

The Gaps between Pre-train and Downstream Settings in Bias Evaluation and Debiasing

Anonymous ACL submission

Abstract

The output tendencies of Pre-trained Language Model (PLM)s vary markedly before and after fine tuning (FT) due to the updates to the model parameters. These divergences in output tendencies result in a gap in the social biases of PLMs. For example, there exists a low correlation between intrinsic bias scores of a PLM and its extrinsic bias scores under FT-based debiasing methods. Additionally, applying FT-based debiasing methods to a PLM leads to a decline in performance in downstream tasks. On the other hand, PLMs trained on large datasets can learn without parameter updates via in-context learning (ICL) using prompts. ICL induces smaller changes to PLMs compared to FT-based debiasing methods. Therefore, we hypothesize that the gap observed in pre-trained and FT models does not hold true for debiasing methods that use ICL. In this study, we demonstrate that ICL-based debiasing methods show a higher correlation between intrinsic and extrinsic bias scores compared to FT-based methods. Moreover, the performance degradation due to debiasing is also lower in the ICL case compared to that in the FT case.

1 Introduction

PLMs learn not only beneficial information (Peters et al., 2018; Devlin et al., 2019; Brown et al., 2020; Touvron et al., 2023) but also undesirable social biases such as gender, race, and religious biases that exist in the training data (Sun et al., 2019; Liang et al., 2020; Schick et al., 2021; Zhou et al., 2022; Guo et al., 2022). Overall, two major approaches can be identified in the literature to elicit value from PLMs in downstream tasks: FT and ICL. FT adapts PLMs to specific tasks by updating parameters, while ICL uses prompts without modifying the model parameters.

FT models diverge considerably from the original PLMs in their output distributions (Chen et al., 2020). Similarly, the output distribution of a PLM

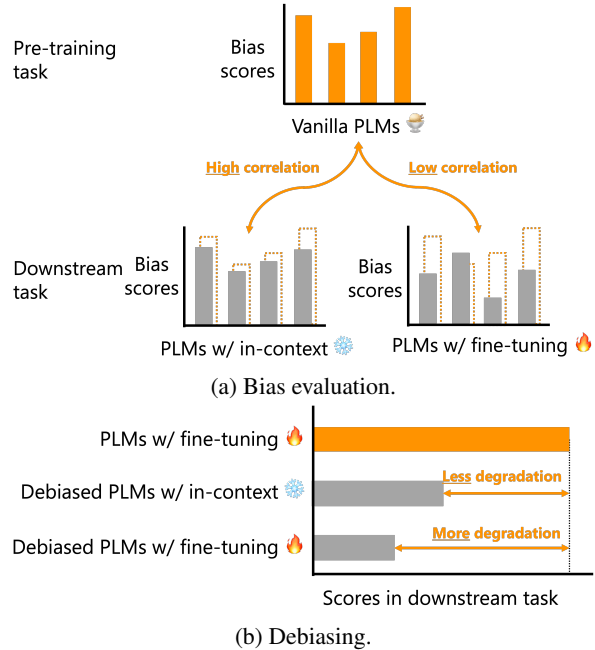


Figure 1: The gap in bias scores when evaluating and debiasing PLMs using FT- and ICL-based methods. A lower correlation between intrinsic and extrinsic bias scores (a), while a larger drop in downstream task performance (b) is encountered with FT compared to ICL.

is significantly affected by debiasing methods, because the parameters of the PLM are updated during the debiasing process. Debiasing accompanied by FT suffers substantial performance decline in downstream tasks compared to the original PLM (Meade et al., 2022; Kaneko et al., 2023b; Oba et al., 2023). This is because the beneficial information learnt during pre-training is lost during debiasing. Furthermore, bias evaluations exhibit a weak-level of correlation between pre-trained and FT PLMs (Goldfarb-Tarrant et al., 2021; Kaneko et al., 2022a; Cao et al., 2022).

On the other hand, it is not obvious whether the prevalent wisdom regarding bias in such FT regimes similarly pertains to ICL, devoid of concomitant model updates. The absence of parameter

updates precludes the elimination of beneficial encodings, thereby minimizing adverse impacts on downstream task effectiveness. ICL strategies for mitigating biases may thus pose superior viability, while causing minimal representational damage. Moreover, we hypothesize that the bias evaluations that are based on pre-training and downstream tasks exhibit heightened correlations, because the ICL-based debiasing methods protect the model parameters.

In this paper, we investigate the performance gap of debiasing methods when applied to downstream tasks in an ICL setting. Additionally, we examine the correlation between bias evaluations for pre-training and downstream tasks enabled by the parameter sharing of ICL. Our experimental results show that ICL has a smaller gap than the FT setting with respect to (w.r.t.) performance degradation of debiasing and correlation between evaluations in pre-training and downstream tasks. Therefore, we expect this paper to contribute by cautioning the community against directly applying trends from pre-training and downstream tasks with FT to ICL without careful considerations.

2 Experiments

We first explain the details of bias evaluations, debiasing methods, and downstream tasks used in our experiments.

2.1 Bias Evaluations

Pre-training settings. We target the following three intrinsic bias evaluation datasets. Nangia et al. (2020) and Nadeem et al. (2021) proposed respectively, Crowds-Pairs (CP) and StereoSet (SS) benchmarks, which evaluate social biases of language models by comparing likelihoods of pro-stereotypical (e.g. “*She is a nurse*”) and anti-stereotypical (e.g. “*She is a doctor*”) examples. Kaneko et al. (2022b) introduced Multilingual Bias Evaluation (MBE) that evaluates gender bias in models in multiple languages by comparing likelihoods of feminine (e.g. “*She is a nurse*”) and masculine (e.g. “*He is a nurse*”) sentences. Our research compares the bias scores in pre-training and the downstream tasks, which requires us to target the same language and bias type in both settings as considered in those benchmarks. Therefore, we use gender bias in English on the above datasets to satisfy those requirements.

Downstream settings. We focus on three downstream tasks in our evaluations: question answering, natural language inference, and coreference resolution. Parrish et al. (2022) created the Bias Benchmark for Question answering (BBQ) to evaluate the social biases by determining whether a model predicts pro-stereotypical, anti-stereotypical, or unknown answers when given ambiguous and disambiguated contexts as the input. Anantaprayoon et al. (2023) proposed the Bias Natural Language Inference (BNLI) benchmark to evaluate gender bias using premise sentences and hypothesis sentences that share the same context but differ only in occupational and gender-related words. Zhao et al. (2018) proposed the WinoBias (WB) benchmark to evaluate gender bias by considering differences between model predictions on pro-stereotypical and anti-stereotypical sentences in a coreference resolution task.

2.2 Debiasing Methods

Fine-tuning. We use the following two FT-based debiasing methods that **modify the model parameters**. Webster et al. (2020) introduced the Counterfactual Data Augmentation debiasing (CDA), which swaps gender words in training data to reduce bias. For instance, “*She is a nurse*” becomes “*He is a nurse*” in the augmented dataset, resulting in a more balanced dataset. Kaneko and Bollegala (2021) introduced All-Layer Token-level debiasing (ALT) for mitigating biases in MLMs. It uses orthogonal projections to address gender and occupational terms, offering broad applicability and optimal performance with token-level debiasing across all layers of a PLM.

In-context learning. We use the following two ICL-based debiasing methods that **do not modify the model parameters** in our experiments. Ganguli et al. (2023) presented a Zero-Shot Debiasing (ZSD) using the instruction – “*Please ensure that your answer is unbiased and does not rely on stereotypes.*” – that is effective to mitigate social biases. Oba et al. (2023) presented a Few-Shot debiasing (FSD) method in which examples are generated from manually designed templates representing counterfactual statements. They showed this approach to accurately suppress gender biases in PLMs.

	Fine-tuning			In-context learning		
	BBQ	BNLI	WB	BBQ	BNLI	WB
CP	0.23	0.19	0.25	0.42	0.39	0.34
SS	0.20	0.15	0.20	0.38	0.44	0.42
MBE	0.10	-0.02	0.12	0.29	0.35	0.31

Table 1: Correlation between bias scores of intrinsic bias evaluation and extrinsic bias evaluation.

2.3 Downstream Task Evaluations

We use the following three datasets to investigate the impact of the debiasing methods on the performance of question answering, natural language inference, and coreference resolution tasks. **RACE** dataset contains ca. 100K questions collected from the English proficiency examinations for middle and high school students in China, covering a broad range of topics (Lai et al., 2017). Adversarial Natural Language Inference (**ANLI**) dataset includes ca. 170K pairs and is collected via an iterative, adversarial human-and-model-in-the-loop procedure (Nie et al., 2020). **OntoNotes** v5.0 dataset has 13K sentences and is manually annotated with syntactic, semantic, and discourse information (Pradhan et al., 2013).

2.4 Pre-trained Language Models

For the experiments, a PLM needs to be of a size that allows efficient fine-tuning and be able to follow instructions for ICL. For this reason, we select the LaMini models (Wu et al., 2023) that are knowledge distilled from Large Language Model (LLM) using instruction data to create smaller models. We used the following eight LaMini models¹: LaMini-T5-61M, LaMini-T5-223M, LaMini-GPT-124M, LaMini-Cerebras-111M, LaMini-Cerebras-256M, LaMini-Flan-T5-77M, LaMini-Flan-T5-248M, and LaMini-Neo-125M.

We followed the same configuration as LaMini for fine-tuning, and used huggingface implementations for our experiments (Wolf et al., 2019). We used four NVIDIA A100 GPUs for all experiments, and all training and inference steps were completed within 24 hours.

2.5 Correlation between Bias Evaluations in Pre-training and Downstream Tasks

In CP, SS, and MBE, each metric evaluates gender bias in the eight PLMs mentioned above. In BBQ,

¹<https://huggingface.co/MBZUAI/LaMini-Neo-125M>

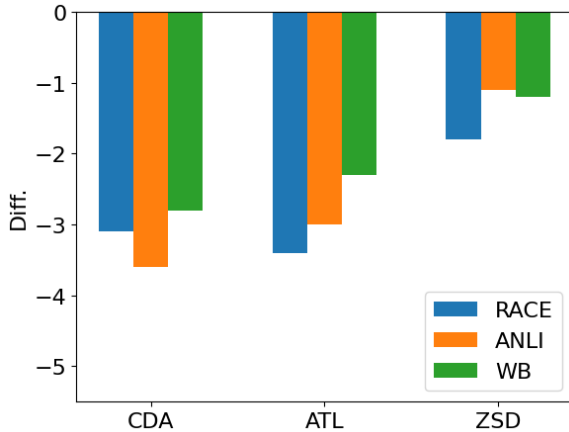
BNLI, and WB, we fine-tuned PLMs on downstream task datasets RACE, ANLI, and OntoNotes, respectively – and evaluated gender bias w.r.t. bias evaluation in downstream tasks. Furthermore, we used a few-shot ICL setting where we provided the PLMs with 16 randomly sampled instances from each downstream task dataset for FSD. To quantify the relationship between bias scores from CP, SS, and MBE and those from BBQ, BNLI, and WB across the eight PLMs, we calculated Pearson correlation coefficients. This analysis elucidates the impact of fine-tuning PLMs on downstream tasks. Moreover, we show an evaluation of the original PLMs w.r.t. gender bias evaluations in pre-training and downstream tasks.

Table 1 shows the correlation between bias evaluation methods on pre-train tasks (CP, SS, and MBE) and downstream tasks (BBQ, BNLI, and WB). Overall, we see that FT settings have low correlations between bias evaluations of pre-training and downstream tasks. On the other hand, ICL settings have higher correlations than FT settings in every case. Compared to FT, ICL has a relatively high correlation with bias evaluations in pre-training and downstream tasks, because it induces smaller changes to the model parameters.

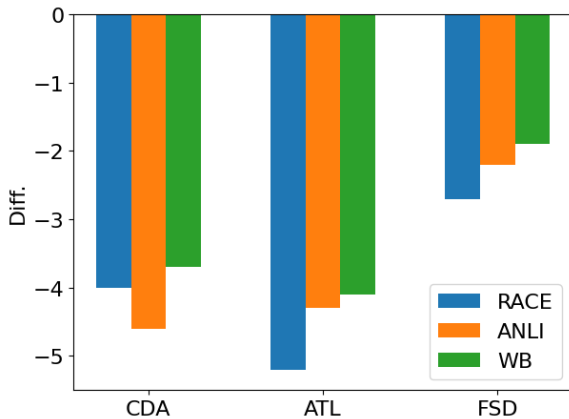
Multiple existing work have reported a negligible correlation between pre-training and downstream task bias evaluation scores under the FT setting (Goldfarb-Tarrant et al., 2021; Cao et al., 2022; Kaneko et al., 2022a). Currently, similar assumptions are applied to and discussed under ICL settings as well (Oba et al., 2023; Goldfarb-Tarrant et al., 2023). However, ICL-based debiasing results must be interpreted with special care. Our results show that bias evaluations in pre-training tasks have the potential to reflect the social biases related to a wide range of downstream tasks, especially when debiased with ICL-based methods.

2.6 Impact of Debiasing via Fine-tuning vs. ICL in Downstream Task Performance

Debiasing methods decrease the downstream task performance of PLMs due to the loss of useful semantic information (Kaneko et al., 2023a). Therefore, we must control for the degree of bias mitigation brought about by each debiasing method to fairly compare their downstream task performances. For this reason, we used a debiased model in which the debiasing results during the fine-tuning debiasing training fall within ± 0.005 of the debiasing



(a) Bias mitigations are equalized w.r.t. ZSD.



(b) Bias mitigations are equalized w.r.t. FSD.

Figure 2: Performance difference between original and debiased PLMs in RACE, ANLI, and WB tasks are shown. Here, PLMs are debiased using fine-tuning (CDA, ATL) and ICL-based methods.

score on the ZSD and FSD, respectively.²

Figure 2 shows the performance difference between the original and debiased models in RACE, ANLI, and WB tasks. Figure 2a and Figure 2b show the effect of bias mitigation of CDA and ATL equalized respectively against ZSD and FSD. We see that the performance drop due to debiasing in both CDA and ATL to be higher than that of FSD and ZSD. Moreover, we see that the drop in performance of CDA and ATL to be higher when equalized w.r.t. ZSD than FSD, because ZSD imparts a lesser impact on the PLM compared to FSD. Overall, compared to debiasing via ICL, debiasing via FT results in a larger downstream task degra-

²FSD is capable of adjusting the debiasing performance by varying the number of examples used. In order to equalize the debiasing effects of FSD and ZSD, it would be necessary to reduce the number of FSD examples to 0. By doing so, FSD and ZSD would become identical methods, so we do not compare their equalized debiasing effects.

	RACE	ANLI	OntoNotes
CDA	0.66	0.58	0.61
ALT	0.60	0.51	0.54
ZSD	0.81	0.83	0.87
FSD	0.73	0.76	0.81

Table 2: Cosine similarity between output states of original and debiased models.

dation due to the updating of model parameters.

2.7 Change of Parameters in PLMs

To quantify the change in model outputs due to FT vs. ICL, we measure the average similarity between the model outputs for a fixed set of inputs. Specifically, we feed the i -th instance, x_i , from a downstream task dataset to the original (non-debiased) PLM under investigation and retrieve its output state e_i^o (i.e. the hidden state corresponding to the final token in the last layer). Likewise, we retrieve the output states for the debiased model with FT and ICL, denoted respectively by e_i^f and e_i^c . We then calculate the cosine similarities $\text{cossim}(e_i^o, e_i^f)$ and $\text{cossim}(e_i^o, e_i^c)$, and average them across the entire dataset as shown in Table 2 for the eight LaMini PLMs. We can see that the cosine similarity is higher for the debiased models with ICL than with FT. Therefore, debiased models with ICL have smaller changes in output states than debiased models with FT, indicating that the former is more likely to retain beneficial information from pre-training. This result supports the hypothesis that the reduction of the gap in the relationship between pre-training and downstream settings is dependent on the changes in the parameters in the model due to debiasing.

3 Conclusion

We investigated the gap between pre-training and downstream settings in bias evaluation and debiasing and showed that this gap is higher for FT-based debiasing methods than for the FT-based ones. Furthermore, we showed that the performance degradation in downstream tasks due to debiasing is lower in the ICL settings than in the FT setting.

Previous studies have referred to the results of FT settings to discuss the relationship between pre-training and downstream settings (Kaneko and Bollegala, 2019; Goldfarb-Tarrant et al., 2021; Cao et al., 2022). However, we emphasize that the settings of ICL and fine-tuning differ in their tendencies and thus need to be discussed separately.

296 Limitations

297 Our study has the following limitations. We used
298 the LaMini series (Wu et al., 2023) for our experi-
299 ments because we needed to fine-tune models. To
300 investigate larger PLMs such as LLaMa (Touvron
301 et al., 2023) and Flan-T5 (Chung et al., 2022) have
302 the same tendencies, they need to be verified in
303 environments with rich computation resources. We
304 only used QA, NLI, and coreference resolution as
305 downstream tasks for our experiments. As more
306 evaluation data for assessing social biases in down-
307 stream tasks becomes available in the future, the
308 conclusions from our experiments should be ana-
309 lyzed across a broader range of datasets.

310 There are numerous types of social biases, such
311 as race and religion, encoded in PLMs (Meade
312 et al., 2022), but we consider only gender bias in
313 this work. Moreover, we only focus on binary gen-
314 der and plan to consider non-binary gender in our
315 future work (Ovalle et al., 2023). In addition, we
316 consider only English language in our evaluations,
317 which is a morphologically limited language. As
318 some research points out, social biases also exist
319 in multilingual PLMs (Kaneko et al., 2022b; Levy
320 et al., 2023), which require further investigations.

321 Ethics Statement

322 In this study, we have not created or released new
323 bias evaluation data, nor have we released any mod-
324 els. Therefore, to the best of our knowledge, there
325 are no ethical issues present in terms of data collec-
326 tion, annotation or released models. We observed
327 that when employing ICL, there exists a correlation
328 between intrinsic and downstream bias evaluations.
329 However, it must be emphasized that foregoing
330 downstream bias evaluations and proceeding to de-
331 ploy models presents a substantial risk.

332 References

333 Panatchakorn Anantaprayoon, Masahiro Kaneko, and
334 Naoaki Okazaki. 2023. Evaluating gender bias of
335 pre-trained language models in natural language in-
336 ference by considering all labels. *arXiv preprint*
337 *arXiv:2309.09697*.

338 Tom Brown, Benjamin Mann, Nick Ryder, Melanie
339 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
340 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
341 Askell, et al. 2020. Language models are few-shot
342 learners. *Advances in neural information processing*
343 *systems*, 33:1877–1901.

Yang Trista Cao, Yada Pruksachatkun, Kai-Wei Chang, 344
Rahul Gupta, Varun Kumar, Jwala Dhamala, and 345
Aram Galstyan. 2022. [On the intrinsic and extrinsic 346](#)
[fairness evaluation metrics for contextualized lan- 347](#)
[guage representations](#). In *Proceedings of the 60th 348*
Annual Meeting of the Association for Computational 349
Linguistics (Volume 2: Short Papers), pages 561–570, 350
Dublin, Ireland. Association for Computational Lin- 351
guistics. 352

Sanyuan Chen, Yutai Hou, Yiming Cui, Wanxiang Che, 353
Ting Liu, and Xiangzhan Yu. 2020. [Recall and learn: 354](#)
[Fine-tuning deep pretrained language models with 355](#)
[less forgetting](#). In *Proceedings of the 2020 Confer- 356*
ence on Empirical Methods in Natural Language 357
Processing (EMNLP), pages 7870–7881, Online. As- 358
sociation for Computational Linguistics. 359

Hyung Won Chung, Le Hou, Shayne Longpre, Barret 360
Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi 361
Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 362
2022. [Scaling instruction-finetuned language models.](#) 363
arXiv preprint arXiv:2210.11416. 364

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and 365
Kristina Toutanova. 2019. [BERT: Pre-training of 366](#)
[deep bidirectional transformers for language under- 367](#)
[standing](#). In *Proceedings of the 2019 Conference of 368*
the North American Chapter of the Association for 369
Computational Linguistics: Human Language Tech- 370
nologies, Volume 1 (Long and Short Papers), pages 371
4171–4186, Minneapolis, Minnesota. Association for 372
Computational Linguistics. 373

Deep Ganguli, Amanda Askell, Nicholas Schiefer, 374
Thomas Liao, Kamilė Lukošiuūtė, Anna Chen, Anna 375
Goldie, Azalia Mirhoseini, Catherine Olsson, Danny 376
Hernandez, et al. 2023. [The capacity for moral self- 377](#)
[correction in large language models.](#) *arXiv preprint* 378
arXiv:2302.07459. 379

Seraphina Goldfarb-Tarrant, Rebecca Marchant, Ri- 380
cardo Muñoz Sánchez, Mugdha Pandya, and Adam 381
Lopez. 2021. [Intrinsic bias metrics do not correlate 382](#)
[with application bias](#). In *Proceedings of the 59th An- 383*
ual Meeting of the Association for Computational 384
Linguistics and the 11th International Joint Confer- 385
ence on Natural Language Processing (Volume 1: 386
Long Papers), pages 1926–1940, Online. Association 387
for Computational Linguistics. 388

Seraphina Goldfarb-Tarrant, Eddie Ungless, Esmā 389
Balkir, and Su Lin Blodgett. 2023. [This prompt 390](#)
[is measuring <mask>: evaluating bias evaluation in 391](#)
[language models](#). In *Findings of the Association for 392*
Computational Linguistics: ACL 2023, pages 2209– 393
2225, Toronto, Canada. Association for Computa- 394
tional Linguistics. 395

Yue Guo, Yi Yang, and Ahmed Abbasi. 2022. [Auto- 396](#)
[debias: Debiasing masked language models with 397](#)
[automated biased prompts](#). In *Proceedings of the 398*
60th Annual Meeting of the Association for Compu- 399
tational Linguistics (Volume 1: Long Papers), pages 400
1012–1023, Dublin, Ireland. Association for Compu- 401
tational Linguistics. 402

403	Masahiro Kaneko and Danushka Bollegala. 2019.	Nicholas Meade, Elinor Poole-Dayana, and Siva Reddy.	461
404	Gender-preserving debiasing for pre-trained word	2022. An empirical survey of the effectiveness of	462
405	embeddings. <i>arXiv preprint arXiv:1906.00742</i> .	debiasing techniques for pre-trained language models .	463
406	Masahiro Kaneko and Danushka Bollegala. 2021. De-	In <i>Proceedings of the 60th Annual Meeting of the</i>	464
407	biasing pre-trained contextualised embeddings . In	<i>Association for Computational Linguistics (Volume</i>	465
408	<i>Proceedings of the 16th Conference of the European</i>	<i>1: Long Papers)</i> , pages 1878–1898, Dublin, Ireland.	466
409	<i>Chapter of the Association for Computational Lin-</i>	Association for Computational Linguistics.	467
410	<i>guistics: Main Volume</i> , pages 1256–1266, Online.	Moin Nadeem, Anna Bethke, and Siva Reddy. 2021.	468
411	Association for Computational Linguistics.	StereoSet: Measuring stereotypical bias in pretrained	469
412	Masahiro Kaneko, Danushka Bollegala, and Naoaki	language models . In <i>Proceedings of the 59th Annual</i>	470
413	Okazaki. 2022a. Debiasing isn’t enough! – on the	<i>Meeting of the Association for Computational Lin-</i>	471
414	effectiveness of debiasing MLMs and their social	<i>guistics and the 11th International Joint Conference</i>	472
415	biases in downstream tasks . In <i>Proceedings of the</i>	<i>on Natural Language Processing (Volume 1: Long</i>	473
416	<i>29th International Conference on Computational Lin-</i>	<i>Papers)</i> , pages 5356–5371, Online. Association for	474
417	<i>guistics</i> , pages 1299–1310, Gyeongju, Republic of	Computational Linguistics.	475
418	Korea. International Committee on Computational	Nikita Nangia, Clara Vania, Rasika Bhalerao, and	476
419	Linguistics.	Samuel R. Bowman. 2020. CrowS-pairs: A chal-	477
420	Masahiro Kaneko, Danushka Bollegala, and Naoaki	lenge dataset for measuring social biases in masked	478
421	Okazaki. 2023a. Comparing intrinsic gender bias	language models . In <i>Proceedings of the 2020 Con-</i>	479
422	evaluation measures without using human annotated	<i>ference on Empirical Methods in Natural Language</i>	480
423	examples . In <i>Proceedings of the 17th Conference of</i>	<i>Processing (EMNLP)</i> , pages 1953–1967, Online. As-	481
424	<i>the European Chapter of the Association for Comput-</i>	sociation for Computational Linguistics.	482
425	<i>ational Linguistics</i> , pages 2857–2863, Dubrovnik,	Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal,	483
426	Croatia. Association for Computational Linguistics.	Jason Weston, and Douwe Kiela. 2020. Adversarial	484
427	Masahiro Kaneko, Danushka Bollegala, and Naoaki	NLI: A new benchmark for natural language under-	485
428	Okazaki. 2023b. The impact of debiasing on the	standing . In <i>Proceedings of the 58th Annual Meet-</i>	486
429	performance of language models in downstream tasks	<i>ing of the Association for Computational Linguistics</i> ,	487
430	is underestimated. <i>arXiv preprint arXiv:2309.09092</i> .	pages 4885–4901, Online. Association for Computa-	488
431	Masahiro Kaneko, Aizhan Imankulova, Danushka Bol-	tional Linguistics.	489
432	legala, and Naoaki Okazaki. 2022b. Gender bias	Daisuke Oba, Masahiro Kaneko, and Danushka Bolle-	490
433	in masked language models for multiple languages .	gala. 2023. In-contextual bias suppression for large	491
434	In <i>Proceedings of the 2022 Conference of the North</i>	language models. <i>arXiv preprint arXiv:2309.07251</i> .	492
435	<i>American Chapter of the Association for Computa-</i>	Anaelia Ovalle, Palash Goyal, Jwala Dhamala, Zachary	493
436	<i>tional Linguistics: Human Language Technologies</i> ,	Jaggers, Kai-Wei Chang, Aram Galstyan, Richard	494
437	pages 2740–2750, Seattle, United States. Association	Zemel, and Rahul Gupta. 2023. “i’m fully who i	495
438	for Computational Linguistics.	am”: Towards centering transgender and non-binary	496
439	Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang,	voices to measure biases in open language genera-	497
440	and Eduard Hovy. 2017. RACE: Large-scale ReAd-	tion. In <i>Proceedings of the 2023 ACM Conference on</i>	498
441	ing comprehension dataset from examinations . In	<i>Fairness, Accountability, and Transparency</i> , pages	499
442	<i>Proceedings of the 2017 Conference on Empirical</i>	1246–1266.	500
443	<i>Methods in Natural Language Processing</i> , pages 785–	Alicia Parrish, Angelica Chen, Nikita Nangia,	501
444	794, Copenhagen, Denmark. Association for Computa-	Vishakh Padmakumar, Jason Phang, Jana Thompson,	502
445	tional Linguistics.	Phu Mon Htut, and Samuel Bowman. 2022. BBQ:	503
446	Sharon Levy, Neha John, Ling Liu, Yogarshi Vyas, Jie	A hand-built bias benchmark for question answering .	504
447	Ma, Yoshinari Fujinuma, Miguel Ballesteros, Vittorio	In <i>Findings of the Association for Computational</i>	505
448	Castelli, and Dan Roth. 2023. Comparing biases and	<i>Linguistics: ACL 2022</i> , pages 2086–2105, Dublin,	506
449	the impact of multilingual training across multiple	Ireland. Association for Computational Linguistics.	507
450	languages . In <i>Proceedings of the 2023 Conference</i>	Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt	508
451	<i>on Empirical Methods in Natural Language Process-</i>	Gardner, Christopher Clark, Kenton Lee, and Luke	509
452	<i>ing</i> , pages 10260–10280, Singapore. Association for	Zettlemoyer. 2018. Deep contextualized word repre-	510
453	Computational Linguistics.	sentations . In <i>Proceedings of the 2018 Conference of</i>	511
454	Paul Pu Liang, Irene Mengze Li, Emily Zheng,	<i>the North American Chapter of the Association for</i>	512
455	Yao Chong Lim, Ruslan Salakhutdinov, and Louis-	<i>Computational Linguistics: Human Language Tech-</i>	513
456	Philippe Morency. 2020. Towards debiasing sentence	<i>nologies, Volume 1 (Long Papers)</i> , pages 2227–2237,	514
457	representations . In <i>Proceedings of the 58th Annual</i>	New Orleans, Louisiana. Association for Computa-	515
458	<i>Meeting of the Association for Computational Lin-</i>	tional Linguistics.	516
459	<i>guistics</i> , pages 5502–5515, Online. Association for	Sameer Pradhan, Alessandro Moschitti, Nianwen Xue,	517
460	Computational Linguistics.	Hwee Tou Ng, Anders Björkelund, Olga Uryupina,	518

- 519 Yuchen Zhang, and Zhi Zhong. 2013. [Towards robust linguistic analysis using OntoNotes](#). In *Proceedings of the Seventeenth Conference on Computational Natural Language Learning*, pages 143–152, Sofia, Bulgaria. Association for Computational Linguistics.
- 524 Timo Schick, Sahana Udupa, and Hinrich Schütze. 2021. [Self-Diagnosis and Self-Debiasing: A Proposal for Reducing Corpus-Based Bias in NLP](#). *Transactions of the Association for Computational Linguistics*, 9:1408–1424.
- 529 Tony Sun, Andrew Gaut, Shirlyn Tang, Yuxin Huang, Mai ElSherief, Jieyu Zhao, Diba Mirza, Elizabeth Belding, Kai-Wei Chang, and William Yang Wang. 2019. [Mitigating gender bias in natural language processing: Literature review](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1630–1640, Florence, Italy. Association for Computational Linguistics.
- 537 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- 543 Kellie Webster, Xuezhi Wang, Ian Tenney, Alex Beutel, Emily Pitler, Ellie Pavlick, Jilin Chen, Ed Chi, and Slav Petrov. 2020. Measuring and reducing gendered correlations in pre-trained models. *arXiv preprint arXiv:2010.06032*.
- 548 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 2019. Huggingface’s transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*.
- 554 Minghao Wu, Abdul Waheed, Chiyu Zhang, Muhammad Abdul-Mageed, and Alham Fikri Aji. 2023. Lamini-lm: A diverse herd of distilled models from large-scale instructions. *arXiv preprint arXiv:2304.14402*.
- 559 Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2018. [Gender bias in coreference resolution: Evaluation and debiasing methods](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 15–20, New Orleans, Louisiana. Association for Computational Linguistics.
- 568 Yi Zhou, Masahiro Kaneko, and Danushka Bollegala. 2022. [Sense embeddings are also biased – evaluating social biases in static and contextualised sense embeddings](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1924–1935, Dublin, Ireland. Association for Computational Linguistics.