

Is Machine Learning Speaking My Language? Gender Bias and Under-Representation in Natural Language Processing Across Human Languages

Jeanna Neeffe Matthews, Yan Chen, Christopher Mahoney, Isabella Grasso, Esma Wali, Abigail Matthews, Thomas Middleton, Mariama Njie

Expanded Tech Report (May 4 2021)

ABSTRACT

Natural Language Processing (NLP) systems are at the heart of many critical automated decision-making systems making crucial recommendations about our future world. However, these systems reflect a wide range of bias, from gender bias to a bias in which voices they represent. In this paper, a team including speakers of 9 languages - Chinese, Spanish, English, Arabic, German, French, Farsi, Urdu, and Wolof - reports and analyzes measurements of gender bias in the Wikipedia corpora for these 9 languages. In the process, we also document how our work exposes crucial gaps in the NLP-pipeline for many languages. Despite substantial investments in multilingual support, the modern NLP-pipeline still systematically and dramatically under-represents the majority of human voices in the NLP-guided decisions that are shaping our collective future. We develop extensions to profession-level and corpus-level gender bias metric calculations originally designed for English and apply them to 8 other languages, including languages like Spanish, Arabic, German, French and Urdu that have grammatically gendered nouns including different feminine, masculine and neuter profession words. We compare these gender bias measurements across the Wikipedia corpora in different languages as well as across some corpora of more traditional literature.

KEYWORDS

gender bias, natural language processing, Wikipedia

1. INTRODUCTION

Corpora of human language are regularly fed into machine learning systems as a key way to learn about the world. Natural Language Processing plays a significant role in many powerful applications such as speech recognition, text translation, and autocomplete and is at the heart of many critical automated decision systems making crucial recommendations about our future world (Yordanov 2018)(Banerjee 2020)(Garbade 2018). Systems are taught to identify spam email, suggest medical articles or diagnoses related to a patient's symptoms, sort resumes based on relevance for a given position, and many other tasks that form key components of critical decision making systems in areas such as criminal justice, credit, housing, allocation of public resources and more. Much like facial recognition systems are often trained to represent white men more than black women (Buolamwini 2018), machine learning systems are often trained to represent human expression in languages such as English and Chinese more than in languages such as Urdu or Wolof.

The degree to which some languages are under-represented in commonly used text-based corpora is well-recognized, but the ways in which this effect is magnified throughout the NLP-tool chain is less discussed. Despite huge and admirable investments in multilingual support in projects like Wikipedia (Wikipedia 2020C), BERT (Devlin et al. 2018), Word2Vec (Mikolov et al. 2013), Wikipedia2Vec (Yamada et al. 2018)(Ousia 2016), Natural Language Toolkit (NLTK 2005), MultiNLI (Williams et al. 2020), many NLP tools are only developed for and tested on one or at most a handful of human languages and important advancements in NLP research are rarely extended to or evaluated for multiple languages. For some languages, the NLP-pipeline is streamlined: large publicly available corpora and even pre-trained models exist, tools run without errors and there is a rich set of research results applied to that language (Wali, 2020). However, for the vast majority of human languages, there is hurdle after hurdle. Even when a tool technically does support a given language, that support often comes with substantial caveats such as higher error rates and surprising problems. Also lack of representation at early stages of the pipeline (e.g. small corpora) adds to the lack of representation in later stages of the pipeline (e.g lack of tool support or research results not applied to that language).

In a highly influential paper “Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings”, Bolukbasi et al. (2016) developed a way to measure gender bias using word embedding systems like Word2vec. Specifically, they defined a set of gendered word pairs such as (“he”, “she”) and used the difference between these word pairs to define a gendered vector space. They then evaluated the relationship of profession words like doctor, nurse or teacher relative to this gendered vector space. They demonstrated that word embedding software trained on a corpus of Google news could associate men with the profession computer programmer and women with the profession homemaker. Systems based on such models, trained even with “representative text” like Google news, could lead to biased hiring practices if used to, for example, parse resumes and suggest matches for a computer programming job. However, as with many results in NLP research, this influential result has not been applied beyond English.

In some earlier work from our own team, “Quantifying Gender Bias in Different Corpora”, we applied Bolukbasi et al.’s methodology to computing and comparing corpus-level gender bias metrics across different corpora of the English text (Babaeianjelodar 2020). There we measured the gender bias in pre-trained models based on a “representative” Wikipedia and Book Corpus in English and compared it to models that had been fine-tuned with various smaller corpora including the General Language Understanding Evaluation (GLUE) benchmarks and two collections of toxic speech, RtGender and IdentityToxic. We found that, as might be expected, the RtGender corpora produced the highest gender bias score. However, we also found that the hate speech corpus, IdentityToxic, had lower gender bias scores than some of more representative corpora found in the GLUE benchmarks. By examining the contents of the IdentityToxic corpus, they found that most of the text in Identity Toxic reflected bias towards race or sexual orientation, rather than gender. These results confirmed the use of a corpus-level gender bias metric as a way of measuring gender bias in an unknown corpus and comparing across corpora, but again was only applied in English.

Those results also demonstrated the difficulty in predicting the degree of gender bias in unknown corpora without actually measuring it (Babaeianjelodar 2020). It is an increasingly common practice for application developers to start with pre-trained models and then add in a small amount of fine-tuning customized to their application. However, when the amount of gender bias learned from these “off-the-shelf” ingredients occurs unexpectedly, it can introduce unexpected learned gender bias in deployed applications in unpredictable ways, leading to significant problems when used to make critical decisions impacting the lives of individuals.

This common practice of application developers starting with a pre-trained model and then adding in a small amount of fine-tuning to customize their application has another important consequence. Pre-training from scratch using the large corpora necessary for meaningful NLP-results is expensive (i.e. days on a dozen CPUs). When a team can download a pre-trained model, they avoid this substantial overhead. Fine-tuning is much less expensive (i.e. hours on a single CPU). This makes NLP-based results accessible to a wider range of people, but only if such a pre-trained model is available for their language. When these easy to use pre-trained models exist for only a few languages (BERT, 2020), it steers more teams to languages with these pre-trained models and away from other languages and thus can further exacerbate the disparity in representation and participation among human languages.

In this paper, we build on the work of Bolukbasi et al. and our own earlier work to extend these important techniques in gender bias measurement and analysis beyond English. This is challenging because unlike English, many languages like Spanish, Arabic, German, French and Urdu, have grammatically gendered nouns including feminine, masculine and neuter or neutral profession words. We translate and modify Bolukbasi et al.’s defining sets and profession sets in English for 8 additional languages and develop extensions to the profession-level and corpus-level gender bias metric calculations for languages with grammatically gendered nouns. We use this methodology to analyze the gender bias in Wikipedia corpora for Mandarin Chinese, Spanish, English, Arabic, German, French, Farsi, Urdu, and Wolof.

In the process, we document many ways in which other sources of bias, beyond gender bias, are also being introduced by the modern NLP pipeline. Starting with the input data sets, Wikipedia is often used to train or test NLP-systems and there are substantial differences in the size and quality of the Wikipedia corpora across languages, even when adjusted for the number of speakers of each language. We demonstrate how the modern NLP pipeline not only reflects gender bias, but also leads to substantially over-representing some (especially English voices recorded in the digital text) and under-representing most others (speakers of most of the 7000 human languages and even writers of classic works that have not been digitized). As speakers of 9 languages who have recently examined the modern NLP toolchain, we highlight the difficulties that speakers of many languages face in having their thoughts and expressions included in the NLP-derived conclusions that are being used to direct the future for all of us.

In Section 2, we describe modifications that we made to the defining set and profession set proposed by Bolukbasi et al. in order to extend the methodology beyond English. In Section 3, we discuss the Wikipedia corpora and the occurrence of words in the modified defining and profession sets for 9 languages in Wikipedia. In Section 4, we extend Bolukbasi’s gender bias calculation to languages, like Spanish, Arabic, German, French and Urdu, with grammatically gendered nouns. We apply this to calculate and compare profession-level and corpus-level gender bias metrics for Wikipedia corpora in the 9 languages. In Section 5, we discuss the need to assess corpora beyond Wikipedia. We conclude in Section 6.

Specific aspects of this ongoing work have been described in Wali et al. (2020), Chen et al. (2021) and Matthews et al. (2021).

2. DEFINING SETS AND PROFESSION SETS

Word embedding is a powerful NLP technique that represents words in the form of numeric vectors. It is used for semantic parsing, representing the relationship between words, and capturing the context of a word in a document (Karani 2018). For example, Word2vec is a system used to efficiently create word embeddings by using a two-layer neural network that efficiently processes huge data sets with billions of words, and with millions of words in the vocabulary (Mikolov 2013).

Bolukbasi et al. developed a method for measuring gender bias using word embedding systems like Word2vec. Specifically, they defined a set of highly gendered word pairs such as (“he”, “she”) and used the difference between these word pairs to define a gendered vector space. They then evaluated the relationship of profession words like doctor, nurse or teacher relative to this gendered vector space. Ideally, profession words would not reflect a strong gender bias. However, in practice, they often do. According to such a metric, doctor might be male biased or nurse female biased based on how these words are used in the corpora from which the word embedding model was produced. Thus, this gender bias metric of profession words as calculated from the Word2Vec model can be used as a measure of the gender bias learned from corpora of natural language.

In this section, we describe the modifications we made to the defining set and profession set proposed by Bolukbasi et al. in order to extend the methodology beyond English. Before applying these changes to other languages, we evaluate the impact of the changes on calculations in English. In this section, we also describe the Wikipedia corpora we used across 9 languages and analyze the occurrences of our defining set and profession set words in these corpora.

2.1 Defining Set

We call the list of gendered word pairs used to define what a gendered relationship looks like a defining set. Bolukbasi et al’s original defining set contained 10 English word pairs (she-he, daughter-son, her-his, mother-father, woman-man, gal-guy, Mary-John, girl-boy, herself-himself, and female-male) (Boblusi et al. 2016). We began with this set, but made substantial changes in order to compute gender bias effectively across 9 languages.

Specifically, we removed 6 of the 10 pairs, added 3 new pairs and translated the final set into 8 additional languages. Table 1 summarizes the reasons for the 6 removals. We also added 3 new pairs (queen-king, wife-husband, and madam-sir) for which more consistent translations were available across languages. Our final defining set thus contains 7 word pairs. Table 2 shows our translations of this final defining set across the 9 languages included in our study.

Table 1: Removed word pairs

Word pair	Reasons
she-he	They are the same in some languages like Wolof, Farsi, Urdu, and German relying on context for disambiguation. For example, in German, “sie” without additional context would be very difficult to determine whether it correlates to she, them, they, or you.

her-his	In some languages like French and Spanish, the gender of the possessive word refers to the object rather than to the person to whom the object belongs. In German without context, “ihrer” could mean her, your, theirs or yours.
gal-guy	They are the same words in some languages like Wolof. There are no translations for these words in some languages like Urdu and Arabic.
Mary-John	Simply translating Mary and John to other languages is problematic, but so is trying to identify some alternate "typical" male-female names.
herself-himself	They are the same words in some languages like Wolof, Farsi, Urdu, and German relying on context for disambiguation.
female-male	They can be nouns or adjectives in many languages which introduces ambiguity.

Table 2: Final defining set translated across languages. Note: Wolof is primarily a spoken language and is often written as it would be pronounced in English, French and Arabic. This table shows it written as it would be pronounced in French.

English	Chinese	Spanish	Arabic	German	French	Farsi	Urdu	Wolof
woman	女人	mujer	النساء	Frau	femme	زن	عورت	Jigéen
man	男人	hombre	رجل	Mann	homme	مرد	آدمی	Góor
daughter	女儿	hija	ابنة	Tochter	filles	دختر	بیٹی	Doom ju jigéen
son	儿子	hijo	ولد	Sohn	fil	پسر	بیٹا	Doom ju góor
mother	母亲	madre	ام	Mutter	mère	مادر	مان	Yaay
father	父亲	padre	اب	Vater	père	پدر	باپ	Baay
girl	女孩	niña	ابنة	Mädchen	filles	دختر	لڑکی	Janxa
boy	男孩	niño	صبي	Junge	garçon	پسر	لڑکا	Xale bu góor
queen	女王	reina	ملكة	Königin	reine	ملکہ	ملکہ	Jabari buur
king	国王	rey	ملك	König	roi	پادشاه	بادشاه	Buur
wife	妻子	esposa	زوجة	Ehefrau	épouse	همسر	بیوی	Jabar
husband	丈夫	esposo	الزوج	Ehemann	mari	شوهر	شوہر	jëkkër
madam	女士	señora	سیدتی	Dame	madame	خانم	محترمہ	Ndawsi
sir	男士	señor	سیدی	Herr	monsieur	آقا	جناب	Góorgui

2.2 Profession Set

We began with Bolukbasi et al.'s profession word set in English, but again made substantial changes in order to compute gender bias effectively across 9 languages. Bolukbasi et al. had an original list of 327 profession words (Bolukbasi 2016), including some words that would not technically be classified as professions like saint or drug addict. We narrowed this list down to 32 words including: nurse, teacher, writer, engineer, scientist, manager, driver, banker, musician, artist, chef, filmmaker, judge, comedian, inventor, worker, soldier, journalist, student, athlete, actor, governor, farmer, person, lawyer, adventurer, aide, ambassador, analyst, astronaut, astronomer, and biologist. We tried to choose a diverse set of professions from creative to scientific, from high-paying to lower-paying, etc. that occurred in as many of the 9 languages as we could. As with Bolukbasi et al.'s profession set, one of our profession words, person, is not technically a profession, but we kept it because, unlike many professions, it is especially likely to have a native word in most human languages.

The primary motivation for reducing the profession set from 327 to 32 was to reduce the work needed to translate and validate all of them in 9 languages. Even with 32 words, there were substantial complexities in translation. As we mentioned, languages with grammatically gendered nouns can have feminine, masculine and neuter words for the same profession. For instance, in Spanish, the profession “writer” will be translated as “escritora” for women and “escritor” for men, but the word for journalist, “periodista”, is used for both women and men.

Profession words are often borrowed from other languages. In this study, we found that Urdu and Wolof speakers often use the English word for a profession when speaking in Urdu or Wolof. In some cases, there is a word for that profession in the language as well and in some cases, there is not. For example, in Urdu, it is more common to use the English word “manager” when speaking even though there are Urdu words for the profession manager. In written Urdu, manager could be written directly in English characters (manager) or written phonetically as the representation of the word manager using Urdu/Arabic characters (مينيجر) or written as an Urdu word for manager (منتظم/منتظمه). However, for driver, the English word is almost always used and written either in English or Arabic characters (driver or ڈرائيور).

A similar pattern occurs in Wolof and there are some additional complicating factors as well. Wolof is primarily a spoken language that when written is always transcribed phonetically. This may be done using English, French or Arabic character sets and pronunciation rules. Thus, for the same pronunciation, spelling can vary substantially and this complicates NLP processing such as with Word2Vec significantly. For example the word father could be written as Bai or Baay.

After making these substantial changes to the defining sets and profession sets, the first thing we did was analyze their impact on gender bias measurements in English. Using both Bolukbasi et al.'s original defining and professions sets and our modified sets, we computed the gender bias scores on the English Wikipedia corpus. We computed gender bias results using both our 7 defining set pairs and 32 profession words and Bolukbasi et al.'s 10 defining set pairs and 327 profession words. We conducted a T-test and even with these substantial changes the T-test results were insignificant, inferring that the resulting gender bias scores in both instances have no statistically significant difference for the English Wikipedia corpus. This result was an encouraging validation that our method was measuring the same gender bias effects as in Bolukbasi et al. even with the modified and reduced defining and profession sets.

While our goal in this study was to identify a defining set and profession set that could more easily be used across many languages and for which the T-test results indicated no statistically significant difference in results over the English Wikipedia corpus, it would be interesting to repeat this analysis with additional variations in the defining set and profession set. For example, we considered adding additional pairs like sister-brother or grandmother-grandfather. In some languages like Chinese, Arabic and Wolof, there are different words for younger and older sister or brother. In Chinese, there are different words for paternal and maternal grandmother and grandfather. We also considered and discarded many other profession words such as bartender, policeman, celebrity, and electrician. For example, we discarded bartender because it is not a legal profession in some countries.

3. WIKIPEDIA CORPORA ACROSS LANGUAGES

Bolukbasi et al. applied their gender bias calculations to a Word2Vec model trained with a corpus of Google news in English. In Babaeianjelodar et al. (2020), we used the same defining and profession sets as Bolukbasi et al. to compute gender bias metrics for a BERT model trained with Wikipedia and a BookCorpus also in English. In this paper, we train Word2Vec models using our modified defining and profession sets and the Wikipedia corpora for 9 languages. Specifically, we use the Chinese, Spanish, Arabic, German, French, Farsi, Urdu and Wolof corpora downloaded from Wikipedia on 2020-06-20.

3.1 Differences in Wikipedia across Languages

While there are Wikipedia corpora for all 9 of our languages, they differ substantially in size and quality. Table 3 illustrates that even when accounting for differences in the number of speakers of each language worldwide, some languages have substantially more representation than others. It is interesting to note that the number of speakers of each language does not track evenly with the number of articles. French has the highest ratio of articles per 1000 Speakers at 29.15 and Wolof the lowest at 0.26. English has the largest number of articles, even including languages not in our list, but its ratio of articles to speakers is lower than some other languages. In addition to the difference in the number of articles and articles/1000 speakers, Wikipedia corpora for different languages also vary widely along many other dimensions including the total size of corpora in MB, total pages, percentage of articles that are simply stub articles with no content, number of edits, number of admins working in that language, total number of users and total number of active users.

Table 3: Comparing the number of speakers of a language to the size of the Wikipedia Corpora for that language. For the number of articles and estimates of the number of speakers of Chinese, Spanish, English, Arabic, German, French, Farsi, Urdu, and Wolof. All of the numbers of speakers only accounted for native speakers. These statistics on the number of articles and the number of speakers are taken from Wikipedia itself (WikipediaA)(WikipediaB)(WikipediaC)(WikipediaD).

Language	Number of Articles	Number of Speakers (thousand)	Articles/1000 Speakers
Chinese	1,149,477	921,500	1.25
Spanish	1,629,888	463,000	3.52
English	6,167,101	369,700	16.68
Arabic	1,067,664	310,000	3.44
German	2,485,274	95,000	26.16
French	2,253,331	77,300	29.15
Farsi	747,551	70,000	10.68
Urdu	157,475	69,000	2.28
Wolof	1,422	5,500	0.26

Wikipedia is a very commonly used dataset for testing NLP tools and even for building pre-trained models. However, for many reasons, a checkmark simply saying that a Wikipedia corpus exists for a language hides many caveats to full representation and participation. In addition to variation in size and quality across languages, not all speakers of a language have equal access to contributing to Wikipedia. For example, in the case of Chinese, Chinese speakers in mainland China have little access to Wikipedia because it is banned by the Chinese government (Siegel 2019). Thus, Chinese articles in Wikipedia are more likely to have been contributed by the 40 million Chinese speakers in Taiwan, Hong Kong, Singapore and elsewhere (Su 2019). In other

cases, the percentage of speakers with access to Wikipedia may vary for other reasons such as access to computing devices and Internet access.

Using Wikipedia as the basis of pre-trained models and testing of NLP tools also means that the voices of those producing digital text are prioritized. Even authors of classic works of literature that fundamentally shaped cultures are under-represented in favor of writers typing Wikipedia articles on their computer or even translating text written in other languages with automated tools.

3.2 Word count results

One critical aspect of our process was to examine the number of times each word in our defining set (7 pairs) and 32 profession words occurs in the Wikipedia corpus for each language. This proved an invaluable step in refining our defining and profession sets, understanding the nature of the Wikipedia corpora themselves, catching additional instances where NLP tools were not designed to handle the complexities of some languages and even catching simple errors in our own translations and process. For example, when our original word count results for German showed a count of zero for all words, we discovered that even though all nouns in German are capitalized, in the Word2vec processed Wikipedia corpus for German, all words were in lowercase. This was an easy problem to fix, but illustrates the kind of “death by a thousand cuts” list of surprising errors that can occur for many languages throughout the NLP pipeline.

One important limitation to note is that for many languages, if a word is expressed with a multi-word phrase (e.g. astronomer(علم الفلك) in Arabic). The word count reported by Word2Vec for this phrase will be zero. For each language, there is a tokenizer that identifies the words or phrases to be tracked. In many cases, the tokenizer identifies words as being separated by a space. The Chinese tokenizer however attempts to recognize when multiple characters that are separated with spaces should be tracked as a multi-character word or concept. This involves looking up a string of characters in a dictionary. Once again this demonstrates the types of surprising errors that can occur for many languages throughout the NLP pipeline. It is also possible to add the word vectors for component words together as a measure of the multi-word pair, but this is not always ideal. In this study, we did not attempt this, but it would be interesting future work.

Another important factor is that, as we described earlier, the Wikipedia corpora for some languages are quite small. In Wolof, for example, only two of our profession words occurred (“nit”, the word for person, occurred 1401 times and waykat, the word for musician, occurred 5 times). This is partly because of multi-word pairs and partly because of variants in spelling. However, the percentage of profession words amongst the total words for Wolof is similar to that of other languages suggesting that it is the small size of the Wolof corpus that is the primary problem. In fact, the percentage of profession words varied from 0.014% and 0.037% across the 9 languages and Wolof had one of the higher percentages at 0.026%. On the other hand, Wolof's overall Wikipedia corpus is tiny (1422 articles or less than 1% of the number of articles even in Urdu, the next smallest corpora) and that simply isn't a lot of text with which to work. Even so, Wolof is still much better represented in Wikipedia than the vast majority of the over 7000 human languages spoken today! This is another clear illustration of how the gap in support for so many languages leads directly to the under-representation of many voices in NLP-guided decision-making.

In the supplementary material, we will include tables with the word counts for all languages - defining sets in Appendix 1 and profession sets in Appendix 2.

4. PROFESSION AND CORPORA LEVEL GENDER BIAS METRICS

We have already described how we established a modified defining set and profession set for use across 9 languages and then evaluated the use of these sets of words in Wikipedia. We also described how we used the Wikipedia corpora of these 9 languages to train Word2Vec models for each language. In this section, we describe how we extend Bolukbasi et al.'s method for computing the gender bias of each word.

We begin with Bolukbasi et al.'s method for computing a gender bias metric for each word. Specifically, each word is expressed as a vector by Word2Vec and we calculate the center of the vectors for each definitional pair. For example, to calculate the center of the definitional pair she/he, we average the vector for “she” with the vector for “he”. Then, we calculate the distance of each word in the definitional pair from the center by subtracting the center from each word in the pair (e.g. “she” - center). We then apply Principal Component Analysis (PCA) to the matrix of these distances. PCA is an approach that compresses multiple dimensions into fewer dimensions, ideally in a way that the information within the original data is not lost. Usually the number of

reduced dimensions is 1-3 as it allows for easier visualization of a dataset. Bolukbasi et al. used the first eigenvalue from the PCA matrix (i.e. the one that is larger than the rest). Because the defining set pairs were chosen to be highly gendered, they expected this dimension to be related primarily to gender and therefore called it the gender direction or the g direction. (Note: The effectiveness of this compression can vary and in some cases, the first eigenvalue may not actually be that much larger than the second. We see cases of this in our study as we will discuss.)

We use Bolukbasi et al.’s formula for direct gender bias:

$$\text{DirectBias}_c = \frac{1}{|N|} \sum_{w \in N} |\cos(w, g)|^c$$

where N represents the list of profession words, N_count represents the total occurrences of the profession words in the corpora, g represents the gender direction calculated, w represents each profession word, w_count is the occurrences of w in the corpora and c is a parameter to measure the strictness of the bias. In this paper, we used $c = 1$, however c values and their effects are explained in more detail in Bolukbasi et al. We examine this gender bias score both for the individual words as well as an average gender bias across profession words as a measure of gender bias in a corpus.

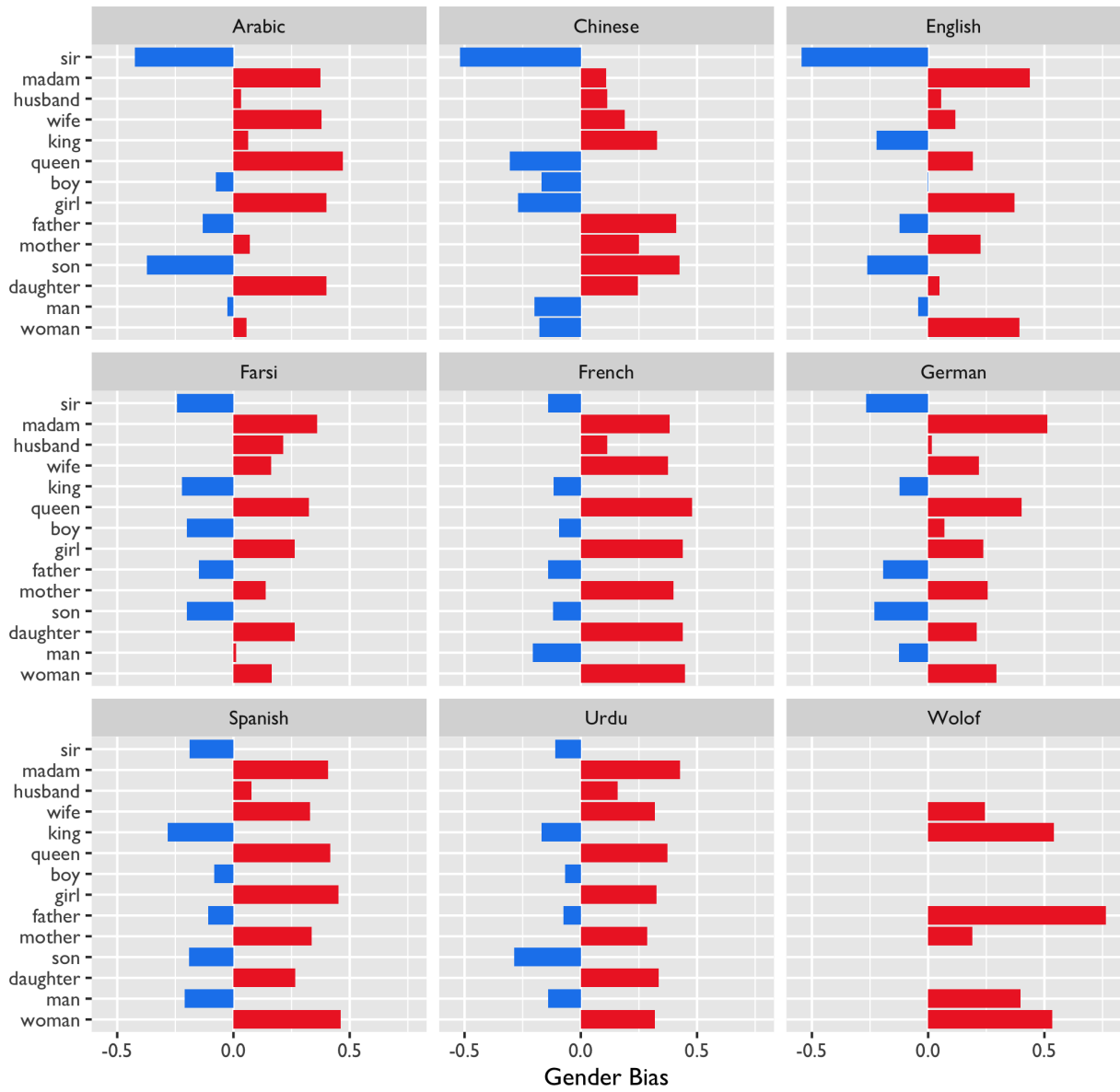
To our knowledge, this is the first paper to apply this methodology across languages and some important modifications and extensions were required, especially to handle languages, like Spanish, Arabic, German, French and Urdu, that have grammatically gendered nouns. In this section, we describe our modifications and apply them to computing and comparing both profession-level and corpus-level gender bias metrics across the Wikipedia corpora for 9 languages.

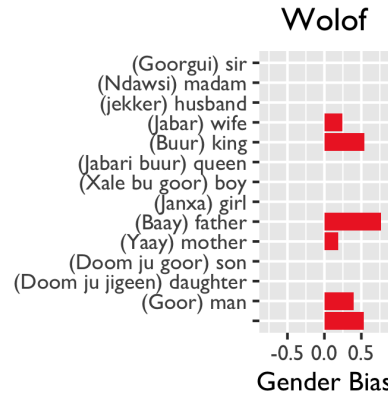
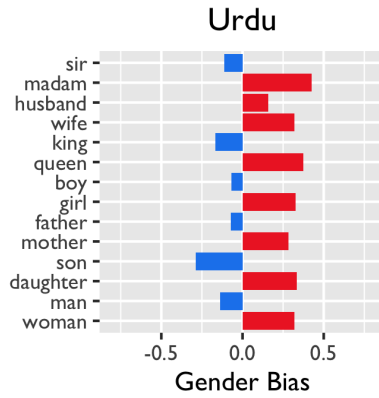
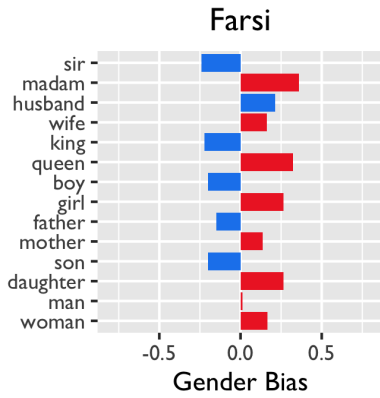
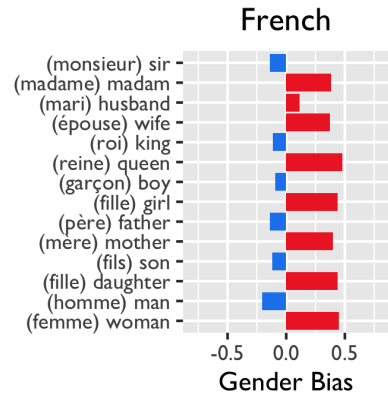
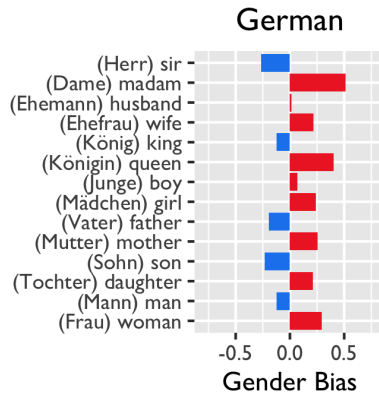
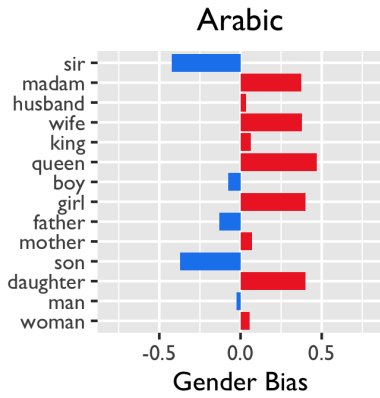
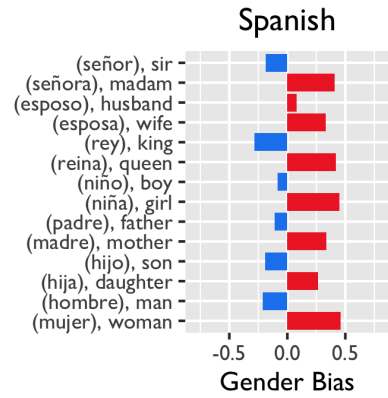
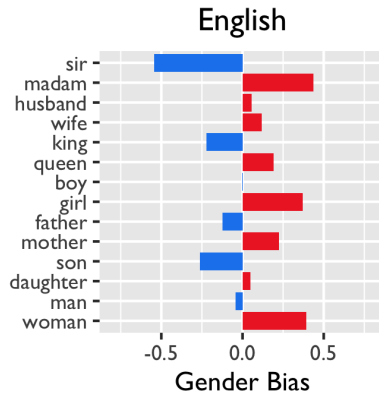
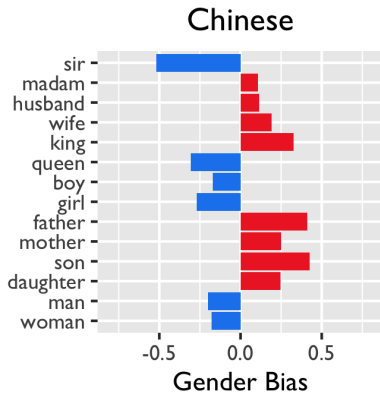
4.1 Comparing the Gender Bias of Defining Sets

To begin, in Figure 1, we present the gender bias scores, calculated as described above according to Bolukbasi et al.’s methodology, for each of our 14 defining set words (7 pairs) across 9 languages. Female gender bias is represented as a positive number (red bar) and male gender bias is represented as a negative number (blue bar). As we have discussed, not all defining set words occur in the Wikipedia corpus for Wolof. As another illustration of the hurdles encountered in various languages, for one variant of Figure 1, our process of graphing CSV files in R failed to display Chinese, Arabic, Farsi and Urdu words using the proper character sets. In those cases, the y-axis contains the English word.

The defining set pairs were specifically chosen because we expect them to be highly gendered. So not surprisingly, in most cases, the defining set words indicated male or female bias as expected, but there were some exceptions. More surprisingly, a common exception was the word husband. Husband has a female bias in every language except Wolof where it did not occur in the corpora. We hypothesize that “husband” may more often be used in relationship to women (e.g. “her husband”). One might guess that the same pattern would happen for wife then but it does not appear to be the case. We hypothesize that it may be less likely for a man to be defined as a husband outside of a female context, where women may often be defined by their role as a wife even when not in the context of the husband. This is an interesting effect we saw across many languages and it would be interesting to explore it further in future work. Husband is part of the set of 3 pairs (queen-king, wife-husband, and madam-sir) that were added in this study and not used in Bolukbasi et al. nor in Babaeianjelodar et al. In our ongoing work, we are repeating this analysis without the pair husband-wife in the defining set.

Defining Sets Across All Languages (Weighted By Word Count)





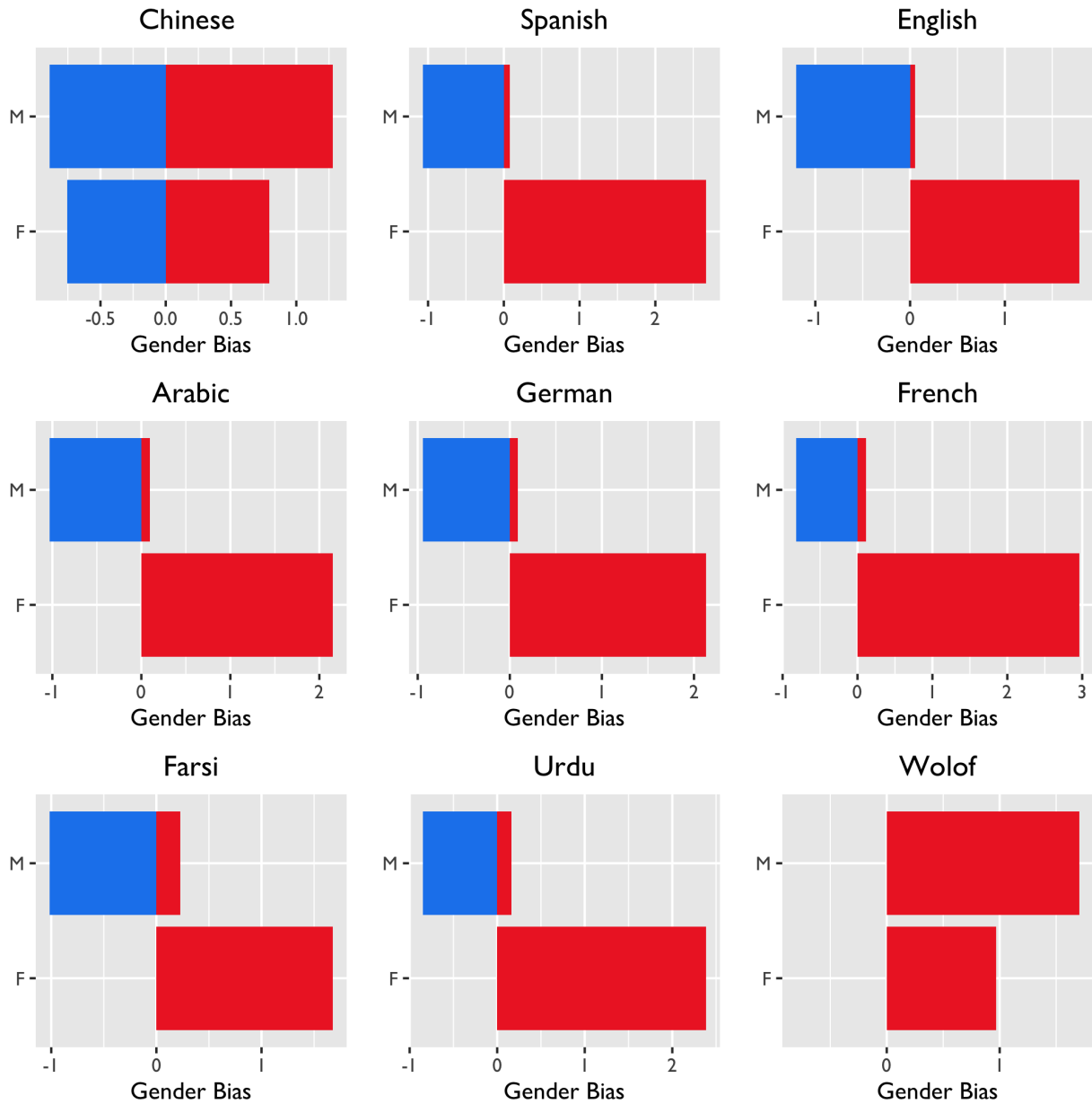


Figure 1: Defining Sets Across Languages The x-axis represents per-word gender bias scores as proposed by Bolukbasi et al. Female gender bias is represented as a positive number (red bar) and male gender bias is represented as a negative number (blue bar). Not all defining set words occur in the small Wikipedia corpus for Wolof. As another illustration of the hurdles encountered in various languages, our process of graphing CSV files in R failed to display Chinese, Arabic, Farsi and Urdu words using the proper character sets. In those cases, the y-axis contains the English word. We note that boy in English has a gender bias of -0.002 which is such a small blue line that it is difficult to see.

In Chinese, in particular, we saw more surprising results with words like father and son taking on a female bias. After much investigation, we isolated an issue related to the Principal Component Analysis (PCA) in Chinese. As we described at the beginning of this section, Bolukbasi et al.'s methodology calls for using the largest eigenvalue and in their experience the first eigenvalue was much larger than the second and they analyzed their results using only this dominant dimension. However, we found that this was not always the case. In particular for the Chinese Wikipedia corpus, the largest eigenvalue of the PCA matrix

is not much larger than the second. In Figure 2, we report the difference in PCA scores between the dominant component and the next most dominant component across 9 languages in our study. We also add a bar for this reported difference that Bolukbasi et al. reported for the Google News Corpora in English that they analyzed. Chinese has the lowest. Wolof has the highest with 1.0, but only because there were not enough defining pairs to meaningfully perform dimension reduction into 2 dimensions.

Thus, for the Chinese Wikipedia corpus, even though the defining set was chosen to be highly gendered, when PCA is used to reduce the number of dimensions, there is not a clearly dominant gender direction. We believe this is the key reason that the gender bias of the defining set words is not as intuitive for Chinese as it is for other languages. If the primary PCA dimension does not capture gender it may suggest the need to add more defining set word pairs and it would be interesting to repeat this analysis with an expanded defining set in Chinese.

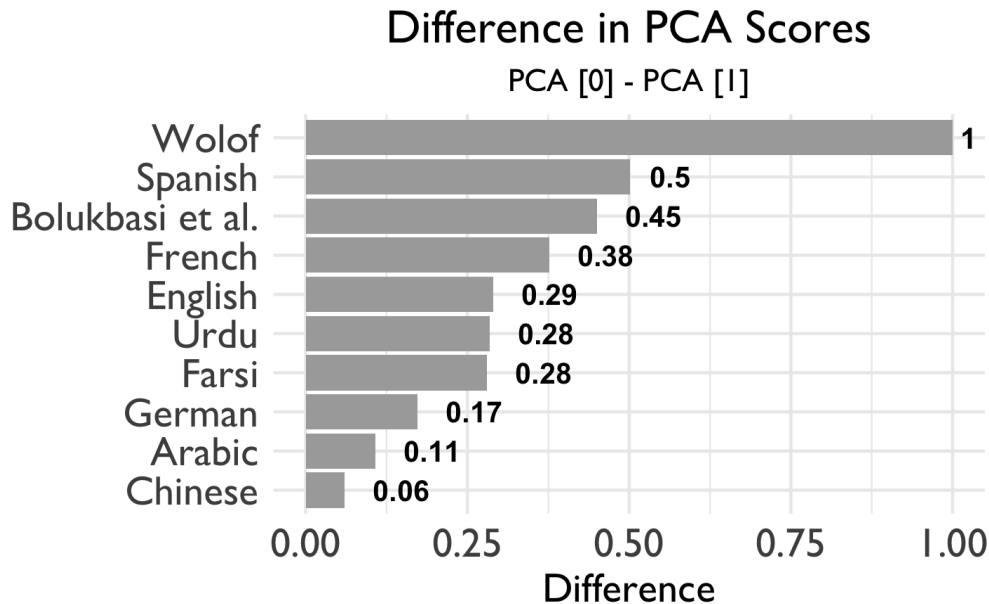


Figure 2: Difference in importance between first and second principal components by language. A larger difference increases the confidence we’ve isolated the gender direction. Note Wolof is 1.0 because there were not enough defining pairs to meaningfully perform dimension reduction into 2 dimensions.

The word boy in German, Junge, is also an exception and highlights some important issues. Junge can also be used as an adjective such as in “junge Leute” (young people) and it is also a common surname. Since these different uses of the word are not disambiguated, it is likely that the token “junge” encompasses more meanings than simply boy. We saw this with the defining set word “fille” in French which means both girl and daughter. Also, we have mentioned the issue of Word2Vec changing all words to lowercase and this also contributes to combining words in German that should be considered different parts of speech. Since all nouns in German are capitalized, maintaining capitalization would have provided some level of term separation. We suspect that this may contribute to Junge having a feminine gender bias. This problem of disambiguation is not unique to German and multiple meanings for words should be considered when selecting terms. For example, in English we included doctor as a profession, however had concerns of ambiguity with the verb to doctor. Such disambiguation of terms warrants further investigation in all languages.

4.2 Comparing the Gender Bias of Professions and Profession Sets

Having analyzed the defining set results where there is a clearly expected gender for each word, we move on to the question of computing the gender bias scores for each of our 32 profession words. Bolukbasi et al.’s methodology can be applied directly in English and also in other languages which, like English, do not have many grammatically gendered nouns. Of the 9 languages, we studied, Chinese, Farsi and Wolof are also in this category. Figure 3 shows the gender bias scores for the 32 profession words for the English Wikipedia corpus. Not surprisingly, we see many of the same patterns as documented by Bolukbasi et al. In this

figure, nurse is the profession with the largest absolute value of bias having a female gender bias of 0.32. Engineer is the profession with the largest male gender bias at -0.26.

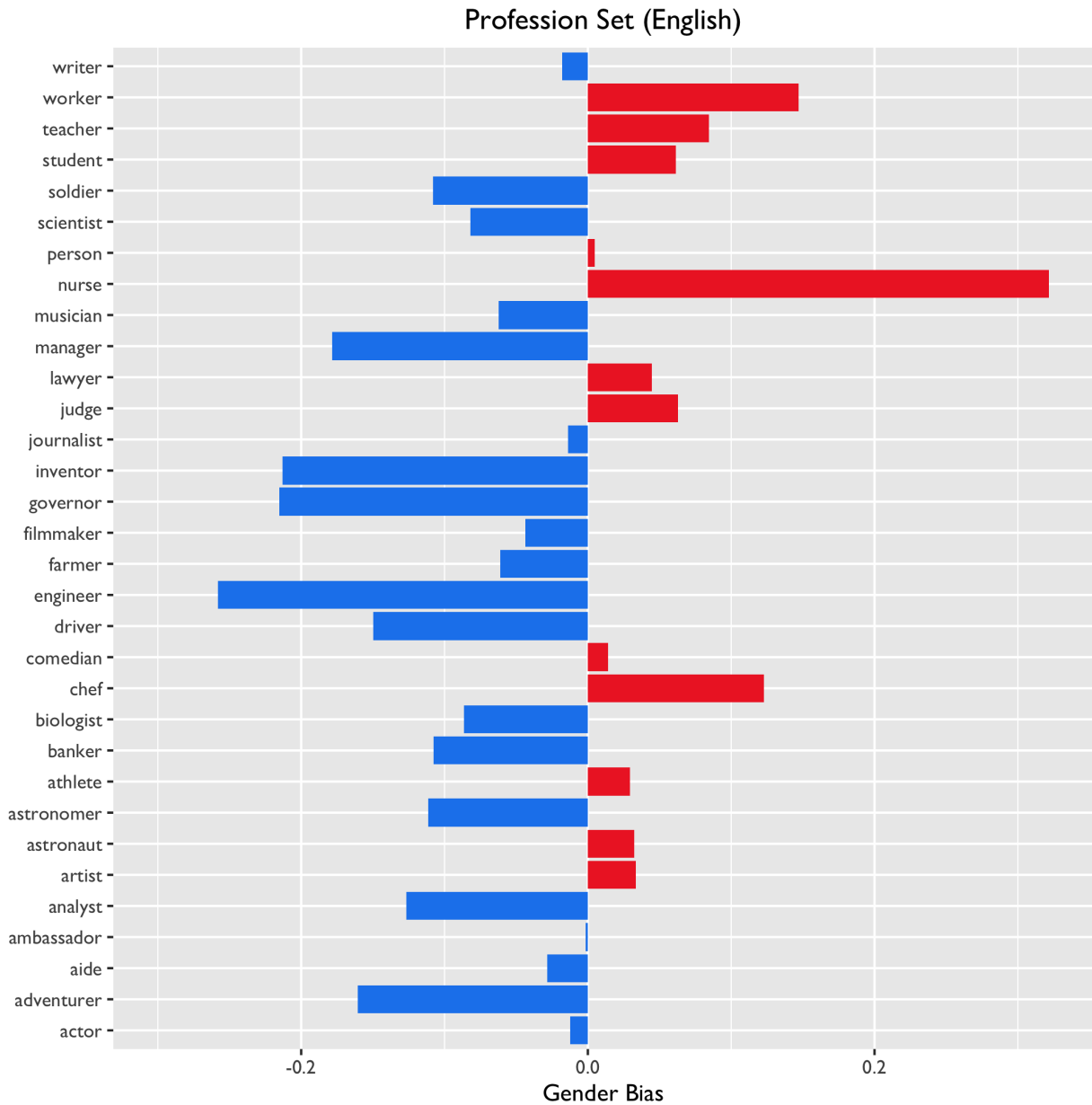


Figure 3: Per Profession Gender Bias for English Wikipedia.

The situation is more complicated in languages with grammatically gendered nouns. Five of the languages we are studying fall into this category: Spanish, Arabic, German, French, and Urdu. In these languages, many professions have both a feminine and masculine form. In some cases, there is also a neutral form and in some cases there is only a neutral form. In Section 2.2, we discussed how Urdu also often uses English words directly. Thus there are neutral Urdu words and neutral English words used in Urdu. To form a per-profession bias metric, we averaged the bias metrics of these various forms in several different ways. First, we averaged them, weighting each different form of a profession equally. However, we found that this overestimated the female bias in many cases. For example, in German the male form of scientist, Wissenschaftler, has a slight male gender bias (-0.06) and the female form, Wissenschaftlerin, has a strong female gender bias (0.32). When averaged together evenly, we would get an

overall female gender bias of 0.13. However, the male form occurs 32,467 times in the German Wikipedia corpus while the female form occurs only 1354 times. To take this difference into account, we also computed a weighted average resulting in an overall male gender bias of -0.04. With this weighted average, we could observe intuitive patterns across languages with grammatically gendered nouns and languages without. This increases our confidence in the usefulness of these profession level metrics and in particular the weighted average.

In Figure 4, we show the breakdown of the gender bias scores for the Spanish profession words. We show female only variants, male only variants and neutral only variants. In Appendix 3, we show a breakdown like this for all 5 of the gendered languages in our study. Notice in Figure 4, that the gender bias for all female words is indeed female and that the gender bias for all male words is indeed male. This is an intuitive and encouraging result that further supports the use of per-word gender bias calculations across languages. There are some exceptions in some languages as shown in detail in Appendix 3, but it is generally the case. Neutral words show a mix of male and female bias.

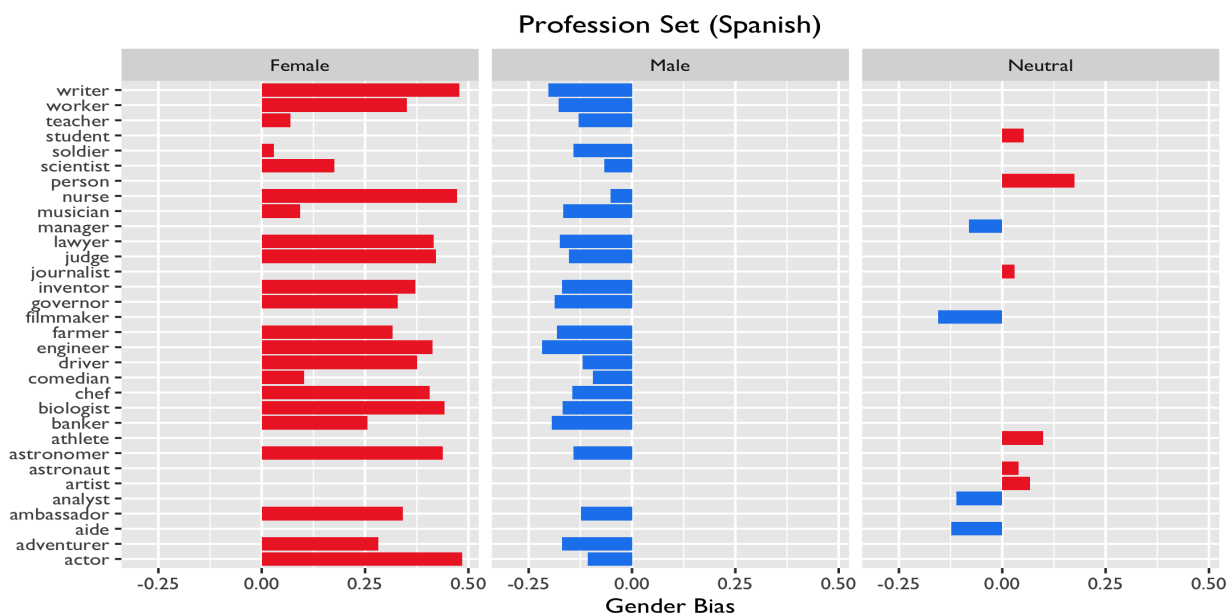


Figure 4: Per Profession Gender Bias for Spanish. Broken down into female only variants, male only variants, and neutral variants.

In Figures 5 and 6, we compare these profession-level gender bias scores across languages. In Figure 5, we show results for the languages without grammatically gendered nouns. In Figure 6, we show results across all languages using the weighted average (weighted by word count). We have removed the y-axis labels in Figure 6 for readability, but the order is the same as in the previous figures and here our emphasis is on observing the patterns across languages rather than on a drill down into specific words. Earlier in the paper, we note some problems with the Chinese results (lack of a non-dominant PCA dimension) and Wolof (a very small corpus). It is interesting to note how similar English and French are. Of these 4 languages, we believe that English and French are the most similar. Notice also the similarities in patterns between Spanish, English, Arabic, German, French, Farsi and Urdu. In Appendix 4, we include a similar graph using an evenly weighted average. In that alternate graph, the languages with grammatically gendered nouns are similar to each other, but not to English and Farsi. Based on these results, our recommendation is to use the weighted average as shown in Figure 6.

Languages Without Grammatically Gendered Nouns

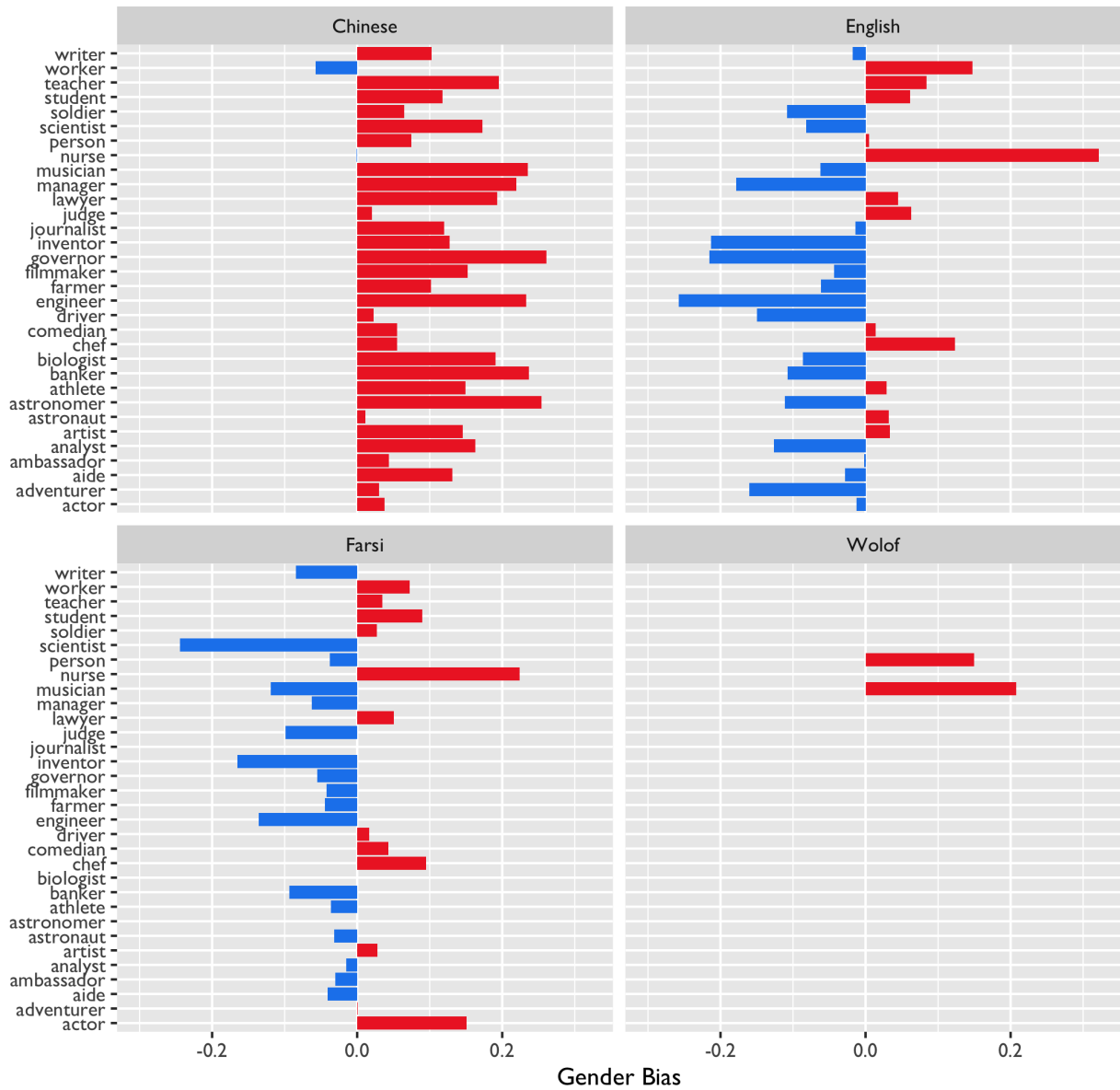


Figure 5 – Per-Profession Gender Bias Metrics for Languages Without Grammatically Gendered Nouns

Profession Sets Across All Languages (Weighted By Word Count)

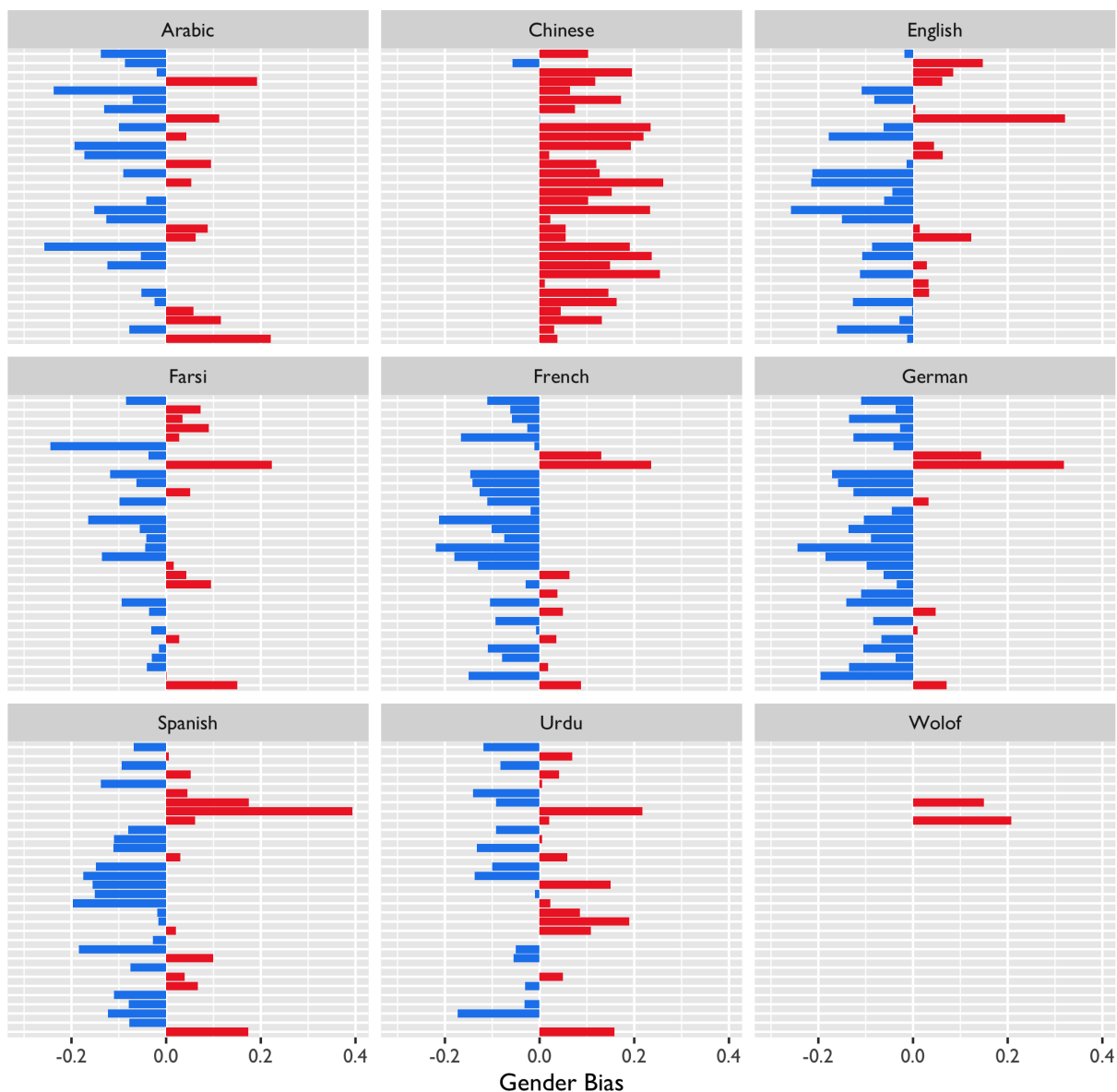


Figure 6 – Per-Profession Gender Bias Metrics for All Languages Weighting by Word Count

4.3 Comparing the Gender Bias of Corpora

Having analyzed the gender bias of our defining set and profession words, in this section, we discuss how to combine them into an overall gender bias score for each corpus. Bolukbasi et al. did not use their gender bias metric to compare different corpora. In Babaeianjelodar et al., we used an evenly weighted average across 327 profession words in English. We observed that the corpus-level gender bias score for a hate speech corpora like RtGender was substantially higher than more “representative” corpora like the contents of the GLUE benchmarks. We were also able to use their calculated gender bias metric to diagnose a surprisingly low level of gender bias in a hate speech corpora like IdentityToxic and when we investigated this surprising result, we were indeed able to confirm that IdentityToxic contained mostly hate speech targeting race and sexual orientation, rather than gender. Our earlier work comparing corpora in English demonstrated the ability to meaningfully compare the gender bias across corpora in English, but the study described here is the first to attempt to do such a comparison across languages.

Wikipedia offers an interesting basis for an initial cross-language comparison. As we have discussed the Wikipedia corpora in various languages vary substantially in size and quality, however they do have the same goal of offering an open-collaborative online encyclopedia. They have similar patterns of authorship (i.e. maintained by a community of volunteer editors using a wiki-based editing system). The comparison is certainly not exactly apples-to-apples, but it is an interesting and meaningful one, especially given how often Wikipedia is used in NLP research.

In both Bolukbasi et al. and Babaeianjelodar et al. results across the profession set words were averaged evenly. Here, inspired by our findings with the weighted average in the last section, we compute a weighted average over entire documents or corpora as well. Even in English, this produces a difference. Instead of summing the gender bias scores for each profession word and dividing by the number of profession words, the weighted average adds the gender bias for each profession word each time it occurs and then divides by the total number of times a profession word occurs.

In Figure 7, we show both the average and weighted average. They are close for all languages, but interesting, in Spanish the simple average is male biased and the weighted average is female biased. Also interestingly, the gender bias metric is negative or male for most of the languages. Besides the weighted average in Spanish, only Chinese and Wolof for which we have previously described some substantial concerns are female biased.

In future work, we would like to gain more experience using these metrics to compare more corpora for which we have more intuition of a ground truth- first within one language and then between languages. We expect there are important lessons to be learned from comparison within a language as we did Babaeianjelodar et al. with English and then perhaps between gendered languages separately from non-gendered languages. Even with this in mind, our methodology and the results presented here substantially improve on the state of the art in calculating and comparing gender bias across human languages. We would love to apply these metrics to the corpora we used in Babaeianjelodar et al. for which there is a ground truth understanding of their contents, but have not had an opportunity to do so. It is notable that those results were done with BERT rather than Word2Vec which although similar would introduce an additional level of variation in the results.

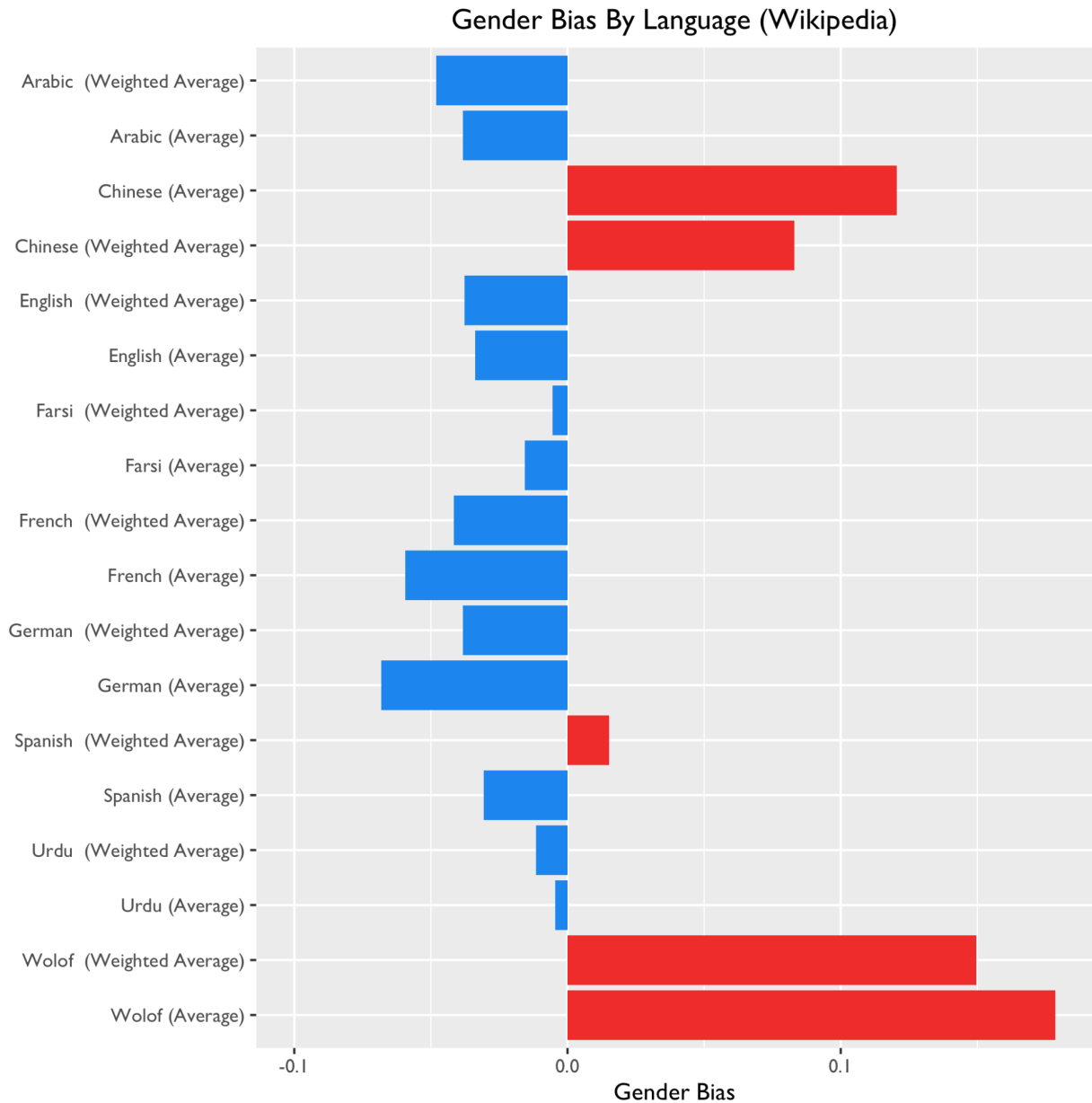


Figure 7: Per-Corpora Gender Bias Metrics: Overall gender bias across Languages for Wikipedia.

5. Beyond Wikipedia

Given how often Wikipedia is used in NLP research, we felt that documenting the difference in gender bias across Wikipedia corpora in different languages was a good first step. However, beyond Wikipedia, we would also like to experiment with additional corpora.

A corpus of documents that are widely translated into many languages would be one interesting type of apples-to-apples comparison. According to ITC Translations, the Bible, the Universal Declaration of Human Rights and the Adventures of Pinocchio are among the most widely translated documents. However, clearly those documents would not be equally reflective of different cultures and may very well not even contain language generated by native speakers (e.g. they may be produced with

translation software). We see similar problems with Wikipedia. Wikipedia prioritizes the voices of those creating content digitally and does not necessarily reflect the views of many speakers of those languages.

We would like to work with scholars in other disciplines such as sociology, linguistics and literature to compare gender bias across different corpora of culturally important texts written by native speakers. Even within one language, it would be interesting to examine collections with different emphasis such as gender of author, different time periods, different genres of text, different country of origin, etc. We could also work to produce versions of our tools that scholars in these disciplines could use to analyze different corpora for themselves.

As an initial experiment in this space, we assembled a small collection of works that might be considered “classics” in English, Chinese and Spanish to provide an additional datapoint beyond just Wikipedia. Some example texts include *Pride and Prejudice* by Jane Austin and *The Great Gatsby* by F. Scott Fitzgerald in English, *Cien años de soledad* by Gabriel García Márquez and *Ficciones* by Jorge Luis Borges in Spanish and *司马迁 (Records of the Grand Historian)* by Qian Sima and *萧红 (Tales of Hulan River)* by Hong Xiao in Chinese. Unlike with Wikipedia, these are pieces that have had a profound impact on the culture and were written by native speakers of the language. Where Wikipedia focuses on recently produced digital writing, many of these texts are older, as far back as 475 BC in the case of *论语 (The Analects)* in Chinese.

The Chinese corpus in particular spans a wide range of years (475 BC-1992) and this allowed us to observe some important nuances that we did not see with Wikipedia. Classical Chinese (before 1900) uses a vocabulary and grammar that differs significantly from modern Chinese. It may not be surprising that there are changes in professions over that length of time, but we even found changes when it comes to some of the most fundamental defining set words. For example, the word woman can be translated in many ways, including “女子”, “女人”, and “妇女”. “女子” was popularly used in ancient times, but its usage has decreased in modern writing. This shows us that as languages evolve over time, defining sets and profession sets may also have to evolve to measure gender bias.

6. CONCLUSION

In this paper, we extended an influential method for computing gender bias from Bolukbasi et al. It had only been applied in English and we made key modifications that allowed us to extend the methodology to 8 additional languages. Specifically, we modified and translated the original defining sets and profession sets. We extended the methodology to include languages with grammatically gendered nouns. With this, we quantified how gender bias varies across 9 languages within Wikipedia. We also assembled an initial classics corpora in 3 of the 9 languages and applied our methodology to it as well.

As speakers of 9 languages, we also used this process as an opportunity to shed light on the ways in which the modern NLP-pipeline does not reflect the voices of much of the world. For most languages, corpora are small and tool support is weak. Many published research methods, like Bolukbasi et al.’s gender bias metric calculations, are designed without consideration of the complexities of the multiple languages. This highlights the difficulties that speakers of many languages still face in having their thoughts and expressions fully included in the NLP-derived conclusions that are being used to direct the future. Despite substantial and admirable investments in multilingual support in projects like Wikipedia and Word2vec, we are still making NLP-guided decisions that systematically and dramatically under-represents many voices.

This work is an important step toward quantifying and comparing gender bias across languages - what we can measure, we can more easily begin to track and improve, but it is only a start. We focused on 9 languages from approximately 7,000 languages in the world. The majority of human languages need more useful tools and resources to overcome the barriers such that we can build NLP tools with less gender bias and such that NLP can deliver more value to every part of the world.

ACKNOWLEDGMENTS

We gratefully acknowledge the input and assistance of our wider research team including Hunter Bashaw, Marzieh Babaeianjelodar, William Smialek, Josh Gordon, Stephen Lorenz, Graham Northup and Cameron Weinfurt. We appreciate the insight of Clarkson faculty Golshan Madraki and Yu Liu. Support from The Clarkson Open Source Institute was essential to this work.

REFERENCES

- Babaeianjelodar, M.; Lorenz, S.; Gordon, J.; Matthews, J.; and Freitag, E. 2020. Quantifying gender bias in different corpora. In Companion Proceedings of the Web Conference 2020, WWW '20, page 752–759, New York, NY, USA, 2020. Association for Computing Machinery. DOI: <https://doi.org/10.1145/3366424.3383559>.
- Banerjee, D. 2020. Natural Language Processing (NLP) Simplified: A Step-by-step Guide. Datascience foundation. Retrieved from <https://datascience.foundation/sciencewhitepaper/natural-language-processing-nlp-simplified-a-step-by-step-guide>
- BERT. 2020. BERT Pretrained models. Github. Retrieved from <https://github.com/google-research/bert#bert>.
- Bolukbasi, T.; Chang, K.; Zou, J. Y.; Saligrama, V.; and Kalai, A.T. 2016. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. In Advances in neural information processing systems, pp. 4349-4357.
- Buolamwini, J; and GebruGender T. 2018. Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification. Proceedings of the 1st Conference on Fairness, Accountability and Transparency, PMLR 81:77-91, 2018.
- Bussieck, J. 2017. Demystifying Word2Vec. Retrieved from <https://www.deeplearningweekly.com/blog/demystifying-word2vec/>
- Chen, Y., Mahoney, C., Grasso, I., Wali, E., Matthews, A., Middleton, T., Njie, M., and Matthews, J. Gender Bias and Under-Representation in Natural Language Processing Across Human Languages Proceedings of the 2021 AAAI/ACM Conference on Artificial Intelligence, Ethics and Society (AIES), May 19-21 2021.
- Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805.
- Garbade, M. J. 2018. *A Simple Introduction To Natural Language Processing*. Retrieved from <https://becominghuman.ai/a-simple-introduction-to-natural-language-processing-ea66a1747b32>
- Holley, R. 2009. How good can it get? Analysing and improving OCR accuracy in large scale historic newspaper digitisation programs. D-Lib Magazine, March/April 2009, Volume 15 Number 3/4 ISSN 1082-9873. Retrieved from <http://www.dlib.org/dlib/march09/holley/03holley.html>
- Karani, D. 2018. Introduction to Word Embedding and Word2Vec. Retrieved from <https://towardsdatascience.com/introduction-to-word-embedding-and-word2vec-652d0c2060fa>
- Karch, M. 2020. How to Use the Ngram Viewer Tool in Google Books. Retrieved from <https://www.lifewire.com/google-books-ngram-viewer-1616701>
- Matthews, A., Grasso, I., Mahoney, C., Chen, Y., Wali, E., Middleton, T., Njie, M., and Matthews, J. Gender Bias in Natural Language Processing Across Human Languages. TrustNLP: First Workshop on Trustworthy Natural Language Processing Colocated with the Annual Conference of the North American Chapter of the Association for Computational Linguistics, June 10 2021.
- Mikolov, T.; Chen, K.; Corrado, G.; and Dean, J. 2013. Efficient estimation of word representations in vector space. *arXiv:1301.3781*. Retrieved from <https://arxiv.org/abs/1301.3781>
- Mithe, R.; Indalkar, S.; and Divekar, N. 2013. Optical character recognition. International journal of recent technology and engineering. ISSN: 2277-3878, Volume-2, Issue-1, March 2013.
- NLTK (2005). Natural Language Toolkit. Retrieved from <http://www.nltk.org/>
- Nosek, B. A.; Banaji, M. R.; and Greenwald, A. G. 2002. Harvesting implicit group attitudes and beliefs from a demonstration web site. *Group Dynamics: Theory, Research, and Practice*, 6(1):101, 2002.

- Ohio University. Wolof Language. Retrieved from <https://www.ohio.edu/cis/african/languages/wolof>
- Rong, X. 2016. word2vec Parameter Learning Explained. arXiv:1411.2738. Retrieved from <https://arxiv.org/abs/1411.2738>
- Siegel, R. 2019. Search result not found: China bans Wikipedia in all languages. Retrieved from <https://www.washingtonpost.com/business/2019/05/15/china-bans-wikipedia-all-languages/>
- Su, Q, G. 2019. Which Parts of the World Speaks Mandarin Chinese?. Retrieved from [https://www.thoughtco.com/where-is-mandarin-spoken-2278443#:~:text=Mandarin%20Spoken%20Here&text=There%20are%20an%20estimated%2040.countries%20\(about%2030%20million\)](https://www.thoughtco.com/where-is-mandarin-spoken-2278443#:~:text=Mandarin%20Spoken%20Here&text=There%20are%20an%20estimated%2040.countries%20(about%2030%20million))
- Wali, E., Chen, Y., Mahoney, C., Middleton, T., Babaeian-jelodar, M., Njie, M., and Matthews, J. Is Machine Learning Speaking my Language? A Critical Look at the NLP-Pipeline Across 8 Human Languages Participatory ML Workshop, Thirty-seventh International Conference on Machine Learning (ICML2020), July 17 2020.
- WikipediaA. 2020. Wikipedia: *German language*. Retrieved from https://en.wikipedia.org/wiki/German_language
- WikipediaB. 2020. Wikipedia: *List of languages by total number of speakers*. Retrieved from https://en.wikipedia.org/wiki/List_of_languages_by_total_number_of_speakers
- WikipediaC. 2020. Wikipedia: *List of Wikipedias*. Retrieved from https://en.wikipedia.org/wiki/List_of_Wikipedias
- WikipediaD. 2020. Wikipedia: *Wolof Wikipedia*. Retrieved from https://en.wikipedia.org/wiki/Wolof_Wikipedia
- Williams, A., Nangia, N., Bowma, S. (2020). MultiNLI, The Multi-Genre NLI Corpus. Retrieved from <https://cims.nyu.edu/~sbowman/multinli>
- Yordanov, V. 2018. *Introduction To Natural Language Processing For Text*. Medium. Retrieved from <https://towardsdatascience.com/introduction-to-natural-language-processing-for-text-df845750fb63>
- Yamada, I.; Asai, A.; Sakuma, J.; Shindo, H.; Takeda, H.; Takefuji, Y.; and Matsumoto Y. 2018. Wikipedia2Vec: An Efficient Toolkit for Learning and Visualizing the Embeddings of Words and Entities from Wikipedia. arXiv:1812.06280. Retrieved from <https://arxiv.org/abs/1812.06280>

1. APPENDIX 1: Word Count Data for Defining Sets in All Languages (Wikipedia)

These tables contain the number of times each of our 14 defining set words (7 pairs) occurs in the Wikipedia Corpora for 9 languages as downloaded from Wikipedia on 2020-06-20.

Table 1.1: Defining set word counts in Chinese

Defining word	Word counts
woman(女人)	14,276
man(男人)	12,948
daughter(女儿)	14,014
son(儿子)	21,605
mother(母亲)	17,775
father(父亲)	29,291
girl(女孩)	15,286
boy(男孩)	8,105
queen(女王)	15,637
king(国王)	28,110
wife(妻子)	34,673
husband(丈夫)	17,443
madam(女士)	6,865
sir(男士)	1,583

Table 1.2: Defining set word counts in Spanish

Defining word	Word counts
woman(mujer)	136,418
man(hombre)	128,558
daughter(hija)	122,514
son(hijo)	228,001
mother(madre)	137,180
father(padre)	213,878
girl(niña)	25,307
boy(niño)	60,857
queen(reina)	75,833

king(rey)	257,774
wife(esposa)	90,521
husband(esposo)	22,445
madam(señora)	56,254
sir(señor)	65,051

Table 1.3: Defining set word counts in English

Defining word	Word counts
woman	391,942
man	853,753
daughter	488,907
son	869,322
mother	532,677
father	819,271
girl	259,817
boy	222,415
queen	346,530
king	918,788
wife	564,337
husband	272,398
madam	5,946
sir	369,141

Table 1.4: Defining set word counts in Arabic

Defining word	Word counts
woman(النساء)	38,964
man(رجل)	29,707
daughter(ابنة)	11,379
son(ولد)	263,029
mother(ام)	13,534
father(اب)	10,584
girl(ابنة)	11,379
boy(صبي)	1,297
queen(ملكة)	10,595
king(ملك)	30,198
wife(زوجة)	8,844
husband(الزوج)	2,998
madam(سيدتي)	345

sir(سيدي)	25,456
-----------	--------

Table 1.5: Defining set word counts in German

Defining word	Word counts
woman(Frau)	243,752
man(Mann)	169,581
daughter(Tochter)	211,261
son(Sohn)	369,234
mother(Mutter)	148,965
father(Vater)	214,064
girl(Mädchen)	56,891
boy(Junge)	81,498
queen(Königin)	52,631
king(König)	279,997
wife(Ehefrau)	65,893
husband(Ehemann)	24,979
madam(Dame)	32,372
sir(Herr)	49,193

Table 1.6: Defining set word counts in French

Defining word	Word counts
woman(femme)	258,208
man(homme)	412,795
daughter(fille)	239,968
son(fils)	408,486
mother(mère)	198,220
father(père)	349,479
girl(fille)	239,968
boy(garçon)	27,253
queen(reine)	84,742
king(roi)	373,417
wife(épouse)	162,337
husband(mari)	58,311
madam(madame)	38,384
sir(monsieur)	40,306

Table 1.7: Defining set word counts in Farsi

Defining word	Word counts
woman(زن)	149,008

man(مرد)	167,114
daughter(دختر)	21,599
son(پسر)	24,438
mother(مادر)	16,497
father(پدر)	23,139
girl(دختر)	21,599
boy(پسر)	24,438
queen(ملکه)	7,087
king(پادشاه)	17,678
wife(همسر)	13,229
husband(شوهر)	2,520
madam(خانم)	8,446
sir(آقا)	6,668

Table 1.8: Defining set word counts in Urdu

Defining word	Word counts
woman(عورت)	2,427
man(آدمی)	2,139
daughter(بیٹی)	2,263
son(بیٹا)	2,433
mother(مان)	2,775
father(باپ)	2,773
girl(لڑکی)	1,265
boy(لڑکا)	509
queen(ملکہ)	1,678
king(پادشاه)	5,825
wife(بیوی)	2,474
husband(شوہر)	1,376
madam(محترمہ)	308
sir(جناب)	1,420

Table 1.9: Defining set word counts in Wolof

Defining word	Word counts
woman(Jigéen)	53
man(Góor)	36
daughter(Doom ju jigéen)	0
son(Doom ju góor)	0
mother(Yaay)	0

father(Baay)	0
girl(Janxa)	0
boy(Xale bu góor)	0
queen(Jabari buur)	0
king(Buur)	0
wife(Jabar)	13
husband(jëkkër)	0
madam(Ndawsì)	0
sir(Góorgui)	0

2. APPENDIX 2: Word Count Data for Profession Sets in All Languages (Wikipedia)

These tables contain the number of times each of our 32 profession set words occurs in the Wikipedia Corpora for 9 languages as downloaded from Wikipedia on 2020-06-20. We break down the total count of words in a number of ways. N represents a neutral word variant, M represents a male word variant and F represents a female variant. E stands for phonetically-written English word used in another language. These are also neutral, but they are not a native part of the language as designated with an N. For example, some Urdu words are spelled based on English pronunciation, some Wolof words use English words directly. * stands for multi-word phrases and they are not found with Word2Vec.

N= Neutral Variant

M= Male Variant

F = Female Variant

E = English word used in another language

*= multi-word phrases

Table 2.1: Profession words and word counts in Chinese

Profession	Word counts	Profession	Word counts
nurse(护士)	N: 1,187	soldier(战士)	N: 5,234
teacher(教师)	N: 10,494	journalist(记者)	N: 11,304
writer(作家)	N: 45,936	student(学生)	N: 38,716
engineer(工程师)	N: 7,827	athlete(运动员)	N: 52,954
scientist(科学家)	N: 8,267	actor(演员)	N: 17,395
manager(经理)	N: 3,379	governor(总督)	N: 8,627
driver(司机)	N: 4,802	farmer(农民)	N: 8,212
banker(银行家)	N: 788	person(人)	N: 1,021,664
musician(音乐家)	N: 2,796	lawyer(律师)	N: 7,554
artist(艺术家)	N: 7,514	adventurer(冒险家)	N: 195
chef(厨师)	N: 702	aide(助手)	N: 6,306
filmmaker(制片人)	N: 3,274	ambassador(大使)	N: 26,723
judge(法官)	N: 12,397	analyst(分析员)	N: 78
comedian(喜剧演员)	N: 650	astronaut(宇航员)	N: 3,052

inventor(发明家)	N: 667	astronomer(天文学家)	N: 3,545
worker(工人)	N: 19,004	biologist(生物学家)	N: 1,108

Table 2.2: Profession words and word counts in Spanish

Profession	Word counts	Profession	Word counts
nurse(enfermera/enfermero)	F: 5,043 M: 900	soldier(soldada/soldado)	F: 303 M: 18,859
teacher(profesora/profesor)	F: 20,683 M: 97,104	journalist(periodista)	N: 53,537
writer(escritora/escritor)	F: 20,822 M: 85,364	student(estudiante)	N: 25,939
engineer(ingeniera/ingeniero)	F: 1,366 M: 39,743	athlete(atleta)	N: 13,679
scientist(científica/científico)	F: 30,364 M: 35,632	actor(actriz/actor)	F: 96,016 M: 107,015
manager(gerente)	N: 11,667	governor(gobernadora/gobernador)	F: 2,107 M: 88,821
driver(conductora/conductor)	F: 3,826 M: 14,895	farmer(granjera/granjero)	F: 123 M: 1,812
banker(banquera/banquero)	F: 65 M: 3,127	person(persona)	N: 90,789
musician(música/músico)	F: 273,196 M: 38,549	lawyer(abogada/abogado)	F: 4,879 M: 39,118
artist(artista)	N: 75,012	adventurer(aventurera/aventurero)	F: 569 M: 2,225
chef(cocinera/cocinero)	F: 1,066 M: 2,481	aide(ayudante)	N: 13,296
filmmaker(cineasta)	N: 8,904	ambassador(ayudante/ayudante)	F: 2,633 M: 24,332
judge(jueza/juez)	F: 2,087 M: 26,833	analyst(analista)	N: 3,149
comedian(cómica/cómico)	F: 5,169 M: 7,881	astronaut(astronauta)	N: 2,590
inventor(inventora/inventor)	F: 286 M: 6,851	astronomer(astrónoma/astrónomo)	F: 1,167 M: 9,134
worker(trabajadora/trabajador)	F: 3,918 M: 7,389	biologist(bióloga/biólogo)	F: 848 M: 2,849

Table 2.3: Profession words and word counts in English

Profession	Word counts	Profession	Word counts
nurse	N: 46,823	soldier	N: 106,010
teacher	N: 204,870	journalist	N: 170,129
writer	N: 410,919	student	N: 355,461
engineer	N: 188,339	athlete	N: 69,242
scientist	N: 83,568	actor	N: 386,534
manager	N: 341,512	governor	N: 413,356

driver	N: 150,006	farmer	N: 61,820
banker	N: 25,603	person	N: 365,444
musician	N: 160,519	lawyer	N: 127,866
artist	N: 438,972	adventurer	N: 7,454
chef	N: 36,865	aide	N: 18,832
filmmaker	N: 35,809	ambassador	N: 107,442
judge	N: 238,220	analyst	N: 26,908
comedian	N: 51,806	astronaut	N: 15,615
inventor	N: 37,636	astronomer	N: 30,370
worker	N: 66,113	biologist	N: 15,689

Table 2.4: Profession words and word counts in Arabic

Profession (M/ F) or (N)	Word counts	Profession	Word counts
nurse(ممرضة/ممرض)	F: 1,087 M: 340	soldier(مجندة/جندي)	F: 45 M: 8,311
teacher(مدرسة/مدرس)	F: 43,592 M: 2,771	journalist(صحفية/صحافي)	F: 4,730 M: 514
writer(كاتبة/كاتب)	F: 5,875 M: 27,020	student(طالبة/طالب علم)	F: 1,741 *M: 0
engineer(مهندسة/مهندس)	F: 728 M: 9,472	athlete(رياضية/رياضي)	F: 39,706 M: 5,957
scientist(عالمة/عالم)	F: 3,072 M: 39,450	actor(ممثلة/الممثل)	F: 25,418 M: 7,118
manager(مديرة/مدير)	F: 1,953 M: 17,986	governor(حاكمة/حاكم)	F: 1,213 M: 10,658
driver(سائقة/سائق)	F: 61 M: 2,788	farmer(مزارعة/مزارع)	F: 35 M: 2,930
banker(مصرفي)	N: 474	person(شخص)	N: 61,332
musician(موسيقار)	N: 325	lawyer(محامية/محامي)	F: 948 M: 8,856
artist(فنانة/فنان)	F: 2,205 M: 5,585	adventurer(مغامر)	N: 179
chef(طاه)	N: 66	aide(مساعدة/مساعد)	F: 10,917 M: 9,266
filmmaker(صانع أفلام)	*N: 0	ambassador(سفيرة/سفير)	F: 618 M: 3,949
judge(قاضية/قاضي)	F: 651 M: 6,491	analyst(محللة/المحلل)	F: 130 M: 619
comedian(المضحك)	N: 126	astronaut(رائدة فضاء/رائد فضاء)	*N: 0
inventor(مخترعة/مخترع)	F: 71 M: 1,228	astronomer(عالم الفلك)	*N: 0
worker(عاملة/عامل)	F: 1,229 M: 11,197	biologist(أحيائي)	N: 290

Table 2.5: Profession words and word counts in German

Profession	Word counts	Profession	Word counts
------------	-------------	------------	-------------

nurse(Krankenschwester/Krankenpfleger)	F: 4,667 M: 1,112	soldier(Soldatin/Soldat)	F: 407 M: 21,717
teacher(Lehrerin/Lehrer)	F: 12,498 M: 82,150	journalist(Journalistin/Journalist)	F: 18,801 M: 63,234
writer(Autorin/Autor)	F: 32,392 M: 122,808	student(Studentin/Student)	F: 2,743 M: 13,122
engineer(Ingenieurin/Ingenieur)	F: 804 M: 35,107	athlete(Athletin/Athlet)	F: 1,792 M: 3,383
scientist(Wissenschaftlerin/Wissenschaftler)	F: 1,354 M: 32,467	actor(Schauspielerin/Schauspieler)	F: 118,359 M: 160,639
manager(Managerin/Manager)	F: 2,181 M: 31,172	governor(Gouverneurin/Gouverneur)	F: 821 M: 44,984
driver(Fahrerin/Fahrer)	F: 519 M: 35,433	farmer(Landwirtin/Landwirt)	F: 225 M: 10,325
banker(Bankerin/Banker)	F: 62 M: 1,212	person(Person)	N: 157,970
musician(Musikerin/Musiker)	F: 6,134 M: 83,175	lawyer(Anwältin/Anwalt)	F: 1,799 M: 19,853
artist(Künstlerin/Künstler)	F: 19,452 M: 117,587	adventurer(Abenteurer)	N: 3,194
chef(Köchin/Koch)	F: 1,613 M: 34,361	aide(Helferin/Helfer)	F: 462 M: 7,233
filmmaker(Filmemacherin/Filmemacher)	F: 1,824 M: 8,124	ambassador(Botschafterin/Botschafter)	F: 2,895 M: 47,101
judge(Richterin/Richter)	F: 3,704 M: 69,650	analyst(Analytikerin/Analytiker)	F: 50 M: 635
comedian(Komikerin/Komiker)	F: 986 M: 6,899	astronaut(Astronautin/Astronaut)	F: 514 M: 3,006
inventor(Erfinderin/Erfinder)	F: 544 M: 19,512	astronomer(Astronomin/Astronom)	F: 1,312 M: 11,797
worker(Arbeiterin/Arbeiter)	F: 676 M: 36,085	biologist(Biologin/Biologe)	F: 1,206 M: 5,275

Table 2.6: Profession words and word counts in French

Profession	Word counts	Profession	Word counts
nurse(infirmière/infirmier)	F: 8,763 M: 2,702	soldier(soldate/soldat)	F: 108 M: 27,846
teacher(professeure/professeur/prof)	F: 6,918 M: 177,007 N: 5,645	journalist/journaliste)	N: 91,261
writer(écrivaine/écrivain/écrivain)	F: 11,232 M: 159,846	student(étudiante/étudiant)	F: 10,492 M: 103,387
engineer(ingénieure/ingénieur)	F: 755 M: 60,176	athlete(athlète)	N: 46,526
scientist(scientifique)	N: 152,849	actor(actrice/acteur)	F: 119,910 M: 148,262
manager(directrice/directeur)	F: 20,510 M: 187,723	governor(gouverneure/gouverneur)	F: 760 M: 84,407
driver(chauffeuse/chauffeur)	F: 72	farmer(agricultrice/agriculteur)	F: 132

	M: 8,227		M: 3,825
banker(banquière/banquier)	F: 267 M: 9,536	person(personne)	N: 173,805
musician(musicienne/musicien)	F: 4,002 M: 38,976	lawyer(avocate/avocat)	F: 4,576 M: 55,779
artist(artiste)	N: 124,066	adventurer(aventurière/aventurier)	F: 802 M: 4,272
chef(cuisinière/cuisinier/chef)	F: 1,683 M: 4,908 N:314,644	aide(aide)	N: 188,941
filmmaker(réalisatrice/réalisateur/cinéaste)	F: 10,388 M: 84,907 N:12,263	ambassador(ambassadrice/ambassadeur)	F: 2,427 M: 30,600
judge(juge)	N: 59,559	analyst(analyste)	N: 2,945
comedian(comédienne/comédien)	F: 10,111 M: 14,965	astronaut(astronaute)	N: 5,051
inventor(inventrice/inventeur)	F: 224 M: 14,832	astronomer(astronome)	N: 23,352
worker(ouvrière/ouvrier)	F: 12,749 M: 20,204	biologist(biologiste)	N: 5,842

Table 2.7: Profession words and word counts in Farsi

Profession	Word counts	Profession	Word counts
nurse(پرستار)	N: 1,118	soldier(سرباز)	N: 7,127
teacher(معلم)	N: 4,771	journalist(روزنامه نگار)	*N: 0
writer(نویسنده)	N: 38,012	student(دانشجو)	N: 4,020
engineer(مهندس)	N: 5,577	athlete(ورزشکار)	N: 5,895
scientist(دانشمند)	N: 5,788	actor(بازیگر)	N: 38,749
manager(مدیر)	N: 29,687	governor(فرماندار)	N: 3,375
driver(راننده)	N: 3,722	farmer(کشاورز)	N: 1,783
banker(بانکدار)	N: 527	person(شخص)	N: 16,400
musician(موسیقیدان)	N: 840	lawyer(وکیل)	N: 4,821
artist(هنرمند)	N: 11,193	adventurer(ماجراجو)	N: 137
chef(آشپز)	N: 361	aide(مشاور)	N: 5,755
filmmaker(فیلمساز)	N: 928	ambassador(سفیر)	N: 4,665
judge(قاضی)	N: 5	analyst(تحلیلگر)	N: 461
comedian(کمدین)	N: 6,828	astronaut(فضانورد)	N: 1,596
inventor(مخترع)	N: 1,545	astronomer(ستاره شناس)	*N: 0
worker(کارگر)	N: 4,176	biologist(زیست شناس)	*N: 0

Table 2.8: Profession words and word counts in Urdu

Profession (E/F/M/N)	Word counts	Profession (E/F/M/N)	Word counts
nurse E: نرس	E: 64	soldier N: فوجی	N:4,626
teacher E: ٹیچر	E: 140 F: 33	journalist	N:1,828

F: استاتی M: استاد	M: 2,453	N: صحافی	
writer F: مصنفہ M: مصنف	F: 0 M: 4,045	student E: اسٹوڈنٹ F: طالبة M: طالب علم	E: 47 F: 0 *M: 0
engineer E: انجینئر N: معمار	E: 384 N: 317	athlete N: کھلاڑی	N: 15,518
scientist N: سائنسدان	N: 290	actor F: اداکارہ M: اداکار	F: 4,740 M: 10,280
manager E: مینجر F: منتظمہ M: منتظم	E: 553 F: 0 M: 1,128	governor E: گورنر N: حاکم	E: 2,963 N: 1,290
driver E: ڈرائیور F: گاڑی چلانے والی M: گاڑی چلانے والا	E: 266 *F: 0 *M: 0	farmer N: کسان	N: 330
banker E: بینکار	E: 71	person N: شخص	N: 6,418
musician N: موسیقار	N: 2,090	lawyer N: وکیل	N: 843
artist E: آرٹسٹ F: فنکارہ M: فنکار	E: 207 F: 0 M: 976	adventurer N: مہم جو	*N: 0
chef E: شیف F: باورچن M: باورچی	E: 24 F: 0 M: 118	aide N: مددگار	N: 282
filmmaker N: فلمساز	N: 234	ambassador N: سفیر	N: 710
judge E: جج N: قاضی	E: 573 N: 2,381	analyst N: تجزیہ کار	*N: 0
comedian E: کامیڈین N: مزاح نگار	E: 41 *N: 0	astronaut N: خلا باز	N: 34
inventor N: موجد	N: 269	astronomer N: ماہر فلکیات	*N: 0
worker N: کارکن	N: 835	biologist N: ماہر حیاتیات	*N: 0

Table 2.9: Profession words and word counts in Wolof

Profession	Word counts	Profession	Word counts
nurse(nurse)	E:0	soldier(soldar)	N:0
teacher(Jangalehkat)	N:0	journalist(journalist)	E:0

writer(bindakat)	N:0	student(Jangakat)	N:0
engineer(engineer)	E:0	athlete(dawkat)	N:0
scientist(scientist)	E:0	actor(tiyatarkat)	N:0
manager(manager)	E:0	governor(Governor)	E:0
driver(dawalkati moto)	*N:0	farmer(baikat)	N:0
banker(ligeykati bank)	*N:0	person(nit)	N:1,401
musician(waykat)	N:5	lawyer(lawyer)	E:0
artist(artist)	E:0	adventurer(adventurer)	E:0
chef(tougakat)	N:0	aide(dimbalehkat)	N:0
filmmaker(dafarkati film)	N:0	ambassador(ambassador)	E:0
judge(judge)	E:0	analyst(analyst)	E:0
comedian(comedian)	E:0	astronaut(astronaut)	E:0
inventor(inventor)	E:0	astronomer(astronomer)	E:0
worker(ligaikat)	N:0	biologist(biologist)	E:0

3. APPENDIX 3: Gender Bias Calculations for Profession Sets in All Languages (Wikipedia)

These figures contain the gender bias calculated for each of our 32 profession set words occur in the Wikipedia Corpora for 9 languages as downloaded from Wikipedia on 2020-06-20.

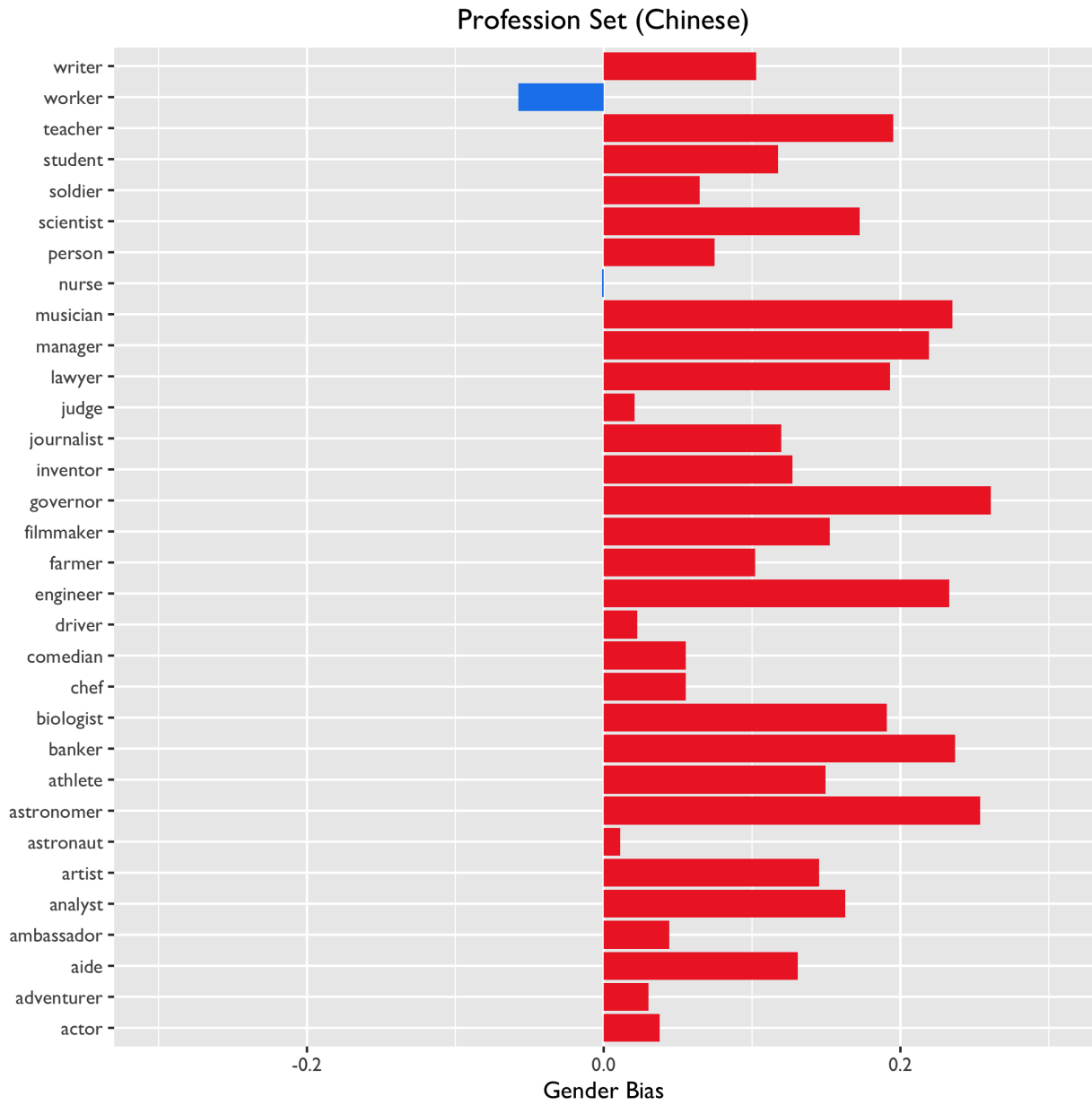


Figure 3.1: Per Profession Gender Bias for Chinese

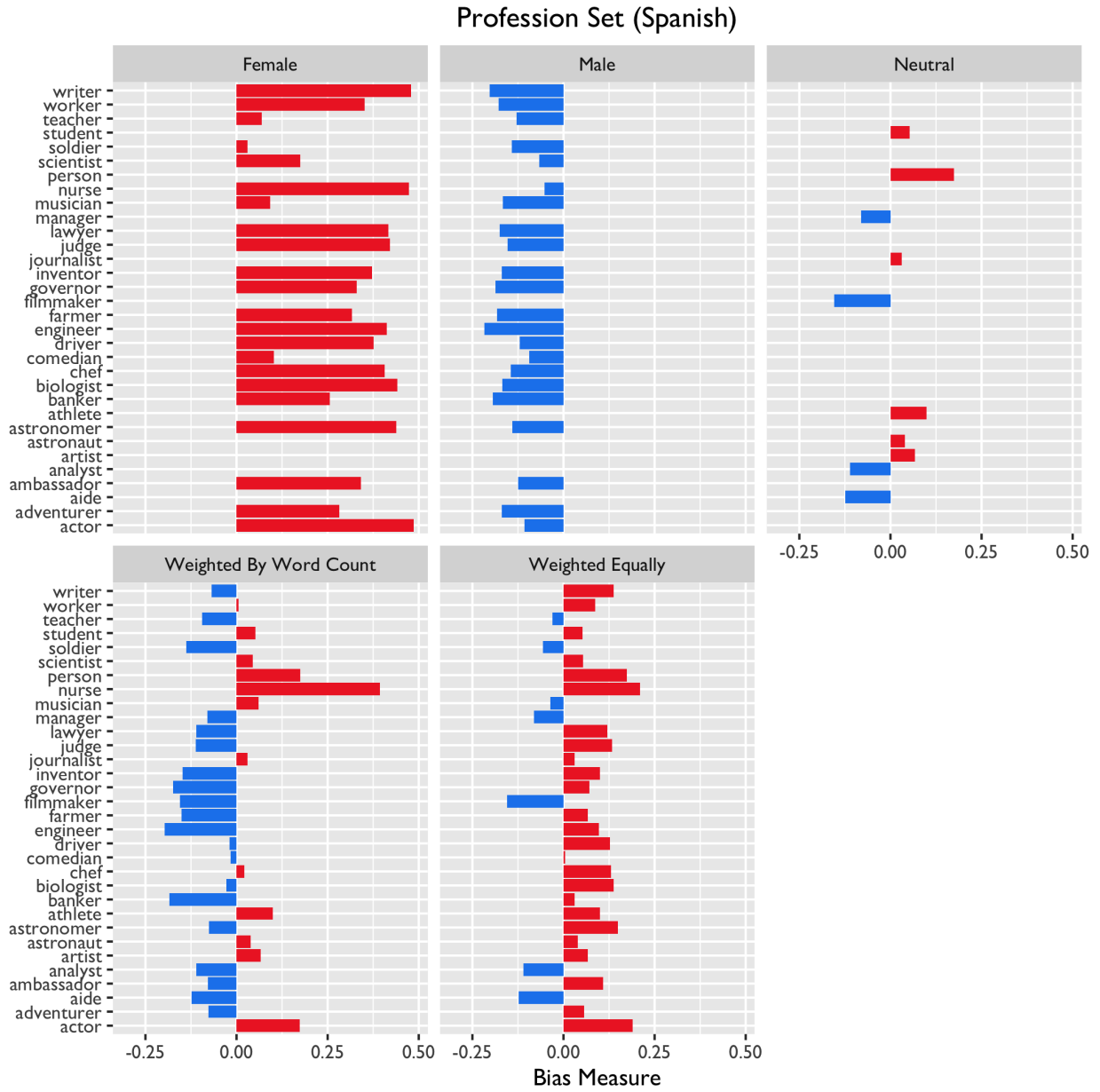


Figure 3.2: Per Profession Gender Bias for Spanish

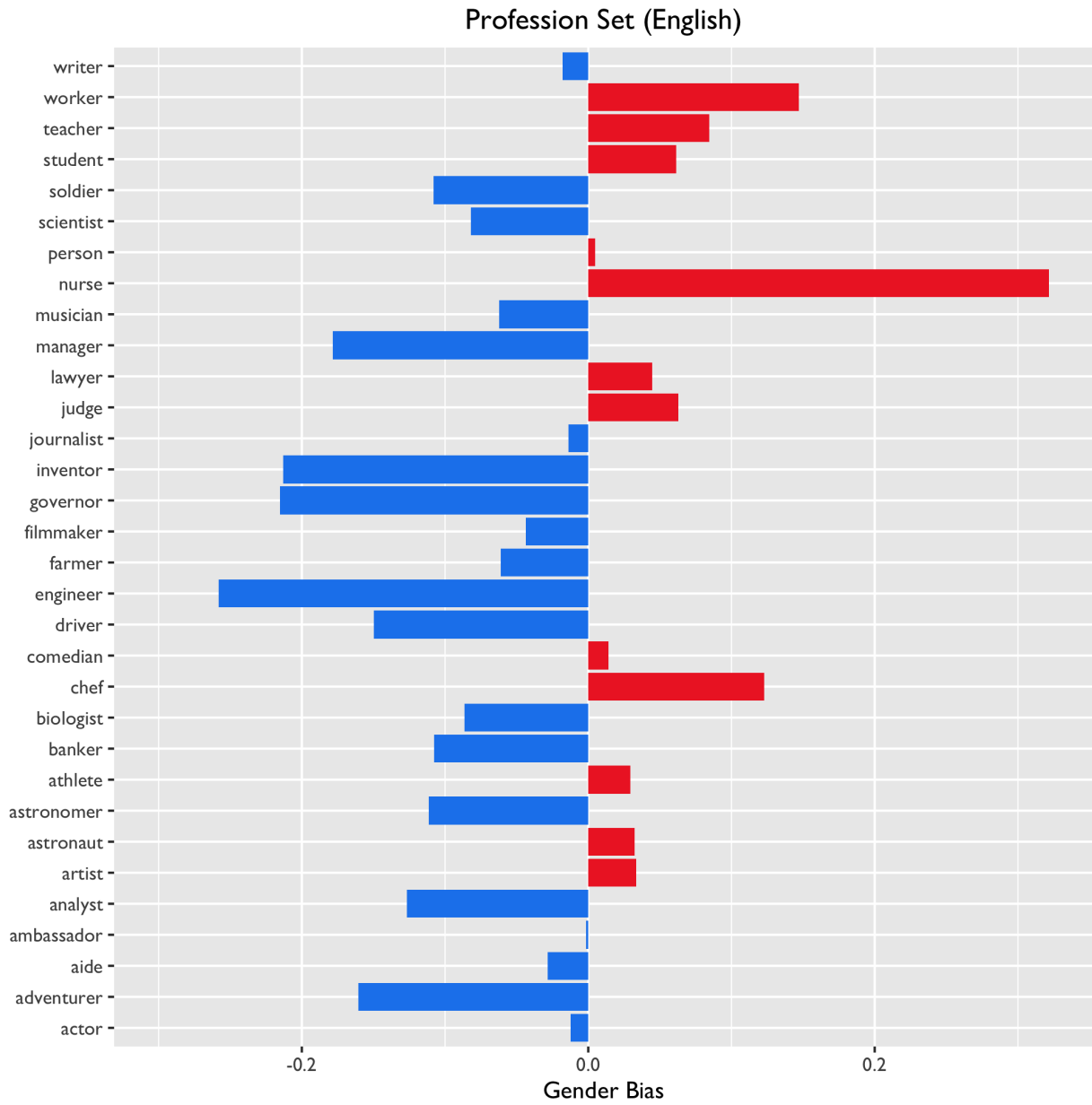


Figure 3.3: Per Profession Gender Bias for English

Profession Set (Arabic)

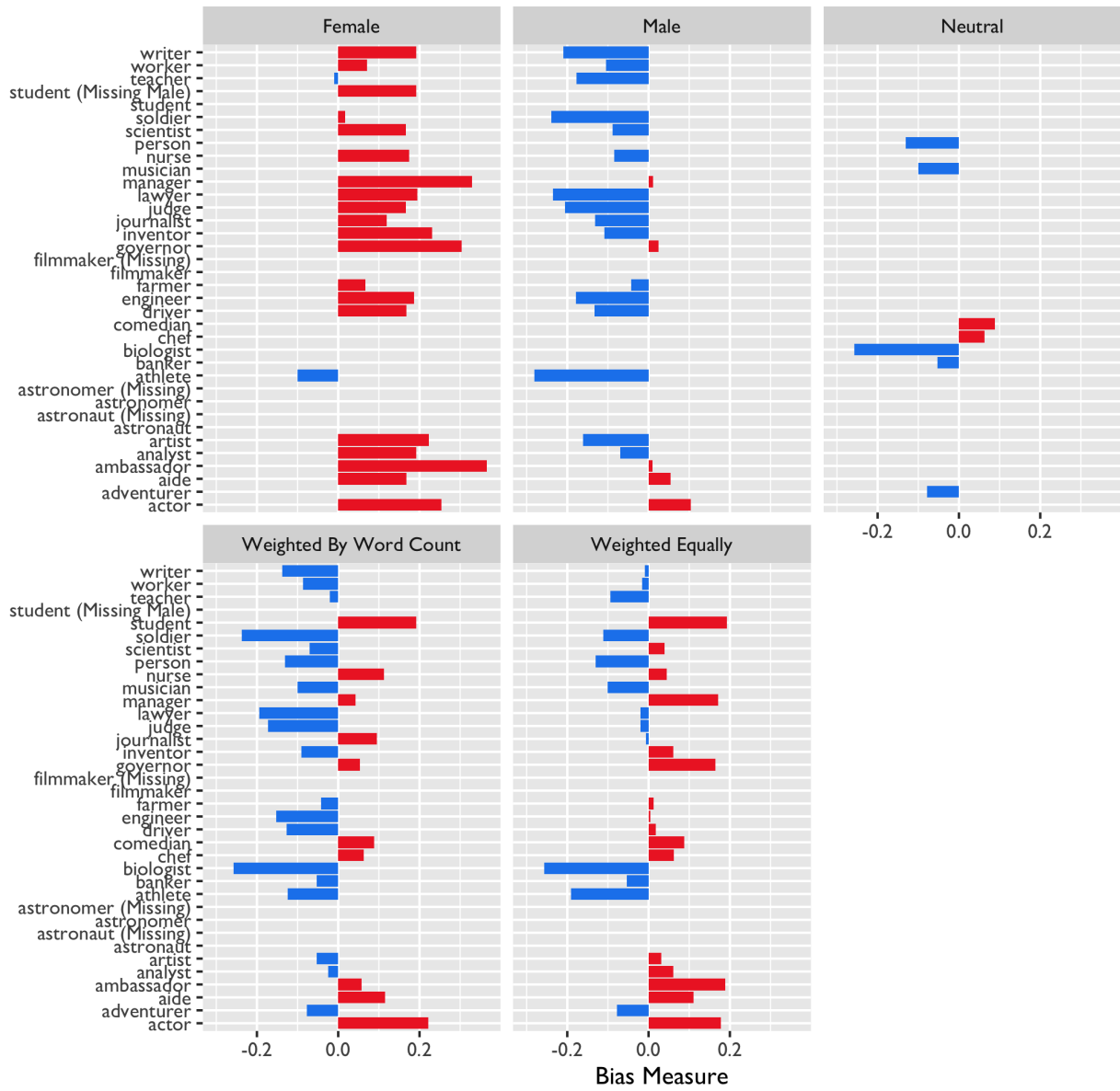


Figure 3.4: Per Profession Gender Bias for Arabic

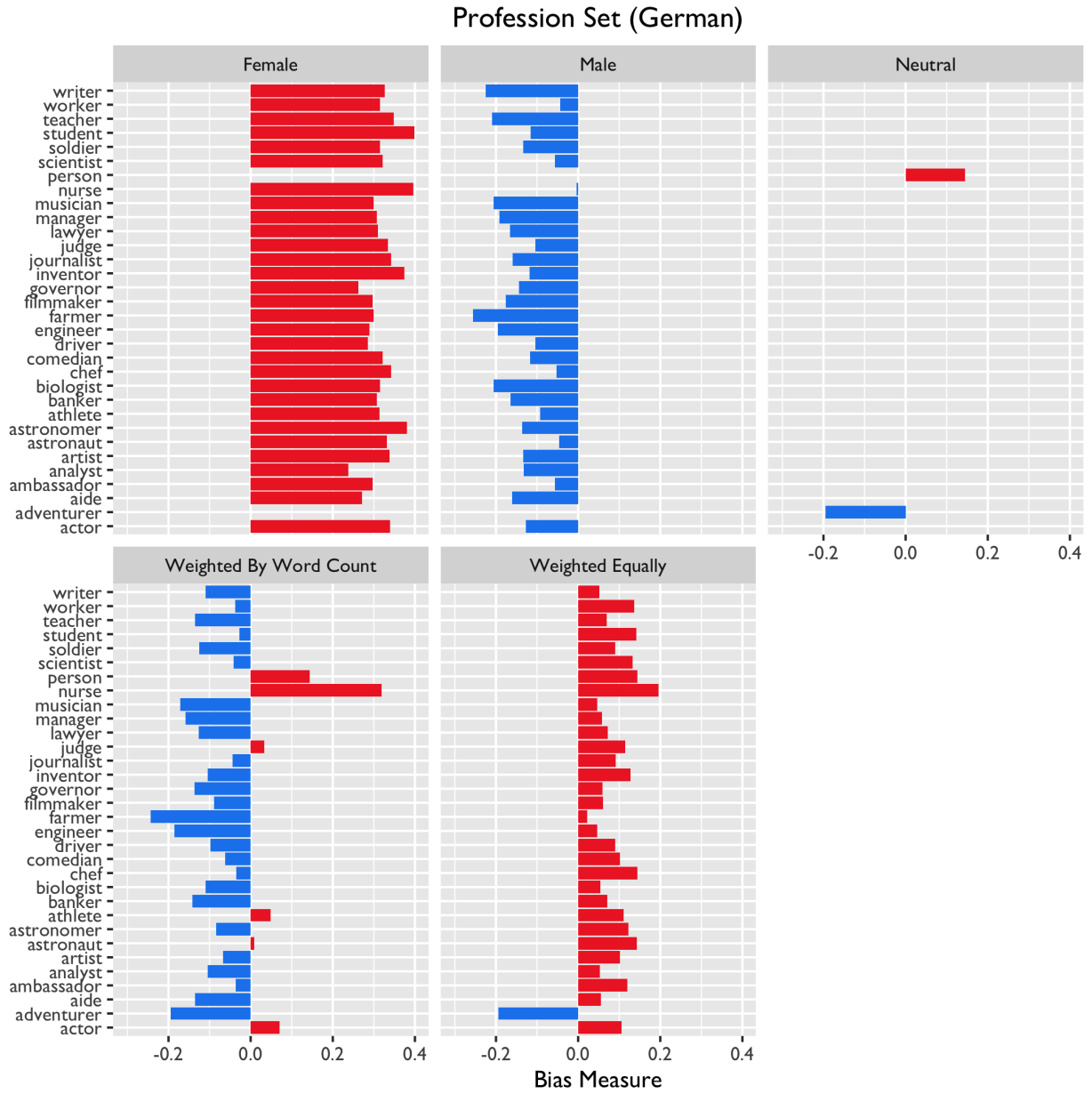


Figure 3.5: Per Profession Gender Bias for German

Profession Set (French)

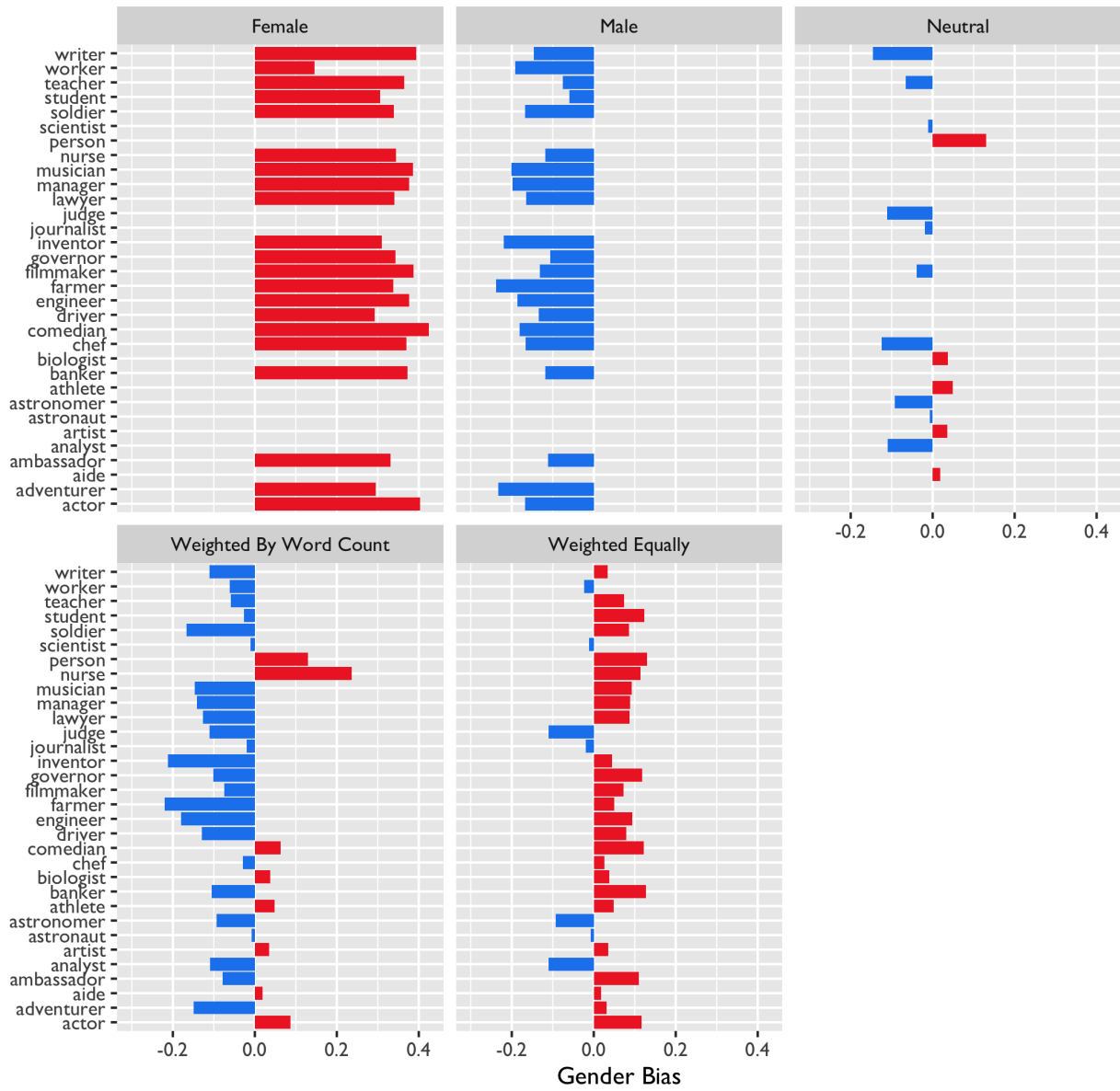


Figure 3.6: Per Profession Gender Bias for French

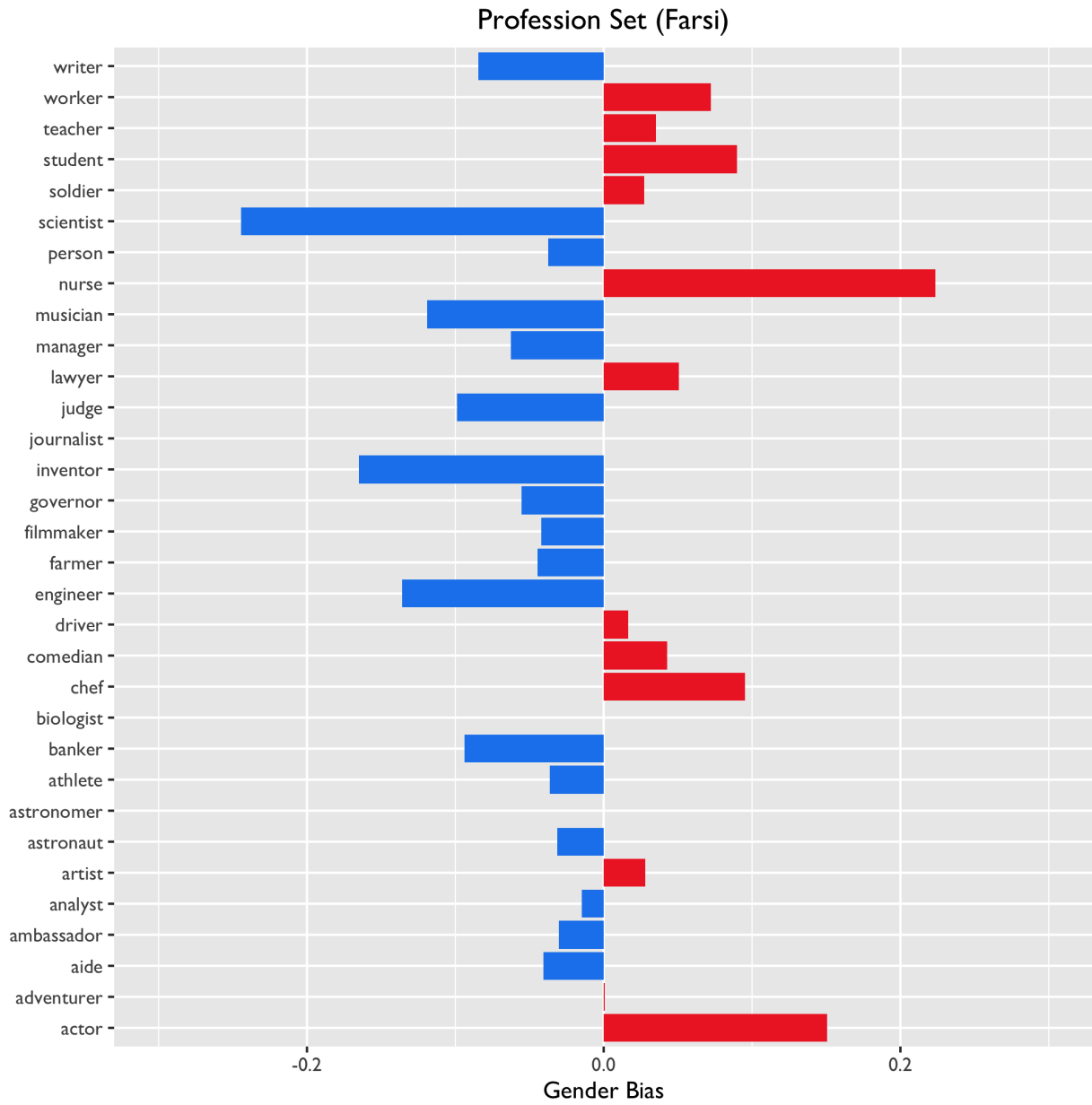


Figure 3.7: Per Profession Gender Bias for Farsi

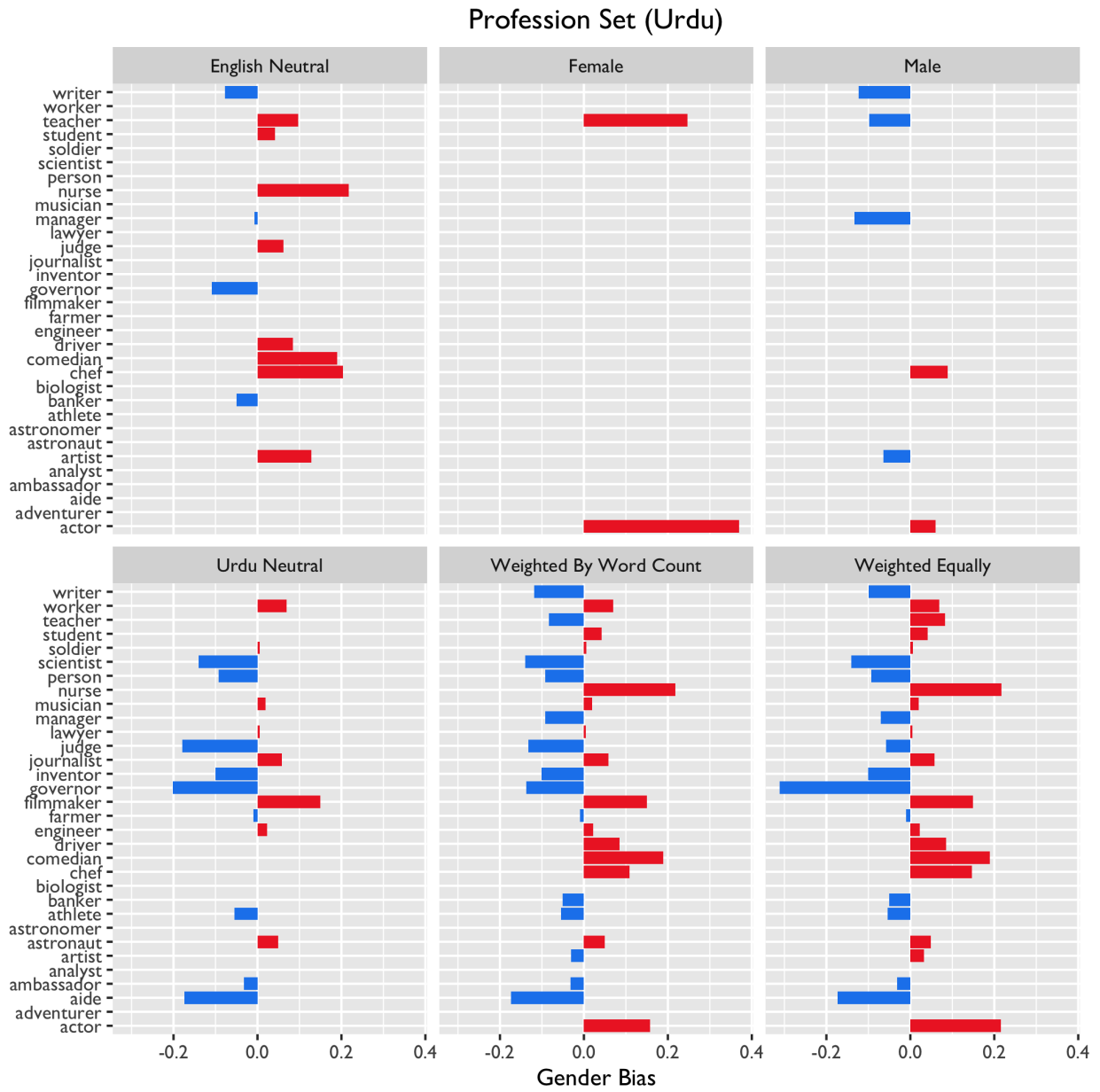


Figure 3.8: Per Profession Gender Bias for Urdu

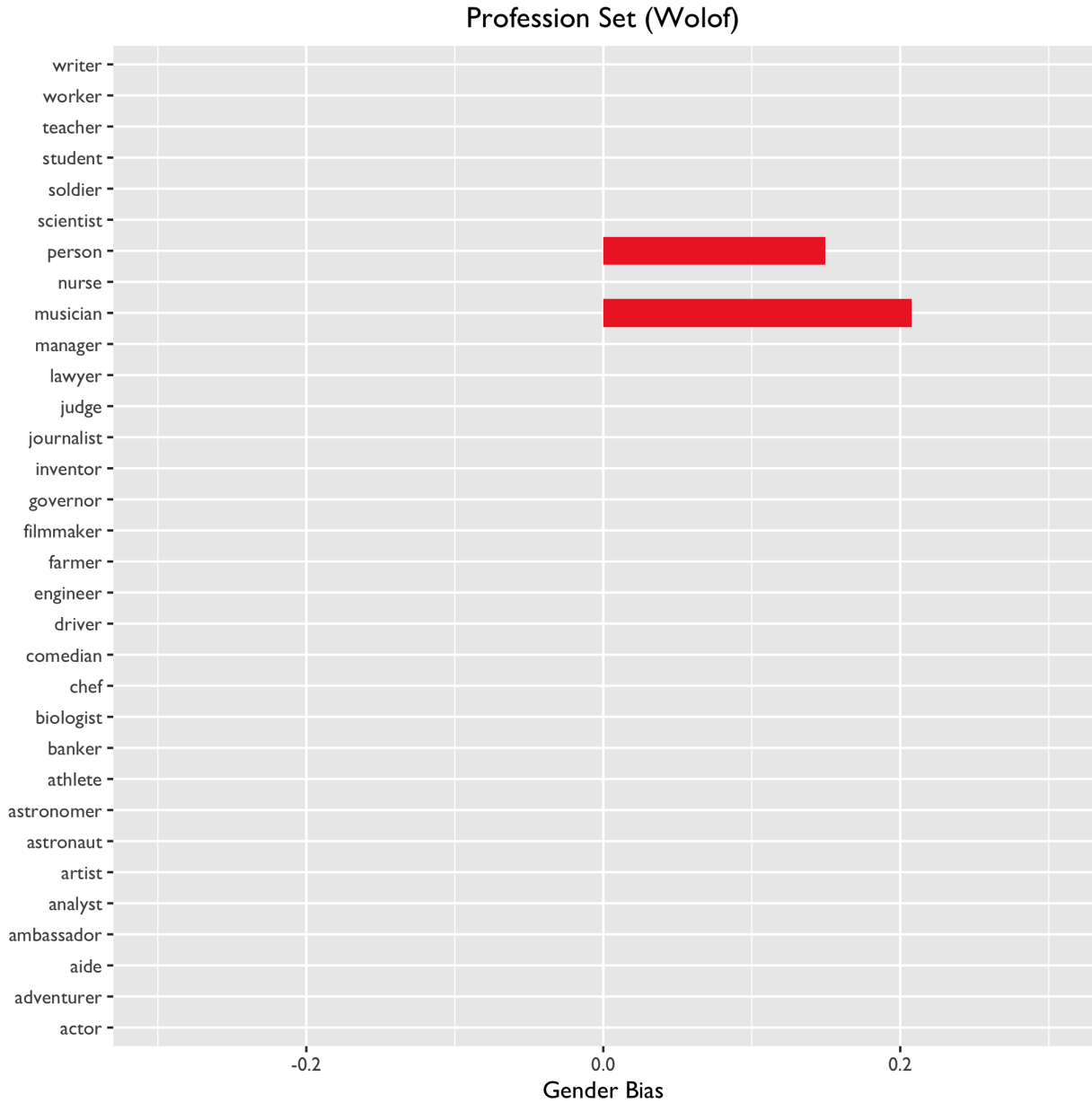


Figure 3.9: Per Profession Gender Bias for Wolof

4. APPENDIX 4: COMPARISON OF THE WEIGHTED AVERAGE TO THE SIMPLE AVERAGE ACROSS LANGUAGES

In Figures 4.1 - 4.2, we compare these profession-level gender bias scores across languages. In Figure 4.1, we show results for the languages without grammatically gendered nouns. It is interesting to note how similar English and French are. Of these 4 languages, we believe that English and French are. Earlier in the paper, we note some problems with the Chinese results (lack of a non-dominant PCA dimension) and Wolof (a very small corpora). In Figure 4.2, we show results across all languages using the weighted average (weighted by word count). We have removed the y-axis labels for readability, but the order is the same as in the previous figures and our emphasis is on observing the patterns across languages rather than on a drill down into specific words. Notice the similarities in patterns between Spanish, English, Arabic, German, French, Farsi and Urdu. Figure 4.3 contains a

similar comparison, but using an evenly weighted average. In that graph, the languages with grammatically gendered nouns are similar to each other, but not to English and Farsi. Based on these results, our recommendation is to use the weighted average.

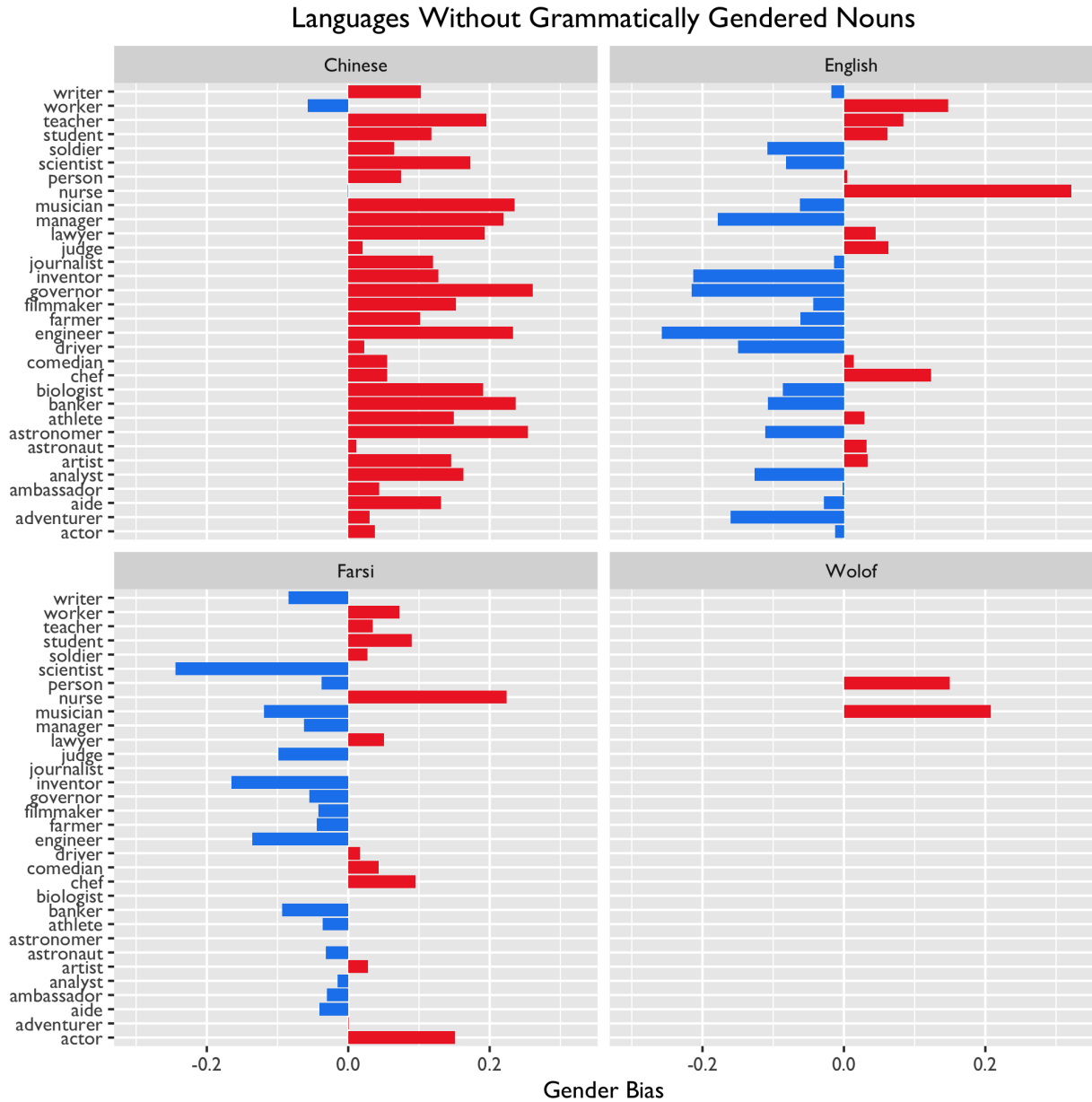


Figure 4.1 – Per-Profession Gender Bias Metrics for Languages Without Grammatically Gendered Nouns

Profession Sets Across All Languages (Weighted By Word Count)

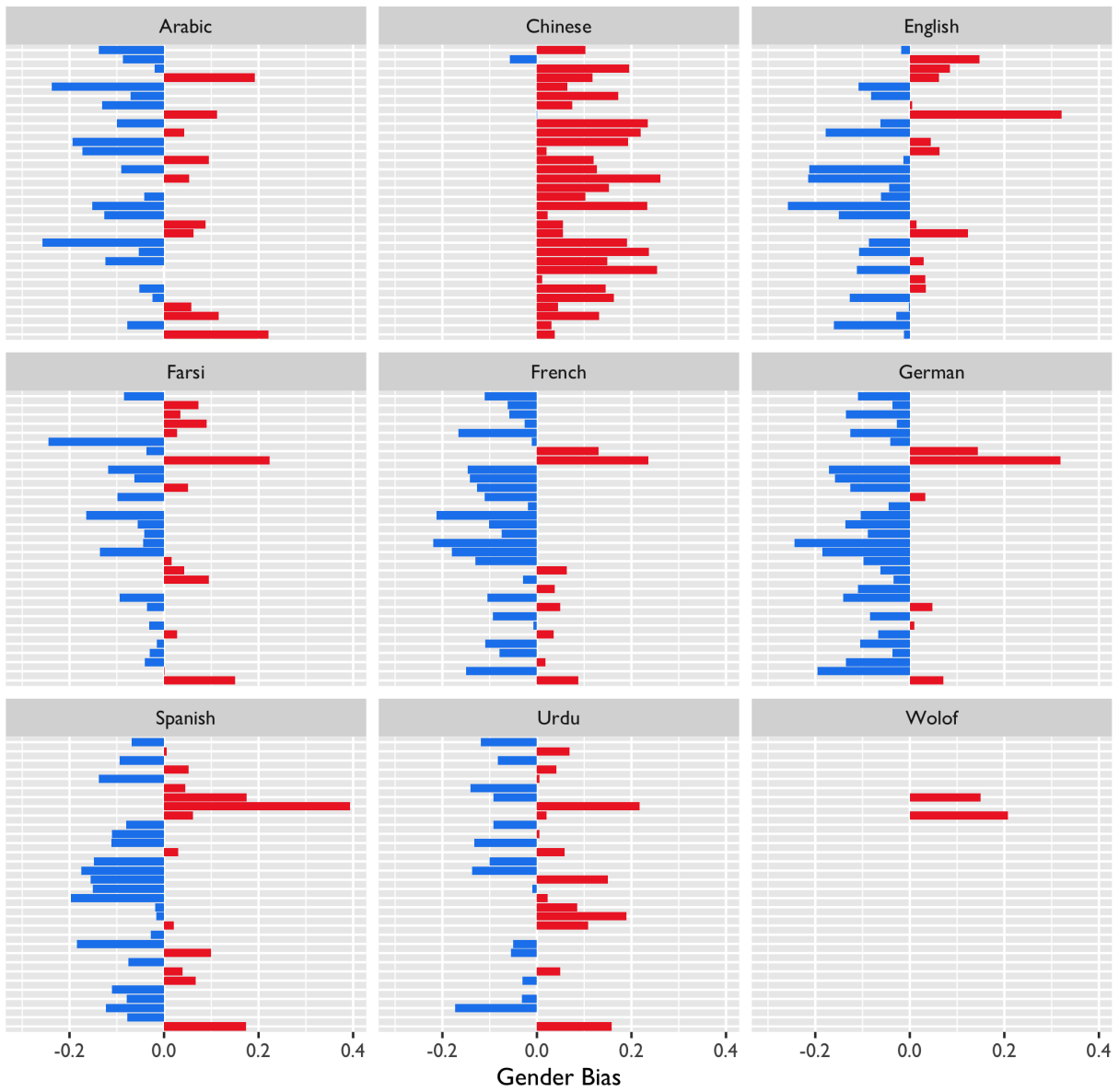


Figure 4.2 – Per-Profession Gender Bias Metrics for Languages With Grammatically Gendered, Weighting by Word Count

Profession Sets Across All Languages (Weighted Evenly)

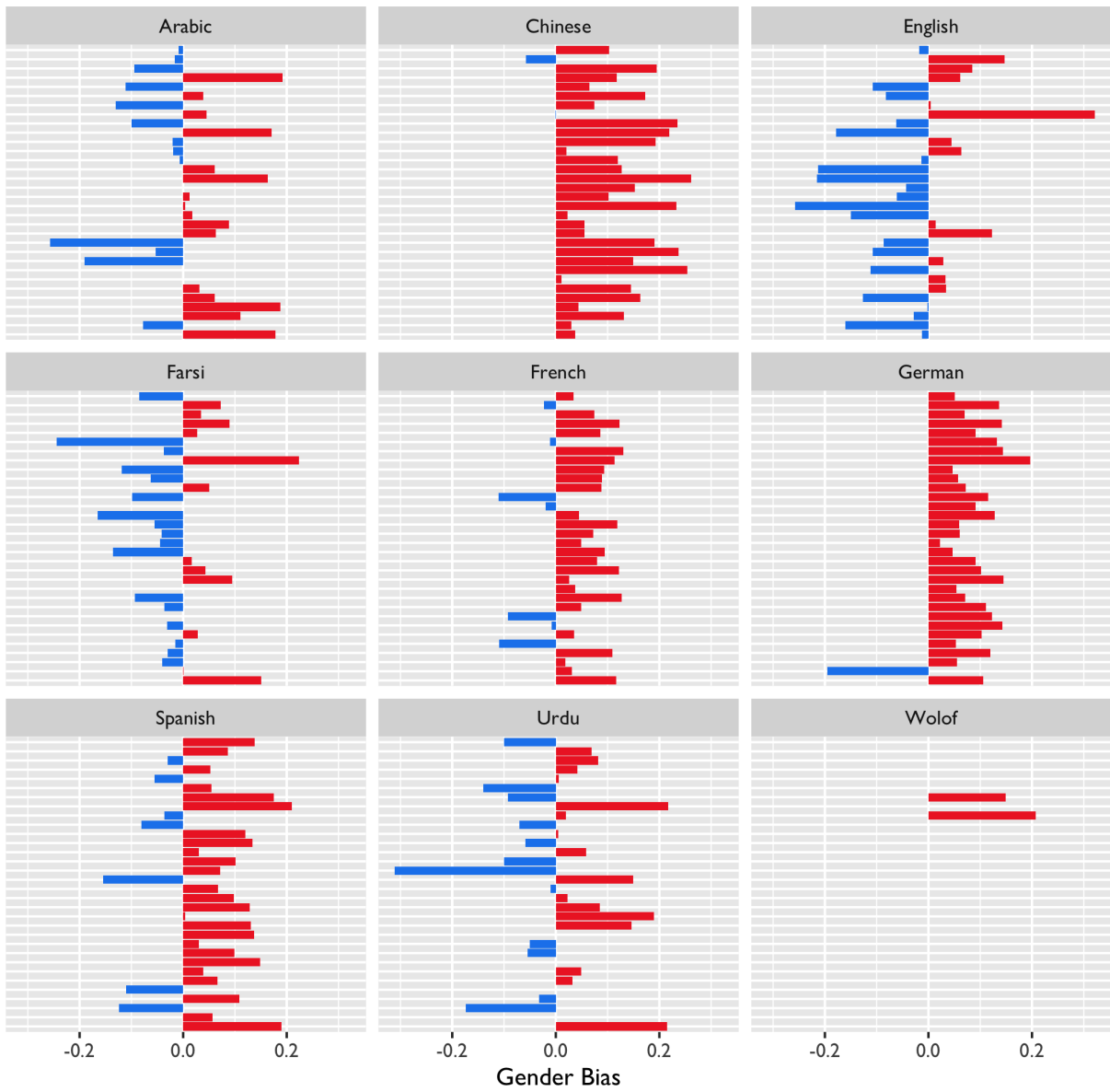


Figure 4.3 Overall Gender Bias Scores for the Wikipedia Corpora

5. APPENDIX 6: Additional Comparison Data for the Wikipedia Corpora Across Languages

Table 1.1: Comparing the number of speakers of a language to the size of the Wikipedia Corpora for that language. For the number of articles and estimates of the number of speakers of Chinese, Spanish, English, Arabic, German, French, Farsi, Urdu, and Wolof. All of the numbers of speakers only accounted for native speakers. Note the * in the Chinese word count reflects that it is the count after being tokenized. To combat ambiguity Chinese tokenizers make an effort to report multiple possible tokenization of a sentence which in some cases can lead to double counting.

Language	Number of Articles (WikipediaC)	Total Word Count (million)	Number of Speakers (thousands)	Articles/ 1000 Speakers	Profession Word Count
Chinese	1,149,477	4011.2*	921,500 (WikipediaB)	1.25	1,342,351
Spanish	1,629,888	7467.85	463,000 (WikipediaB)	3.52	1,436,102
English	6,167,101	27624.37	369,700 (WikipediaB)	16.68	4,795,752
Arabic	1,067,664	2026.73	310,000 (WikipediaA)	3.44	391,323
German	2,485,274	10191.46	95,000 (WikipediaA)	26.16	1,469,592
French	2,253,331	10755.87	77,300 (WikipediaB)	29.15	267,0372
Farsi	747,551	1511.1	70,000 (WikipediaB)	10.68	209,862
Urdu	157,475	264.73	69000 (WikipediaB)	2.28	67,401
Wolof	1,422	5.4	5500 (WikipediaD)	0.26	1,406

Table 6.1: Comparing the number of speakers of a language to the size of the Wikipedia Corpora for that language. For the number of articles and estimates of the number of speakers of Chinese, Spanish, English, Arabic, German, French, Farsi, Urdu, and Wolof. All of the numbers of speakers only accounted for native speakers. These statistics on the number of articles and the number of speakers are taken from Wikipedia itself (WikipediaA)(WikipediaB)(WikipediaC)(WikipediaD).