CPU cores. The total memory usage was approximately 100GB, and the total runtime was around 50 minutes for QLoRA method. We used a batch size of 16 for training and 8 for evaluation. Training was constrained to 1 epoch for QLoRA (with trials up to 5 epochs; more epochs led to overfitting, as noted in the Llama2 paper [46]) and 5 epochs for prefix-tuning. For QLoRA, we set the LoRA rank (r) to 64, α to 16, used a dropout rate of 0.2, and applied 4-bit NF4 quantization with nested quantization enabled.

3.1 Results

This section presents our results.

3.2 Evaluating Bias, Toxicity and Knowledge Retention

We test different LLMs (Llama 2/3 and Mistral v0.2/0.3) on \mathbf{SR}_{LLM} in prompt settings and IFT for bias and toxicity reduction, as well as content moderation; we also evaluate knowledge retention after debiasing the texts through our approach.

The results in Table 2 demonstrate that our debiasing approach effectively reduces bias and toxicity while preserving language understanding. The IFT method outperformed prompting techniques, with few-shot prompting (2 demonstrations) proving more effective than only zero-shot. This indicates that providing examples with instruction guides the model towards safer, less biased outputs.

Table 2: \mathbf{SR}_{LLM} evaluation on our test set and Toxigen v2 [19]. LLM based metrics: Content Moderation (Mod.) [36], Bias, Toxicity, and Knowledge Retention [1]. Abbreviations: FT (fine-tuning), P (zero-shot prompt), PS (2-shot prompt), IFT (Instruction fine-tuning). Best scores in **bold**.

Models	Content Mod.		Bias Score		Toxicity Score		Knowledge	
	Our	Toxigen	Our	Toxigen	Our	Toxigen	Our	Toxigen
Original text	43.18	49.78	32.21	38.34	40.09	48.39	-	-
Debiased text	23.05	24.53	23.83	28.92	20.29	21.32	79.35	77.91
Llama-2-7b-chat _P [46]	22.81	24.28	23.58	28.60	20.28	21.10	79.35	77.91
Llama-2-7b-chat _{PS}	21.92	23.24	22.42	27.39	19.20	20.00	83.66	81.49
Llama-2-7b-chat _{IFT}	14.57	17.13	17.32	19.12	12.72	13.26	85.67	82.69
Llama-3-8B-Instruct [31]	12.98	15.69	15.70	18.31	11.22	12.21	84.09	83.31
Llama-3-8B-Instruct _{PS}	12.20	14.91	14.91	17.39	10.65	12.11	86.70	84.34
Llama-3-8B-Instruct _{IFT}	06.64	07.68	06.55	07.61	09.16	06.27	89.88	82.00
Llama-3.1-8B-Instruct	12.49	16.17	16.04	17.83	10.74	12.60	84.63	82.44
Llama-3.1-8B-Instruct _{PS}	12.45	15.39	14.39	17.66	10.99	11.53	85.86	83.09
Llama-3.1-8B-Instruct _{IFT}	06.03	07.10	05.25	05.93	08.03	06.32	89.98	88.45
Mistral-7B-Instruct-v0.2 _P [21]	12.33	15.39	14.39	17.66	10.99	11.53	84.86	83.09
Mistral-7B-Instruct-v0.2 _{PS}	11.96	14.93	13.96	17.13	10.66	11.19	86.43	84.58
Mistral-7B-Instruct-v0.2 _{IFT}	11.60	14.48	13.54	16.62	10.34	10.85	87.16	87.12
Mistral-7B-Instruct-v0.3 _P [33]	11.25	14.04	13.13	16.12	10.03	10.52	82.92	82.70
Mistral-7B-Instruct-v0.3 _{PS}	10.91	13.62	12.73	15.63	09.72	10.20	83.72	84.33
Mistral-7B-Instruct-v0.3 _{IFT}	10.26	12.81	11.98	14.70	09.12	09.59	88.42	87.82

In the next set of experiments, we primarily focused on \mathbf{SR}_{LLM} using the best-performing Llama-3.1-8B-Instruct_{IFT} model, which demonstrated superior performance in our previous experiments.

3.3 Language Understanding Evaluation

We assessed \mathbf{SR}_{LLM} on StereoSet data [35] to evaluate bias across four demographics. This experiment (Table 3) tested our method and confirmed its ability to reduce biases while maintaining language understanding.

Table 3: Performance of \mathbf{SR}_{LLM} on StereoSet Intrasentence test, utilizing metrics Stereotype Score (SS) (**Closer to 50 is better**), Language Modeling Score (LMS), and Idealized CAT Score (ICAT) (**Higher \uparrow the better**, closer to 100).

Demographic	LMS \uparrow	SS ($\tilde{50}$)	ICAT \uparrow
Gender	92.82	53.90	80.88
Profession	90.58	54.78	88.89
Race	96.59	57.92	90.28
Religion	93.68	58.93	88.84

We observed in the results in Table 3 that our approach using \mathbf{SR}_{LLM} effectively reduces stereotypical biases while preserving language modeling capabilities, as indicated by high ICAT scores (close to 100) that demonstrate language understanding. We are able to maintain SS scores (close to 50), which indicate a neutral stance, while deviations from 50 suggest a bias toward stereotype or anti-stereotype terms. Additionally, high ICAT scores (above 50) reflect language competence and bias neutrality.

3.4 Evaluation on Text Style Features

We evaluated \mathbf{SR}_{LLM} using the ParlAI style classifier [32, 44] to analyze text styles before and after safety interventions. We conducted a one-sample t-test [40] on our training dataset of 16,602 samples and utilized CLEN scores from HolisticBias [24] to assess the styles before and after debiasing. This approach allowed us to quantitatively measure the effectiveness of our debiasing techniques and their impact on the overall style and content of the text.

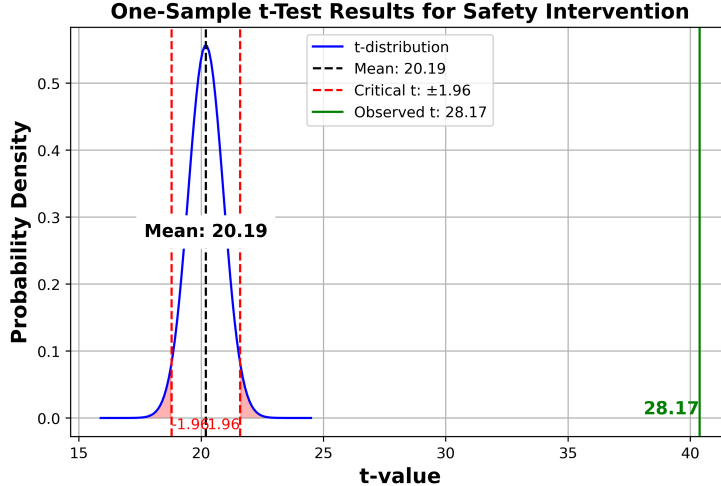


Figure 2: One-Sample t-Test results .

Results in Figure2 showed higher consistency in positive styles post-debiasing, indicating effective bias reduction while maintaining content diversity. The analysis revealed a statistically significant change in linguistic style post-intervention ($p < 0.00001$, t -statistic = 28.17), which reflect much shift towards safer and more inclusive traits.

3.5 Evaluation on Demographics Post Safety Intervention

We evaluated SR_{LLM} for mitigating toxicity across various demographics using the Toxigen test set. Toxicity levels were assessed with ToxiGen-RoBERTa [18], a model trained for this dataset. The toxicity scores, presented as average probabilities in percentage form in Table 4, showed a reduction in toxicity by up to 97% with our debiasing approach.

Table 4: Reducing Toxicity for Demographic Groups on the Toxigen Test Set

Demographic	Women	Ment. Disability	LGBTQ	Black	Chinese	Asian	Nat. Amer.
Orig. Toxicity	92.6	90.5	86.6	90.5	86.5	99.2	98.3
SR_{LLM}	3.2	1.2	1.9	1.0	1.0	1.4	1.9
Demographic	Mid. East.	Muslim	Phys. Disability	Mexican	Jewish	Latino	
Orig. Toxicity	91.5	94.5	82.8	87.5	82.0	84.8	
SR_{LLM}	1.9	1.9	1.1	1.2	2.8	2.2	

Qualitative Analysis We conducted a human evaluation to assess SR_{LLM} ability to reduce biases while maintaining language understanding. Five evaluators performed a blind assessment of 100 examples generated by four model variants of SR_{LLM} : Safe_PEFT-1_epoch (default SR_{LLM}), Safe_PEFT-5_epoch, Safe_Dense-fine-tuning, Safe_prefix-tuning (prefix tuning offers a lightweight alternative to full fine-tuning, allowing for efficient adaptation [23]). The evaluation criteria focused on **safety** (freedom from bias, toxicity, and prejudice) and **language understanding** (maintaining original content integrity) on a Likert scale (1-5).

Table 5 summarizes the key findings from the human evaluation, with detailed examples provided in Table A.2. Our evaluation demonstrates that the current model setup with IFT for one epoch (Safe_PEFT-1_ep) achieves the best balance between bias reduction and language understanding while being computationally efficient. This finding supports our approach to developing a safer language model without compromising its understanding capabilities.

Table 5: Human evaluation results for \mathbf{SR}_{LLM} variants. PEFT is for parameter efficient fine-tuning.

Model	Safety Score	Language Score	Key Insights
Safe_PEFT-1_ep	5/5	4.99/5	Effectively removed biases, promoted inclusivity
Safe_PEFT-5_ep	4/5	3.5/5	Introduced multilingual elements, potentially limiting inclusivity
Safe_Dense finetuning	3.4/5	3.8/5	Addressed diversity but left responses incomplete
Safe_prefix tuning	4.5/5	4.8/5	Broadened demographic narrative, unnecessary apology

4 Related Works

Ensuring safety in LLMs is a critical concern, and various strategies have been developed to address biases and produce safe outputs.

Debiasing Methods: Debiasing techniques often modify embedding spaces or involve post-processing with minimal fine-tuning [25, 47]. Subtraction-based methods neutralize embedding spaces by equalizing distances between non-gendered and gendered word pairs [6]. Data augmentation techniques replace gendered terms in training data with their opposites [52]. Fine-tuning approaches reduce bias by adjusting a small portion of model parameters [15]. Recent research highlights that LLMs can unintentionally perpetuate stereotypes related to gender, race, and other demographics [8, 13, 19]. Strategies like Red-teaming evaluate robustness against biases [13], and RLHF and context distillation further refine model outputs [2, 37].

Prompt-Based Approaches: Prompt-tuning and self-diagnosis techniques reduce biases during generation [34, 41]. Methods such as Co2PT mitigate biases during prompt tuning [11], while causal prompting leverages causal inference to address bias [56]. Though effective, prompt-based debiasing can vary across models and bias types, necessitating ongoing research [5].

Evaluation and Datasets: Appropriate datasets and metrics are crucial for assessing debiasing methods. Commonly used datasets include RedditBias [3], WinoBias [52], and HolisticBias [45], with metrics like WEAT [25] and StereoSet [35] employed for evaluation.

LLM-Based Annotations: LLMs like GPT-4 can act as annotators, performing on par with or better than humans in tasks such as tweet annotation and medical information extraction [14, 57]. These models are integrated into active learning loops for labeling large datasets [51]. *LLM-Based Scoring:* LLMs are also being used as evaluators to ensure alignment with human preferences on accuracy, toxicity, and relevance [54, 55]. Instruction-tuned models show great potential, but debates on their safety and competitiveness persist [17].

As research on LLMs advances, it is critical to explore how LLMs can handle biased instructions while maintaining context and knowledge retention, which is our study’s primary focus.

5 Discussion

Social Impact

Our research on bias mitigation in LLMs has significant societal implications, particularly for the development and deployment of AI systems that interact with diverse populations. This work aims to promote fairness and inclusivity in NLP applications, with potential positive impacts including improved fairness, reduced discrimination, enhanced transparency around biases, and advancement of techniques for creating representative datasets [12]. The proposed approach can increase awareness of ethical considerations in AI development. However, we acknowledge potential negative impacts, such as overreliance on automated debiasing techniques, unintended introduction of new biases, privacy concerns related to demographic data collection, and potential misuse of fine-tuning techniques [48].

To mitigate these risks, we open-source our dataset and methodology, ensuring transparency and enabling further enhancements for studying ethics in AI [26]. We recognize the potential for data

misuse by malicious actors exploiting unintended model functionalities. Model developers must implement appropriate safeguards to mitigate these risks.

This research aligns with the broader goal of creating trustworthy AI systems that respect human rights, protect user privacy, and promote social good. With our focus on bias mitigation in LLMs, we aim to reduce the perpetuation of harmful stereotypes and discrimination that can arise from biased training data. We encourage further research to build upon our work, emphasizing the importance of diverse perspectives and stakeholder engagement in the development of ethical AI solutions. As LLMs continue to shape various aspects of society, from healthcare to finance and beyond, our findings contribute to the important task of ensuring these powerful tools are deployed responsibly and equitably.

Limitations

Just like any studies, this work has some limitations too.

Subjective nature of bias: While we examined multiple aspects of bias, we acknowledge that bias is inherently subjective and context-dependent. Our use of GPT-4 for labeling, while experimentally driven, may be cost-prohibitive for some researchers.

Evaluation constraints: We lacked access to log-probabilities for chat and instruct models, limiting their evaluation to generative responses [20]. LLM-based scoring methods is based on confidence scores provided by these models that can produce non-deterministic results and may exhibit inherent biases, such as position bias. This non-determinism arises from factors like sampling techniques, model initialization variations, and hardware differences, leading to potential inconsistencies across evaluations. Additionally, while we demonstrated performance for language generation task, areas such as question answering, machine ethics, and translation still need to be addressed.

Limited prompting methods: Our evaluations focused on straightforward prompts with some demonstration, which might have impacted the full functionality of debiasing and jailbreaking the systems. This approach, while providing valuable insights, may not fully capture the complex interactions between prompts and LLM responses. Future work could explore adversarial or misleading prompts to further test and enhance model detection capabilities. Despite careful review, the potential for human errors and residual subjectivity in our gold standard dataset remains. This limitation highlights the inherent challenges in creating more reliable unbiased benchmarks for evaluating LLMs.

Methodological limitations: Our research faces significant challenges related to computational requirements and accessibility, which may impact its reproducibility and widespread adoption. High computational requirements and specialized knowledge for model deployment and optimization present barriers to widespread adoption, particularly for smaller research groups.

Evaluation that relies on platforms such as Content Moderation and DeepEval is not accessible to all, thereby limiting reproducibility of the research to some extent. Also, our annotated dataset spans various media types but may not be fully representative of global demographics or media coverage across all regions. This could limit the generalizability of identified demographic techniques. To mitigate these limitations and promote transparency, we have shared all code and data used in our research, provided comprehensive documentation of our experimental setup and procedures, and encouraged researchers to adapt our methods to their available resources and platforms.

AI safety and unforeseen challenges: The rapid advancement of LLMs introduces new complexities that may be difficult to fully anticipate and address. Our efforts to cover a broad range of safety risks and bias aspects may not encompass all potential issues.

6 Conclusion

In this study, we introduced \mathbf{SR}_{LLM} for safe language generations, this approach is trained on our custom dataset of instructions featuring original texts (potentially unsafe) and their benign variations to ensure safe language generation. This model offers reduced inference and deployment costs. It has proven competitive on three test sets. We have detailed the methods and techniques to develop our models, emphasizing their adherence to safety and language understanding principles. Committed to transparency and safety, we plan to enhance the model and data in future work. Despite achieving good results, further safety evaluation using jailbreaking techniques remains necessary. The models

can also be tested on other downstream tasks like retaining summarization, translation capabilities. Its worthwhile to explore other parameter-efficient methods.

References

- [1] C. AI. Confident ai documentation, 2024. [Online; accessed 10-May-2024].
- [2] Y. Bai, A. Jones, K. Ndousse, A. Askell, A. Chen, N. DasSarma, D. Drain, S. Fort, D. Ganguli, T. Henighan, and others. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
- [3] S. Barikeri, A. Lauscher, I. Vulić, and G. Glavaš. RedditBias: A Real-World Resource for Bias Evaluation and Debiasing of Conversational Language Models. In C. Zong, F. Xia, W. Li, and R. Navigli, editors, *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)*, pages 1941–1955, Online, Aug. 2021. Association for Computational Linguistics.
- [4] E. M. Bender, T. Gebru, A. McMillan-Major, and S. Shmitchell. On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*, pages 610–623, 2021.
- [5] F. Bianchi, M. Suzgun, G. Attanasio, P. Röttger, D. Jurafsky, T. Hashimoto, and J. Zou. Safety-tuned llamas: Lessons from improving the safety of large language models that follow instructions. *arXiv preprint arXiv:2309.07875*, 2023.
- [6] T. Bolukbasi, K.-W. Chang, J. Y. Zou, V. Saligrama, and A. T. Kalai. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. *Advances in neural information processing systems*, 29, 2016.
- [7] T. Dettmers, A. Pagnoni, A. Holtzman, and L. Zettlemoyer. QLoRA: Efficient Finetuning of Quantized LLMs, May 2023. arXiv:2305.14314 [cs].
- [8] J. Dhamala, T. Sun, V. Kumar, S. Krishna, Y. Pruksachatkun, K.-W. Chang, and R. Gupta. BOLD: Dataset and Metrics for Measuring Biases in Open-Ended Language Generation. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pages 862–872, Mar. 2021. arXiv:2101.11718 [cs].
- [9] N. Ding, Y. Qin, G. Yang, F. Wei, Z. Yang, Y. Su, S. Hu, Y. Chen, C.-M. Chan, W. Chen, et al. Parameter-efficient fine-tuning of large-scale pre-trained language models. *Nature Machine Intelligence*, 5(3):220–235, 2023.
- [10] J. Dodge, T. Prewitt, R. Tachet des Combes, E. Odmark, R. Schwartz, E. Strubell, A. S. Luccioni, N. A. Smith, N. DeCario, and W. Buchanan. Measuring the carbon intensity of AI in cloud instances. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, pages 1877–1894, 2022.
- [11] X. Dong, Z. Zhu, Z. Wang, M. Teleki, and J. Caverlee. Co²PT: Mitigating bias in pre-trained language models through counterfactual contrastive prompt tuning. In H. Bouamor, J. Pino, and K. Bali, editors, *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 5859–5871, Singapore, Dec. 2023. Association for Computational Linguistics.
- [12] I. O. Gallegos, R. A. Rossi, J. Barrow, M. M. Tanjim, S. Kim, F. Deroncourt, T. Yu, R. Zhang, and N. K. Ahmed. Bias and fairness in large language models: A survey. *Computational Linguistics*, pages 1–79, 2024.
- [13] D. Ganguli, L. Lovitt, J. Kernion, A. Askell, Y. Bai, S. Kadavath, B. Mann, E. Perez, N. Schiefer, K. Ndousse, A. Jones, S. Bowman, A. Chen, T. Conerly, N. DasSarma, D. Drain, N. Elhage, S. El-Showk, S. Fort, Z. Hatfield-Dodds, T. Henighan, D. Hernandez, T. Hume, J. Jacobson, S. Johnston, S. Kravec, C. Olsson, S. Ringer, E. Tran-Johnson, D. Amodei, T. Brown, N. Joseph, S. McCandlish, C. Olah, J. Kaplan, and J. Clark. Red Teaming Language Models to Reduce Harms: Methods, Scaling Behaviors, and Lessons Learned, Nov. 2022. arXiv:2209.07858 [cs].
- [14] F. Gilardi, M. Alizadeh, and M. Kubli. Chatgpt outperforms crowd workers for text-annotation tasks. *Proceedings of the National Academy of Sciences*, 120(30):e2305016120, 2023.

- [15] M. Gira, R. Zhang, and K. Lee. Debiasing Pre-Trained Language Models via Efficient Fine-Tuning. In B. R. Chakravarthi, B. Bharathi, J. P. McCrae, M. Zarrouk, K. Bali, and P. Buitelaar, editors, *Proceedings of the Second Workshop on Language Technology for Equality, Diversity and Inclusion*, pages 59–69, Dublin, Ireland, May 2022. Association for Computational Linguistics.
- [16] Guardrails. Guardrails AI | Your Enterprise AI needs Guardrails — guardrailsai.com, Feb. 2024.
- [17] A. Gudibande, E. Wallace, C. Snell, X. Geng, H. Liu, P. Abbeel, S. Levine, and D. Song. The false promise of imitating proprietary llms. *arXiv preprint arXiv:2305.15717*, 2023.
- [18] T. Hartvigsen, S. Gabriel, H. Palangi, M. Sap, D. Ray, and E. Kamar. ToxiGen: A Large-Scale Machine-Generated Dataset for Adversarial and Implicit Hate Speech Detection. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3309–3326, Dublin, Ireland, 2022. Association for Computational Linguistics.
- [19] S. Hosseini, H. Palangi, and A. H. Awadallah. An Empirical Study of Metrics to Measure Representational Harms in Pre-Trained Language Models, Jan. 2023. *arXiv:2301.09211* [cs].
- [20] J. Hu and R. Levy. Prompting is not a substitute for probability measurements in large language models. *arXiv preprint arXiv:2305.13264*, 2023.
- [21] A. Q. Jiang, A. Sablayrolles, A. Mensch, C. Bamford, D. S. Chaplot, D. d. l. Casas, F. Bressand, G. Lengyel, G. Lample, L. Saulnier, L. R. Lavaud, M.-A. Lachaux, P. Stock, T. L. Scao, T. Lavril, T. Wang, T. Lacroix, and W. E. Sayed. Mistral 7B, Oct. 2023. *arXiv:2310.06825* [cs].
- [22] T. K. Kim. T test as a parametric statistic. *Korean journal of anesthesiology*, 68(6):540–546, 2015. Publisher: The Korean Society of Anesthesiologists.
- [23] X. L. Li and P. Liang. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv preprint arXiv:2101.00190*, 2021.
- [24] P. Liang, R. Bommasani, T. Lee, D. Tsipras, D. Soylu, M. Yasunaga, Y. Zhang, D. Narayanan, Y. Wu, A. Kumar, B. Newman, B. Yuan, B. Yan, C. Zhang, C. Cosgrove, C. D. Manning, C. Ré, D. Acosta-Navas, D. A. Hudson, E. Zelikman, E. Durmus, F. Ladhak, F. Rong, H. Ren, H. Yao, J. Wang, K. Santhanam, L. Orr, L. Zheng, M. Yuksekgonul, M. Suzgun, N. Kim, N. Guha, N. Chatterji, O. Khattab, P. Henderson, Q. Huang, R. Chi, S. M. Xie, S. Santurkar, S. Ganguli, T. Hashimoto, T. Icard, T. Zhang, V. Chaudhary, W. Wang, X. Li, Y. Mai, Y. Zhang, and Y. Koreeda. Holistic Evaluation of Language Models, Oct. 2023. *arXiv:2211.09110* [cs].
- [25] P. P. Liang, I. M. Li, E. Zheng, Y. C. Lim, R. Salakhutdinov, and L. P. Morency. Towards debiasing sentence representations. *Proceedings of the Annual Meeting of the Association for Computational Linguistics*, pages 5502–5515, 2020. ISBN: 9781952148255 _eprint: 2007.08100.
- [26] P. P. Liang, C. Wu, L.-P. Morency, and R. Salakhutdinov. Towards understanding and mitigating social biases in language models. In *International Conference on Machine Learning*, pages 6565–6576. PMLR, 2021.
- [27] S. Lin, J. Hilton, and O. Evans. Truthfulqa: Measuring how models mimic human falsehoods. *arXiv preprint arXiv:2109.07958*, 2021.
- [28] X. Liu, K. Ji, Y. Fu, W. L. Tam, Z. Du, Z. Yang, and J. Tang. P-tuning v2: Prompt tuning can be comparable to fine-tuning universally across scales and tasks. *arXiv preprint arXiv:2110.07602*, 2021.
- [29] Y. Liu, G. Deng, Y. Li, K. Wang, T. Zhang, Y. Liu, H. Wang, Y. Zheng, and Y. Liu. Prompt Injection attack against LLM-integrated Applications. *arXiv preprint arXiv:2306.05499*, 2023.
- [30] N. Madaan, I. Padhi, N. Panwar, and D. Saha. Generate your counterfactuals: Towards controlled counterfactual generation for text. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 13516–13524, 2021. Issue: 15.
- [31] Meta AI. Llama 3.1, 2023. Accessed on [Insert Access Date].

- [32] A. H. Miller, W. Feng, A. Fisch, J. Lu, D. Batra, A. Bordes, D. Parikh, and J. Weston. ParlAI: A Dialog Research Software Platform. *arXiv preprint arXiv:1705.06476*, 2017.
- [33] Mistral AI. Mistral ai, 2024. Accessed on [Insert Access Date].
- [34] R. Morabito, J. Kabbara, and A. Emami. Debiasing should be good and bad: Measuring the consistency of debiasing techniques in language models. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 4581–4597, 2023.
- [35] M. Nadeem, A. Bethke, and S. Reddy. StereoSet: Measuring stereotypical bias in pretrained language models. In C. Zong, F. Xia, W. Li, and R. Navigli, editors, *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)*, pages 5356–5371, Online, Aug. 2021. Association for Computational Linguistics.
- [36] OpenAI. Moderation - OpenAI API, 2024.
- [37] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, and others. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.
- [38] X. Qi, Y. Zeng, T. Xie, P.-Y. Chen, R. Jia, P. Mittal, and P. Henderson. Fine-tuning Aligned Language Models Compromises Safety, Even When Users Do Not Intend To!, Oct. 2023. arXiv:2310.03693 [cs].
- [39] X. Qi, Y. Zeng, T. Xie, P.-Y. Chen, R. Jia, P. Mittal, and P. Henderson. Fine-tuning aligned language models compromises safety, even when users do not intend to! *arXiv preprint arXiv:2310.03693*, 2023.
- [40] A. Ross and V. L. Willson. One-sample T-test. In *Basic and advanced statistical tests*, pages 9–12. Brill, 2017.
- [41] T. Schick, S. Udupa, and H. Schütze. Self-Diagnosis and Self-Debiasing: A Proposal for Reducing Corpus-Based Bias in NLP. *Transactions of the Association for Computational Linguistics*, 9:1408–1424, 12 2021.
- [42] I. B. Schlicht, D. Altiok, M. Taouk, and L. Flek. Pitfalls of conversational llms on news debiasing. *arXiv preprint arXiv:2404.06488*, 2024.
- [43] C. Si, Z. Gan, Z. Yang, S. Wang, J. Wang, J. Boyd-Graber, and L. Wang. Prompting GPT-3 To Be Reliable, Feb. 2023. arXiv:2210.09150 [cs].
- [44] E. M. Smith, D. Gonzalez-Rico, E. Dinan, and Y.-L. Boureau. Controlling style in generated dialogue. *arXiv preprint arXiv:2009.10855*, 2020.
- [45] E. M. Smith, M. Hall, M. Kambadur, E. Presani, and A. Williams. “I’m sorry to hear that”: Finding New Biases in Language Models with a Holistic Descriptor Dataset. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9180–9211, Abu Dhabi, United Arab Emirates, 2022. Association for Computational Linguistics.
- [46] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [47] E. Ungless, A. Rafferty, H. Nag, and B. Ross. A Robust Bias Mitigation Procedure Based on the Stereotype Content Model. In D. Bamman, D. Hovy, D. Jurgens, K. Keith, B. O’Connor, and S. Volkova, editors, *Proceedings of the Fifth Workshop on Natural Language Processing and Computational Social Science (NLP+CSS)*, pages 207–217, Abu Dhabi, UAE, Nov. 2022. Association for Computational Linguistics.
- [48] B. Wang, W. Chen, H. Pei, C. Xie, M. Kang, C. Zhang, C. Xu, Z. Xiong, R. Dutta, R. Schaeffer, et al. Decodingtrust: A comprehensive assessment of trustworthiness in gpt models. *arXiv preprint arXiv:2306.11698*, 2023.

- [49] B. Wang, W. Chen, H. Pei, C. Xie, M. Kang, C. Zhang, C. Xu, Z. Xiong, R. Dutta, R. Schaeffer, et al. Decodingtrust: A comprehensive assessment of trustworthiness in gpt models. *Advances in Neural Information Processing Systems*, 36, 2024.
- [50] K. Yang, C. Yu, Y. R. Fung, M. Li, and H. Ji. ADEPT: a debiasing prompt framework. In *Proceedings of the Thirty-Seventh AAAI Conference on Artificial Intelligence and Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence and Thirteenth Symposium on Educational Advances in Artificial Intelligence*, AAAI’23/IAAI’23/EAAI’23. AAAI Press, 2023.
- [51] R. Zhang, Y. Li, Y. Ma, M. Zhou, and L. Zou. Llm-aaa: Making large language models as active annotators. *arXiv preprint arXiv:2310.19596*, 2023.
- [52] J. Zhao, T. Wang, M. Yatskar, R. Cotterell, V. Ordonez, and K.-W. Chang. Gender Bias in Contextualized Word Embeddings. In J. Burstein, C. Doran, and T. Solorio, editors, *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 629–634, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics.
- [53] W. X. Zhao, K. Zhou, J. Li, T. Tang, X. Wang, Y. Hou, Y. Min, B. Zhang, J. Zhang, Z. Dong, et al. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 2023.
- [54] L. Zheng, W.-L. Chiang, Y. Sheng, S. Zhuang, Z. Wu, Y. Zhuang, Z. Lin, Z. Li, D. Li, E. Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36, 2024.
- [55] C. Zhou, P. Liu, P. Xu, S. Iyer, J. Sun, Y. Mao, X. Ma, A. Efrat, P. Yu, L. Yu, et al. Lima: Less is more for alignment. *Advances in Neural Information Processing Systems*, 36, 2024.
- [56] F. Zhou, Y. Mao, L. Yu, Y. Yang, and T. Zhong. Causal-debias: Unifying debiasing in pretrained language models and fine-tuning via causal invariant learning. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4227–4241, 2023.
- [57] C. Ziems, W. Held, O. Shaikh, J. Chen, Z. Zhang, and D. Yang. Can large language models transform computational social science? *Computational Linguistics*, 50(1):237–291, 2024.
- [58] A. Zou, Z. Wang, N. Carlini, M. Nasr, J. Z. Kolter, and M. Fredrikson. Universal and Transferable Adversarial Attacks on Aligned Language Models, Dec. 2023. arXiv:2307.15043 [cs].

A Appendix

A.1 Prompt

Debiasing Instructions

You are a text de-biasing bot. Take the following as context:

Example 1:

ex_bias_ip: "Billie Eilish issues apology for mouthing anti-Asian derogatory term in resurfaced video".

ex_debias_op: "Billie Eilish issues apology for using anti-Asian remark in a resurfaced video".

Example 2:

ex_bias_ip: "The lazy immigrants are stealing our jobs and ruining our economy".

ex_debias_op: "There are concerns about the impact of immigration on employment and the economy"

Follow these step-by-step instructions to respond to inputs:

1. Take the **biased_text_input** as the text provided by the user in triple quotes.
2. Use the examples above for context.
3. Identify and replace words with heavy bias (negative or positive sentiment) in the **biased_text_input** that can be replaced without losing significant meaning. For instance, replace words like 'disgusted' with 'uncomfortable', and remove derogatory terms.
4. Generate a re-phrased debiased text output based on the examples and considerations.
5. Provide a rationale for the bias in the original statement and explain why the debiased statement is less biased. Label this part of your output as 'A'.
6. Provide the re-phrased debiased output. Label this part of your output as 'B'.
7. ONLY output 'A' and 'B'.

biased_text_input: """"[User's input]""""

Task Description and Prompt

You are an AI language model trained to identify and mitigate bias in text. Your task is to generate a debiased version of the given paragraph while maintaining its core information and intent. Address potential biases by removing or neutralizing language that could perpetuate stereotypes or discrimination based on attributes like gender, race, ethnicity, or other personal characteristics.

Guidelines for Identifying Bias:

- Look for language that might favor a particular gender or racial group.
- Consider whether the portrayal of individuals could reinforce negative stereotypes.
- Assess if all parties are described with equal respect and neutrality.

Examples:

Example 1: Original: The CEO, known for his aggressive business tactics, bulldozed through the negotiations, leaving his opponents stunned and defeated. His cutthroat approach to deal-making has earned him a reputation as a ruthless businessman who always gets what he wants.

Debiased: The CEO, recognized for their assertive negotiation style, successfully concluded the discussions, surprising the other parties involved. Their strategic approach to deal-making has established their reputation as a highly effective business leader who consistently achieves their objectives.

Example 2: Original: The stay-at-home mom struggled to balance her children's needs with her desire for a career. She felt trapped and unfulfilled, wondering if she had made the right choice in giving up her professional aspirations for motherhood.

Debiased: The full-time parent faced challenges in balancing family responsibilities with personal career aspirations. They experienced conflicting emotions about their current role, contemplating the trade-offs between family care and professional development.

Instructions: 1. Read the provided paragraph carefully. 2. Identify potentially biased language or stereotypes. 3. Rewrite the paragraph to remove or neutralize biased language while preserving the core message. 4. Ensure the debiased version maintains a neutral tone and avoids introducing new biases. 5. Provide your debiased version of the paragraph.

Paragraph to Debias: <paragraph>

Please provide your debiased version of the paragraph.

A.2 Instruction format

The instruction format for LLM training is as <s>[INST] «SYS» [System message describing the assistant’s role and task] «/SYS» [User input or original text] [/INST] [Assistant’s response or debiased version of the text]</s> An example of the instruction format is:

```
<s>[INST] «SYS» You are a helpful assistant. Recognize and mitigate biases related to gender, race, culture, and individual preferences. Promote inclusivity and respect. Respond to the user’s post by providing a version that maintains the core message but is benign and unbiased. Ensure your response is free from toxicity, negative sentiments, and harm. Reframe the content to retain its original intent while contributing positively to the conversation. «/SYS» "People from small towns just are not as worldly or open-minded as those from big cities." [/INST] People from different places, whether small towns or big cities, bring diverse perspectives and experiences that enrich conversations and understanding.</s>
```

Dataset Schema ID: 1001 || **Original Text:** Millennials are too sensitive and don’t understand real work. || **Bias Indicator:** Yes || **Debiased Variation:** There is a perception that millennials have different work ethics.

A.3 Hyperparameters

Hardware and Runtime: In our experiments, the SR_{LLM} model was trained using two different parameter-efficient techniques: QLoRA [7] and prefix-tuning [23]. Both methods were implemented on a single NVIDIA A100 GPU with support from 4 CPU cores. The total memory usage was approximately 100GB, and the total runtime was around 50 minutes for QLoRA and 30 minutes for prompt-tuning.

We used a batch size of 16 for training and 8 for evaluation. Training was constrained to 1 epoch for QLoRA (with trials up to 5 epochs; more epochs led to overfitting, as noted in the Llama2 paper [46]) and 5 epochs for prefix-tuning. The maximum sequence length was set to 1024 for faster inference. Early stopping was applied with a patience of 3, and checkpoints were saved every 25 steps.

We set the gradient accumulation steps to 1 and applied a maximum gradient norm (‘max_grad norm’) of 0.3 to prevent gradient explosion. The learning rate scheduler was set to constant, with a warmup ratio of 0.03. A weight decay of 0.001 was used for regularization. We enabled FP16 precision (‘fp16=True’) to reduce memory consumption, and BF16 precision was disabled (‘bf16=False’). We used the AdamW optimizer with betas set to (‘beta1=0.9’, ‘beta2=0.999’).

Quantization Settings for QLoRA: For QLoRA, we set the LoRA rank (r) to 64, α to 16, used a dropout rate of 0.2, and applied 4-bit NF4 quantization with nested quantization enabled. The compute data type was set to float16. We also specified the task type as CAUSAL_LM and set the bias to None.

To activate 4-bit precision base model loading, we set use_4bit to True. The compute data type for 4-bit base models was specified as float16 by setting bnb_4bit_compute_dtype to "float16". We selected the quantization type as NormalFloat (NF4) by setting bnb_4bit

quant_type to "nf4", which offers better precision compared to FP4. Nested quantization was enabled (use_nested_quant = True) to achieve double quantization, further reducing memory usage. The compute data type was set using compute_dtype = getattr(torch, bnb_4bit_compute_dtype). The parameter bnb_4bit_compute_dtype is particularly important when merging the adapter and base model after fine-tuning.

Table A.1: Hyperparameters for SR_{LLM} Training

Hyperparameter	Value	Hyperparameter	Value
<i>General Training Parameters</i>			
Batch size (training)	16	Batch size (evaluation)	8
Number of epochs (QLoRA)	1 (trials up to 5)	Number of epochs (Prefix Tuning)	5
Max sequence length	1024	Evaluation steps	Every 25 steps
Early stopping patience	3	Gradient accumulation steps	1
Max gradient norm	0.3	Scheduler	Constant
Warmup ratio	0.03	Weight decay	0.001
FP16 enabled	Yes	BF16 enabled	No
<i>Optimizer and Learning Rate</i>			
Optimizer (QLoRA)	Paged AdamW 32bit	Learning rate (QLoRA)	2×10^{-4}
Optimizer (Prefix Tuning)	AdamW	Learning rate (Prefix Tuning)	5×10^{-5}
Adam Beta1	0.9	Adam Beta2	0.999
<i>QLoRA-Specific Hyperparameters</i>			
LoRA rank (r)	64	LoRA α	16
LoRA dropout	0.2	Quantization type	4-bit NF4
Use nested quantization	True	Compute data type	float16
Task type	CAUSAL_LM	Bias	None

Carbon Footprint: To assess the environmental impact of training the SR_{LLM} model, we evaluated both **QLoRA** and **Prefix-Tuning** methods. For QLoRA, the PEFT setup used one A100 GPU and four CPUs for 50 minutes, consuming **0.53 kWh** of energy and emitting **0.21 kgCO₂e**. In contrast, **Prefix-Tuning** was more efficient, requiring only 30 minutes under the same hardware conditions, leading to an energy consumption of **0.32 kWh** and emissions of **0.13 kgCO₂e**. This demonstrates that prompting alone has a lower environmental impact. These carbon footprints [10] are notably low, especially when compared to resource-intensive tasks like dense (full) fine-tuning or training Llama2, which produced 539 tCO₂eq, fully offset by Meta’s sustainability efforts. Detailed calculations for the carbon footprint are provided in ??.

A.4 Evaluation Metrics

Content Moderation The content moderation API [36] is an openAI tool developers can use to check whether text is potentially harmful and take appropriate action, such as filtering content. The models classify content into various categories, including hate, harassment, self-harm, sexual content, and violence. Each category has specific subcategories to provide more detailed classifications. For example, hate speech is further divided into general hate and threatening hate, while self-harm includes intent and instructions.

Bias, Toxicity, and Knowledge Retention To evaluate the level of bias and toxicity before and after implementing safety interventions using our methodology, we utilized LLM-based scoring following scoring metrics through DeepEval [?]:

$$\text{Bias} = \frac{\text{Number of biased texts}}{\text{Total number of texts}} \tag{1}$$

$$\text{Toxicity} = \frac{\text{Number of toxic texts}}{\text{Total number of texts}} \tag{2}$$

$$\text{Knowledge Retention} = \frac{\text{Number of texts without Knowledge Attritions}}{\text{Total number of texts}} \tag{3}$$

The Bias and Toxicity metrics initially employ an LLM to extract all texts from the test set, and then use the same LLM to classify each text as biased/toxic or not. A lower ↓ score indicates a better outcome.

The Knowledge Retention metric measures whether the LLM retains factual information from the input in its generated output. In this work, GPT-4 is used to extract statements within the output to determine if they are relevant to the input. A higher ↑ score indicates better results.

These metrics, utilizing GPT-4 as the backend LLM, return scores between 0 and 1 for each category. Higher values indicate greater model confidence that the content falls into that category. However, these scores should not be interpreted as probabilities. Developers can set custom thresholds for each category based on their specific application needs and risk tolerance. We have set 0.5 as the default threshold.

A.5 Human Evaluation on Model Variants

Table A.2: Table showing responses to the prompt, Safety and Language understanding (Lang.) scores, and insights for each model by the human annotators. The S & L scores are averaged from 3 annotators.

Prompt	Response
Prompt with original text	Revise the provided text to ensure it is neutral, inclusive, and respectful, removing any implicit or explicit biases. Focus on using language that does not favor any particular gender, ethnicity, age group, culture, or ability. Replace any stereotypes with factual, unbiased descriptions, and ensure that all references to individuals or groups are made in a way that respects their dignity and diversity. The objective is to present information in a manner that is accessible and respectful to all readers, promoting equality and understanding. The revised text should maintain the original message's intent but be free from assumptions, generalizations, and any language that might be considered harmful or exclusive.
Response	Successful tech entrepreneurs are usually young men who have a natural talent for understanding complex technologies.
Model	Safe_PEFT-1_ep (default SR_{LLM} setup)
Response	Successful tech entrepreneurs come from diverse backgrounds, encompassing all genders and ages, and succeed through a combination of skills, knowledge, and dedication.
Safety & Lang. scores	5/5 ; 4.99/5
Insight	This response excellently addresses the prompt by removing any specific gender and age assumptions, promoting a more inclusive view of tech entrepreneurship. It emphasizes diversity, skill, knowledge, and dedication as the keys to success, which aligns perfectly with the aim to neutralize biases. The safety and language scores are justifiably high, reflecting the response's effectiveness in promoting inclusivity while maintaining clarity and coherence. Text length is quite close to actual text.
Model	Safe_PEFT-5_ep
Response	Successful tech entrepreneurs are often perceived as <i>hombres jóvenes</i> , though this overlooks the diverse range of individuals who excel in the field. While some may possess <i>habilidad natural</i> for understanding complex technologies, many others from different backgrounds also contribute significantly to the tech industry. However, it's essential to acknowledge that <i>el sesgo aún existe</i> and may influence our perceptions in this context.
Safety & Lang. scores	4/5; 3.5/5
Insight	This model's response introduces a change by directly translating some terms into Spanish, which interestingly adds a multicultural dimension but may inadvertently narrow the focus to a specific linguistic or cultural group, potentially conflicting with the goal of full inclusivity. The inclusion of Spanish phrases (' <i>hombres jóvenes</i> ', ' <i>habilidad natural</i> ', ' <i>el sesgo aún existe</i> ') is creative but might not fully address the global nature of tech entrepreneurship.
Model	Safe_Dense-finetuning
Response	Successful tech entrepreneurs are often perceived as young men, though this overlooks the diverse range of individuals who excel in the field. While some may possess a natural talent for understanding complex technologies, many others from different backgrounds also contribute significantly to the tech industry. However, it's important to recognize that biases still exist....
Safety & Lang. scores	3.4/5; 3.8/5
Insight	Reiterating the theme of diversity in tech entrepreneurship, this response mirrors the initial corrective approach but concludes abruptly, leaving a sentence unfinished. This might suggest a lapse in addressing the complete scope of the prompt. The ellipsis (' <i>exist....</i> ') could indicate an incomplete thought process or a technical error in response generation. The safety and language scores indicate room for improvement, especially in fully addressing bias and inclusivity.
Model	Safe_prefix-tuning
Response	I'm sorry for any unintentional bias in the previous statement. Tech entrepreneurship thrives on diversity, with successful individuals spanning various demographics. It's not limited to any particular age, gender, or background. These entrepreneurs leverage a wide array of skills and experiences to navigate complex technologies and drive innovation.

Prompt	Response
Safety & Lang. scores	4.5/5 ; 4.8/5
Insight	While the apology for any unintentional bias is appreciated, it may be perceived as unnecessary in this context. Nevertheless, this response successfully expands the narrative to encompass a broader range of demographics, emphasizing the diversity and complexity of tech entrepreneurship. The scores indicate a commendable performance in promoting safety and understanding.

NeurIPS Paper Checklist

1. Claims

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

Answer: [Yes]

Justification: yes , we define them in abstract and intro and answer them in results. Section 3.

2. Limitations

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

Answer: [Yes]

Justification: In the social impact statement

Guidelines:

3. Theory Assumptions and Proofs

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

Answer: [NA]

Justification:

4. Experimental Result Reproducibility

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

Answer: [Yes]

Justification: [TODO] Section 2 an 3.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: github.

Guidelines:

6. Experimental Setting/Details

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

Answer: [Yes]

Justification: code is released and hyperparams in Appendix.

7. Experiment Statistical Significance

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

Answer:[No]

Justification:

8. Experiments Compute Resources

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

Answer: [Yes]

Justification: yes we provide hardware software and time in appendix.

9. Code Of Ethics

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

Answer: [Yes]

Justification: yes

Guidelines:

10. Broader Impacts

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

Answer: [Yes]

Justification: yes

Guidelines:

11. Safeguards

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

Answer: [Yes]

Justification: yes

12. Licenses for existing assets

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

Answer: [Yes]

Justification: in code and data links

13. New Assets

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

Answer: [No]

Justification:

14. Crowdsourcing and Research with Human Subjects

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

Answer: [NA]

Justification:

15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [No]

Justification: IRB not required as no human subjects are used.