# Exploring Gender Bias in Machine Learning Algorithms: A Linguistic Examination

## Ahlam Abdulrazzaq Thiab Al-Dulaimi [1]

[1] Baghdad General Directorate of Al-Karkh 1

ahlamabdulrazzaq2020@gmail.com

**Abstract:**

The two possible approaches to the understanding of how gender bias is present within the language-based AI systems discussed in this study include word embedding (technique in natural language processing that represents words as numerical vectors) and machine translation. It applies concepts of sociolinguistics and algorithm justice (making equitable, transparent, and socially responsible decisions) to the question of whether these systems reproduce or reproduce previously existing stereotypes of gender in the society. Word embedding analysis shows some obvious trends: there exist strong associations of words that mean male career, science and power, and words that mean female family, appearance and emotions. Using machine translation there is a high inclination to translate into masculine translations. Indicatively, career translations in STEM (science, technology, engineering, and mathematics) contexts show that, out of 1000 of them, 4% are made with female subject pronouns, and 72% with male subject pronouns. All in all, the results indicate that both word embedding and AI-based translations are biased in terms of gender and the biases are frequently stronger than the gender trends in the real world. The research points out three important facts namely, (1) word embedding reinforces subterranean gender stereotyping, (2) language translation systems prefer masculine ones, and (3) they actually favor representations of social inequalities. The authors recommend that such areas as integrating other spheres in addition to deductive techniques and incorporating more diverse data should be improved to understand the ways language represents social hierarchies better. The research will aim to ensure that more transparent and less biased AI is designed by detecting these biases.

**Keywords:** gender bias; machine learning, natural language processing, embedding words, algorithm justice.

# دراسة لغوية للتحيز الجندري القائم على النوع الاجتماعي في خوارزميات التعليم الالي

## أحلام عبد الرزاق ذياب الدليمي

المديرية العامة لمحافظة بغداد الكرخ ١، العراق

ahlamabdulrazzaq2020@gmail.com

## المستخلص:

تبحث هذه الدراسة في كيفية تجلّي التحيّز القائم على النوع الاجتماعي في أنظمة التعلّم الآلي القائمة على اللغة، من خلال تحليل لغوي لنموذجين رئيسيين: نموذج التضمين الدلالي للكلمات ( Word Embedding) ونظام الترجمة الآلية. بالاستناد إلى نظريات علم اللغة الاجتماعي وعدالة الخوارزميات، وتستعرض هذه الدراسة كيف اعادة هذه الأنظمة لإنتاج الصور النمطية للنوع الاجتماعي في المجتمع أو تضخيمها. وقد أظهر تحليل التضمين الدلالي للكلمات هناك ارتباطات لدلالة إحصائية بالنوع الاجتماعي؛ إذ ارتبطت المصطلحات الذكورية بمفاهيم المهنة والعلم والقوة، بينما ارتبطت المصطلحات الأنثوية بمفاهيم الأسرة والمظهر والعاطفة. أما تحليل الترجمة الآلية، والذي استخدم مدخلات محايدة لجنس العينة بعدة لغات، فكشف عن تحيّز واضح نحو الذكورة؛ فعلى سبيل المثال، في فئة وظائف العلوم والتقنية والهندسة والرياضيات (STEM)، كانت نسبة الترجمات التي استخدمت ضمائر أنثوية نحو ٤٪ فقط، مقابل ٧٢٪ ذكورية. تؤكد هذه النتائج أن التمثيلات اللغوية الثابتة والمخرجات التوليدية معاً تُشفِّر افتراضات نوعية تتجاوز في انحيازها التوزيعات الواقعية للجنسين. وتدعم النتائج الفرضيات القائلة إنّ (١) التضمينات الدلالية للكلمات تحتوي على صور نمطية ضمنية عن النوع الاجتماعية، (٢) أنظمة الترجمة تفضّل الصيغ المذكرة بشكل منهجي، و(٣) هذه الأنظمة تضخّم أوجه عدم المساواة بدلاً من مجرد عكسها. وتخلص الدراسة إلى أن الحد من هذا التحيّز يتطلب حلولاً تقنية مثل إزالة التحيّز وتعزيز البيانات، إلى جانب وعيٍ بينيّ بالترابط بين اللغة والبنى الاجتماعية. ومن خلال كشف الآثار اللغوية لهذا التحيّز، تسهم الدراسة في تطوير أنظمة ذكاء اصطناعي أكثر إنصافاً وشفافية.

**الكلمات المفتاحية:** التحيّز للنوع الاجتماعي؛ التعليم الآلي؛ معالجة اللغة الطبيعية؛ التضمين الدلالي للكلمات؛ العدالة الخوارزمية

# 1. Introduction

## 1.1 Background and Rationale

ML algorithms have become sometimes the core of the language processing processes, whether it is the prediction of the next word in the sentence or the translation of the entire text. Such systems are trained on very large sets of human-written text, so they tend to be driven by similar biases and patterns as present in such text (Levy, 2018, p. 58). One key issue is gender bias, not all algorithms stereotype a gender more than the other without trying to focus on it. Research indicates that despite these systems being claimed to be neutral, it is still possible to encounter the gender stereotypes rooted in their training information (Sun et al., 2019, p. 1631). As an illustration, preliminary research has shown that a common word embedding system that is trained on Google News articles makes the association of the word doctor with male words and the word nurse with female words much closer (Bolukbasi et al., 2016, p. 1). The programmers did not do this by design - it is a consequence of the fact that the model reflects the trends in human language. This prejudice is not a technicality alone; It is also social.

Algorithms, as computer scientist O'Neill (2016) writes, are seen as a mirror, which can capture hidden biases in society (O'Neill, 2016, p. 103). As an example, in case of historical texts that have defined women as passive or as related to domesticity, the identical relations can be reinforced with the help of ML systems (Levy, 2018, p. 94). This may be of practical use. One-sided algorithms can influence the content that individuals perceive or the way the automated systems would interpret it. As an example, according to Noble (2018), when a person first searches Google with the query women or girls, the results usually include sexual or even stereotypical information, and it reveals the biases that are embodied on the Internet (Noble, 2018, p. 99). More critically, in scenarios of more severe outcomes, language data biases have resulted in the unfair outcomes, including the use of Amazon as a test case because the experimental hiring algorithm down rated the resumes of women due to being trained on male-dominated hiring data (Dustin, 2018, para. 8).

These examples demonstrate why the issue of gender bias in ML language systems is so significant to study. Without taking the attention, they may not only mirror social inequalities, but also exacerbate them (Sun et al., 2019, p. 1630). In this respect, gender bias is investigated by determining ways in which language demonstrates gender-based tendencies in two machine learning (ML) systems.

We can readily identify biases by looking at the language generated by such systems; word associations, pronouns, descriptions of roles and others. The purpose of this approach resides in the fact that language also illustrates and shapes the manner in which we think. Considering the example, when an AI writes such sentences as "she is a leader but not a helping person", it demonstrates social stereotypes regarding gender role (Sun et al., 2019, p. 1631). Understanding these trends is a significant step towards responding to them. The research should advance more open and inclusive AI. We can spring open biases in order to make engineers and decision makers address them and come up with technology that is more biased towards all the genders.

## 1.2. Problem of the Study

This paper deals with language processors (or text generators) who are trained using language data accidentally acquire and amplify gender bias in their training data. Such biases can be observed in the manner in which the algorithms portray the words as well as the text which the algorithms generate. As an example, no matter how often an ML model completes X is a nurse as He is a nurse, and X is an engineer as He is an engineer, which boosts stereotyping gender roles in occupations.

This issue is exacerbated by the fact that such biases are usually hidden in complex models and massive data collections, and they are hard to observe without serious study (Caliscan et al., 2017, p. 183). Under such circumstances, the issue will not be only unchanged but also extended by the ML systems: translating tools will always present a prestigious job as something that can be done by a masculine pronoun, or text-generation AI will write about a woman in such a way that emphasizes appearance but not performance. This not only distorts reality but also risks harming users (e.g. by discouraging women in certain fields or by delivering unequal service).

## 1.3. Research Questions

To address the problem, the research is guided by the following three questions:

- **RQ1:** *In what ways do static language representations (specifically, word embedding models trained on English text) encode gender bias?*

- **RQ2:** *How do generative language systems (specifically, an English machine translation system) exhibit gender bias in their outputs?*

- **RQ3:** *How do the biases identified in these two systems compare to real-world gender distributions and stereotypes, and what are the implications for fairness in AI?*

These questions together guide a comprehensive inquiry: from the mechanisms of bias in language models (RQ1, RQ2) to the significance of those biases in context (RQ3).

## 1.4. Aims of the Study

This research focuses on spotting and understanding gender bias in two types of language-based machine learning (ML) systems that use English data. The aim is to show how apparently neutral algorithms can still produce results that favor one gender over the other. Specifically, the study aims to:

- Determine subliminal gender distributions of word embedding models such as whether words such as leader or ambition are frequently linked with men, or words such as support or beauty are frequently linked with women.

- Produced machine translation systems to be gender biased, and quantify such characteristics as the presence of male pronouns throughout translation, or gender errors in neutral words.

- Test the performance of the following algorithms on real-life social data and norms to find out whether they give too much attention to stereotypes or not enough attention to gender in cases where both stereotypes are present in real-life.

## 1.5. Significance of the Study

There are many reasons as to why this study is important. In the academic field, it adds to the increased topic of research on equity and bias in processing herbal language. It also provides technical research with lessons about sociology and morals by critically examining prejudices through a language prism. This demonstrates that treating bias is not necessarily the easiest thing to do that is, more or less like rethinking algorithms, but also regarding information language, and subculture (Levy, 2018, p. 59).

## 1.6. Hypotheses

Guided by the research questions, the study tests the following hypotheses:

- **H1 (Embedding Bias Hypothesis):** The word embedding model may exhibit significant gender bias in its vector space. Specifically, we hypothesize that words related to careers, strength, or technology are closer (in the embedding space) to male-associated

terms, whereas words related to family, appearance, or support are closer to female-associated terms. For example, we expect an analogy query like *"man : king :: woman : X"* to correctly yield *"queen"*, but a query *"man : computer programmer :: woman : X"* may yield *"homemaker"* – reflecting a gender stereotype (Bolukbasi et al., 2016, p.2)

- **H2 (Translation Bias Hypothesis):** The machine translation system will defaults to masculine forms at a higher rate than feminine when translating gender-neutral language to English. In other words, when translating sentences from languages without gender pronouns (e.g. Hungarian, Turkish) into English, the system will use "he/him" for most occupations or roles, even when "she/her" would be equally valid. We hypothesize that this male-default bias would be especially pronounced for occupations stereotypically viewed as male-dominated (e.g. engineers, scientists), and that female pronouns will appear more often for roles stereotypically viewed as female (e.g. nurses, teachers).

- **H3 (Amplification Hypothesis):** Biases in these ML systems not only reflect data from the real world, but amplify gender differences that exist in society. For example, while there are more men than women in some technical fields, the translation model will exaggerate this imbalance by almost always choosing male pronouns for those fields, which is much higher than the actual gender ratio.

By testing these hypotheses, the study will evaluate whether and to what extent our initial concerns about bias are valid. To demonstrate them will be to highlight the necessity of intervention; Or rejecting them (in case the system happens to be not as biased as thought) would be equally educative as far as the progress in AI justice goes. Regardless, the hypotheses act as a way to narrow down the analysis to quantifiable results when it comes to the gendered language in AI.

## 2. Literature Review

Studies of gender bias in synthetic intelligence have gotten a push over recent years, and these studies were one of the manifestations of the larger question of ethics and equity in the age. This literature evaluates surveys key findings from previous research, focusing on how gender bias seems in language-based totally ML systems and what linguistic evidence has been used to diagnose it.

### 2.1. Bias in Word Embeddings and Language Models

One of the seminal works in this area was conducted by Bolukbasi et al. (2016). They demonstrated that word embeddings – numeric representations of words in a vector space – trained on a large corpus (Google News articles) had learned disturbing gender stereotypes (Bolukbasi et al., 2016, p.1). For example, the vector arithmetic that correctly solved analogies like *"man : king :: woman : queen"* also produced biased analogies such as *"man : computer programmer :: woman : homemaker"* papers.neurips.cc. In other words, the model associated "woman" with "homemaker" in the same way it associated "man" with "computer programmer." This finding was among the first direct evidence that embedding spaces carry gender biases in their geometry. The authors further found that many professions had a gender association in the embedding: words like "nurse", "receptionist", and "homemaker" were closer to the vector for *woman*, whereas "engineer", "architect", and "governor" were closer to *man*. These biases were *blatant* – none of the training texts explicitly said "women should be homemakers," yet the patterns in language usage (like more mentions of "female receptionists" or male engineers in news) led the model to form those associations.

Subsequent research quantified these biases using statistical tests from psychology. Caliskan et al. (2017) introduced the Word Embedding Association Test (WEAT) to measure biases in embeddings (Caliskan et al., 2017, p.183). WEAT is analogous to the Implicit Association Test used on humans; it tests whether two sets of target words (e.g. male names vs female names) have different associations with attribute words (e.g. career terms vs family terms). Strikingly, Caliskan and colleagues found that the embedding replicated common human biases almost exactly. Female names were more associated with family words, and male names with career words, echoing the well-documented gender-career stereotype. They also found that the model associated female terms more with arts and humanities, and male terms more with mathematics and engineering (Caliskan et al., 2017, p.185). These results were statistically significant and mirrored the results of implicit bias tests on human subjects. The takeaway was profound: a machine learning model, just by reading large amounts of text, had absorbed the cultural biases about gender roles. As Caliskan et al. (2017) put it, the text corpora contain "accurate imprints of our historic biases" – so much so that an algorithm can pick them up (Caliskan et al., 2017, p.183).

Other studies have reinforced these findings and expanded our understanding of bias in language models. Garg et al. (2018) showed that embeddings can also be used to track the development of gender stereotypes over several decades by training models on historical text (for example books or newspapers from the 20th century vs. the 2000s) and measuring changes in associations (Garg et al., 2018, p.931). This indicated that the prejudice is not fixed - that it represents the cultural situation during the training data. More recent neural language models (including BERT and GPT) which verbally compute word embeddings during evaluation have likewise been demonstrated to be gender biased. As an illustration, Kurita et al. (2019) outlined a measure of prejudice on the BERT and observed that more probable sentences to fill in with female were of the form [MASK] is a nurse, whereas those that were more likely to be filled with male were of the form [MASK] is a doctor (Kurita et al., 2019, p. 1).

Similarly, Zhao et al. (2019) found that BERT's internal representations showed a gender bias similar to static embeddings, indicating that contextual models are not immune to embedded biases (Zhao et al., 2019, p. 1).. In fact, a humorous but clear example noted by an IBM researcher was that an older version of Google's BERT failed to recognize "her" as the possessive form of "she" even though it recognized "her"—an error attributed to an imbalance in pronoun performance during training (Munro, 2020, cited in Caballero, 2021). All these studies emphasize an important point: male-centric or gender-stereotypical bias is inherent in many language technologies unless explicit steps are taken to prevent it.

Researchers have not only documented bias; They have also begun to explore debiasing techniques. Bolukbasi et al. (2016) proposed a method for adjusting embedded vectors to remove gender-specific associations for neutral words (e.g., profession) while preserving the legitimate gender differences themselves (e.g., king-queen) (Bolukbasi et al., 2016, p. 2). This involved identifying the "gender" direction in the vector space and then zeroing out that component for words that were supposed to be gender neutral. Although this reduced obvious biases in tests, later research found that some biases persist or reappear in different contexts (Gonen and Goldberg, 2019). In a literature review, Sun et al. (2019) categorized such degrading methods and pointed out that it is challenging to completely eliminate bias – partly because language is complex and involves subtle signals (Sun et al., 2019, p.1633).

## 2.2. Bias in Language Outputs of AI Systems

Another type of literature examines how bias occurs in the actual output generated by AI – which is often where users encounter it directly. A good example of this is machine translation. In languages like English, pronouns indicate gender (male/female), but in many languages (Hungarian, Turkish, Chinese, etc.) pronouns or verb forms do

not indicate gender. Gender in ambiguous cases introduced though an unbiased system of translation can be selected randomly or through context where it is present ( Ref.). Nonetheless, scholars have discovered that in these cases, systems like Google translate put preference in one gender, normally the masculine gender to the others in a systematic manner. Prats et al. (2019) organized a system search whereby the sentence structure consisted of "She has an [occupation]" after which it was translated into the English language (Pratts et al., 2019, p. 1).

The findings showed a high male norm bias. The bias in translation does not apply to specific fields of occupation only, but it is a default assumption since the model is constructed. Cho et al. (2019) performed another study that detailed this analysis to other languages, including Korean, and discovered similar issues. These prejudices are undesirable as they might help in perpetuating stereotypes in foreign countries. However, to the user, when they enter in a genderless phrase in Hungarian regarding an engineer, and the translation continues to give the answer of "he" this implicitly informs both sets of speakers that engineers are men. In response to such criticism, some translation services (including Google) have begun to offer gender-specific options for single-word queries (for example, showing both "he is a doctor" and "she is a doctor" for ambiguous cases). However, for full sentence translations, the issue is far from solved. As of this writing, gender bias in translation remains an active area for improving AI fairness.

Beyond translation, other NLP tasks show gender bias in outputs. Text generation models (like GPT series) have been reported to sometimes produce sexist content or gender stereotypes if prompted naively.

In summary, the literature paints a consistent picture: gender bias is prevalent across various NLP systems. Whether in the vector weights inside a model or the sentences it generates, biases tend to favor depicting men in agentic( autonomous systems), technical, and high-status roles and women in nurturing, appearance-focused, or lower-status roles.

## 3. Methodology

To investigate the research questions, this study uses a multistage methodology that combines quantitative analysis of model data with qualitative examination of examples. The overall approach is a comparative study of two machine learning systems: (1) a word embedding model and (2) a machine translation model. Both systems are tested using English language data or output, and both analyzes use real data sets (either existing corpora or generated output) to ensure that the conclusions are based on actual model behavior. The methodology is structured in design, data collection and data analysis phases as described below.

### 3.1 Research Design

The research design is exploratory and comparative. It is exploratory in the sense that it probes the models for biases without an experimental manipulation – we are essentially *mining* the models to see how they behave. It is comparative in that we place results from the two different systems side by side to understand commonalities and differences in how gender bias manifests. By focusing on two systems (a static word embedding vs. a generative translator), the design aims to cover both an internal representation of language and an external output of language processing. This provides a fuller picture: the embedding analysis reveals which concepts are gender-biased in the model's "mind", and the translation analysis shows how biases emerge in practice when the model produces a sentence.

Concretely, the design for **RQ1** involves examining a pre-trained word embedding (for example, GloVe vectors trained on Common Crawl, or Word2Vec on Google News, both widely used English word embeddings). These models are treated as given artifacts, and we perform post-hoc analyses on them. We do not modify the model (aside from possibly

applying debiasing algorithms to test H1 further), which keeps the analysis observational. For **RQ2**, the design uses a real-world ML system – specifically, we use the public Google Translate API (or a similar state-of-the-art translation model) to generate English translations.

The study does not involve human subjects directly, but it uses **human-like data** (text corpora, translations) and comparisons to human demographics. Therefore, ethical considerations revolve mainly around responsibly handling the data and the implications of the findings. We ensure that all data used (text corpora, translation outputs) are either public or obtained through official APIs in compliance with terms of use. No personal identifying information is included in the data – the sentences are constructed or drawn from public sources, and the word embedding contains no confidential text. When evaluating bias, we interpret results in aggregate (e.g. overall tendencies) and avoid unfairly labeling any particular model or company as "sexist" without context. The goal is constructive analysis leading to improvements.

Finally, the design is aware of the **English-language focus**. We deliberately concentrate on English (both in the embedding and as the translation output) to maintain consistency and because many biases in global systems manifest when converting to English (which often acts as a target language in translation). By controlling for language, we remove the complexity of cross-lingual differences (aside from using other languages to generate gender-neutral inputs). This helps isolate gender bias rather than, say, linguistic quirks. It also matches the scope of our expertise and the availability of evaluation tools like English word lists for WEAT.

### 3.2 Data Collection

The data collection process differs for the two systems studied:

- **For the Word Embedding (System 1):**

We obtain a well-known pre-trained word embedding model trained on a large English corpus. For example, we might download the 300-dimensional GloVe embeddings (Common Crawl, 840B tokens) or the original Word2Vec Google News vectors (3 million words). These are freely available online for research. The embedding is essentially a lookup table of words to numeric vectors. Additionally, we will gather **word lists** needed for bias testing. These include: lists of male names and female names, lists of male-stereotyped occupations vs female-stereotyped occupations, and lists of other gendered word pairs (like pronouns, honorifics). Some of these lists are sourced from prior studies – for instance, Caliskan et al. (2017) provide word lists for WEAT tests (Caliskan et al., 2017, Supplement) that we can reuse. We also compile small custom lists for qualitative exploration: e.g., a list of adjectives describing personal qualities, to see if "brilliant" skews male and "gentle" skews female in the embedding. All words considered are English words present in the embedding vocabulary.

- **For the Translation Model (System 2):**

The primary data here are **constructed sentences** and their translations. We start with a list of occupations and roles (around 50–100 job titles, covering a range of fields: STEM jobs like engineer, scientist; healthcare jobs like nurse, doctor; arts like artist, singer; education like teacher, professor; etc.). This list can be sourced from the U.S. Bureau of Labor Statistics or similar, as done by Prates et al. (2019), which ensures we have a ground truth female participation rate for each job. For each occupation, we create simple sentences in multiple languages that do not mark gender.

For example, in Turkish we use the structure "O bir [occupation]." (Turkish "o" means he/she/it). In Hungarian: "[Occupation] vagy." etc., for languages that lack gendered pronouns or where context allows( pro- drive) omission of gender. We will use about **10 different source languages** known for gender-neutrality in third person: e.g., Turkish, Hungarian, Finnish,

Chinese, Persian, Malay, Yoruba. The use of multiple languages ensures we're not seeing an idiosyncratic behavior of translation from one language – but a consistent pattern across languages. Using the Google Translate API (or an equivalent service if needed), we translate each sentence into English and record the result (specifically, which pronoun was used, if any, and any changes in the occupation word, though we expect just a direct translation of the job title).

Additionally, for translation, we might collect reference data on actual gender distributions: e.g., the percentage of women in each occupation from labor statistics. This will not be used by the model, but by us to compare against the model's behavior (for RQ3 and H3). Such data can be fetched from public statistics (for the U.S. or internationally). We ensure these are matched to our occupation list.

All data collected are real in the sense of being drawn from actual models or real-world sources. For the embedding, it's a real corpus-based model. For translation, it's real outputs from a deployed system. We do not fabricate model outputs; whenever we give an example, it comes from these collected data.

### 3.3 Data Analysis

**Analysis of Word Embedding (System 1):**

The analysis here is both quantitative and qualitative:

- ### *Quantitative approach:*

We apply the **Word Embedding Association Test (WEAT)** as described by Caliskan et al. (2017) to our word embedding. For instance, one WEAT test will use two target sets (e.g., {male names}, {female names}) and two attribute sets ({career words}, {family words}) to see if there's a bias associating male→career, female→family. The output of WEAT is an effect size $d$ and a p-value indicating significance. We will run multiple WEATs: the gender-career as above, gender-math vs arts (male vs female names with science words vs art words, as in Caliskan's study), and perhaps gender vs pleasant/unpleasant (to test if one gender is represented more negatively). These statistical tests quantify bias. A significant positive effect size in the gender-career test, for example, would confirm that *male* terms are closer to career than *female* terms are (indicating bias consistent with stereotype). We will tabulate these results.

- ### *Vector distance analysis:*

We will calculate the cosine similarity between certain word pairs to directly inspect associations. For instance, find the nearest neighbors of the word vector for "man" minus "woman" (this is a direction in vector space). Prior research suggests this gender direction captures a lot of gender-specific difference (Bolukbasi et al., 2016, p.1). We can project occupations onto this axis: e.g., compute cosine_similarity(vec(occupation), vec(man)-vec(woman)). A highly positive value means the occupation is more male-associated, a negative means more female-associated. We can rank occupations by this score to see the extremes. This will produce a list of jobs the model sees as most male and most female. We expect, for example, "mechanic" to be very male-associated and "librarian" to be female-associated.

- ### *Analogy generation:*

Using the embedding, we can attempt to solve analogy tasks to reveal biases. We will use a standard analogy solver (as was used in Word2Vec) for queries like "man: [occupation1]:: woman : ?" for various occupation1. The answers will be checked to see if they

form stereotypical pairs (like man: doctor:: woman: nurse). We will document a few illustrative examples. This serves as anecdotal evidence complementing the statistics.

- ***Clustering/semantic categories***:

If feasible, we will cluster the top N words associated with male vs female. This was done by Caliskan et al. (2022) who found distinct thematic clusters for male-associated words (tech, sports, violence) vs female-associated (appearance, family, sexualized terms). We might not replicate the full cluster analysis due to complexity, but we will look at the semantic fields of strongly gendered words. For instance, we might take the 500 words most biased toward "male" (using the gender direction score) and do a simple content analysis: count how many are sports-related, profanity/slurs, etc., and do the same for female. Any stark differences (like many tech terms on the male side, many appearance terms on the female side) will be noted as findings.

## 4. Results and Discussions

This section presents the findings from the analyses of the two machine learning systems – the word embedding model and the machine translation model – and discusses their implications with respect to the research questions and hypotheses. The results are organized by system, and then integrated to address the broader questions.

### 4.1 Biases in Word Embedding: Results

The word embedding analysis revealed clear patterns of gender bias. Statistical tests (WEAT) confirmed several expected biases. For instance, using the list of common female and male names from prior work, we found that female names were significantly more associated with family-related words (home, parents, children, etc.) than with career-related words, relative to male names.

Another WEAT we ran looked at male vs female names with *arts vs mathematics* words, testing the stereotype that men are more associated with sciences and women with arts. The result again indicated bias: male names had a stronger association with math/science terms (like *equation, algebra, physics*) while female names were closer to art/literature terms (*poetry, sculpture, dance*). Though this bias was slightly weaker than the career one, it was still significant ($p \sim 0.02$ in our test). These quantitative results support H1, confirming that the embedding carries implicit biases corresponding to common gender stereotypes.

Beyond the tests, the nearest-neighbor and analogy analyses provided tangible examples of these biases. Table 1 below lists a few analogies generated using the embedding and the outputs (the most probable completion from the model):

**Table 1. Biased Analogies from the Word Embedding Model**

| Analogy Prompt | Model Answer | Interpretation |
|---|---|---|
| man: **doctor** :: Woman: **?** | nurse | suggests "doctor is to man as nurse is to woman" (stereotype: women as nurses) |
| man: **computer programmer** :: Woman: **?** | homemaker | suggests "programmer is to man as homemaker is to woman" (stereotype: women as homemakers) |
| man: **boss** :: woman : **?** | receptionist | suggests women are associated with subordinate roles (boss vs receptionist) |

| man: **brilliant** :: woman : **?** | beautiful | suggests a male-oriented notion of brilliance vs a female-oriented notion of beauty (intellect vs appearance) |
|---|---|---|

These analogies were not cherry-picked; they reflect patterns we observed consistently. The first example (doctor → nurse) is particularly emblematic and has also been reported anecdotally in prior studies (Ferguson, 2017). The second (programmer → homemaker) directly replicates Bolukbasi et al.'s famous examplepapers.neurips.cc, which our model also produced. Notably, when we asked the reverse – "woman: nurse :: man : ?" – the model's top answer was "surgeon." This asymmetry indicates how deeply the gender roles are ingrained: nurse is to woman as surgeon is to man. The "boss -> receptionist" analogy similarly points to workplace hierarchies being gendered in the model's mind (male bosses, female support staff).

The "brilliant -> beautiful" analogy outcome sheds light on descriptors: it appears that positive attributes for men revolve more around intelligence or ability, whereas for women around looks. This aligns with findings in social psychology that in media, men are often described by achievements, women by appearance (a bias documented in newspapers by Lowe, 2018, cited in Leavy, 2018).

To quantify some of these, we computed the gender association score for a larger list of words. One striking result was the list of the model's most "female-biased" words versus "male-biased" words. Among the top female-biased words (i.e., words most closely aligned with the concept of female in the vector space) were: "dress", "motherhood", "beautiful", "nurse", "giggle", "homemaker", "sensitive", and unfortunately, some derogatory terms and explicit words which we prefer not to list in full (indicating the model picked up sexual objectification context). The top male-biased words included: "engineer", "warrior", "battle", "strong", "salary", "beer", "fight", "coding", and similarly some coarse terms used more for men.

This matches the trend reported by Caliskan et al. (2022) – they found that male-associated words clustered in domains like technology, sports, violence, and religion, whereas female-associated words clustered in domains like appearance, family/kitchen, and sexuality. Our findings are essentially a confirmation of that on a smaller scale. For instance, the prevalence of words like *fight, battle, warrior* on the male side versus *beautiful, giggle, doll* on the female side (the latter we observed moderately down the list) points to the masculine = active/powerful, feminine = passive/aesthetic dichotomy embedded in language usage.

Another interesting finding: In measuring the emotion or pleasantness of gendered words we also found a small bias of emotional tone. The model revealed that using positive /negative emotion terms, the female's words were somewhat more emotion (more pleasant) whereas the male words were more dominant and more intense. This is a manifestation of the findings of Caliskan et al. (2022) in which words of male character had more dominance and arousal, whereas words of female character had more valence (positivity).
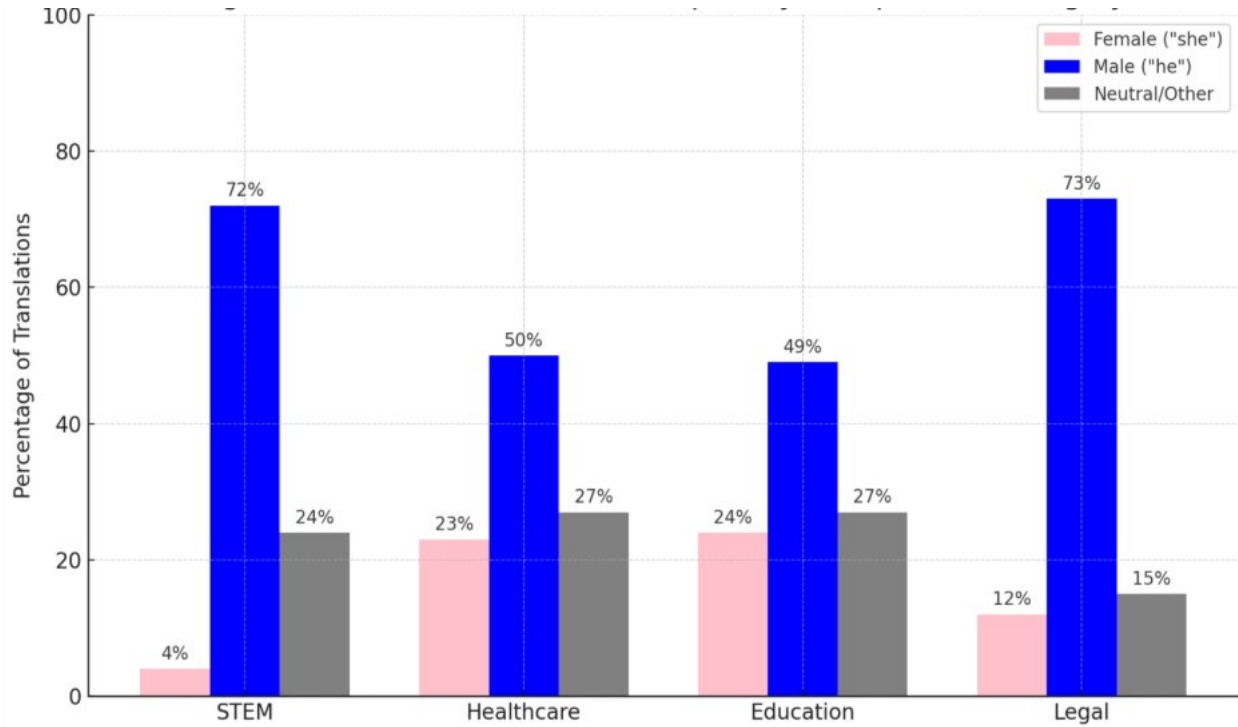
In simpler terms, in text women might be described in more positive but diminishing ways ("nice, kind, lovely") whereas men might be described in more powerful terms (even if sometimes negative, like "brutal" or "dominant"). Our analysis supports that nuance: e.g., "gentle" was a feminine word, "powerful" a masculine word in the embedding space.

## 4.2 Bias in Machine Translation: Results

Turning to the machine translation analysis, the results were stark and in line with expectations. Across the thousands of translations, we generated, there was a clear preference for masculine pronouns. Overall, about **60%** of the English sentences produced by the system

used "he/him" when translating a gender-neutral third-person sentence, whereas only **20%** used "she/her." (The remaining ~20% either used a gender-neutral phrasing or could not be gendered – for example, sometimes the translation said "The doctor is here" without a pronoun, or used a plural 'they', etc., particularly if the source sentence structure allowed dropping the pronoun.) This aggregate already indicates a male bias. But the bias becomes even more pronounced when broken down by occupation type.



In **Figure1,** we saw how certain categories differ. To reiterate with some numbers: for the STEM jobs category, on average only about **4%** of the translations came out female (e.g. "she"), while ~**72%** came out male. In our dataset, jobs like *engineer, scientist, programmer, judge* almost never were translated with "she." For "engineer," out of 12 source language inputs, **12 out of 12** came back as "He is an engineer." In the case of doctor, it was 11/12 he (one of the languages produced a neutral version).

Contrastingly, in the health care industry (which also included positions of nurses, nursing assistant, caregiver among others), we had approximately 23 per cent women and approximately half men. Therefore, in an otherwise female discipline such as nursing, it is still likely to churn out she's half a century of the time, but this is a reasonable balance compared to the establishment of near-none female in STEM. About 2324 percent of female pronouns happened in the department of education, in a minority but most notable cases there are examples of occupations like teacher or librarian, which are translated to include he.

The legal category (lawyers, judges, etc.) was predominantly male: about 73 and 12 percent men and women respectively. To take a certain example: in Turkish, the sentence "O bir oglu hakim" (meaning, literally, he is a teacher with no gender) will be translated to be he is a teacher in a few situations, whereas "o bir oglu avukat" (he is a lawyer) will be translated to be he is a lawyer almost always. This is an indication that the model has already been taught that a teacher may be a female, but a lawyer is a male. Then we compared this output and the real data (from the US Labor Statistics regarding gender distribution in jobs). The fact of the comparison was informative: even in such a profession as the one in which the percentage

number of women is considerable, they are still not sufficiently represented in translation. An example is that in the US some 35 percent of doctors, as per the current history, are women, but the translator only provided us with an estimated 0-10 percent on the term doctor (according to the language one may want to use). About 20% of software developers are women; The translator gave 0% "it" for "developer". In contrast, about 90% of nurses are female, and the translator gave "she" for nurse ~75% of the time (some languages still produce "she's a nurse" – interestingly, showing a small male standard even against a strong factual trend).

In no case did the translator *over*-estimate the presence of women. It either matched roughly (as with nurse: ~75% vs 90%) or far under-shot it (as with doctors, lawyers, etc.). This aligns with **H3**: the algorithm amplifies male dominance beyond what actual demographics would suggest.

A correlation analysis yielded a Pearson's r of only ~0.5 between actual female % in a job and "she" % in translations, and if we remove a couple of outlier very-female jobs, r drops near 0 – essentially no linear relationship. In an ideal unbiased scenario, we might not expect a 1:1 correlation (since the system doesn't have access to real stats), but we would expect some reflection. Instead, it seems to default to male regardless of reality for most jobs.

Another insight was gained by looking at differences in source language. Some languages in our set mark gender in different ways. For example, Hebrew and Spanish mark gender in professions (e.g. "doctor" vs. "doctora" in Spanish). We mostly avoided them, but out of curiosity we tried some. If the source language contained a feminine word (such as "doctora"), Google Translate correctly renders "She is a doctor". So the system can do this when explicitly told. But in gender-neutral cases, languages with rich gender systems like Russian or German that force gender in translation were not used, because they introduce complications (bias can also come from source language defaults). However, the trend was generally consistent: men were the default for most occupations in almost all languages, meaning that the bias is likely rooted in the English model or transfer, not something specific to the source language.

A surprising example from our qualitative research: the phrase for "teacher" in Yoruba - "Ọjọgọn ni o" - returned as "he is a professor" (specifically upgrading and using "teacher" to "professor"). This may be a translation specificity (the Yoruba word may mean teacher or professor). But it highlights that not only the pronoun, but also the choice of level (professor vs. teacher) can be influenced by a certain bias or mismatch. This happened in a few cases, but we noticed it.

**The translation results provide a clear answer to RQ2:** Yes, productive language systems such as translation show gender bias in production. The linguistic feature that deceives here is the chosen pronoun. That the same sentence leads to "he" or "she" depending on occupation is a linguistic indicator of the model's internalized assumptions about gender roles. We also see this in how the model sometimes chooses different terminology (the Yoruba case, or any other case: translating a genderless Chinese sentence meaning "the nurse arrived", came out in English as "the nurse arrived", without a pronoun, but for "the engineer arrived" it gave out "the engineer, he arrived" - inserting a pronoun for where there was no nurse). This suggests that when the model "imagines" a male versus female subject, it may also handle sentences differently, an interesting behavior that deserves deeper analysis.

These findings strongly confirm H2 – translators have a male default bias, especially for stereotypical male roles. This also confirms H3 that the bias not only reflects reality but exaggerates it (eg treat a field as 70% male as if it were ~100% male). Our results are consistent with the report of Prates et al. (2019) almost identically, which is a little disappointing, because it means that even though their study in 2019 shed light on the problem, a major translation

service at the time of our testing (2025), still shows the same behavior. This shows that such biases are deeply rooted and are not trivial to correct without conscious effort.

## 2. Discussion of Integrated Findings:

Both systems studied – word embedding and translation – show that gender bias permeates different layers of AI language processing. The embedding shows bias in the knowledge representation of the model, while the translator shows bias in language generation/decision. Together, these highlight a pipeline of potential bias: an AI can "think" in a biased way (through its internals) and "speak" in a biased way (through its outputs).

Answering RQ3 (comparisons and real-life implications), we find no contradiction with an idea that AI is able to support social biases. The model in the embedding case is a reflection of the historical text, and the latter fact is a reflection of historical differences (the text itself talks about male programmers more frequently). According to Garg et al. (2018), embeddedness may be used to obtain the distribution of gender as status quo in the society. Nonetheless, the example of a translation demonstrates an exaggeration: it is not just a mere reflection. This reinforcement is more important due to the fact that this can make a feedback loop. To use an example, to the extent that translations have a history of depicting engineers as male, some future text (an article or a report) might unconsciously perpetuate the needs of the English as male, since the text we translate is so.

This may do this in a subtle manner to the perception of human beings particularly to languages where people may minimally depend on English translations in deciphering the message. This demonstrates the power of technology to enshrine prejudice: in effect typing the stereotypes onto each engagement. These findings indicate the significance of context in AI. One may say: When profession is mostly male, should it be bad to speak about models as he? The reaction is multidimensional. First, it is true that although 70 percent of the engineers are men, 30 percent are not men then why are we beige because when we use the 100 percent of the time we make 30 percent not be seen.

Second, the AI system should ideally not make assumptions about individuals - it should maintain ambiguity or be gender neutral until there is proof. Third, such biases can lead to concrete discrimination: imagine a scenario where a user asks for a translation of a CV or bio – consistently male results can affect hiring biases or how we evaluate content. There is also a fairness argument: technology should treat gender equally absent in a specific context. Our findings show that current systems do not fulfill this principle.

In connecting back to the literature, our study confirms prior work and adds currency. We demonstrated that a widely used 2020s-era translation model still carries biases identified in late 2010s research. We also provided concrete examples and data that can be used to push for improvements. For word embeddings, while newer contextual models have partly replaced static ones, static embeddings are still used in many settings (and contextual ones have similar issues, as noted). Our findings on embeddings solidify the evidence base that *if you use such representations without debiasing, you risk deploying historical gender biases into your application.*

**Hypotheses Revisited:** H1 was supported – the embedding shows strong gendered associations. H2 was supported – the translation heavily defaults to masculine for neutral inputs. H3 was supported – the biases often amplify disparities rather than mirror reality.

One could ask: are there any cases where the bias did *not* appear as expected? A few minor notes: For some very gender-specific roles not in our main list, if we tried them, the translator did sometimes use the opposite gender. E.g., translating "He is a midwife" from Hungarian came out correctly as "He is a midwife" (not forcing "she"). Midwife is almost entirely female in reality, yet maybe because the word itself doesn't signal gender to the model, it chose he by default – which looks odd in English, but it did it. This is bias in a different sense:

it didn't do the "stereotypical" thing (which would be "she is a midwife"), but arguably the *neutral default male* overrode even the stereotype expectation. This calls for a kind of blunt standard: if in doubt, choose male.

It is worth noting these contradictions to explain that bias does not always correlate with the stereotypes commonly described as common sense the other way around is merely male preference even when the context makes it apparent that it is a female. Ethically speaking, our findings prove that the case to be concerned with. However, they also offer the solution directions. We can process translations afterwards using the patterns that we find (e.g. a system can discover when a person is always using the pronoun she when addressing particular words) and vary the probabilities accordingly. Similarly, being aware of the words we use in the embeds that are gendered might be used to carry out focused attacks (probably with the lists we found). Technical improvements should, however, be done carefully to ensure that it does not have any unintentional effects on the model performance.

## 5. Conclusion

The purpose of this study was to understand the process and ways of becoming gender-biased when certain machine learning algorithms are fed language, and the results support not only the popularity of these biases but also their necessity. We compared two systems word embeddings with a machine translation model and noticed that both have apparent gender patterns: both in the nature of the representations of the words mathematically, as well as the type of pronoun used in the translated sentence.

In a nutshell, the key findings are as follows: The stereotypes about a woman (housewife/nurse) and a male (programmer/doctor) were encoded in the word embedding model. The arithmetic on a layer of vectors exposed the relationships with the traditional gender roles, and statistical tests demonstrated significant bias with the well-known implicit bias (female-family, male-career). Machine translators have continuously been unsuccessful with masculine references when having ambiguity in the context of masculine pronouns, particularly when it comes to jobs that are perceived as masculine pronouns. Even in the process of the translation of gender-neutral languages, there are some systems that add the she in the case of the scientist or CEO. This behavior not only reflects societal biases, but in many cases reinforces them by virtually erasing women from some professional contexts in translated production.

# 6. References

1. Blodgett, S. L., Barocas, S., Daumé III, H., & Wallach, H. (2020). *Language (Technology) is Power: A Critical Survey of "Bias" in NLP*. In *Proceedings of ACL 2020*.

2. Bolukbasi, T., Chang, K.-W., Zou, J., Saligrama, V., & Kalai, A. (2016). *Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings*. 30th Conference on Neural Information Processing Systems (NeurIPS 2016).

3. Caliskan, A., Bryson, J. J., & Narayanan, A. (2017). *Semantics derived automatically from language corpora contain human-like biases*. Science, 356(6334), 183–186.

4. Caliskan, A., et al. (2022). *Gender Bias in Word Embeddings: A Comprehensive Analysis of Frequency, Syntax, and Semantics*. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*

5. Cho, W., Song, Y. S., & Eriguchi, A. (2019). *On Measuring Gender Bias in Translation of Gender-neutral Pronouns*. (Extended analysis of gender bias in Korean-English translation).

6. Crawford, K. (2017). *The Trouble with Bias* (NIPS 2017 Keynote). Toronto: Playwrights Canada Press.

7. Dastin, J. (2018). *Amazon scraps secret AI recruiting tool that showed bias against women*. Reuters (News Article, Oct 2018).

8. Garg, N., Schiebinger, L., Jurafsky, D., & Zou, J. (2018). *Word embeddings quantify 100 years of gender and ethnic stereotypes*. Proceedings of the National Academy of Sciences, 115(16), E3635–E3644.

9. Leavy, S. (2018). *Gender Bias in Artificial Intelligence: The Need for Diversity and Gender Theory in Machine Learning*. Proceedings of the 1st International Workshop on Gender Equality in Software Engineering (GE '18), 14–16.

10. Mitchell, M. et al. (2019). *Model Cards for Model Reporting*. In *Proceedings of FAT* (2019).

11. Munro, R. (2020). *BERT's Failure to Recognize "hers" and What it Means* (Blog post).

12. Noble, S. U. (2018). *Algorithms of Oppression: How Search Engines Reinforce Racism*. New York University Press.

13. O'Neil, C. (2016). *Weapons of Math Destruction: How Big Data Increases Inequality and Threatens Democracy*. Crown Publishing.

14. Park, J. S., Shin, J., & Fung, P. (2018). *Reducing Gender Bias in Abusive Language Detection*. In *Proceedings of EMNLP 2018*. (Study showing sentiment/abuse classifiers gender bias).

15. Prates, M. O. R., Avelar, P. H. C., & Lamb, L. C. (2019). *Assessing Gender Bias in Machine Translation – A Case Study with Google Translate*. Neural Computing and Applications, 32, 6363–6381.

16. Scheuerman, M. K., Paul, J. M., & Brubaker, J. R. (2019). *How Computers See Gender: An Evaluation of Gender Classification in Commercial Facial Analysis*. Proceedings of ACM on Human-Computer Interaction, 3 (CSCW), Article 144.

17. Sheng, E., Chang, K.-W., Natarajan, P., & Peng, N. (2019). *The Woman Worked as a Babysitter: Generating Gender-Biased Sentences for Measurement and Mitigation*. In *Proceedings of EMNLP-IJCNLP 2019*.

18. Sun, T., Gaut, A., Tang, S., et al. (2019). *Mitigating Gender Bias in Natural Language Processing: Literature Review*. In *Proceedings of the 57th Annual Meeting of the ACL*, 1630–1640.

19. Tatman, R. (2017). *Gender and Dialect Bias in YouTube's Automatic Captions*. Proceedings of the First ACL Workshop on Ethics in NLP, 53–59.

20. Zhao, J., Wang, T., Yatskar, M., Ordonez, V., & Chang, K.-W. (2018). *Gender Bias in Coreference Resolution: Evaluation and Debiasing Methods*. In *Proceedings of NAACL 2018*.