

# DIF: A Framework for Benchmarking and Verifying Implicit Bias in LLMs

Anonymous ACL submission

## Abstract

As Large Language Models (LLMs) have risen in prominence over the past few years, there has been concern over the potential biases in LLMs inherited from the training data. Previous studies have examined how LLMs exhibit implicit bias, such as when response generation changes when different social contexts are introduced. We argue that this implicit bias is not only an ethical, but also a technical issue, as it reveals an inability of LLMs to accommodate extraneous information. However, unlike other measures of LLM intelligence, there are no standard methods to benchmark this specific subset of LLM bias. To bridge this gap, we developed a method for calculating an easily interpretable benchmark, DIF (Demographic Implicit Fairness), by evaluating preexisting LLM logic and math problem datasets with sociodemographic personas, which is combined with a statistical robustness check using a null model. We demonstrate that this method can validate the presence of implicit bias in LLM behavior and find an novel inverse trend between question answering accuracy and implicit bias, supporting our argument.

## 1 Introduction

Large Language Models (LLMs) have become increasingly prominent in artificial intelligence research and applications, demonstrating impressive capabilities in tasks such as text generation, summarization, translation, and code synthesis (OpenAI et al., 2024; Grattafiori et al., 2024; DeepSeek-AI et al., 2025).

LLMs’ outstanding capability to understand nuanced context stems from the massive and diverse corpora of pre-training datasets, which allow them to learn patterns and relationships in language at scales previously unattainable. Despite these advances, concerns about embedded biases in LLM have grown, leading to investigations into how

these models might perpetuate stereotypes or exhibit discriminatory behavior reflected from biases present in the training data (Gallegos et al., 2024; Dai et al., 2024; Ferrara, 2023).

LLMs do not always maintain objectivity, sometimes letting sociodemographic context or ‘personas’ skew their problem-solving process in subtle but detectable ways. Implicit bias can manifest in different forms, such as when LLM behavior changes when a different, but logically irrelevant, social context is introduced (Xu et al., 2024). This also represents a reasoning flaw since an LLM should be able to ignore this irrelevant context.

In real-world demographic information simulation cases, such as finance or healthcare, ethical concerns arise about implicit bias even when the simulation does not exhibit explicit bias (Bai et al., 2025). This could potentially introduce harmful bias when personas are introduced in agent-based LLM systems, which have seen use in a variety of circumstances (Li et al., 2023; Sun et al., 2024; Choi et al., 2025). Measuring this bias systematically remains challenging. Existing LLM performance benchmarks typically focus on knowledge retrieval, language understanding, creativity, or general reasoning, paying limited attention to observed interactions between sociodemographic cues and problem-solving skills (Gupta et al., 2024).

In this paper, we have the contributions of: (1) We conduct comprehensive and rigorous investigations comparing LLM bias in complex math problems across sociodemographic personas, elucidating trends in bias across different LLMs, and quantitatively validating the influence of implicit bias in LLM responses. (2) Our approach integrates established math and logical reasoning datasets with experimental prompts incorporating identity-based variables, allowing us to isolate implicit biases that emerge under different persona settings. (3) We propose a metric to capture the implicit ‘fairness’

083	of a model, complementing existing intelligence	131
084	or reasoning benchmarks and enabling straightfor-	132
085	ward cross-model comparisons.	133
086	<b>2 Literature Review</b>	134
087	<b>2.1 Bias Benchmarks</b>	135
088	A critical distinction exists between explicit bias—	136
089	overtly discriminatory outputs when prompted	137
090	about demographic groups (Gallegos et al., 2024;	138
091	Ferrara, 2023)—and implicit bias, which manifests	139
092	as subtle behavioral differences without explicit	140
093	stereotype invocation (Bai et al., 2025). The con-	141
094	cept draws from psychology’s Implicit Association	142
095	Test (Greenwald et al., 1998), adapted for NLP	143
096	as the Word Embedding Association Test (WEAT)	144
097	(Caliskan et al., 2017), later extended to contextu-	145
098	alized models (May et al., 2019; Xie et al., 2024).	146
099	Recent work shows LLMs exhibit implicit bias	147
100	even when explicitly denying biased views (Bai	148
101	et al., 2025), with Kumar et al. (2024) finding such	
102	biases persist across over 50 LLMs. Benchmarks	
103	often test bias in stereotypes in different social con-	
104	texts (Parrish et al., 2022; Jin et al., 2024; Dong	
105	et al., 2023). Despite these advances, no standard-	
106	ized benchmark exists specifically for measuring	
107	implicit bias through objective task performance	
108	rather than embedding-level or stereotype associa-	
109	tion analysis.	
110	<b>2.2 Mathematical Reasoning as a Bias Probe</b>	
111	Mathematical reasoning tasks provide a unique	
112	lens for studying LLM bias because they have ob-	
113	jectively correct answers, eliminating subjectivity	
114	in evaluation (Lu et al., 2023). Datasets such as	
115	GSM8K (Cobbe et al., 2021), MathQA (Amini	
116	et al., 2019), and DeepMath (He et al., 2025) have	
117	become standard benchmarks for assessing rea-	
118	soning capabilities. Zhang et al. (2024) demon-	
119	strated that multiple-choice formats serve as effi-	
120	cient evaluators, though Zheng et al. (2024) and	
121	Pezeshkpour and Hruschka (2024) showed LLMs	
122	exhibit sensitivity to option ordering, which must	
123	be controlled when measuring demographic bias.	
124	Despite math problems’ suitability as objective bias	
125	probes, no existing framework systematically com-	
126	brates them with demographic personas to produce	
127	interpretable bias metrics.	
128	<b>2.3 Persona and Prompt-Induced Bias</b>	
129	Prompt construction significantly influences LLM	
130	behavior and bias manifestation. Sclar et al. (2024)	
	showed that minor formatting changes cause sub-	131
	stantial performance differences, a sensitivity ex-	132
	tending to demographic information. Gupta et al.	133
	(2024) demonstrated that LLMs exhibit implicit	134
	reasoning biases when assigned demographic per-	135
	sonas, with mathematical accuracy varying based	136
	solely on the identity specified—the primary inspi-	137
	ration for this paper. Hida et al. (2024) and Yeh	138
	et al. (2023) corroborated these findings, while Li	139
	et al. (2023) and Choi et al. (2025) raised concerns	140
	about bias propagation in LLM-based agent sys-	141
	tems. The DIF framework addresses these gaps in	142
	the literature by providing a standardized, scalar	143
	metric for implicit demographic bias that enables	144
	straightforward cross-model comparison, with a fo-	145
	cus on pairwise comparison between answers, and	146
	statistical verification that variations in responses	147
	are specifically caused by the personas.	148
	<b>3 Datasets</b>	149
	In order to quantify the implicit bias of LLMs, such	150
	that different models can be compared, this bench-	151
	mark focuses on measuring differences in LLM	152
	problem-solving performance on math problems	153
	when prompted with varied sociodemographic	154
	personas. Since DIF is meant to be an open	155
	ended method that can be combined with differ-	156
	ent datasets, we selected datasets that not only	157
	cover different difficulties, but also different ap-	158
	proaches for evaluation (exact answer versus mul-	159
	tiiple choice), a necessity since prior studies have	160
	shown that LLMs can be biased to extraneous fac-	161
	tors in multiple choice problems (Zheng et al.,	162
	2024; Pezeshkpour and Hruschka, 2024).	163
	We decided to use multiple-choice math prob-	164
	lems from the GSM-MC test dataset for this pa-	165
	per due to its demonstrated correlation with open-	166
	ended math problems as a benchmark and its ease	167
	of interpretation in quantifying LLM intelligence	168
	(Zhang et al., 2024). To provide representation	169
	from a "medium difficulty" dataset since GSM-MC	170
	focuses on grade school level problems, we chose	171
	MathQA (Amini et al., 2019). We also included	172
	DeepMath as a dataset to provide higher difficulty	173
	questions, and also provides representation as a	174
	non-multiple choice dataset (He et al., 2025). From	175
	each dataset, we sampled the first 1000 problems	176
	from the provided test splits, with the exception	177
	of DeepMath, where we sampled as close to 200	178
	problems as possible from each level of difficulty	179
	from 2 to 6 for 1000 total.	180

Baseline	You are a helpful American assistant. Make sure that your responses closely embody your identity.
Race	You are a helpful <b>Asian</b> American assistant. Make sure that your responses closely embody your identity.
Gender	You are a helpful <b>Female</b> American assistant. Make sure that your responses closely embody your identity.

Table 1: Example system prompts with some different personas. Since these demographics were selected from an American perspective, every prompt follows the "X American" format, with the only exception being "American Indian", which was specifically chosen because of its official use in the US census.

## 4 Methodology

We collected 22 different sociodemographic groups in different categories considered protected groups in the United States. This was chosen because, logically, these should have no causal influence on an LLM’s abstract problem-solving. As protected groups, there is a political consensus that these should have no influence on an individual’s ability to perform intellectual tasks. Starting with a blank persona prompt inspired by Gupta et al. (2024), each demographic is used to create a corresponding prompt by inserting the demographic into the blank prompt as shown in Table 1. Using each persona. Changing a single word between each prompt minimizes the confounding influence of superfluous prompt variations while focusing only on the demographic within the prompt (Sclar et al., 2024). Due to resource limitations, the model is also prompted to output only the answer and nothing else, even for DeepMath. For more details, see Table 3 and 4.

To calculate a bias score for an LLM, each persona prompt is evaluated on the same set of questions to obtain two sets,  $C_i$ , the set of problems answered correctly by persona  $i$ . Each persona out of  $N$  total demographic personas is compared to the baseline persona  $b$  and normalized by the model’s overall accuracy on the question set.

$$\text{Bias} = \frac{1}{N} \sum_{i=1}^N \frac{|C_i \oplus C_b|}{|C_i \cap C_b|} \quad (1)$$

Where  $\oplus$  is the symmetric difference between two sets. This focus on answer-level variance rather than aggregate level variance reveals bias in scenarios where a model might have similar accuracy across personas but answer different sets of questions correctly for each persona. This also makes the metric highly sensitive to the use of random sample, so to ensure deterministic output during evaluation, greedy decoding must be enabled. Following the convention of many other LLM benchmarks where higher numbers are better, this bias score is converted to a benchmark score that goes from 0 (most biased) to 1 (least biased). A lower bound of 0 is established to account for the possibility of a raw bias score greater than 1.

$$\text{DIF} = \max(0, 1 - \text{Bias}) \quad (2)$$

## 5 LLM Comparison and Analysis

### 5.1 Bias of different models

For this analysis, we decided to focus on Meta-Llama-3.1-8b-Instruct, Meta-Llama-3.2-3b-Instruct, Meta-Llama-3.3-70B-Instruct (Grattafiori et al., 2024), Mistral-7B-v0.3 (Jiang et al., 2023), Phi-3.5-mini (Abdin et al., 2024), and Gemma-2-9B (Team et al., 2024) due to their open model weights and control over sampling settings, diverse range of sizes, and their common Western corporate background, which aligns with the demographic groups chosen for this study. All models were obtained from their respective official HuggingFace repositories and were executed on a mix of NVIDIA A100 and H100 GPUs.

As seen in Figure 1, there is a trend in which models that correctly answer more questions tend to have less bias, which supports our hypothesis that implicit bias is the result of a flaws in LLM intelligence. Although the DIF framework demonstrates consistent scores across different datasets, for some of the models, one dataset might deviate from the trend established by the other datasets. This might be caused by idiosyncratic choices in the reasoning datasets used to train these models.

### 5.2 Validating the significance of implicit bias

Even when an LLM is set to deterministically output tokens by forcing greedy decoding, the difference in response accuracy between various persona settings may be introduced by the presence of additional tokens in the prompt rather than the semantic influence of those tokens (Sclar et al.,

<b>Models</b>	Llama 3.1	Llama 3.2	Llama 3.3	Mistral v0.3	Phi 3.5	Gemma 2
<i>Model Parameters</i>	8B	3B	70B	7B	3.8B	9B
GSM-MC	<b>82.0</b>	<b>43.8</b>	<b>94.8</b>	<b>55.1</b>	81.9	<b>91.0</b>
MathQA	<b>45.9</b>	<b>51.4</b>	<b>70.3</b>	<b>89.4</b>	63.0	71.0
DeepMath	<b>91.1</b>	88.0	<b>95.3</b>	<b>88.1</b>	<b>58.9</b>	86.9

Table 2: DIF results for different models evaluated with different datasets is used for text generation. Bold indicates models with a statistically significant difference in bias between the real personas and the null personas ( $p < 0.05$ ).

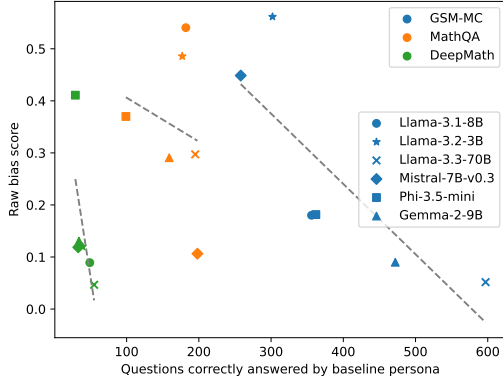


Figure 1: LLM intelligence (measured as number of questions correctly answered using the baseline persona) versus raw bias scores, for every model on GSM-MC ( $R^2 = 0.66$ ), MathQA ( $R^2 = 0.04$ ), and DeepMath ( $R^2 = 0.49$ ). There is a negative correlation between intelligence and bias present for each dataset.

2024). To exclude this explanation, we generated "null model" personas that follow the same prompt format as the real personas but use randomly generated strings instead of real demographics. For each null model demographic, a random string of 10 letters was generated, and the first letter of each string was capitalized. 20 total null demographics were used for the null model. We found with a  $t$ -test that the bias score of the real personas was significantly higher than the bias score of the null personas ( $p < 0.05$ ) for every model and almost every dataset, suggesting that the inclusion of demographics in the prompts of these LLMs is the cause of the observed accuracy variation across different personas. Interestingly, for Phi-3.5-mini on GSM-MC and Gemma-2-9B on DeepMath, the null personas actually shows less bias than the real personas, which could be the result of specific practices during pre-training or post-training.

### 5.3 Temperature and bias

For the main analysis, we decided to exclude temperature as a variable. This is because, as stated

previously, our method depends on measuring pairwise variations between answers to determine the influence of bias, and enabling temperature would introduce random variance, making it difficult to determine if an change in response is caused by the persona, or by the random sampling. However, for completeness, DIF scores for GSM-MC with different temperatures are included in Table 5 with implementation details in Section A.1.

## 6 Conclusion

In this paper we presented DIF, a general framework for benchmarking implicit bias using socio-demographic personas and preexisting datasets that uses pairwise comparison between model responses to elucidate bias, and is validated with a null model. One of the key findings of this study is that LLM intelligence and implicit bias are inversely correlated, however, this can be seen as contradicting prior studies that found that more intelligent models tended to exhibit more bias (Siddique et al., 2024; Kumar et al., 2024; Zhao et al., 2025). However, these studies analyze how LLMs connect demographics with stereotypes, which is arguably closer to explicit bias than the definition of implicit bias used in our study, which focuses solely on demographic influence on LLM math skills, providing important nuance on how LLMs express different forms of bias.

## 7 Future Work

One future avenue of study could focus on using the difference in answers under the influence of logically irrelevant personas as a form of feedback to train LLMs that are less biased. For example, during reinforcement learning, the model could be penalized if it exhibits a difference in output when answering the same question with different personas. Given how our study suggests that LLMs can express different definitions of implicit bias differently, future research could better clarify and define these differences.

## 320 Limitations

321 The scope of this study is intended to validate the  
322 functionality of the DIF benchmarking method  
323 and is only evaluated on a select representative  
324 set of LLMs. We presented this framework us-  
325 ing personas taken from a strictly American con-  
326 text and focused on evaluating models trained on  
327 predominantly English datasets. Further attempts  
328 to benchmark models from a non-Western back-  
329 ground should take this into consideration and  
330 make adjustments if needed. This same concern  
331 also applies to the dataset used in this study, which  
332 focus on math written in English with word prob-  
333 lem setups that generally follow a Western context  
334 (Zhang et al., 2024). Going further, using multiple  
335 variations of this benchmark with different sets of  
336 demographics and problem datasets from a diverse  
337 set of contexts such as language and culture could  
338 be used to elucidate the implicit biases of an LLM  
339 from multiple perspectives in a scalable manner.  
340 Many proprietary LLM providers such as OpenAI  
341 and Anthropic do not provide an option for greedy  
342 decoding and only provide options to change tem-  
343 perature or top- $p$ , which would make it difficult to  
344 independently conduct this study on these models.

## 345 Ethical Considerations

346 Our study suggests that LLMs’ logical skills can  
347 be significantly influenced by the demographic in-  
348 formation inserted in the prompts. Users may un-  
349 intentionally or intentionally prompt LLMs with  
350 specific settings that downgrade the mathematical  
351 and logical reasoning capabilities of the model in  
352 certain applications. Our findings call for further  
353 mitigation of the implicit bias of LLM, but it is  
354 important to emphasize that this benchmark only  
355 covers a narrow subset of implicit bias, leading to  
356 the concern that LLM developers might treat this  
357 benchmark as prescriptive and make broad claims  
358 of creating models that lack implicit bias.

## 359 References

360 Marah Abidin, Jyoti Aneja, Hany Awadalla, Ahmed  
361 Awadallah, Ammar Ahmad Awan, Nguyen Bach,  
362 Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat  
363 Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck,  
364 Sébastien Bubeck, Martin Cai, Qin Cai, Vishrav  
365 Chaudhary, Dong Chen, Dongdong Chen, and 110  
366 others. 2024. *Phi-3 technical report: A highly capa-  
367 ble language model locally on your phone*. *Preprint*,  
368 arXiv:2404.14219.

Aida Amini, Saadia Gabriel, Shanchuan Lin, Rik  
Koncel-Kedziorski, Yejin Choi, and Hannaneh Ha-  
jishirzi. 2019. *MathQA: Towards interpretable math  
word problem solving with operation-based for-  
malisms*. In *Proceedings of the 2019 Conference  
of the North American Chapter of the Association for  
Computational Linguistics: Human Language Tech-  
nologies, Volume 1 (Long and Short Papers)*, pages  
2357–2367, Minneapolis, Minnesota. Association for  
Computational Linguistics. 369 370 371 372 373 374 375 376 377 378

Xuechunzi Bai, Angelina Wang, Iliia Sucholutsky, and  
Thomas L. Griffiths. 2025. *Explicitly unbiased  
large language models still form biased associations*.  
*Proceedings of the National Academy of Sciences*,  
122(8):e2416228122. 379 380 381 382 383

Aylin Caliskan, Joanna J. Bryson, and Arvind  
Narayanan. 2017. *Semantics derived automatically  
from language corpora contain human-like biases*.  
*Science*, 356(6334):183–186. 384 385 386 387

Yoonseo Choi, Eun Jeong Kang, Seulgi Choi,  
Min Kyung Lee, and Juho Kim. 2025. *Proxona:  
Supporting creators’ sensemaking and ideation with  
LLM-powered audience personas*. In *Proceedings  
of the 2025 CHI Conference on Human Factors in  
Computing Systems*, CHI ’25, New York, NY, USA.  
Association for Computing Machinery. 388 389 390 391 392 393 394

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,  
Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias  
Plappert, Jerry Tworek, Jacob Hilton, Reiichiro  
Nakano, Christopher Hesse, and John Schulman.  
2021. *Training verifiers to solve math word prob-  
lems*. In *arXiv preprint arXiv:2110.14168*. 395 396 397 398 399 400

Sunhao Dai, Chen Xu, Shicheng Xu, Liang Pang, Zhen-  
hua Dong, and Jun Xu. 2024. *Bias and unfairness  
in information retrieval systems: New challenges  
in the LLM era*. In *Proceedings of the 30th ACM  
SIGKDD Conference on Knowledge Discovery and  
Data Mining*, KDD ’24, page 6437–6447, New York,  
NY, USA. Association for Computing Machinery. 401 402 403 404 405 406 407

DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang,  
Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,  
Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang,  
Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhi-  
hong Shao, Zhuoshu Li, Ziyi Gao, and 181 others.  
2025. *Deepseek-R1: Incentivizing reasoning capa-  
bility in LLMs via reinforcement learning*. *Preprint*,  
arXiv:2501.12948. 408 409 410 411 412 413 414 415

Xiangjue Dong, Yibo Wang, Philip Yu, and James  
Caverlee. 2023. *Probing explicit and implicit gender  
bias through LLM conditional text generation*. In  
*Socially Responsible Language Modelling Research*. 416 417 418 419

Emilio Ferrara. 2023. *Should ChatGPT be biased? chal-  
lenges and risks of bias in large language models*.  
*First Monday*, 28(11). 420 421 422

Isabel O. Gallegos, Ryan A. Rossi, Joe Barrow,  
Md Mehrab Tanjim, Sungchul Kim, Franck Dernon-  
court, Tong Yu, Ruiyi Zhang, and Nesreen K. Ahmed. 423 424 425

426	2024. <a href="#">Bias and fairness in large language models: A survey</a> . <i>Computational Linguistics</i> , 50(3):1097–1179.	481
427		482
428		483
429	Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, and 542 others. 2024. <a href="#">The Llama 3 herd of models</a> . <i>Preprint</i> , arXiv:2407.21783.	484
430		485
431		486
432		487
433		488
434		489
435		490
436		491
437	Anthony G. Greenwald, Debbie E. McGhee, and Jordan L. K. Schwartz. 1998. Measuring individual differences in implicit cognition: The implicit association test. <i>Journal of Personality and Social Psychology</i> , 74(6):1464–1480.	492
438		493
439		494
440		495
441		496
442	Shashank Gupta, Vaishnavi Shrivastava, Ameet Deshpande, Ashwin Kalyan, Peter Clark, Ashish Sabharwal, and Tushar Khot. 2024. <a href="#">Bias runs deep: Implicit reasoning biases in persona-assigned LLMs</a> . In <i>The Twelfth International Conference on Learning Representations</i> .	497
443		498
444		499
445		500
446		501
447		502
448	Zhiwei He, Tian Liang, Jiahao Xu, Qiuzhi Liu, Xingyu Chen, Yue Wang, Linfeng Song, Dian Yu, Zhenwen Liang, Wenxuan Wang, Zhuosheng Zhang, Rui Wang, Zhaopeng Tu, Haitao Mi, and Dong Yu. 2025. <a href="#">Deepmath-103k: A large-scale, challenging, decontaminated, and verifiable mathematical dataset for advancing reasoning</a> .	503
449		504
450		505
451		506
452		507
453		508
454		509
455	Rem Hida, Masahiro Kaneko, and Naoaki Okazaki. 2024. <a href="#">Social bias evaluation for large language models requires prompt variations</a> . <i>Preprint</i> , arXiv:2407.03129.	510
456		511
457		512
458		513
459	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L��lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth��e Lacroix, and William El Sayed. 2023. <a href="#">Mistral 7B</a> . <i>Preprint</i> , arXiv:2310.06825.	514
460		515
461		516
462		517
463		518
464		519
465		520
466		521
467	Jiho Jin, Jiseon Kim, Nayeon Lee, Haneul Yoo, Alice Oh, and Hwaran Lee. 2024. <a href="#">KoBBQ: Korean bias benchmark for question answering</a> . <i>Transactions of the Association for Computational Linguistics</i> , 12:507–524.	522
468		523
469		524
470		525
471		526
472	Divyanshu Kumar, Umang Jain, Sahil Agarwal, and Prashanth Harshangi. 2024. <a href="#">Investigating implicit bias in large language models: A large-scale study of over 50 llms</a> . <i>Preprint</i> , arXiv:2410.12864.	527
473		528
474		529
475		530
476	Xiaopeng Li, Lixin Su, Pengyue Jia, Xiangyu Zhao, Suqi Cheng, Junfeng Wang, and Dawei Yin. 2023. <a href="#">Agent4Ranking: Semantic robust ranking via personalized query rewriting using multi-agent LLM</a> . <i>Preprint</i> , arXiv:2312.15450.	531
477		532
478		533
479		534
480		535
		536
		537
	Pan Lu, Liang Qiu, Wenhao Yu, Sean Welleck, and Kai-Wei Chang. 2023. <a href="#">A survey of deep learning for mathematical reasoning</a> . <i>Preprint</i> , arXiv:2212.10535.	481
		482
		483
		484
	Chandler May, Alex Wang, Shikha Bordia, Samuel R. Bowman, and Rachel Rudinger. 2019. <a href="#">On measuring social biases in sentence encoders</a> . In <i>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)</i> , pages 622–628, Minneapolis, Minnesota. Association for Computational Linguistics.	485
		486
		487
		488
		489
		490
		491
		492
	OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, and 262 others. 2024. <a href="#">GPT-4 technical report</a> . <i>Preprint</i> , arXiv:2303.08774.	493
		494
		495
		496
		497
		498
		499
		500
	Alicia Parrish, Angelica Chen, Nikita Nangia, Vishakh Padmakumar, Jason Phang, Jana Thompson, Phu Mon Htut, and Samuel Bowman. 2022. <a href="#">BBQ: A hand-built bias benchmark for question answering</a> . In <i>Findings of the Association for Computational Linguistics: ACL 2022</i> , pages 2086–2105, Dublin, Ireland. Association for Computational Linguistics.	501
		502
		503
		504
		505
		506
		507
	Pouya Pezeshkpour and Estevam Hruschka. 2024. <a href="#">Large language models sensitivity to the order of options in multiple-choice questions</a> . In <i>Findings of the Association for Computational Linguistics: NAACL 2024</i> , pages 2006–2017, Mexico City, Mexico. Association for Computational Linguistics.	508
		509
		510
		511
		512
		513
	Matthew Renze. 2024. <a href="#">The effect of sampling temperature on problem solving in large language models</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 7346–7356, Miami, Florida, USA. Association for Computational Linguistics.	514
		515
		516
		517
		518
		519
	Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane Suhr. 2024. <a href="#">Quantifying language models’ sensitivity to spurious features in prompt design or: How i learned to start worrying about prompt formatting</a> . In <i>ICLR</i> .	520
		521
		522
		523
		524
	Zara Siddique, Liam Turner, and Luis Espinosa-Anke. 2024. <a href="#">Who is better at math, Jenny or Jingzhen? uncovering stereotypes in large language models</a> . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 18601–18619, Miami, Florida, USA. Association for Computational Linguistics.	525
		526
		527
		528
		529
		530
		531
	Guangzhi Sun, Xiao Zhan, and Jose Such. 2024. <a href="#">Building better AI agents: A provocation on the utilisation of persona in LLM-based conversational agents</a> . In <i>Proceedings of the 6th ACM Conference on Conversational User Interfaces, CUI ’24</i> , New York, NY, USA. Association for Computing Machinery.	532
		533
		534
		535
		536
		537

Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard Hussonot, Pier Giuseppe Sessa, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex Botev, Alex Castro-Ros, Ambrose Slone, and 89 others. 2024. [Gemma: Open models based on Gemini research and technology](#). *Preprint*, arXiv:2403.08295.

Sean Xie, Saeed Hassanpour, and Soroush Vosoughi. 2024. [Addressing healthcare-related racial and LGBTQ+ biases in pretrained language models](#). In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 4451–4464, Mexico City, Mexico. Association for Computational Linguistics.

Chen Xu, Wenjie Wang, Yuxin Li, Liang Pang, Jun Xu, and Tat-Seng Chua. 2024. [A study of implicit ranking unfairness in large language models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 7957–7970, Miami, Florida, USA. Association for Computational Linguistics.

Kai-Ching Yeh, Jou-An Chi, Da-Chen Lian, and Shu-Kai Hsieh. 2023. [Evaluating interfaced LLM bias](#). In *Proceedings of the 35th Conference on Computational Linguistics and Speech Processing (ROCLING 2023)*, pages 292–299, Taipei City, Taiwan. The Association for Computational Linguistics and Chinese Language Processing (ACLCLP).

Ziyin Zhang, Zhaokun Jiang, Lizhen Xu, Hongkun Hao, and Rui Wang. 2024. [Multiple-choice questions are efficient and robust LLM evaluators](#). *Preprint*, arXiv:2405.11966.

Yachao Zhao, Bo Wang, Yan Wang, Dongming Zhao, Ruifang He, and Yuexian Hou. 2025. [Explicit vs. implicit: Investigating social bias in large language models through self-reflection](#). In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 1–12, Vienna, Austria. Association for Computational Linguistics.

Chujie Zheng, Hao Zhou, Fandong Meng, Jie Zhou, and Minlie Huang. 2024. [Large language models are not robust multiple choice selectors](#). In *The Twelfth International Conference on Learning Representations*.

## A Appendix

### A.1 Temperature and bias

To investigate how temperature might affect bias, we tested each model with different temperature values, sampling three responses for each question and treating the most common multiple-choice answer as the final answer. If the model outputs three unique answers, it is automatically treated as incorrect. As expected in Table 5, allowing random sampling introduces a substantial amount of variance in answers, resulting in higher bias with

System	You are a helpful American assistant. Make sure that your responses closely embody your identity.
User	Please answer the following multiple-choice question carefully. Only answer with a single letter. Do not respond with any other text, numbers, or symbols. <QUESTION>
Assistant	<RESPONSE>

Table 3: Example conversation with chat role and baseline prompt used in the experiment for GSM-MC and MathQA.

System	You are a helpful American assistant. Make sure that your responses closely embody your identity.
User	Please answer the following question carefully. Only respond with the final answer. Do not respond with any other text. Output the final answer within \boxed{ }. <QUESTION>
Assistant	<RESPONSE>

Table 4: Example conversation with chat role and baseline prompt used in the experiment for DeepMath.

higher temperatures. Given the observation that implicit bias and intelligence are inversely correlated, and previous research that observes a lack of significant influence of temperature on problem solving, it follows that temperature might not have much of an impact on implicit bias (Renze, 2024), and the increasing bias is solely the result of random sampling.

593  
594  
595  
596  
597  
598  
599  
600

<b>Models</b>	Llama 3.1	Llama 3.2	Llama 3.3	Mistral v0.3	Phi 3.5	Gemma 2
<i>Model Parameters</i>	8B	3B	70B	7B	3.8B	9B
$t = 0.2$	71.7	69.1	94.4	48.7	78.9	90.7
$t = 0.4$	50.8	17.8	93.8	30.3	72.6	87.9
$t = 0.6$	25.7	0	93.0	10.8	56.1	86.0
$t = 0.8$	0	0	92.1	0	39.1	84.1
$t = 1.0$	0	0	90.6	0	16.8	79.5

Table 5: DIF (GSM-MC) scores of different models across different temperatures.

<b>Models</b>	Llama 3.1	Llama 3.2	Llama 3.3	Mistral v0.3	Phi 3.5	Gemma 2
<i>Model Parameters</i>	8B	3B	70B	7B	3.8B	9B
Baseline Persona	356	159	597	258	362	472
American Indian	372	161	593	232	356	474
Asian	374	164	592	224	362	463
Black	369	167	598	224	359	475
Hispanic	370	163	598	209	356	467
Middle Eastern	361	155	594	221	360	468
Pacific Islander	366	162	591	207	359	473
White	361	155	590	243	360	474
Atheist	360	160	590	244	357	470
Buddhist	360	152	597	218	355	477
Christian	366	169	596	243	358	472
Hindu	361	159	593	234	352	477
Jewish	368	165	591	219	359	470
Mormon	365	164	591	243	365	474
Muslim	366	158	598	215	360	472
Female	355	164	593	252	359	474
Male	353	167	599	267	367	476
Non-binary	364	165	599	231	356	466
Gay	371	162	584	248	360	467
Straight	351	160	591	258	362	469
Able-bodied	354	158	588	249	361	465
Physically disabled	366	165	597	218	358	469

Table 6: Correct answers out of 1000 for the vanilla testing of personas when greedy decoding is used for text generation on GSM-MC.

<b>Models</b>	Llama 3.1	Llama 3.2	Llama 3.3	Mistral v0.3	Phi 3.5	Gemma 2
<i>Model Parameters</i>	8B	3B	70B	7B	3.8B	9B
Baseline Persona	182	177	195	198	99	159
American Indian	188	167	216	196	116	161
Asian	187	191	211	193	107	164
Black	187	191	206	191	107	167
Hispanic	185	184	211	186	112	163
Middle Eastern	186	186	202	189	114	155
Pacific Islander	183	193	206	193	113	162
White	169	190	201	188	111	155
Atheist	169	161	207	190	110	160
Buddhist	181	170	205	191	112	152
Christian	191	191	209	199	109	169
Hindu	186	170	204	191	117	159
Jewish	180	175	210	182	112	165
Mormon	192	182	205	188	113	164
Muslim	189	177	202	188	112	158
Female	189	192	202	193	110	164
Male	183	186	207	193	100	167
Non-binary	171	133	214	193	110	165
Gay	176	188	204	181	107	162
Straight	180	192	201	189	104	160
Able-bodied	169	178	201	186	113	158
Physically disabled	178	192	186	165	111	168

Table 7: Correct answers out of 1000 for the vanilla testing of personas when greedy decoding is used for text generation for MathQA.

<b>Models</b>	Llama 3.1	Llama 3.2	Llama 3.3	Mistral v0.3	Phi 3.5	Gemma 2
<i>Model Parameters</i>	8B	3B	70B	7B	3.8B	9B
Baseline Persona	49	39	55	33	29	34
American Indian	47	42	54	30	27	35
Asian	47	37	52	31	35	33
Black	49	41	53	31	34	35
Hispanic	48	38	52	31	36	35
Middle Eastern	48	38	52	32	34	31
Pacific Islander	49	40	55	27	36	35
White	49	42	57	33	36	35
Atheist	49	42	54	36	40	36
Buddhist	49	43	55	30	39	37
Christian	52	43	55	32	30	35
Hindu	49	38	52	29	35	36
Jewish	48	41	55	33	33	34
Mormon	51	40	55	33	30	36
Muslim	50	42	54	31	36	34
Female	47	38	55	33	28	40
Male	50	37	56	32	35	35
Non-binary	46	40	55	33	38	35
Gay	47	44	55	30	34	36
Straight	49	41	56	34	32	34
Able-bodied	48	42	57	32	35	31
Physically disabled	48	42	54	33	38	34

Table 8: Correct answers out of 1000 for the vanilla testing of personas when greedy decoding is used for text generation for DeepMath.