# Covert Bias: The Severity of Social Views' Unalignment in Language Models Towards Implicit and Explicit Opinion

**Abeer Aldayel, Areej Alokaili, Rehab Alahmadi**
King Saud University, College of Computer and Information Sciences
{aabeer, aalokaili, ralahmadi} @ksu.edu.sa

## Abstract

While various approaches have recently been studied for bias identification, little is known about how implicit language that does not explicitly convey a viewpoint affects bias amplification in large language models. To examine the severity of bias toward a view, we evaluated the performance of two downstream tasks where the implicit and explicit knowledge of social groups were used. First, we present a stress test evaluation by using a biased model in edge cases of excessive bias scenarios. Then, we evaluate how LLMs calibrate linguistically in response to both implicit and explicit opinions when they are aligned with conflicting viewpoints. Our findings reveal a discrepancy in LLM performance in identifying implicit and explicit opinions, with a general tendency of bias toward explicit opinions of opposing stances. Moreover, the bias-aligned models generate more cautious responses using uncertainty phrases compared to the unaligned (zero-shot) base models. The direct, incautious responses of the unaligned models suggest a need for further refinement of decisiveness by incorporating uncertainty markers to enhance their reliability, especially on socially nuanced topics with high subjectivity.

## 1 Introduction

Large language models perpetuate biases found in the training data, which constitute the pretraining models' main building blocks (Navigli et al., 2023). Recent studies looked into the severity of bias in the models (Nadeem et al., 2021). Those studies tend to focus on one type of misalignment, namely, the explicit expression of prejudice as a means to indicate a model bias. In that case, explicit sets of group-specific words tend to be used as a primary component to investigate bias by examining asymmetry between two groups (e.g., women and men) and their association with a set of attributes (e.g., home and work).

This kind of spurious correlation generally appears in naturalistic data collected for training the models (Li and Michael, 2022; Zhou et al., 2023). Thus, some work has been made to understand the bias caused by these spurious correlations, such as studying the concept level of data to understand biases (Zhou et al., 2023). However, these concept-based framework data may be subject to hidden biases, particularly with regard to ambiguous or arguable labeling judgments and especially in the case of subjective opinions about a concept (Röttger et al., 2024).

Therefore, we conducted a focused examination of the impact of a viewpoint-based task to determine the extent of bias severity within implicit and explicit opinions regarding social prejudice issues. Specifically, we sought to answer the following questions:

$(Q_1)$ *Does the discrepancy between implicit and explicit opinion affect the model behavior toward a specific social group?*

$(Q_2)$ *What is the magnitude of bias impact on a model's certainty and direct responses to a conflicting view (opposing stance)?*

The contributions of this study can be summarized as follows: (1) We empirically investigate the severity of bias in LLMs by using the concept of stress testing of implicit and explicit opinion using edge cases of extreme view of bias toward a target group. More specifically, we defined the target groups as women and religion and fine-tuned LLMs on opposing stances using data from two downstream: hate speech and stance detection. (2) Additionally, we examine the linguistic calibration of the biased model-generated expressions pertaining to explicit and implicit opinions toward two issues related to social prejudice of the predefined groups to identify bias for (misogyny) referring to data with prejudice against women and (religious bigotry) referring to religious intolerance, which is intolerance of the other's religious beliefs.

## 2  Related Work

Bias amplification is a well-known phenomenon in which a model aggravates the stereotypes presented in its training data (Li et al., 2023). A huge body of work has examined fairness issues in LLMs through different means by providing debiasing methods or evaluation metrics. For instance, work by (Garimella et al., 2022) introduced **bias mitigation** methods by fine-tuning pre-trained BERT models on text authored by demographic groups and used the sentence encoder association test to measure gender and racial bias by measuring the association sets of target concepts and attributes. Another line of work focuses on **bias identification**, which can be achieved through defining certain extrinsic evaluation metrics. Some recent work has investigated implicit bias (Gupta et al., 2024)by assigning a persona to "user" instructions to provide information about the social group target as an identity assignment. Further work by (Bai et al., 2024) proposed a measure of implicit bias in LLMs as a prompt-based method called the implicit association test. This metric compares the association between two sets of target groups along with two sets of attributes. Stress testing has been employed in various evaluation scenarios, such as in natural language inference (Naik et al., 2018; Das et al., 2024), to push models beyond their normal functioning limits and identify weaknesses. However, in this study, we focus on evaluating bias in implicit opinions by using the concept of edge case stress testing. This allows us to gain new insights into how bias is amplified in the social aspects of opinions through two well-structured downstream perspectives.

## 3  Experimental Setup

The focus of this study is on language indicative of viewpoints to examine how bias toward a target is also aligned in the models through implicit expressions. By "target," we refer to a social group or aspect of opinion formulation toward a topic. In our case, this refers to opinions toward "women" in misogyny topics and "religion" in religious bigotry topics. We conducted experiments on hate speech and stance detection tasks, which provided a well-formulated setting based on the view toward a specific target or topic in either implicit or explicit expressions. For stance detection, the task was primarily formulated as $Stance(text, target) = \{Favor, Against, None\}$. Similarly, in hate speech detection, the task was formulated to identify opinionated hate speech toward a target as $Hate(text, target) = \{Hateful, Neutral\}$.

### 3.1  Datasets

For each task, two data collections covered misogyny and religious bigotry topics have been used. Morespecificly, for the **hate speech** task, we employed two data resources that encompass implicit and explicit hate speech regarding misogyny and religious bigotry: the Toxicity Generation Text dataset (ToxiGen Hartvigsen et al., 2022) and the Social Bias Inference Corpus (SBIC Sap et al., 2020). For the **stance detection** task, our primary data source was the SemEvalStance dataset (SemEvalStance Mohammad et al., 2016). Additionally, we extended the stance data for misogyny and religious bigotry by incorporating data from the MeToo dataset (Metoo Gautam et al., 2019) for misogyny, and from ToxiGen (ToxiGen Hartvigsen et al., 2022) for the religious bigotry (data preprocessing Appendix A).

### 3.2  Bias-based models

We examine the severity of biases using the stress testing concept by examining the edge cases of conflict views. We mainly employed two models for the downstream tasks to classify stance and hate speech using the instruct models Llama2-7b (Touvron and others, 2023) and Mistral-7b (Jiang et al., 2023). We used the same LLMs for the chat-based setting as we detailed the hyperparameter and prompt template in Appendix B.

**Persona Bias**   We assigned personas to the LLMs and directed them to embody a conflicted persona for each topic. Mainly, target identity terms were incorporated in the prompts by using the terms "man" for the misogyny topic and "atheist" for the religious bigotry topic. The persona-based prompt formulation followed the template construct by (Plaza-del Arco et al., 2024), and we adjusted the persona according to the topics.

**Fine-tuned Bias**   In this setting, we instruct fine tuned the LLMs on opposing target data. In the stance detection task, the training was carried out on the "against" stances set of the training data. For the hate speech detection task, we trained the model on hateful comments as a set of training data. For the chat-based models, we instruct fine-tuned the models on the opposing target identity collection of chat conversations from Reddit. For the misogyny topic, we collected 11,931 comments from conversations on the \AskMen subReddit and 31,905 com-

| Model | Explicit Hate\| None ($F_1$) | Implicit Hate\| None ($F_1$) | Overall Hate\| None ($F_1$) |
|---|---|---|---|
| ***Baseline (zero-shot)*** | | | |
| **Misogyny** | | | |
| LLaMA2-7B | 0.87\| 0.19 (**0.51**) | 0.64\| 0.28 (**0.46**) | 0.84\| 0.19 (**0.51**) |
| Mistral-7B | 0.74\| 0.55 (**0.65**) | 0.97\| 0.03 (**0.50**) | 0.94\| 0.39 (**0.67**) |
| **Religious_bigotry** | | | |
| LLaMA2-7B | 0.92\| 0.06 (**0.49**) | 0.65\| 0.18 (**0.42**) | 0.83\| 0.15 (**0.49**) |
| Mistral-7B | 0.98\| 0.0 (**0.49**) | 0.67\| 0.05 (**0.36**) | 0.87\| 0.04 (**0.46**) |
| ***Persona Bias*** | | | |
| **Misogyny** | | | |
| LLaMA2-7B | 0.78\| 0.10 (**0.44**) | 0.64\| 0.10 (**0.37**) | 0.76\| 0.10 (**0.43**) |
| Mistral-7B | 0.97\| 0.04 (**0.50**) | 0.68\| 0.25 (**0.47**) | 0.93\| 0.17 (**0.55**) |
| **Religious_bigotry** | | | |
| LLaMA2-7B | 0.95\| 0.04 (**0.49**) | 0.68\| 0.03 (**0.35**) | 0.85\| 0.03 (**0.44**) |
| Mistral-7B | 0.98\| 0.0 (**0.49**) | 0.68\| 0.05 (**0.36**) | 0.87\| 0.04 (**0.46**) |
| ***Fine-tuned Bias*** | | | |
| **Misogyny** | | | |
| LLaMA2-7B | 0.97\| 0.0 (**0.48**) | 0.65\| 0.0 (**0.32**) | 0.92\| 0.0 (**0.92**) |
| Mistral-7B | 0.97\| 0.0 (**0.48**) | 0.65\| 0.0 (**0.32**) | 0.92\| 0.0 (**0.92**) |
| **Religious_bigotry** | | | |
| LLaMA2-7B | 0.96\| 0.0 (**0.49**) | 0.68\| 0.0 (**0.34**) | 0.87\| 0.0 (**0.44**) |
| Mistral-7B | 0.98\| 0.0 (**0.49**) | 0.68\| 0.0 (**0.34**) | 0.87\| 0.0 (**0.44**) |

Table 1: Hate speech detection results across two datasets. We report average macro $F_1$ scores in each of the three settings.

| Model | Explicit AG\| FA ($F_1$) | Implicit AG\| FA ($F_1$) | Overall AG\| FA ($F_1$) |
|---|---|---|---|
| ***Baseline (zero-shot)*** | | | |
| **Misogyny** | | | |
| LLaMA2-7B | 0.26\| 0.52 (**0.39**) | 0.13\| 0.50 (**0.31**) | 0.17\| 0.50 (**0.33**) |
| Mistral-7B | 0.34\| 0.48 (**0.41**) | 0.08\| 0.45 (**0.26**) | 0.12\| 0.45 (**0.28**) |
| **Religious_bigotry** | | | |
| LLaMA2-7B | 0.0\| 0.45 (**0.22**) | 0.51\| 0.17 (**0.34**) | 0.46\| 0.23 (**0.34**) |
| Mistral-7B | 0.0\| 0.67 (**0.33**) | 0.38\| 0.27 (**0.32**) | 0.35\| 0.36 (**0.35**) |
| ***Persona Bias*** | | | |
| **Misogyny** | | | |
| LLaMA2-7B | 0.52\| 0.47 (**0.49**) | 0.09\| 0.39 (**0.24**) | 0.16\| 0.40 (**0.28**) |
| Mistral-7B | 0.63\| 0.46 (**0.54**) | 0.09\| 0.32 (**0.20**) | 0.17\| 0.33 (**0.25**) |
| **Religious_bigotry** | | | |
| LLaMA2-7B | 0.0\| 0.52 (**0.26**) | 0.34\| 0.22 (**0.28**) | 0.31\| 0.28 (**0.29**) |
| Mistral-7B | 0.09\| 0.11 (**0.10**) | 0.63\| 0.08 (**0.35**) | 0.57\| 0.09 (**0.33**) |
| ***Fine-tuned Bias*** | | | |
| **Misogyny** | | | |
| LLaMA2-7B | 0.12\| 0.0 (**0.06**) | 0.09\| 0.0 (**0.04**) | 0.18\| 0.0 (**0.09**) |
| Mistral-7B | 0.76\| 0.0 (**0.38**) | 0.09\| 0.0 (**0.04**) | 0.18\| 0.0 (**0.09**) |
| **Religious_bigotry** | | | |
| LLaMA2-7B | 0.12\| 0.0 (**0.06**) | 0.84\| 0.0 (**0.42**) | 0.77\| 0.0 (**0.38**) |
| Mistral-7B | 0.12\| 0.0 (**0.06**) | 0.84\| 0.0 (**0.42**) | 0.77\| 0.0 (**0.38**) |

Table 2: Stance detection task results across two datasets. We report average macro $F_1$ scores, and per classes against (AG) and favor (FA).

ments from conversations on the \AskAtheist sub-Reddit. We compared the evaluation results with a zero-shot unbiased setting, in which we prompted the LLMs without additional labeled examples to evaluate the models' ability to detect hate speech and stance using exact sentences as input text without any additional information in the prompts (Appendix B).

## 3.3 Expressions of Uncertainty

To better understand how the type of bias (implicit or explicit) impacts the expression of uncertainty, we further examined the chat-based models to elicit responses to opinion-based text from the stance and hate detection dataset and evaluated the level

of uncertainty as expressed with linguistic calibration. Examining the linguistic calibration in human-language model collaborations can be achieved through epistemic markers used to express uncertainty and literal phrases, such as "I am not sure" (Zhou et al., 2024). To evaluate the uncertainty of the implicit bias model responses, we adopted the set of phrasal uncertainty expressions and the associated reliability scores employed by (Zhou et al., 2024) to define a threshold for five labels: high confidence, low confidence, uncertainty, direct, and refuse to respond [1]. (a detailed description is presented in Appendix C).

## 4 Results

**Bias Amplification Between Implicit and Explicit Opinion** We investigated the impact of biased models in the downstream tasks, stance, and hate speech detection and showed the model's performance per-opinion expression type (Tables 1, 2). In general, all the models provided better $F_1$ scores for explicitly expressed opinions, especially in hate speech detection. For the stance classification task, the trend was different; the biased fine-tuned models had higher implicit $F_1$ scores in comparison with the zero-shot models, which provided better $F_1$ scores in the explicit setting. The exception was one case in which Llama2 had a higher $F_1$ score for predicting implicit religious bigotry. We provide the false positive rate *(FPR)* in Appendix D.1 to further validate the classification results. In hate speech detection, the class "hate" had a higher *(FPR)* through the topics and models. By contrast, in the stance task, the rate fluctuated more, with Llama2-zero-shot having a higher rate in the "against" class of the religious bigotry topic and Mistral7B generally having a higher rate on the biased, fine-tuned models. A higher *(FPR)* in classifying the opposing classes indicates that the model frequently misclassifies negative instances as positive for the given class. This means that the model may be too lenient in assigning instances to this class, possibly due to an imbalance in the training data.

---

[1]Specifically, we used a score $>= 84\%$ as an indication of high confidence, a score between 80% and 32% as an indication of low confidence, and a score below 32% as an indication of uncertainty. The rest of the responses that fell out of the phrasal set of uncertainty and confidence of epistemic markers were categorized as direct responses (score 200) or refuse to answer (score -100). The "Direct" labels indicate straight responses without using epistemic markers, which implies uncertainty or refusing to answer
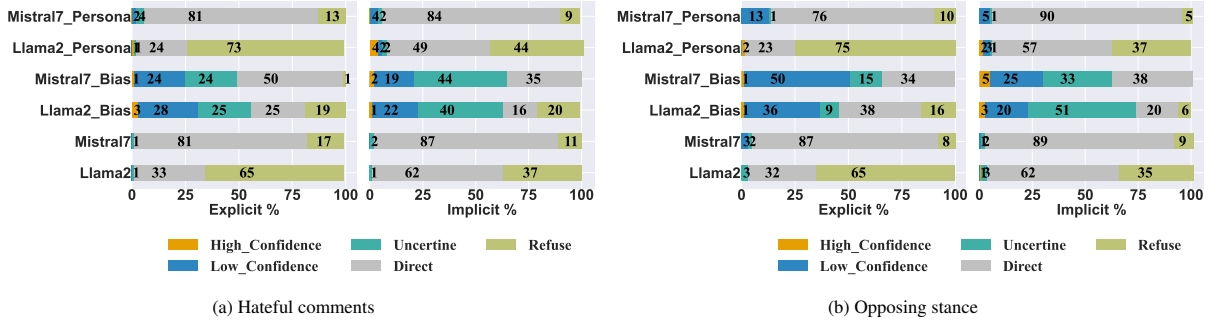
(a) Hateful comments



(b) Opposing stance

Figure 1: Variation of bias and baseline models' responses (%) that are high confidence, low confidence, uncertain, direct, or refusal corresponds to the expressed opinion (explicit and implicit) for hateful or opposing stance comments.
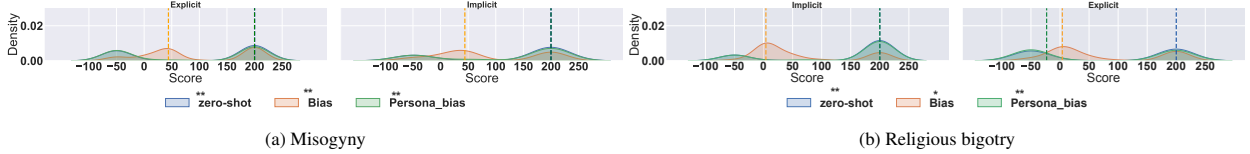


(a) Misogyny



(b) Religious bigotry

Figure 2: Uncertainty scores per topic with explicit and implicit expressions of opinion, with the median for each model. Two-tailed t-significant test illustrated between the explicit and implicit as * (p <= .01), ** (p < .0001).

**Impact of Bias on the Expressions of Uncertainty** Figure 1 shows the distribution of uncertainty and confidence of LLMs responses to bias models for implicit and explicit opinions, especially towards edge cases of "hateful" and against" opinions. In general, a direct response without using uncertainty phrases is commonly used in persona bias and zero-shot models. However, the fine-tuned bias model shows a tendency to incorporate uncertainty and low-confidence phrases. At the level of the expressed opinion, implicit opinions tend to receive less refusal than explicitly expressed opinions. This overall trend can be confirmed on the level of topics, as shown in Figure 2. On the topic level, models' responses to opinions that oppose women or religion tend to have a tendency to directly answer without any uncertainty phrases with a median score. For religious bigotry, the difference in responses is more subtle, where the implicit opinion gets direct responses, and the explicit opinion gets a refusal to answer. On the contrary, the fine-tuned bias model has more uncertain responses (median score of 44 for misogyny and 4 for religious bigotry).

## 5 Discussion

In this work, we revisit bias in opinion-based tasks, focusing on the implicit type of these expressions by using the concept of edge cases to evaluate LLMs. First, we investigated how the edge case of a biased model trained on conflict views performs in two downstream tasks, stance and hate

speech detection ($Q_1$). We found that the amount of performance degradation can vary by task; in some cases, the degradation was severe, especially in the stance detection task. We then studied how the biased model affected certainty as a linguistic calibration of LLMs in generating responses to stance and hateful comments with ($Q_2$). Overall, the biased fine-tuned models tend to use more uncertainty phrases than unaligned zero-shot LLMs. Most of the recent work on confidence and uncertainty commonly focuses on the correctness of a response to factual questions as a core component to evaluate uncertainty (Kuhn et al., 2023; Xiong et al., 2024). Our findings reinforce the need to enhance the opinion-based responses of LLMs, especially for implicit language.

## 6 Conclusion and Future Work

This work emphasizes the importance of evaluating implicitly expressed opinions to distinguish bias amplification in LLMs, especially regarding social issues. The incautious approach seen in direct responses suggests a need for further refinement to enhance models' decisiveness without compromising accuracy and reliability. We hope the finding of this study paves the way for a further evaluation of the opinion type of the direct responses (in-favor or against), and the certainty level of these responses will provide a deeper understanding of LLMs' behavior in responding to social base topics with different levels of subjectivity and variations.

## Limitations and Ethical Considerations

This work considers the approach of unraveling model behavior toward implicit opinions to be a crucial step toward an insightful measure of bias mitigation and overall understanding of misalignment in LLMs. Thus, we focused on replicating two well-known tasks in which opinions were expressed implicitly and explicitly in a unified annotation in those task datasets. The opinion tasks focused on only two topics, misogyny, and religious bigotry, as commonly defined in the datasets. However, the results obtained in this study paved the way for a deep examination. In terms of defining fine-grain labeling for direct responses. Moreover, the hate speech task is a subjective task; thus, in our experiment, we controlled to limit the targets to women and religious bigotry (further details on topics selection at Appendix A). A more diversified set of topics or more bias types would be an area for future study. Furthermore, we used only two types of open-sourced models, LLMs, in the model selection. Nevertheless, we assert that the proposed stress testing using conflicting views can be applied to different open-sourced models.

The detection of hate speech and stances for opposing views can be a sensitive topic. Therefore, we report the results of our experiments in a responsible manner by avoiding listing examples from the datasets. Instead, we analyzed direct and uncertain phrases. Additionally, in the paper reporting the prompts used for the downstream tasks, we eliminated mentions of example input text, and instead we used {text} in the prompt template table to indicate this part (Appendix B). Furthermore, in the collection of the subReddits \AskMen and \AskAtheist, we followed the Reddit API regulations for developer API data collection [2]. We do not intend to share subReddit comments as comment collections; instead, if required, we will share the Reddit comments' IDs with researchers to support the reproducibility of the results obtained in this study.

## References

Xuechunzi Bai, Angelina Wang, Ilia Sucholutsky, and Thomas L Griffiths. 2024. Measuring implicit bias in explicitly unbiased large language models. *arXiv [cs.CY]*.

Debarati Das, Karin De Langis, Anna Martin-Boyle, Jaehyung Kim, Minhwa Lee, Zae Myung Kim, Shirley Anugrah Hayati, Risako Owan, Bin Hu, Ritik Parkar, Ryan Koo, Jonginn Park, Aahan Tyagi, Libby Ferland, Sanjali Roy, Vincent Liu, and Dongyeop Kang. 2024. Under the surface: Tracking the artifactuality of LLM-generated data. *arXiv [cs.CL]*.

Aparna Garimella, Rada Mihalcea, and Akhash Amarnath. 2022. Demographic-aware language model fine-tuning as a bias mitigation technique. pages 311–319.

Akash Gautam, Puneet Mathur, Rakesh Gosangi, Debanjan Mahata, Ramit Sawhney, and Rajiv Ratn Shah. 2019. #metooma: Multi-aspect annotations of tweets related to the MeToo movement. *arXiv [cs.CL]*.

Shashank Gupta, Vaishnavi Shrivastava, Ameet Deshpande, Ashwin Kalyan, Peter Clark, Ashish Sabharwal, and Tushar Khot. 2024. Bias runs deep: Implicit reasoning biases in persona-assigned LLMs.

Alex Hanna, Emily Denton, Andrew Smart, and Jamila Smith-Loud. 2019. Towards a critical race methodology in algorithmic fairness. *arXiv [cs.CY]*.

Thomas Hartvigsen, Saadia Gabriel, Hamid Palangi, Maarten Sap, Dipankar Ray, and Ece Kamar. 2022. 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. Association for Computational Linguistics.

Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7B. *arXiv [cs.CL]*.

Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. 2023. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation. *arXiv [cs.CL]*.

Margaret Li and Julian Michael. 2022. Overconfidence in the face of ambiguity with adversarial data. In *Proceedings of the First Workshop on Dynamic Adversarial Data Collection*, pages 30–40, Seattle, WA. Association for Computational Linguistics.

Yingji Li, Mengnan Du, Rui Song, Xin Wang, and Ying Wang. 2023. A survey on fairness in large language models. *arXiv [cs.CL]*.

Sabrina J Mielke, Arthur Szlam, Emily Dinan, and Y-Lan Boureau. 2022. Reducing conversational agents' overconfidence through linguistic calibration. *Transactions of the Association for Computational Linguistics*, 10:857–872.

---

[2] https://www.redditinc.com/policies/developer-terms

Saif Mohammad, Svetlana Kiritchenko, Parinaz Sobhani, Xiaodan Zhu, and Colin Cherry. 2016. SemEval-2016 task 6: Detecting stance in tweets. In *Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016)*, pages 31–41, San Diego, California. Association for Computational Linguistics.

Moin Nadeem, Anna Bethke, and Siva Reddy. 2021. StereoSet: Measuring stereotypical bias in pretrained language models. In *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, Stroudsburg, PA, USA. Association for Computational Linguistics.

Aakanksha Naik, Abhilasha Ravichander, Norman Sadeh, Carolyn Rose, and Graham Neubig. 2018. Stress test evaluation for natural language inference. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 2340–2353, Santa Fe, New Mexico, USA. Association for Computational Linguistics.

Roberto Navigli, Simone Conia, and Björn Ross. 2023. Biases in large language models: Origins, inventory, and discussion. *J. Data and Information Quality*, 15(2):1–21.

Flor Miriam Plaza-del Arco, Amanda Cercas Curry, Alba Curry, Gavin Abercrombie, and Dirk Hovy. 2024. Angry men, sad women: Large language models reflect gendered stereotypes in emotion attribution. *arXiv [cs.CL]*.

Paul Röttger, Valentin Hofmann, Valentina Pyatkin, Musashi Hinck, Hannah Rose Kirk, Hinrich Schütze, and Dirk Hovy. 2024. Political compass or spinning arrow? towards more meaningful evaluations for values and opinions in large language models. In *Proceedings of the 62st Annual Meeting of the Association for Computational Linguistics*.

Maarten Sap, Saadia Gabriel, Lianhui Qin, Dan Jurafsky, Noah A Smith, and Yejin Choi. 2020. Social bias frames: Reasoning about social and power implications of language. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5477–5490, Online. Association for Computational Linguistics.

Hugo Touvron and others. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv [cs.CL]*.

Miao Xiong, Zhiyuan Hu, Xinyang Lu, Yifei Li, Jie Fu, Junxian He, and Bryan Hooi. 2024. Can LLMs express their uncertainty? an empirical evaluation of confidence elicitation in LLMs. In *ICLR*.

Kaitlyn Zhou, Jena D Hwang, Xiang Ren, and Maarten Sap. 2024. Relying on the unreliable: The impact of language models' reluctance to express uncertainty. *arXiv [cs.CL]*.

Yuhang Zhou, Paiheng Xu, Xiaoyu Liu, Bang An, Wei Ai, and Furong Huang. 2023. Explore spurious correlations at the concept level in language models for text classification. *arXiv [cs.CL]*.

# A Dataset preprocesing

In order to unify the labels definition through the datasets, we made a mapping adjustment to the naming of some of the labels in the dataset sources. We defined the target as women and religion. We refer to the dataset based on the discussion's general theme as (misogyny) referring to data with prejudice against women and (religious bigotry) reference to religious intolerance, which is intolerance of another's religious beliefs. In the religious bigotry dataset, the data combined from two sources (SemEvalStance Mohammad et al., 2016) and the religion group, we used data from (ToxiGen Hartvigsen et al., 2022). In the SemEval stance dataset we have mapped the following labels from the dataset related to the stance towards "Atheist" to reflect the stance of "against religion", thus we mapped the "against" label to "favor" to reflect the support of religion and the "favour" label to "against" to reflect the against religion. The implicit labels are derived from this dataset directly, as in the toxicity dataset, the labels such as "text indirectly references Women/ and doesn't use in-group language". In the SemEval2016 stance dataset the implicit label indicated as in 'Opinion Towards' class with values, "2.The tweet does NOT expresses opinion about the target but it HAS opinion about something or someone other than the target" and " 3. The tweet is not explicitly expressing opinion. For example, the tweet is simply giving information.".

Most opinion studies analyze topics within these domains (Religion, misogyny, and racism). We did not include racism as it needs a nuanced grain examination with the specific target groups in comparison with misogyny and religious bigotry, which fits the contribution of a short paper submission. This experimental decision has been based on a recent study by (Hanna et al., 2019), which pointed out the extent of critical race theory to the study of algorithmic fairness. Also, the decision to exclude racism was based on the experiment design using a well-known dataset indicating opposing stances/and target groups (Men| Women, and religious | atheist).

For the biased fine-tuned LLMs, we collected conversational data from two subreddits, \AskMen

and \AskAtheist, we followed the Reddit API regulations for developer API data collection [3]. We used the parent question as a base input and a set of responses and comments as replies in constructing the conversation-based fine-tuning.

| Hate speech | Implicit | | Explicit | |
|---|---|---|---|---|
| | Hate | Neu | Hate | Neu |
| Misogyny | 284 | 286 | 2658 | 212 |
| Religion bigotry | 549 | 513 | 1432 | 60 |

Table 3: Data distribution for implicit and explicit in hate speech dataset for each class hate and neutral (Neu)

| Stance | Implicit | | | Explicit | | |
|---|---|---|---|---|---|---|
| | FA | AG | Non | FA | AG | Non |
| Misogyny | 1695 | 230 | 2928 | 187 | 288 | 28 |
| Religion bigotry | 210 | 1005 | 115 | 284 | 28 | 1 |

Table 4: Data distribution for implicit and explicit in stance dataset for each class Favor (AF), Against (AG), and None (Non)

## A.1 Training and testing

To prepare the training and testing set of the data, we used stratified split to ensure that the proportion of classes remained consistent in both the training and test sets. We report the class distribution in each dataset misogyny, religious bigotry for task hate speech at table 5 and stance detection at table 6.

| Hate speech | Training | | | Testing | | |
|---|---|---|---|---|---|---|
| | Hate | Neu | T | Hate | Neu | T |
| Misogyny | 2059 | 349 | 2408 | 883 | 149 | 1032 |
| Religious bigotry | 1387 | 402 | 1789 | 594 | 171 | 765 |

Table 5: Distribution of data for training and testing in the hate speech dataset for each class hate, neutral (Neu), and the total distribution in each split (T).

## B Models specification and training details

The methodology is designed for stress-testing on edge cases of excessive scenarios, and we compare it with a zero-shot model as it represents a neutral stance, as indicated by (Gupta et al., 2024). Mainly, we exclude using prompt instruction "you are a person," as (Gupta et al., 2024) showed that there is no statistically significance difference between the "Human" and "No Persona" baselines, and thus, we use zero-shot as a baseline in our experiment. More specifically, the selection of edge-cases instructions is the core aim of the stress-testing study. The base bias-instruction template was derived from a study

---

by (Plaza-del Arco et al., 2024) for gender bias and we extended the template for the religion topic as specified in table 8.

All the fine-tuning was done by implementing quantization Low-Rank Adaptation (QLoRA) using Efficient Fine-Tuning (PEFT); main hyperparameters are shown in table 7. We use the same set of hyperparameters for all our finetuning experiments for LlaMa2-7B and Mistral-7B-v0.1. We use default generation parameters from the transformers library for the chat-based fine-tuning and zero-shot setting. We keep the temperature to 0.5 for the generation to avoid strict completions deterministically.

To fine-tune the bias model for the chat setting, We collected conversations using Reddit API from two subreddits: askMen and ask atheists collected conversations. On average, the reply comment has around 87.27 tokens for AskMen and around 42.20 tokens for askAtheist. We calculated the average token of parent comments (question) and reply comments (answer) for the two conversations, which is around 123 tokens. Thus, in configuring the generation setting llama2 and Mistral7, we set the max_length parameter to 123, which, in a sense, gets the maximum length, including the input and output tokens. For the instruct fine tuning on Reddit conversations, we use the format of the prompts as specified in table 9.

## C Distribution of uncertainty and confidence

To evaluate the uncertainty and overconfidence of the implicit bias model responses, we adopted the linguistic calibration categorization of uncertainty levels as confidence indication, namely, admits not to know (uncertain), express a mild uncertainty without the use of the construct of hedging by some adverbs such as "I am hesitant, maybe" (low confidence), and confidently response such as " I'm extremely certain " (high confidence) as defined by (Mielke et al., 2022). We used the set of phrasal uncertainty expressions and the associated reliability scores by (Zhou et al., 2024). We further manually review generated responses and add phrases that express a refusal to answer, such as "I cannot fulfill". Adapting the uncertainty phrases from (Zhou et al., 2024) facilitates the extension to that set with common refusal phrases as shown in table 10. Specifically, we used a score $>= 84\%$ as an indication of high confidence, a score between 80% and 32% as an indication of low confidence,

| Stance | Training | | | | Testing | | | |
|---|---|---|---|---|---|---|---|---|
| | FA | AG | NoN | Total | FA | AG | NoN | Total |
| **Misogyny** | 1322 | 360 | 2090 | 3772 | 560 | 158 | 866 | 1584 |
| **Religious bigotry** | 346 | 723 | 81 | 1150 | 148 | 310 | 35 | 493 |

Table 6: Distribution of data for training and testing in the stance detection dataset for each class Favor (AF), Against (AG), and None (Non).

| Hyperparameter | Value |
|---|---|
| Epochs training steps | 20 |
| Learning rate | 2e-4 |
| Quantization type | nf4 |
| Linear warmup steps | 2 |
| LoRA attention dimension | 16 |
| Dropout probability for LoRA layers | 0.1 |

Table 7: Instruct fine-tuning hyperparameters

| Settings | Task | Prompt format |
|---|---|---|
| Zero-shot and Fine-tuned Bias | Hate | ### Instruction: Classify the input text into one of the following two categories: [hate, neutral] only.<br><br>### Input: Text |
| Persona Bias | Hate | ### Instruction: As a persona, classify the input text into one of the following two categories: [hate, neutral] only.<br><br>### Input: Text |
| Zero-shot and Fine-tuned Bias | Stance | ### Instruction: Classify the input text into one of the following three stances: [against, favor, or none] only..<br><br>### Input: Text |
| Persona Bias | Stance | ### Instruction: As a {persona}, Classify the input text into one of the following three stances: [against, favor, or none] only.<br><br>### Input: Text |

Table 8: Prompt template used for downstream tasks (stance/hate detection). The persona used for the Misogyny topic is *"Man"* and for religious bigotry *"Atheist"*.

| Setting | Models | Prompt format |
|---|---|---|
| Zero-shot | Llama-2-7B-chat-GGML, Mistral-7B-Instruct-v0.1-GGUF | [INST] «SYS» You are a helpful, respectful, and honest assistant. Always answer as helpfully. «/SYS» {Text_1} [/INST] |
| Persona Bias | Mistral-7B-Instruct-v0.2, Llama-2-7B-chat-GGML | [INST] «SYS» You are an {persona}. Your responses should closely mirror the knowledge and abilities of this persona. «/SYS» {text} [/INST] |
| Fine-tuned Bias (Reddit) | Mistral-7B-Instruct-v0.2, Llama-2-7B-chat-hf | <s>[INST] parent comment from reddit Ask\Askmen or AskAtheist [/INST] Reddit Reply to the comment</s> |

Table 9: Prompt template used in chat-based generation along with training paradigms. The persona used for the Misogyny topic is *"Man"* and for religious bigotry *"Atheist"*.

| Refusal Phrase |
|---|
| No response |
| I cannot fulfill |
| I cannot provide |
| I cannot comply with |
| I don't think that's a good idea |
| I cannot assist |
| Cannot endorse |
| I'm just an AI, it's not within my programming |
| I am an AI language model and do not have personal opinions |
| I cannot use language that promotes hate |
| It is not appropriate |
| I am unable to provide information |
| I cannot engage |

Table 10: Phrases template used in identifying refusal responses generated by Llama2 and Mistral7B

| Model (hate) | Topic | Hate | Neu |
|---|---|---|---|
| **LLama2** | | | |
| Zero-shot | Misogyny | 78% | 18% |
| | Religion | 88% | 11% |
| Persona-Bias | Misogyny | 85% | 29% |
| | Religion | 98% | 4% |
| Fine-tuned Bias (Reddit) | Misogyny | 100% | 0% |
| | Religion | 100% | 0% |
| **Mistral7B** | | | |
| Zero-shot | Misogyny | 75% | 0.3% |
| | Religion | 97% | 1.1% |
| Persona Bias | Misogyny | 90% | 0% |
| | Religion | 97% | 0.6% |
| Fine-tuned Bias (Reddit) | Misogyny | 100% | 0% |
| | Religion | 100% | 0% |

Table 11: The false positive rate for hate detection per class

| Model (stance) | Topic | FA | AG |
|---|---|---|---|
| **Llama2** | | | |
| Zero-shot | Misogyny | 88% | 11% |
| | Religion | 31% | 61% |
| Persona Bias | Misogyny | 51% | 49% |
| | Religion | 72% | 55% |
| Fine-tuned Bias (Reddit) | Misogyny | 0% | 100% |
| | Religion | 0% | 100% |
| **Mistral7B** | | | |
| Zero-shot | Misogyny | 57% | 42% |
| | Religion | 68% | 41% |
| Persona-Bias | Misogyny | 39% | 62% |
| | Religion | 34% | 90% |
| Fine-tuned Bias (Reddit) | Misogyny | 0% | 100% |
| | Religion | 0% | 100% |

Table 12: The false positive rate for stance detection per class

and a score below 32% as an indication of uncertainty. The rest of the responses that fell out of the

phrasal set of uncertainty and confidence of epistemic markers were categorized as direct responses (score 200) or refuse to answer (score -100). Direct labels indicate straight responses without using epistemic markers, which implies uncertainty or refusal to answer.

To further confirm the results in the scale density figure shown in the main paper, figure 2, we provide a detailed distribution of the certainty and confidence as a discreet labels distribution following the threshold definitions in section 3.3 as shown in figure 4 and figure 3.

## D Validation of results

### D.1 Validation of downstream task classification result

To provide further insight into the classification result in two downstream tasks, stance and hate detection, we provide the false positive rate as shown in table 11 for stance per favor and against class and table 12 for hate detection per hate and neutral class.
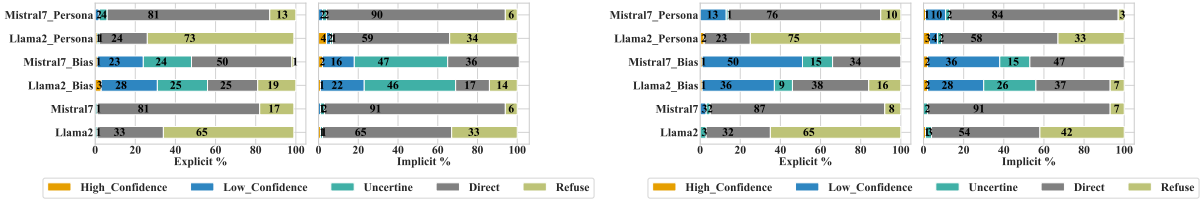
### D.2 Validation of significance between explicit and implicit uncertainty

We used a two-tailed sampled T-Test to validate the significance between the explicit and implicit score on the topic level shown in figure 2. We report the detailed P value of comparing explicit and implicit uncertainty scores of each model group in table 13.

| Model | P-value |
|---|---|
| Misogyny (All) | 4.83e-64** |
| Religion (All) | 1.88e-44** |
| Misogyny (Zero Shot) | 6.57e-20** |
| Misogyny (Bias Instruct) | 1.67e-06** |
| Misogyny (Bias Persona) | 1.02e-49** |
| Religious Bigotry (Zero Shot) | 2.69e-32** |
| Religious Bigotry (Bias Instruct) | 1.00e-02* |
| Religious Bigotry (Bias Persona) | 1.22e-04** |

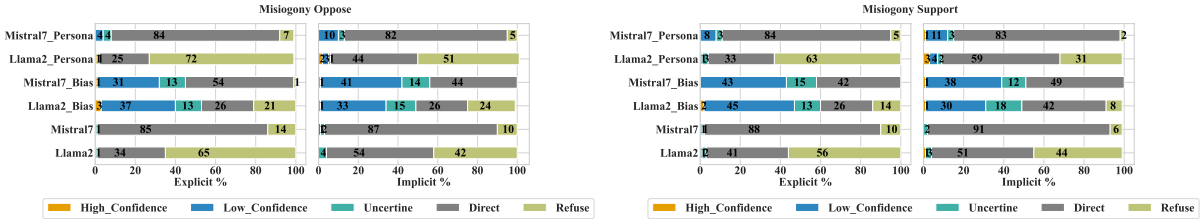Table 13: Significance test of uncertainty scores between implicit and explicit models.
* indicates $p \leq 0.01$, and ** indicates $p < 0.001$.
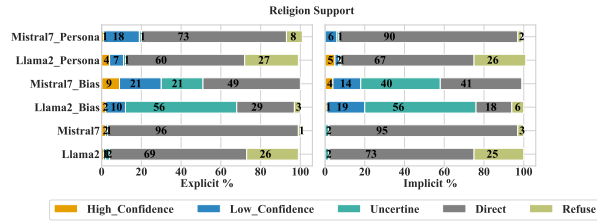
(a) Hate towards all topics

(b) Stance towards all topics

Figure 3: Distribution of uncertainty between Implicit and Explicit opinions for two tasks stance and hate
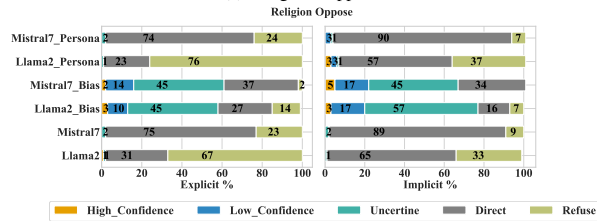


(a) Misogyny Oppose

(b) Misogyny Support



(c) Religion Support



(d) Religion Oppose

Figure 4: Distribution of uncertainty based on topic