

Exploring Gender Bias in Translation Models

Rudy Yuen

zcabcry@ucl.ac.uk

Angel Oyelade

zcabao@ucl.ac.uk

Carol Nghiem

zcabngh@ucl.ac.uk

Katya Piotrovskaya

zcabeap@ucl.ac.uk

Ilya Kosorukov

zcabiko@ucl.ac.uk

Abstract

Modern Large Language Models have been seen to be more competent than ever in their ability to infer languages. Subsequently, this creates a rising concern about mitigating unwanted behaviour from these models, which may have been induced by their training data. These unwanted behaviours are often referred to as “bias”. In this research, we investigate gender bias within 3 modern machine translation models through a prompt engineering approach. We create 3 different hint levels totalling 7 prompts to investigate the probability distribution that the model generates to translate our prompt, which is in English to Spanish. We have seen that all 3 models, NLLB-200, M2M-100 and mBART are still biased towards outputting male acronyms, whereby although the degree of bias has been observed to decrease when the level of hints we provide increases, there are certain occupations where the model still heavily skewed towards using the male acronym, demonstrating an area of improvement in model safety.

1 Introduction

A fundamental concept in Machine Learning is to construct models that spot and replicate behaviours inside the training data with the hope that the distribution of training data simulates that of the wider world (Bishop, 2006; Goodfellow et al., 2016). Yet, in situations where the training data has a skew in distribution towards a particular class, it introduces a possibility of model bias: the model might replicate unwanted or incorrect behaviours to reduce the inaccuracy it has with the training data (Pagano et al., 2022).

An example of such behaviour is Class Imbalance in Classification, where the model might predict all inputs to be from the same class, as there is

a high probability that an unseen input comes from that class (Rawat and Mishra, 2022). On the topic of Large Language Models (LLMs), the model might replicate dangerous behaviours; for example, TayTweets exhibited racist, sexist and politically incorrect behaviours (Kraft, 2016), which was introduced from toxic interactions with users.

Whilst all aforementioned models can infer from unseen data, they tend to choose ethically wrong or biased answers. This is not a preferred behaviour as it introduces a chance of user hostility and decreases the model’s helpfulness.

With the recent explosion of LLMs such as ChatGPT and Gemini, more attention has been paid towards safety. Measures such as Reinforcement Learning through Human Feedback and Supervised Fine-tuning have been used to mitigate bias and promote model safety (Google, 2023). Whilst not as extreme as TayTweets, this biased behaviour is still observed within the realms of Machine Translation (MT), with a particular emphasis on gender. This phenomenon is particularly observed within languages that use different words for different genders, such as most Latin (e.g. Spanish, French) and Slavic (e.g. Polish, Russian) languages (Savoldi et al., 2021).

In this work, we investigate gender bias within MT by observing how carefully designed prompts are translated from a genderless language (English) into a non-genderless one (Spanish). We prompt an LLM with and without gender-identifying hints to observe whether it will comprehend them and provide an accurate translation. We extend our investigation within several state-of-the-art machine translation algorithms such as mBART (Liu et al., 2020), M2M-100 (Fan et al., 2020) and NLLB-200 (Team et al., 2022). The choice of languages was due to both being widely spoken in the world, making the findings of this study highly relevant and impactful.

In this paper, we show that the three investigated LLMs exhibit gender bias towards male acronyms in all three hint levels scenarios (i.e. no hint, subtle hint and obvious hint) and provide a comparison between models’ performance on each occupation basis.

2 Related Work

After a review of existing literature, we have discovered that a comprehensive literature review on this topic has already been conducted in [Savoldi et al. \(2021\)](#). In our review, we will reference this paper and summarise the background knowledge needed for this paper.

2.1 Defining Gender Bias

As discussed in [Savoldi et al. \(2021\)](#), the term Bias can have many partially overlapping definitions in various different disciplines. In our research, we will be using the definition as explained in [Savoldi et al. \(2021\)](#), which states that an MT model is deemed biased if it consistently and unfairly favours specific individuals or groups over others, resulting in discrimination ([Friedman and Nissenbaum, 1996](#)).

Similarly to [Savoldi et al. \(2021\)](#), we have also reached the conclusion that studies on gender bias in machine translation (MT) have commonly operated under an implicit assumption of a binary gender framework. Consequently, we will also be using this definition throughout our work.

2.2 Gender Bias in Machine Translation

[Prates et al. \(2019\)](#) conducted a study on gender bias within Google Translate by testing gender translation statistics against the male and female frequency in a job position taken from the U.S. Bureau of Labor Statistics. They show that Google Translate fails to reproduce a real-world distribution of female workers. We will take a similar approach, as described later. However, we will be investigating other popular translation LLMs to observe if they have the same issues.

2.3 Large Language Models

In this work we investigate gender bias in machine translation (MT) using 3 different LLMs: mBART, NLLB and M2M. These LLMs provide extensive multilingual support and are optimised for translation across a wide range of language pairs including Spanish and English.

mBART (Multilingual Bidirectional and Auto-Regressive Transformers) performs sequence-to-sequence translations using transformer architecture; it was one of the first LLMs to do this instead of only relying on encoders or decoders. It was able to generate text in multiple languages, unlike many models at the time of its introduction, which were either monolingual or bilingual.

NLLB (No Language Left Behind) was developed to reduce language barriers and focuses on achieving language equity by trying to minimise the performance disparity between high and low-resource languages. It is stated to have an improvement of 44% from the previous state-of-the-art in machine translation ([Team et al., 2022](#)).

M2M is aimed at facilitating translations between language pairs without the use of English as an intermediary. It provided a move from English-centric translation.

3 Methodology

To systematically identify and quantify gender bias in LLMs, we have created a dataset of 26 occupations in English alongside their Spanish translations, both male and female versions in the latter case, as that would vary given the context, allowing for a clear evaluation of gender bias. The following experiments are carried out using three distinct LLMs: mBART, M2M-100 and NLLB-200.

3.1 Translation Scenarios

We have designed three separate categories of translation scenarios to incorporate the occupations: no hint, subtle hint, and obvious hint. This has been done to allow us to assess the LLM’s ability to infer gender from context and adjust translations accordingly (see Table 1).

No Hint: Neutral prompt to examine whether the LLM’s translations lean towards a specific gender when no gender information is provided.

Subtle Hint: Includes five prompts designed to assess the LLM’s sensitivity to direct yet non-obvious gender indications.

Obvious Hint: Contains gender-specific prompts to determine the LLM’s translation accuracy in the presence of explicit gender cues.

3.2 Evaluation

For the no hint scenario, we have inserted all 26 occupations into the prompt to examine the out-

Tokens for Correct Translation	Tokens for Occupations							
		Male: analist	Male: a	Male: financier	Male: o	Female: analist	Female: a	Female: financier
	mi	1.69954E-07	1.10153E-05	1.83306E-07	2.74578E-05	1.69954E-07	1.10153E-05	1.83306E-07
	amigo	2.55705E-06	2.68703E-06	1.18563E-08	7.49929E-06	2.55705E-06	2.68703E-06	1.18563E-08
	es	7.28923E-06	1.19383E-05	5.65389E-06	1.1722E-05	7.28923E-06	1.19383E-05	5.65389E-06
	un	0.519850492	9.83368E-07	0.001786403	4.10486E-06	0.519850492	9.83368E-07	0.001786403
	analist	6.02712E-07	0.931759119	3.47089E-06	0.00086118	6.02712E-07	0.931759119	3.47089E-06
	a	1.16286E-05	2.85769E-05	0.926101804	7.64665E-06	1.16286E-05	2.85769E-05	0.926101804
	financier	3.04101E-07	0.316863537	3.51373E-06	0.551219165	3.04101E-07	0.316863537	3.51373E-06
	o	4.47227E-07	3.27988E-06	2.12387E-06	2.37641E-06	4.47227E-07	3.27988E-06	2.12387E-06
	.	3.19549E-07	2.17982E-06	5.73051E-07	2.01175E-06	3.19549E-07	2.17982E-06	5.73051E-07

Figure 1: Example of Misalignment caused by Beam Search

Hint Level	Translation Scenario
No Hint	My friend is a [occupation].
Subtle Hint	I was previously a policeman/policewoman . I am now a [occupation].
	As a busy [occupation] and a father/mother of two, I want to spend as much time with my kids as possible during the holiday.
	My nephews always call me their favourite uncle/aunt . This is because I work as a [occupation].
	I am male/female . I am a [occupation].
	I am the first son/daughter in my family to be a [occupation].
Obvious Hint	He/She is a [occupation].

Table 1: Translation scenarios of No Hint, Subtle Hint, and Obvious Hint Level

put distribution of masculine and feminine translation. A balanced ratio of 50:50 would indicate that the examined model is not biased towards any particular gender. Conversely, a skewed distribution would suggest the presence of gender bias.

For the subtle hint level, we test each scenario for our 26 occupations, with both male and female indicators. When a prompt contains words such as “father” or “uncle”, we expect the model to translate the occupation into its male form. Similarly, words like “mother” or “aunt” should suggest a female translation. To neutralise the potential bias that one hint might provide, we use 5 separate hints to obtain a balanced ratio for the prompting.

Similar to the subtle hint level, we expect each occupation to be translated according to the gender specified in the prompt for the obvious hint level. Thus, we anticipate the highest translation accuracy and output probability in this case as it contains an explicit gender designation “he” or “she”.

To measure the bias, we use binary cross entropy loss (BCEl, formulated as equation 1 (Shannon, 1949)).

$$\text{BCE}(\hat{p}; p) = -\frac{1}{n} \sum_i p_i \ln \hat{p}_i + (1 - p_i) \ln(1 - \hat{p}_i) \quad (1)$$

3.3 Heuristics and Assumptions

Beam search is mostly used in translation models to maximise the overall score (or probability)

of the generated sequence when trying to optimise their output quality (Lemons et al., 2022).

We observe that the algorithm’s focus on maximising the overall sequence score results in the selection of individual tokens that do not align perfectly with the most likely or semantically coherent options in certain contexts. An example of this observation of misalignment is Figure 1. We can see that mBART inserted the token “un” when the highest score is attained by the token “analist”.

To mitigate this issue, we use the whole sequence of token probabilities that are different in male and female occupations to compute the ratio. Using Figure 1 as an example, we sum all probabilities in column “Male: o” and “Female: a” to compute a ratio as these are the only differing tokens in the correct translation of the occupation.

By doing so, we assume that the probabilities of irrelevant tokens are insignificant in the final computed ratio. This is evident from the example “my friend is a teacher”, where both the manually and heuristically computed ratios are approximately 0.999 toward the male gender, with a negligible difference of 1×10^{-5} .

Yet, this assumption is also invalid in some edge cases like Figure 1 where we translate “my friend is a financial analyst”. Token “a” is used in both male and female acronyms of the Spanish translation, which causes the computed ratio to be inaccurate. Yet, as Figure 1 is an edge case, we can

conclude that this change is insignificant to the overall comparison of occupations.

4 Results

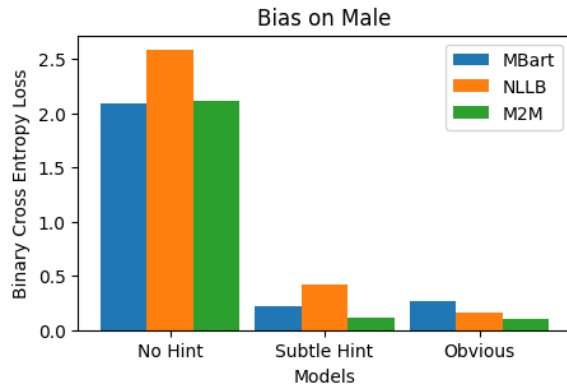


Figure 2: Binary Cross Entropy Loss for mBART, NLLB, and M2M models at varying levels of hints indicating male gender

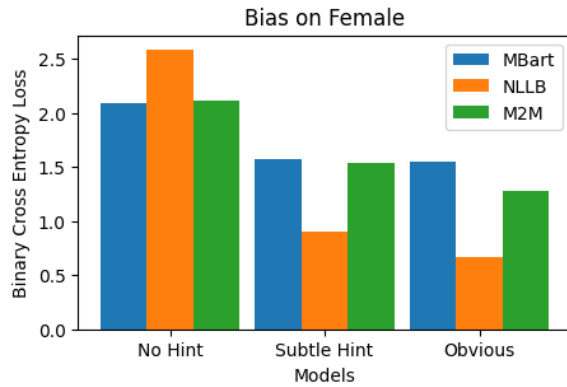


Figure 3: Binary Cross Entropy Loss for mBART, NLLB, and M2M models at varying levels of hints indicating female gender

We will start with a high-level overview of the performance of the models and then go into details for each specific one. Figures 2 and 3 show that in the absence of gender-specific hints, all three models are heavily biased towards male acronyms, with NLLB performing the worst. Once subtle hints are introduced, the results improve but not significantly: the models are still biased towards male connotations; mBART and M2M perform in a similar fashion - both have a high BCEL when faced with female-hinting prompts and a low one with male-hinting ones. NLLB picks up on hints better than its peers and is less biased towards male occupations, yet has a higher BCEL with them; this suggests that it might just be trying to balance out the two genders, irrespective of hints. At the

obvious hint level, the models’ performance improves again. NLLB and M2M improve in regard to both genders, with the former having greater progress in female BCEL and the latter in the male one. mBART’s behaviour almost does not change in the female case and slightly worsens in the male one, which is peculiar since we would expect an LLM to exhibit better behaviour when exposed to a direct statement of gender.

In order to make the graphs and the discussed statistics more tangible, we will use the examples of “nurse” and “musician” occupations consistently throughout this section. As we will further notice, these are also the occupations that all three models are heavily biased towards.

4.1 NLLB

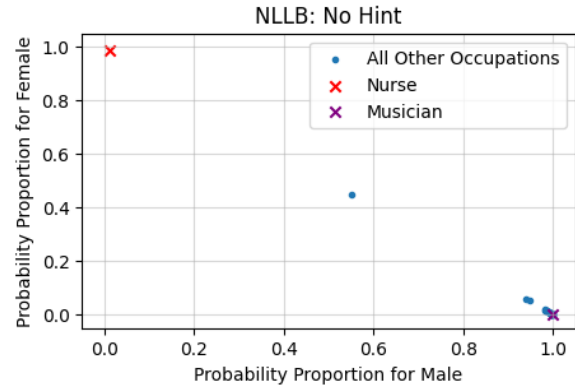


Figure 4: NLLB: Scatter plot of gender probability proportions with no hint

From the above overview, we have observed that NLLB has the best performance with respect to picking up on gender hints. From Figure 4, we can deduce that, as stated above, in the absence of hints, the model is heavily skewed towards the male connotations. Here, we have the probability of “nurse” being translated as a female occupation, which is 0.988, and that of “musician” being translated as a male occupation, which is 0.999.

Now, when subtle hints are introduced (see Figure 5), the distribution changes. When the sentence subject to be translated indicates that the given occupation should be female-gendered, the probability of the correct translation increases. However, one can notice that it is still heavily biased - in the ideal scenario, with a hint present, we expect the translated occupation to reflect it. Yet, in the presented graph, we still see a notable proportion of occupations lacking such coherence. In our recurring example, “musician” has a 0.964

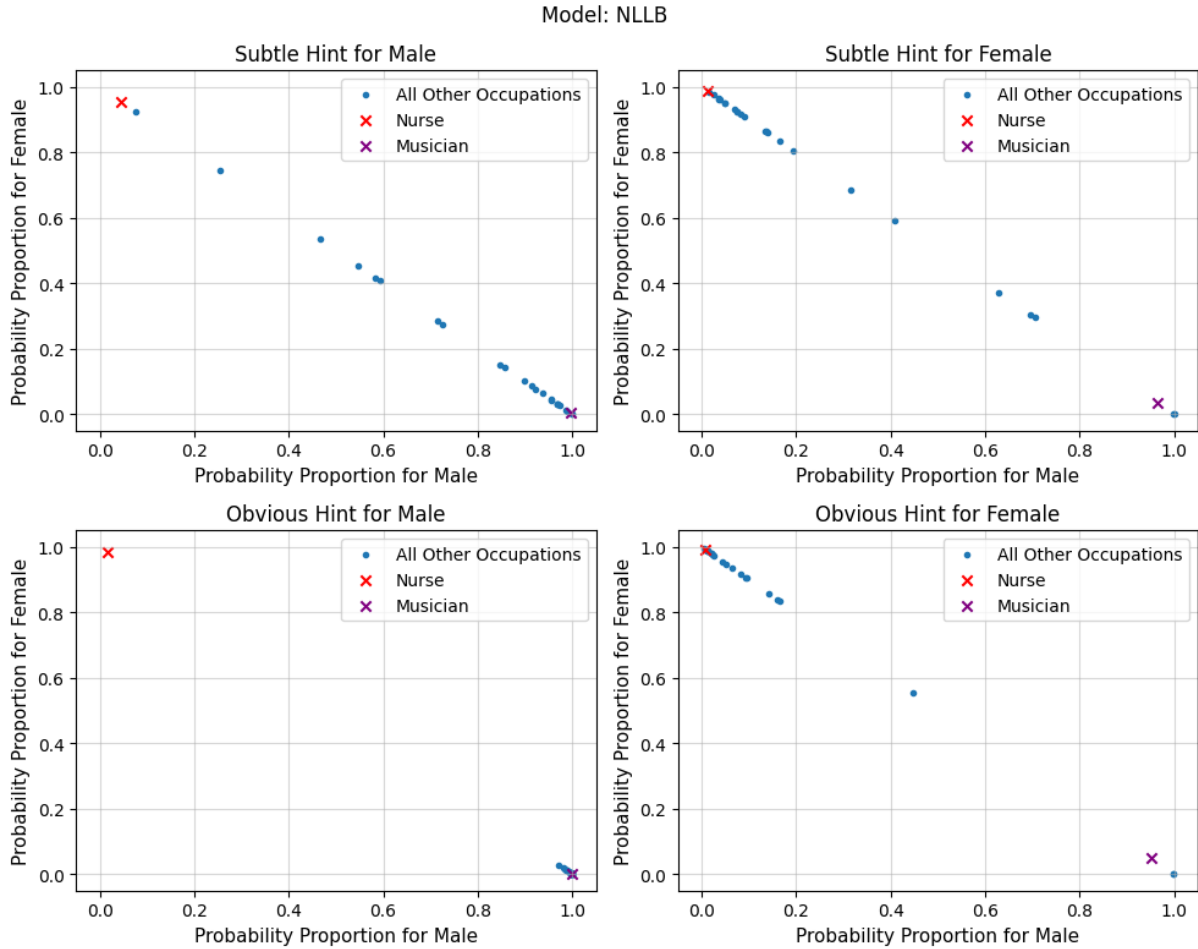


Figure 5: NLLB: Scatter plots of gender probability proportions with subtle hint and obvious hint

chance of being translated as male, even in the presence of subtle female hints. Nonetheless, the improvement is notable compared to the “no hint” scenario.

In the presence of hints suggesting male connotations, we see almost a mirrored graph of that with female aiming hints. Even in this case, “nurse” still has a 0.954 chance of being translated as female. Although the majority of translations have a high probability of being correct, we still see a considerable proportion of occupations appearing to be of an incorrect gender.

With obvious hints present, NLLB performs best out of all three models. When a sentence contains a direct male gender cue (see Figure 5), almost all occupations have a close to 1 probability of being translated correctly, with very few outliers present; “nurse” is an example of such an outlier, with a 0.982 probability to still appear female.

The corresponding performance in the case of obvious female hints is slightly worse than that for male hints yet much better compared to the subtle

hint level. Most occupations have a 0.8 to 1 probability of the correct translation, with some outliers being heavily biased towards male acronyms; one such outlier is “musician” having a 0.951 probability of still being translated as a male occupation.

4.2 M2M

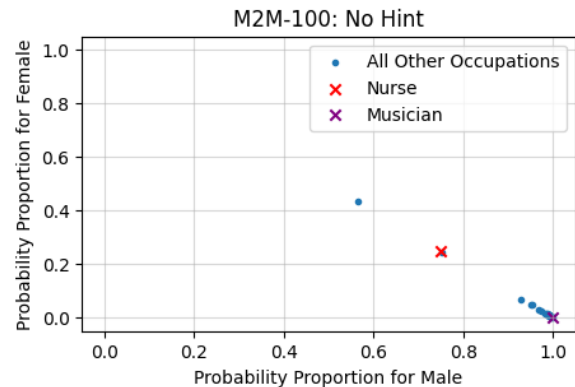


Figure 6: M2M-100: Scatter plot of gender probability proportions with no hint

From Figure 3, we can observe that M2M has one of the worst biases towards translating occupations to their male variant instead of their female variant, even when provided with not only subtle hints but also obvious ones.

Generally M2M had a probability of translating “nurse” as female of 0.339 and “musician” to be translated as a male occupation 0.995.

Figures 6 and 8 similarly show the distributions of predictions across the model having no hints, subtle hints and obvious hints.

As can be seen particularly in Figure 8 when subtle hints are introduced towards a female translation, although the distribution changes slightly from having some context, the translations still mostly come out as the male variant as that is the higher probability. Using our examples of translating “nurse” and “musician”. “nurse” although stereotypical being favoured towards being translated to female, and additionally given a subtle hint, still had a 0.718 probability of incorrectly being translated to male. “musician” had a 0.966 probability of being translated to a male variant, again even with subtle hints. Figure 8 also displays how with more obvious hints, quite a few occupations have successfully moved towards being translated as the female variant. However, as can be seen, a lot of occupations are still being translated as male, which is incorrect.

4.3 mBART

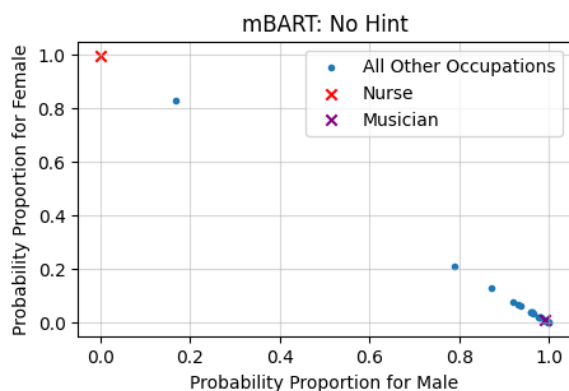


Figure 7: mBART: Scatter plot of gender probability proportions with no hint

Figures 7 and 9 show mBART’s distribution of gender probability proportion with no hints as well as subtle and obvious hints indicating each gender.

When compared to the other model plots, we see that whilst mBART does not perform as well as NLLB, it outperforms M2M. mBART performs

better in the presence of obvious hints indicating either of the genders than with subtle hints showing more points within the 0.8 and 1 range on the target axis in Figure 9. The subtle hints case shows more instances of translations to both target classes for both male and female hints compared to the no hint case. However, there were few points showing a probability of higher than 80% for female translations when the hints indicated female gender. With obvious hints, there was a slight increase in these instances, considering the obvious prompts explicitly stated the target class was female. mBART performed poorly without hints as there were very only two points indicating a high probability of the occupations being translated as female occupations in Figure 9.

The probability of mBART translating “nurse” and “musician” to male with no hints were 0.002 and 0.991, respectively. When a subtle hint indicating male gender is added, the probability of “nurse” being translated as male increases significantly to 0.301, but when an obvious hint indicating male gender is added, the probability decreases to 0.002. The probability in the obvious male gender case is shown as an anomaly in Figure 9, and the increase in errors from subtle hints to obvious hints is also something only seen in the mBART model according to Figure 3. When a subtle hint indicating female gender is added, the probability of “musician” being translated as a male occupation reduces slightly to 0.975, and when an obvious hint indicating female gender is added, the probability reduces further to 0.938, which is still a high probability given the input was a female occupation.

5 Discussion

We shall now interpret the obtained results in the light of both the overall bar charts appearing at the beginning of the previous section (Figures 2 and 3) and the specific more detailed graphs corresponding to each of the three analysed models (Figures 4, 6, 5, 7, 8 and 9). We have noted earlier that from Figures 2 and 3, NLLB seems to be the best-performing model, by a large margin, with respect to picking up on hints. Next, we have M2M performing slightly better than mBART. Now, from the model-specific graphs, we still have that NLLB has the most accurate translations among the three models, even if not exactly ideal. Interestingly enough, we observe an opposite to the overall

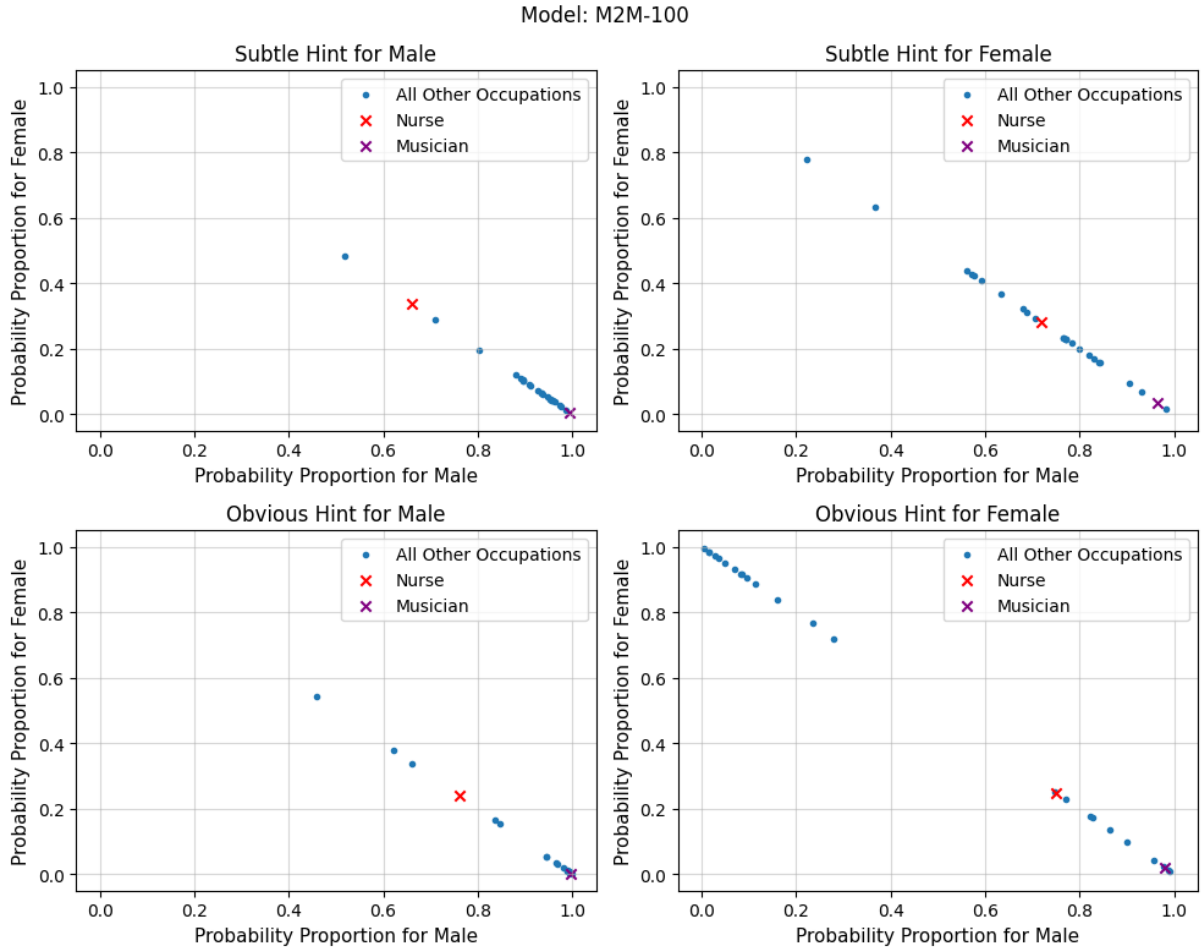


Figure 8: M2M-100: Scatter plots of gender probability proportions with subtle hint and obvious hint

trend in the bar chart for M2M and mBART - according to the detailed graphs, mBART is a better performing model out of these two. When using mBART, more occupations have a higher probability of being translated into a corresponding gender for both subtle and obvious hint levels; this outperformance is preserved for both genders.

A natural question arises: Why do the bar charts and graphs demonstrate seemingly opposing results? Given the nature of the BCEL function used to calculate the BCEL for the models, the “further away” the translation is from the desired one, the more such output is penalised. Hence, when a hinted-to-be-female occupation is translated into its male version with a probability closer to 1, BCEL is much higher than if such probability was, say, slightly above 0.5; in other words, the change is not linear. Hence, the BCEL is good at measuring how biased a particular model is but not necessarily suitable for the model comparison per se.

Thus, we shall refer to the detailed graphs to establish the least and most biased models. As stated

earlier, NLLB exhibits the least bias; it is not perfect, especially in the “no hint” case, but it is good at grasping the context and adapting the occupation’s gender accordingly. We then have mBART as the second-best performing model and M2M as the most biased one. The latter translates most occupations into their male versions with higher than 0.5 probabilities irrespective of context, even when the hint level is “obvious”. The model is heavily biased towards male connotations; further, it is the only model out of the three that has no occupations having a greater than 0.5 probability of appearing female in the absence of hints.

6 Conclusion

In this paper, we have evaluated three LLMs on how gender-biased they are. We have concluded that NLLB exhibits the least bias in the presence of gender-indicating hints and that M2M is the most male acronym-biased model both in the presence and absence of hints. Another contribution of this

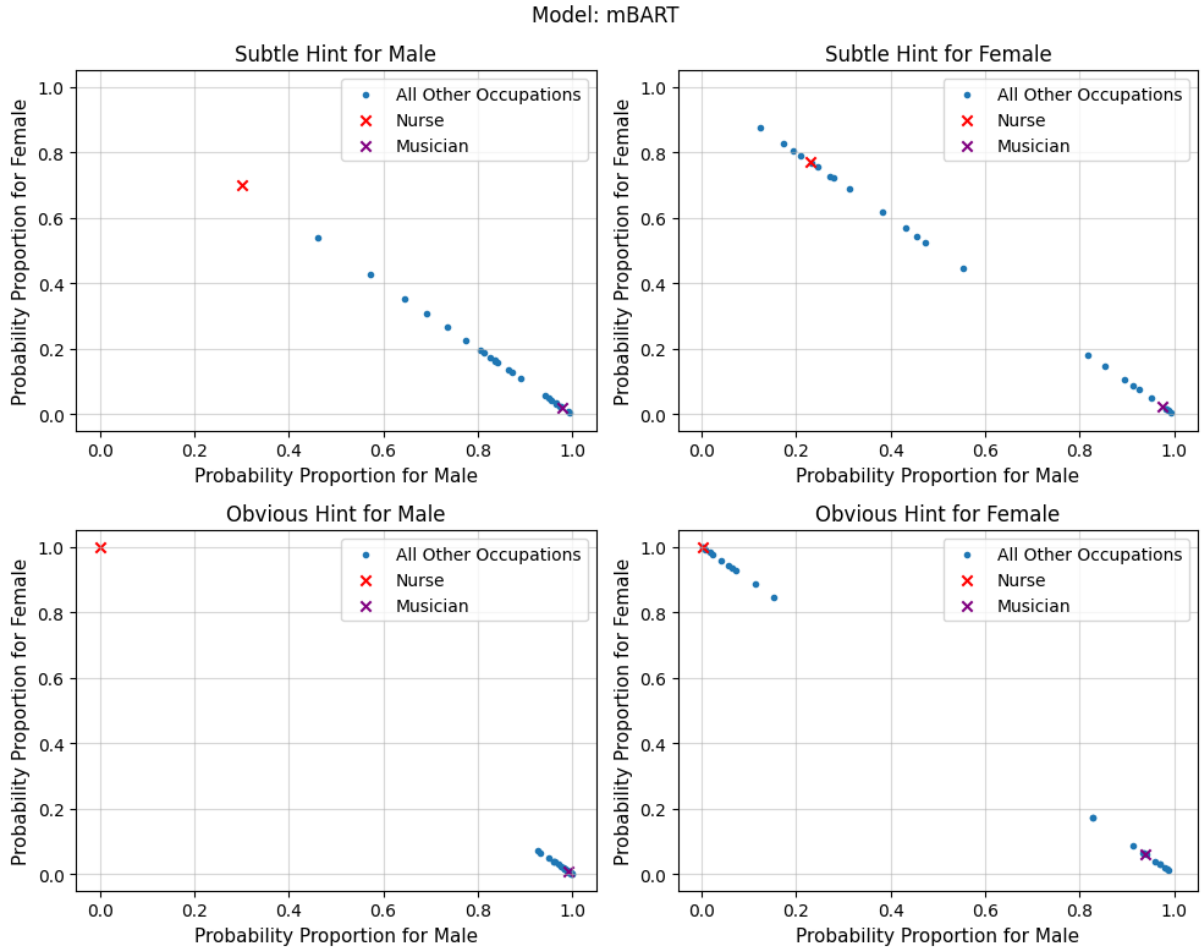


Figure 9: mBART: Scatter plots of gender probability proportions with subtle hint and obvious hint

work is demonstrating that binary cross-entropy loss is not a suitable means of comparison between LLMs. It is a good measure of how biased each separate model is, but due to the function’s nature of assigning weight based on the “distance” from the desired value non-linearly, it is not good at relating models’ behaviours to one another.

7 Limitations

This study required us to create a dataset manually, limiting its size and language pair choices. Exploring gender bias in translation between English and Slavic languages would require time-consuming manual checks due to the complexity of the syntactic and grammatical structure of said languages. However, if one has the resources and availability of people fluent in the desired languages, it is a direction worth exploring in future work. Further on, fully genderless languages such as Finnish or Estonian that do not even have gender pronouns could be investigated to determine whether they

present gender bias when translated into gendered languages.

Our method could benefit from considering that some languages with gendered nouns default to using one gender’s noun irrespective of the subject’s gender. Examples of this are the common use of the word “actor” in English for both actors and actresses and the use of the word “professor” as a gender-neutral term for professors even though it is in the masculine form. We propose that an extension of our study includes both contexts relating to the frequency and normalcy of the use of each gendered noun and the sentiment it’s associated with. Cultural norms and linguistic practices of the speakers of the languages are needed to avoid drawing misleading conclusions, and our future work will look at this.

With limited computing resources in this work, we were only able to run relatively small models. Our results may not be reflective of all LLMs. Future works should consider running this with larger models to increase the validity of the results.

References

Christopher M. Bishop. 2006. *Pattern Recognition and Machine Learning*. Springer.

Angela Fan, Shruti Bhosale, Holger Schwenk, Zhiyi Ma, Ahmed El-Kishky, Siddharth Goyal, Mandeep Baines, Onur Celebi, Guillaume Wenzek, Vishrav Chaudhary, et al. 2020. Beyond english-centric multilingual machine translation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics.

Batya Friedman and Helen Nissenbaum. 1996. *Bias in computer systems*. *ACM Trans. Inf. Syst.*, 14(3):330–347.

Ian Goodfellow, Yoshua Bengio, and Aaron Courville. 2016. *Deep Learning*. The MIT Press.

Gemini Team Google. 2023. *Gemini: A family of highly capable multimodal models*.

Amy Kraft. 2016. *Microsoft shuts down ai chatbot after it turned into a nazi*.

Sofia Lemons, Carlos Linares López, Robert C. Holte, and Wheeler Ruml. 2022. *Beam search: Faster and monotonic*.

Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. 2020. *Multilingual denoising pre-training for neural machine translation*.

Tiago Palma Pagano, Rafael Bessa Loureiro, Fernanda Vitória Nascimento Lisboa, Gustavo Oliveira Ramos Cruz, Rodrigo Matos Peixoto, Guilherme Aragão de Sousa Guimarães, Lucas Lisboa dos Santos, Maira Matos Araujo, Marco Cruz, Ewerton Lopes Silva de Oliveira, Ingrid Winkler, and Erick Giovani Sperandio Nascimento. 2022. *Bias and unfairness in machine learning models: a systematic literature review*.

Marcelo O. R. Prates, Pedro H. C. Avelar, and Luis Lamb. 2019. *Assessing gender bias in machine translation – a case study with google translate*.

Satyendra Singh Rawat and Amit Kumar Mishra. 2022. *Review of methods for handling class-imbalanced in classification problems*.

Beatrice Savoldi, Marco Gaido, Luisa Bentivogli, Matteo Negri, and Marco Turchi. 2021. *Gender Bias in Machine Translation*. *Transactions of the Association for Computational Linguistics*, 9:845–874.

Claude Elwood Shannon. 1949. *A mathematical theory of communication*. University of Illinois Press.

NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semaarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. *No language left behind: Scaling human-centered machine translation*.

Code

All code that is used for this study has been deposited in the following GitHub repository: <https://github.com/RcwYuen/ucl-comp0087-snlp-cw>.