

The paper, “*Gender Bias in Coreference Resolution*” (Rudinger et al., 2018) proposes dataset schemas for measuring gender bias in Coreference Resolution Systems (CRS). To this end, the authors focus on gender bias with respect to occupations, evaluating the accuracy of Rule-based, Statistical, and Neural Coreference Systems in resolving a pronoun (male, female, or neutral) to a coreferent antecedent that is either an occupation or a participant. They constructed a challenge dataset, *Winogender schemas*, in the style of *Winograd schemas*, wherein a pronoun must be resolved to one of two previously mentioned entities in a sentence. The authors followed good practice by validating their hand-crafted dataset on Amazon’s Mechanical Turk (MTurk) with 10-way redundancy, with 94.9% of responses agreeing with their intended answers. This shows that the authors designed test sentences where correct pronoun resolution is not a function of gender. However, they do not report on their MTurk workers’ approval ratings or method for selecting them (e.g., qualification tests). They measure gender bias in coreference resolution systems by varying only the pronoun’s gender and examining the impact of this change on resolution (revealing cases where coreference systems may be more or less likely to recognize a pronoun as coreferent with a particular occupation based on pronoun gender). An unbiased model is expected to exhibit no sensitivity to pronoun gender in its resolution accuracy, resolving a male or female pronoun to an occupation or participant with equal likelihood. They correlate this bias with real-world and textual gender statistics. The models tested were: the Stanford multi-pass sieve system (Lee et al., 2011; Rule-based), Durrett and Klein’s (2013) statistical system, and the Clark and Manning (2016a) deep reinforcement system (neural).

To construct the dataset, the authors used a list of 60 one-word occupations obtained from Caliskan et al. (2017), with corresponding gender percentages available from the U.S. Bureau of Labor Statistics (BLS). For each occupation, there are two similar sentence templates: one in which the pronoun is coreferent with the occupation, and one in which it is coreferent with the participant. For each sentence template, there are two instantiations for the participant referent (a specific participant, e.g., “the passenger,” and a generic participant, “someone.”). Thus, the resulting evaluation set contains 720 sentences: 60 occupations  $\times$  2 sentence templates per occupation  $\times$  2 participants  $\times$  3 pronoun genders.

When each CRS’s predictions diverge based on pronoun gender, they do so in ways that resemble real-world occupational gender disparities. As shown in figure 4 of the paper, the systems’ gender preferences for occupations correlate with BLS and the gender statistics from text (Bergsma and Lin, 2006; B&L), which these systems access directly. All models performed worse in “gotcha” sentences, in which pronoun gender does not match the majority gender (BLS) of the occupation (correct answer). The paper discussed potential bias amplification involving the occupation manager: 38.5% female according to BLS, and mentions of “manager” in the B&L resource are only 5.18%, yet no managers were predicted to be female by any of the coreference systems (percentage-wise differences in real-world statistics may translate into absolute differences in system predictions). A Rule-based system may amplify the biases of its hand-crafted rules (which may amplify biases the dataset(s) and external resources), a Statistical system is vulnerable to the bias of a feature function associating an occupation with a pronoun (which can be informative, yet biased, for occupations occurring less frequently in the data), and a Neural system’s pre-trained embeddings is prone to encoding latent biases from its pre-training data. Gender bias is often introduced into the system as an unintended consequence of task-specific model construction or training. System-level biases can lead to further amplification in society through human-A.I interaction, causing a cycle of bias.

The bias exhibited by all three systems correlates both with real world employment statistics and the text statistics that these systems use. The authors note that while having high positive predictive value, the *Winogender schemas* have low negative predictive value. This follows from the dataset’s focus on gender bias in occupations; the models may be good at coreference resolution in this setting, but exhibit gender bias in different topics (e.g., crime data across genders). The *Winogender schemas* revealed varying degrees of gender bias in all three systems. In particular, 68% of male-female minimal pair test sentences are resolved differently by the Rule-based system; 28% for Statistical; and 13% for Neural. And overall, male pronouns were more likely to be resolved to the occupation antecedent than female or neutral pronouns across all systems.

In conclusion, the paper presents precise schemas for measuring the presence of gender bias in a CRS. Their dataset has gone through rigorous validation through crowdsourcing, and they use appropriate data (BLS and B&L) to compare these systems’ biases with; they are all of North American origin. The authors do not explore or inquire the generalizability of these results across more models, and there is no discussion on the importance of using an evaluation dataset whose national origins are the same as that of the models being evaluated. This is critical due to the varying degrees of gender bias across nations. Furthermore, there is no discussion of how the training method and national origin of training corpora for the CRS models impacts their gender bias in coreference resolution. In the end, the authors measured the presence occupational gender bias in Rule-based, Statistical, and Neural Coreference Resolution Systems successfully, but *Winogender schemas* may be extended broadly to probe for other manifestations of gender bias.

## References

- Rachel Rudinger, Jason Naradowsky, Brian Leonard, and Benjamin Van Durme. 2018. Gender Bias in Coreference Resolution. 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 8–14, New Orleans, Louisiana. Association for Computational Linguistics.
- Aylin Caliskan, Joanna J. Bryson, and Arvind Narayanan. 2017. In *Semantics derived automatically from language corpora contain human-like biases*. *Science*, 356(6334):183–186.
- Shane Bergsma and Dekang Lin. 2006. Bootstrapping path-based pronoun resolution. In *Proceedings of the 21st International Conference on Computational Linguistics and 44th Annual Meeting of the Association for Computational Linguistics*, pages 33–40, Sydney, Australia. Association for Computational Linguistics.
- Heeyoung Lee, Yves Peirsman, Angel Chang, Nathanael Chambers, Mihai Surdeanu, and Dan Jurafsky. 2011. Stanford’s multi-pass sieve coreference resolution system at the conll-2011 shared task. In *Conference on Natural Language Learning (CoNLL) Shared Task*.
- Greg Durrett and Dan Klein. 2013. Easy victories and uphill battles in coreference resolution. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing, Seattle, Washington*. Association for Computational Linguistics.
- Kevin Clark and Christopher D. Manning. 2016a. Deep reinforcement learning for mention-ranking coreference models. In *Empirical Methods on Natural Language Processing (EMNLP)*.
- Crime data. <https://www.ussc.gov/research/quick-facts/individuals-federal-bureau-prisons#:~:text=Individual%20and%20offense%20Characteristics,93.0%25%20are%20men>.