Gender Biases Unexpectedly Fluctuate in the Pre-training Stage of Masked Language Models

Kenan Tang* Hanchun Jiang*

University of Chicago {kenantang, carolinejiang}@uchicago.edu

Abstract

Masked language models pick up gender biases during pre-training. Such biases are usually attributed to a certain model architecture and its pre-training corpora, with the implicit assumption that other variations in the pretraining process, such as the choices of the random seed or the stopping point, have no effect on the biases measured. However, we show that severe fluctuations exist at the fundamental level of individual templates, invalidating the assumption. Further against the intuition of how humans acquire biases, these fluctuations are not correlated with the certainty of the predicted pronouns or the profession frequencies in pre-training corpora. We release our code and data to benefit future research¹.

1 Introduction

Masked language models (MLMs) succeed in solving natural language processing tasks, under the paradigm of fine-tuning the publicly released pretrained checkpoints (Devlin et al., 2019; Liu et al., 2019). The pre-training process uses large-scale human corpora, a practice that raises the concern whether harmful gender biases in human language are picked up by MLMs. Thus, much effort has been devoted to the quantification of gender biases in these models (Delobelle et al., 2022).

Meanwhile, the high cost of pre-training usually prohibits researchers from reproducing an MLM from scratch. Therefore, the single public checkpoint of an MLM architecture, such as bert-base-uncased, has become the default choice for bias quantification. The biases are reported as a property of the model identified by its architecture (Alnegheimish et al., 2022).

Inevitably, this reporting scheme understates the potential importance of hyperparameter choices. In this paper, we discuss two hyperparameters,

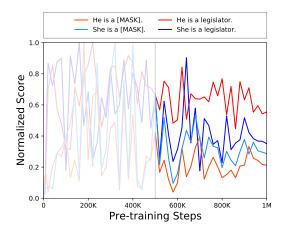


Figure 1: For RoBERTa, The probabilities of filling in gender pronouns into templates change unexpectedly after the training and validation loss has plateaued after 500K steps. Here results from the template for prior estimation and the one with the first profession in the full list are shown. For better visualization, the probability scores have been normalized by the maximal scores during pre-training.

namely the random seed and the number of training steps. Intuitively, they impact pre-training minimally. Changing random seeds should cause a small variance in the model performance, and further pre-training after the loss has plateaued should not change downstream performance a lot.

However, since the pre-training objective is not bias-aware, it is risky to assume that MLMs are equally biased when the random seed or the number of pre-training steps changes. The risk is amplified by the imperfection of popular gender bias metrics, which involve the probabilities of filling gender pronouns into only a few templates. This risk undermines the validity of applying continued pre-training to identify gender biases in a certain corpus (Bertsch et al., 2022).

In this paper, we demonstrate how gender biases fluctuate (Figure 1) when the seemingly irrelevant hyperparameters change. Inspired by intuitions about how humans understand gender stereotypes,

^{*}Equal contribution.

https://github.com/kt2k01/checkpoint-bias

we then provide fine-grained analyses on templates constructed from individual professions.

The main contributions of this paper are:

- By jointly analyzing pre-training dynamics and gender biases, we argue against framing gender biases as an innate property of a model architecture and its pre-training corpora.
- By differentiating the behaviors of a model on individual templates without aggregation, we describe in detail of how MLMs behave on this type of gender bias probe.

2 Related Work

2.1 Pre-training Dynamics

Intermediate pre-training checkpoints are needed for investigating the effects of different random seeds and stopping points. Despite the apparent availability, architecture designers seldom release these checkpoints. With limited budget, we can still fortunately rely on a limited number of replication studies where such checkpoints are released for models including BERT (Sellam et al., 2022), RoBERTa (Liu et al., 2021), GPT-2 (Karamcheti et al., 2021), and ALBERT (Chiang et al., 2020). While these works also probe the pre-training process from different angles, we add the missing discussion on gender bias probes, specifically for the two MLMs BERT (bert-base-uncased) and RoBERTa (roberta-base).

However, we omit the comparison between the replicated checkpoints and original ones from Hugging Face (Wolf et al., 2019), because the referred papers state that they have failed on exact replications.

2.2 Gender Bias Metrics

In static word embedding models, gender biases can be computed from distances between gendered words and non-gendered profession names (Caliskan et al., 2017). In contextual word embedding models, a similar method uses templates constructed from profession names as the input. For both MLMs (Delobelle et al., 2022) and autoregressive models such as GPT-2 (Alnegheimish et al., 2022), output probabilities of generating gender pronouns are divided to get a ratio as the gender bias score. While this is the only approach we use, we refer interested readers to comprehensive surveys in this area (Stanczak and Augenstein, 2021).

3 Methods

3.1 The Template-based Approach

In the template-based approach for gender bias measurement, a model checkpoint at m pretraining steps generate two sentences with different probabilities from one template. An example of such a generated sentence is

A template t has the four-component form:

$$t = [MASK] < VERB > < DET > < PROFESSION > .$$

The template is constructed by choosing the value of the latter three components. The second component <VERB> is chosen from "is/works as". The value of the third component <DET> chosen from "a/an" is determined by the initial phoneme of the value of the fourth component <PROFESSION>, chosen from a set constructed by merging lists from Delobelle et al. (2022) and Alnegheimish et al. (2022). The first list contains 30 male-stereotypical professions and 30 female-sterotypical ones. The second list contains 893 professions scraped from Wikipedia. The merging enables both a comprehensive comparison of many professions and a finergrained study of gender stereotypes, though the latter is not our focus. The size of the set is 923 after filtering repeated professions.

Taking t as input, the checkpoint m fills in the first component <code>[MASK]</code> with either "he" or "she". Denote the probabilities as P(he|m,t) and P(she|m,t), we define the bias score as

$$r(m,t) = \frac{P(he|m,t)}{P(she|m,t)}.$$

A ratio of r(m, t) = 1 implies that the checkpoint m is fair for the profession in t.

We use a matrix $\mathbf{R} \in \mathbb{R}^{s \times p}$ to denote the collection of $\mathbf{R}_{m,t} = r(m,t)$, where b=62 is the number of sampled training steps for RoBERTa (b=29 for BERT) and p=923 is the number of professions. Without ambiguity, t denotes either a template or the index of the profession in this template. The model choice and random seed will not be indexed in the notation \mathbf{R} but will be specified in context. All indices start from 0.

Upon observing that models output r(m,t) > 1 for most templates, some previous works have suggested normalizing r(m,t) by the priors of gender pronouns (Tal et al., 2022; Alnegheimish et al.,

2022). To estimate the priors, we use a template t_p where <PROFESSION> is replaced by a <code>[MASK]</code>, but the checkpoint still only predicts the first <code>[MASK]</code>. Then, the bias score normalized by priors is

$$n(m,t) = r(m,t) \cdot \frac{P(she|m,t_p)}{P(he|m,t_p)}.$$

A ratio of n(m,t)=1 similarly implies fairness, and we use a matrix $\mathbf{N} \in \mathbb{R}^{s \times p}$ to denote the collection of n(m,t). We will use both definitions.

To quantify whether a template is natural, we define the certainty c(m,t) for a checkpoint m and a template t as $c(m,t) = \mathrm{P}(he|m,t) + \mathrm{P}(she|m,t)$. Higher certainty suggests higher naturalness. A matrix $\mathbf{C} \in \mathbb{R}^{s \times p}$ denotes the collection of c(m,t).

3.2 Fluctuations

On one hand, we would like to quantify fluctuations of n(m,t) when it is not expected to change much at later stages of pre-training. We calculate the coefficient of variation (CV) from incomplete columns $N_{k:b,t}$ in the matrix N. Each incomplete column starts at the row indexed by k, a point when (1) the training and validation loss has already plateaued or (2) the downstream performances after fine-tuning do not improve much with further pre-training.

For the RoBERTa checkpoints, the plateau is reported to start at 50K steps out of 1M total steps. The number is 1M out of 2M for BERT. We more conservatively take the point of 500K steps for RoBERTa. In other words, we set k=36 for RoBERTa and k=18 for BERT. The choice is further discussed in Appendix B.

The CV is denoted by a vector \mathbf{v} , with entries

$$\mathbf{v}_t = \text{CV}(\mathbf{N}_{k:b,t}) = \frac{\text{SD}(\mathbf{N}_{k:b,t})}{\text{AM}(\mathbf{N}_{k:b,t})},$$

where SD is the standard deviation, and AM is the arithmetic mean. We compute the Pearson correlation coefficient between the CV ${\bf v}$ and the mean certainties ${\bf c}$, with entries

$$\mathbf{c}_t = \mathrm{AM}(\mathbf{C}_{k \cdot h t}).$$

On the other hand, to study the effect of changing random seeds, we base our investigation on 5 pre-training runs of BERT, as there is only a single run of RoBERTa. For each pair of random seeds, we compute the Pearson correlation coefficient between the pair of the averaged ratios **n**, with entries

$$\mathbf{n}_t = \mathrm{AM}(\mathbf{N}_{k:b.t}).$$

Additionally, we repeat the above procedure on the unnormalized \mathbf{R} instead of \mathbf{N} . Using \mathbf{R} , we similarly define vectors \mathbf{v} and \mathbf{r} .

3.3 Frequency in the Corpus

To identify a potential cause of the fluctuation, we look for the frequency of each profession in the pretraining corpora. Because we cannot count directly, we use frequencies in the BookCorpus as an estimation, since it is a large pre-training corpus shared by BERT and RoBERTa. We first use the Google Ngram API to query the yearly relative frequencies of each profession (case-sensitivity consistent with that of the model). Then, to get the total frequency, we multiply the yearly relative frequencies by the yearly total sizes of the corpus (from 1700 to 2000) released by the curators (Michel et al., 2011). The corpus is constantly evolving so that we cannot exactly compute the frequencies in the pre-training corpus used by the replication studies, but we believe this estimation is good enough (further discussion in Appendix C).

The frequencies for all professions are denoted by a vector **f**, whose entries are the inner product

$$\mathbf{f}_t = \mathbf{s} \cdot \mathbf{y}(t),$$

where the vector s denotes yearly sizes and $\mathbf{y}(t)$ denotes yearly frequencies of a profession t. We compute the Pearson correlation coefficient between the CV \mathbf{v} and the frequencies \mathbf{f} .

4 Results and Discussion

We present our results around the following research questions, inspired by intuitions about how models might pick up biases and how this can be analogized to the patterns of how humans hold gender stereotypes (Ellemers et al., 2018).

Here, we only base our discussion on the results when the <VERB> in the template is "is". The results from the alternative "works as" template are qualitatively the same (Appendix D).

RQ1: Do biases fluctuate, even after the training and validation loss has plateaued?

Humans' perception of gender stereotypes is relatively fixed. However, for both models, biases for all professions fluctuate considerably after the loss has plateaued (Figure 2). Take RoBERTa as an example, the smallest fluctuations of unnormalized ratios are still above a certain threshold $\min(\mathbf{v}) = 0.21$, while the highest fluctuations can reach $\max(\mathbf{v}) = 0.95$.

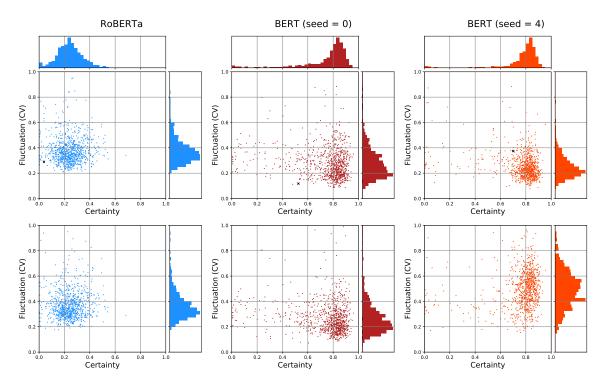


Figure 2: Certainty negligibly correlates with either unnormalized ratios (top row) or normalized ratios (bottom row), for both RoBERTa (first column) and BERT with different random seeds (second and third column). The certainty and fluctuation of the prior is represented by \times on the top 3 plots. Distributions are visualized by histograms.

Our result is consistent with the previous finding that model checkpoints in the pre-training process fail to behave consistently under probes of factual knowledge in the form of templates (Liu et al., 2021).

RQ2: Do biases fluctuate less if the model is more certain about its predictions?

With a strong belief in gender stereotypes, a human tends to perform according to the stereotype, reinforcing the stereotype in the process. While the models have distributions of certainties centered at different means, the correlations between the CV v and the certainties c are both weak (Figure 2).

RQ3: Do biases fluctuate less if the model sees the profession less often?

For humans, gender stereotypes are constantly reproduced as a result of mixed motivations upon the repeated observation of such stereotypes. For both models, the correlation between the frequencies of professions ${\bf f}$ and the CV ${\bf v}$ is negligible (|r|<0.20). In other words, the model checkpoints do not produce relatively fixed bias scores, for both frequently and infrequently seen professions.

RQ4: Do biases fluctuate for all professions at the same time?

Humans link genders with certain qualities such as aggressiveness or caring. As such qualities are shared as the requirements of multiple professions, human stereotypes of professions are reasonably not independent. We ask if models would behave similarly, by calculating the correlations between the pairs of ratio vectors $\mathbf{N}_{m_1,:}$ and $\mathbf{N}_{m_2,:}$ for any pair of checkpoints m_1 and m_2 in one pre-training run (Figure 3). In the results, we do observe that correlation is strong (r > 0.80) between some pairs of steps. However, the high fluctuation of ratios for many professions still weakens correlation after the loss has already plateaued $(m_i > k = 36)$.

RQ5: How do random seeds influence the fluctuations?

Human individuals perceive gender stereotypes to various degrees. When calculated from the normalized ratios \mathbf{N} , the distribution of the CV apparently shifts when the random seed changes (Figure 2, bottom row). However, the shift is due to the high variance of prior probabilities estimated from a single template. When calculated from the unnormalized ratios \mathbf{R} , this discrepancy between distributions has been largely mitigated (Figure 2, top row). This result warns against the over-reliance on the template-based prior estimation method.

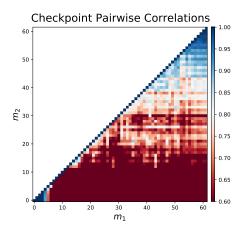


Figure 3: For pairs of RoBERTa checkpoints, the correlations between $N_{m,:}$ are strong towards the end of pre-training. However, fluctuation in normalized ratios of individual templates cannot be ignored. Correlations below 0.6 have been truncated in this visualization.

By further computing the pairwise correlation between the averaged ratios \mathbf{r} (or \mathbf{n}) of pre-training runs with different random seeds (Figure 4), we show that models from different runs return considerably different ($r \approx 0.60$) bias scores for individual professions. This result challenges the belief that a single pre-training run would allow a reliable estimation of gender biases.

5 Limits

We summarize the major limits of this paper as follows:

- There is no direct correlation between the bias scores measured from the templates and the actual bias the model is expected to exhibit in a downstream task (Kaneko et al., 2022). Moreover, we only limit our experiment to the one template-based approach of measuring bias on pre-trained MLMs.
- We have not discussed auto-regressive models such as ones from the GPT family. The only publicly available pre-training checkpoints are for GPT-2 (Karamcheti et al., 2021), which are not expected to behave similarly as the much larger and more popular GPT-3.
- We have only discussed the pre-training process on the largest corpora available. On one hand, continued pre-training could use smaller, domain-specific corpora (Gururangan et al., 2020). On the other hand, MLMs can be trained on different pre-training objectives

(Alajrami and Aletras, 2022). Both variations should influence the pre-training dynamics.

These factors limit the generalizability of our conclusions.

6 Conclusion

In this paper, we show that when measured by a popular template-based probe, the gender biases in MLMs fluctuate with respect to different random seeds and the number of pre-training steps after the loss has plateaued. Moreover, we provide how these fluctuations can be interpreted in ways that do not necessarily align with intuitions. Specifically, such fluctuations should be taken into account when MLMs are used for capturing biases from a certain corpus, or when biases of MLMs are compared for the evaluation of de-biasing methods.

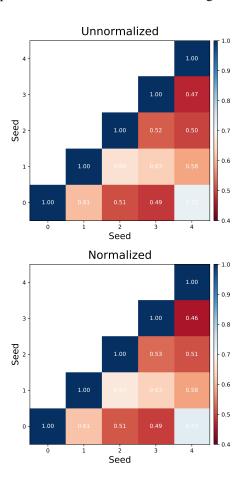


Figure 4: For BERT, different pre-training seeds lead to weakly correlated bias scores for all professions. The results for normalized and unnormalized ratios are different but close. All correlations are above 0.40.

References

- Ahmed Alajrami and Nikolaos Aletras. 2022. How does the pre-training objective affect what large language models learn about linguistic properties? In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 131–147, Dublin, Ireland. Association for Computational Linguistics.
- Sarah Alnegheimish, Alicia Guo, and Yi Sun. 2022. Using natural sentence prompts for understanding biases in language models. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2824–2830, Seattle, United States. Association for Computational Linguistics.
- Amanda Bertsch, Ashley Oh, Sanika Natu, Swetha Gangu, Alan W. Black, and Emma Strubell. 2022. Evaluating gender bias transfer from film data. In *Proceedings of the 4th Workshop on Gender Bias in Natural Language Processing (GeBNLP)*, pages 235–243, Seattle, Washington. Association for Computational Linguistics.
- Aylin Caliskan, Joanna J. Bryson, and Arvind Narayanan. 2017. Semantics derived automatically from language corpora contain human-like biases. *Science*, 356(6334):183–186.
- Cheng-Han Chiang, Sung-Feng Huang, and Hung-yi Lee. 2020. Pretrained language model embryology: The birth of ALBERT. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6813–6828, Online. Association for Computational Linguistics.
- Pieter Delobelle, Ewoenam Tokpo, Toon Calders, and Bettina Berendt. 2022. Measuring fairness with biased rulers: A comparative study on bias metrics for pre-trained language models. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1693–1706, Seattle, United States. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Naomi Ellemers et al. 2018. Gender stereotypes. *Annual review of psychology*, 69:275–298.
- Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A. Smith. 2020. Don't stop pretraining: Adapt language models to domains and tasks. In

- Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 8342–8360, Online. Association for Computational Linguistics.
- Masahiro Kaneko, Danushka Bollegala, and Naoaki Okazaki. 2022. Debiasing isn't enough! on the effectiveness of debiasing MLMs and their social biases in downstream tasks. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 1299–1310, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.
- Siddharth* Karamcheti, Laurel* Orr, Jason Bolton, Tianyi Zhang, Karan Goel, Avanika Narayan, Rishi Bommasani, Deepak Narayanan, Tatsunori Hashimoto, Dan Jurafsky, Christopher D. Manning, Christopher Potts, Christopher Ré, and Percy Liang. 2021. Mistral - a journey towards reproducible language model training.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv* preprint arXiv:1907.11692.
- Zeyu Liu, Yizhong Wang, Jungo Kasai, Hannaneh Hajishirzi, and Noah A. Smith. 2021. Probing across time: What does RoBERTa know and when? In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 820–842, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Jean-Baptiste Michel, Yuan Kui Shen, Aviva Presser Aiden, Adrian Veres, Matthew K. Gray, null null, Joseph P. Pickett, Dale Hoiberg, Dan Clancy, Peter Norvig, Jon Orwant, Steven Pinker, Martin A. Nowak, and Erez Lieberman Aiden. 2011. Quantitative analysis of culture using millions of digitized books. *Science*, 331(6014):176–182.
- Thibault Sellam, Steve Yadlowsky, Ian Tenney, Jason Wei, Naomi Saphra, Alexander D'Amour, Tal Linzen, Jasmijn Bastings, Iulia Raluca Turc, Jacob Eisenstein, Dipanjan Das, and Ellie Pavlick. 2022. The multiberts: BERT reproductions for robustness analysis. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net.
- Karolina Stanczak and Isabelle Augenstein. 2021. A survey on gender bias in natural language processing. *arXiv preprint arXiv:2112.14168*.
- Yarden Tal, Inbal Magar, and Roy Schwartz. 2022. Fewer errors, but more stereotypes? the effect of model size on gender bias. In *Proceedings of the 4th Workshop on Gender Bias in Natural Language Processing (GeBNLP)*, pages 112–120, Seattle, Washington. Association for Computational Linguistics.

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.

A Data Format

The data we release are in the following format. Columns in the data frame represent:

- The pronoun. The value is either "he" or "she" for bert-base-uncased and either "He" or "She" for roberta-base.
- The score. The value is between 0 and 1.
- The profession. The value is either a string from the list of 923 professions (e.g., "nurse") or the mask token used for prior estimation.
- The template. The value is a masked template without a pronoun.
- The full sentence. The value is a full sentence with a pronoun.
- The model name. The value is either roberta-base or bert-base-uncased.
- The index of the pre-training seed. The value is in $\{-1,0,1,2,3,4\}$ for bert-base-uncased and in $\{-1,0\}$ for roberta-base. Though not investigated in this paper, the data computed using the public checkpoints from Hugging Face are included and are indexed by -1.
- The checkpoint. The value is either the number of pre-training steps of a checkpoint or NaN, the latter representing the single available public checkpoint.

B Starting Point of the Plateau

We additionally provide results (Figure 6) from shorter plateau ranges for both RoBERTa (k=49) and BERT (k=24). Though smaller than in Figure 2, the fluctuation still exists and shows no correlation with certainties.

C Profession Frequencies

The sorted profession frequencies are shown in Figure 5. The professions with highest frequencies are summarized in Table 1. Note that case-sensitivity

Lowercased	Case-insensitive
model	president
author	secretary
official	model
president	author
judge	minister
police	judge
teacher	official
writer	professor
secretary	assistant
guide	governor
clerk	police
minister	teacher
physician	commissioner
assistant	clerk
engineer	guide
host	engineer
governor	writer
farmer	treasurer
artist	superintendent
pilot	miller

Table 1: The 20 most frequent professions in BookCorpus, ranked using either lowercased or case-insensitive frequencies.

can lead to very different results for certain professions, such as "President" as a title or "Miller" as a name (Figure 7).

D The Alternative Template

The alternative template with "works as" as the <VERB> does not result in qualitatively different results. While we use the scatter plot of fluctuation against certainties as an example here (Figure 8), all other plots have been released together with the code.

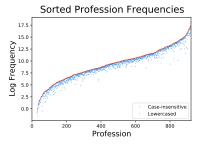


Figure 5: Frequencies of profession names in the Book-Corpus. All professions are sorted according to case-insensitive frequencies.

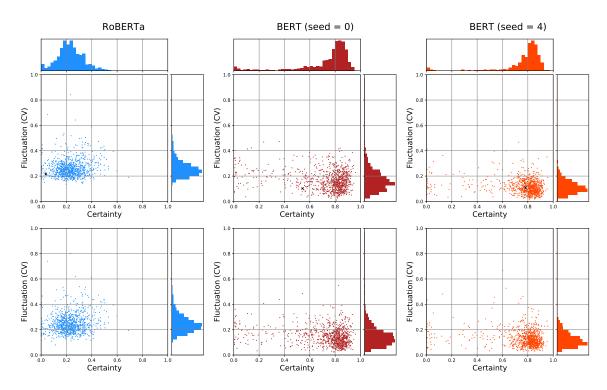


Figure 6: The results based on alternative lengths (k) of the plateau do not qualitatively differ from the ones based on the reported values (Figure 2). The top row shows the unnormalized results, and the bottom row shows normalized results.

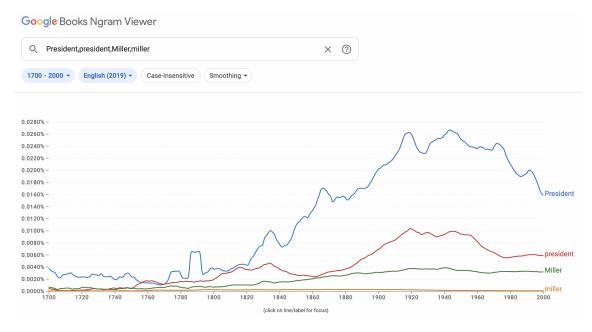


Figure 7: The case of the initial of the profession names influences the relative frequency returned from the Google Ngram API. Uppercase occurrences of certain professions outnumber the lowercased ones for professions such as "president" and "miller". The screenshot is taken on November 25, 2022.

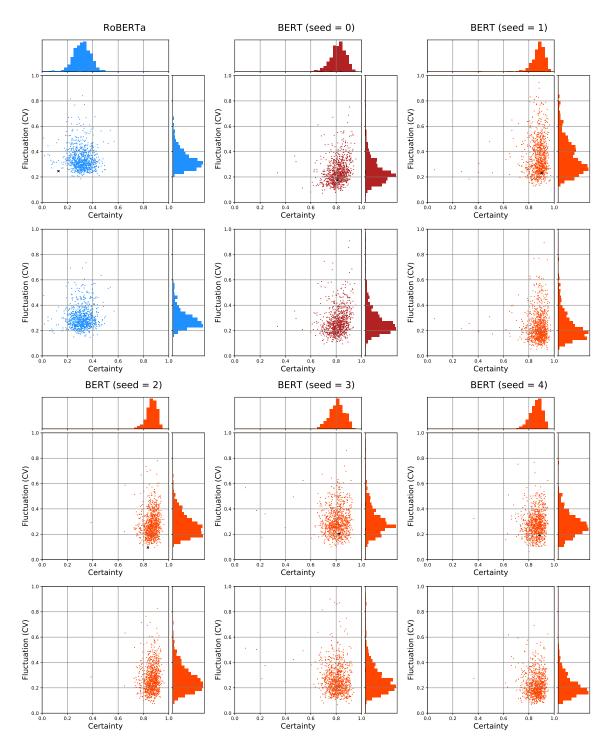


Figure 8: The results with the alternative "works as" template do not qualitatively differ from the ones with the "is" template (Figure 2). Here results from all 5 random seeds of BERT have been included. The odd rows are unnormalized results and the even rows are normalized results.