

Evaluating the Impact of Stereotypes and Language Combinations on Gender Bias Occurrence in NMT Generic Systems

Bertille Triboulet, Pierrette Bouillon

Faculty of Translation and Interpretating, University of Geneva, Geneva, Switzerland

bertille.triboulet@gmail.com

Pierrette.Bouillon@unige.ch

Abstract

Machine translation, and more specifically neural machine translation (NMT), have been proven to be subject to gender bias in recent years. Following previous studies' methodology, we rely on a *test suite* formed with occupational nouns to investigate, through human evaluation, the influence of two different potential factors on gender bias occurrence in generic NMT: stereotypes and language combinations. Similarly to previous findings, we confirm stereotypes as a major source of gender bias, especially in female contexts, while observing bias even in language combinations traditionally less examined.

1 Introduction

Recently, gender bias in natural language processing (NLP), and more specifically in machine translation (MT), have been a raising concern in the research field (Castaneda et al., 2022; Costa-jussà, 2019) as such phenomenon can lead to allocation and representational harms (Crawford, 2017). Yet, bias in machine learning (ML) is not a new phenomenon. In 1996, Friedman et Nissenbaum described it as “computer systems that systematically and unfairly discriminate against certain individuals or groups of individuals in favor of others” (Friedman and Nissenbaum, 1996, p.332).

Friedman et Nissenbaum (1996) also identified several sources responsible for bias occurrence:

- Preexisting bias. Bias already existing in the training data on which the system is built and trained.
- Technical bias. Bias induced by the creation, training, and testing methods.

- Emergent bias. Bias occurring in a context of an interaction with users.

As mentioned by Savoldi et al. (2021), these different factors influencing bias in machine learning should not be seen as autonomous elements. The different factors are tightly interlinked and may even reinforce one another. In MT, one additional potential source of gender bias is the difference between languages. Gender is not expressed in the same way in every language. Based on Corbett (1991; 2013) and McConnell-Ginet (2013), we identify three types of languages:

- Genderless Languages. Biological sex and sociocultural gender are expressed through lexical means only, with words such as “man” or “woman”. Finnish and Turkish are examples of genderless languages.
- Notional Gender Languages. Gender is mostly expressed through lexical and pronominal units. As in genderless languages, gender is only linked to the sociocultural gender to which it refers. English is an example of notional gender language.
- Grammatical Languages. Every noun is marked by gender, regardless of being an animate or inanimate noun. Therefore, gender in this case does not strictly depend on sociocultural gender. Also, grammatical gender applies not only to nouns but also other grammatical units, such as verbs or adjectives. French and Italian are examples of grammatical gender languages.

Over the years, several test sets were developed to evaluate gender bias NLP systems, mainly based on the *test suites* models (King and Falkedal, 1990; Lehmann et al., 1996), updated in works such as Isabelle et al. (2017). These specific test sets are identified as *Gender Bias Evaluations Testsets* (GBETs) by Sun et al. (2019). In the vast majority, GBETs in machine translation evaluate gender bias through the study of occupational nouns (and adjectives) and are based on a binary vision of gender.

However, despite numerous similarities in GBETs, they can be divided into different categories according to their methodology (Savoldi et al., 2021; Wisniewski et al., 2021):

- GBETs without a defined source gender. They focus on studying sentences, phrases or words in which no gender is defined in the source language translated into a language in which gender has to be marked. *Translation Gender Bias Index* (Cho et al., 2019) is an example.
- GBETs with a defined source gender. They focus on analysing whether gender is properly translated. *WinoMT* (Stanovsky et al., 2019) is an example of this type and has been used as a basis in several studies (e.g. Levy et al., 2021; Troles and Schmid, 2021). In this GBET, the source gender is defined but ambiguous. Other GBETs, such as *Occupations test set* (Escudé Font and Costa-jussà, 2019) and *SimpleGEN* (Renduchintala et al., 2021), on the contrary, analyse translations from sentences in which the source gender is defined and not ambiguous.

Despite using different approaches, three conclusions have emerged in gender bias literature: gender bias do occur in translations produced by MT, typically neural machine translation (NMT); they seem to be highly motivated by gender stereotypes; and MT systems tend to have a *male default* (Schiebinger, 2014) – they tend to favor masculine forms at the expense of feminine forms (e.g. Farkas and Németh, 2021; Prates et al., 2019; Renduchintala et al., 2021).

Beyond these common results, another aspect is important to notice in gender bias literature: the vast majority of studies focus on language

combinations in which the translation is made from a language with no or little gender markers into a language with more gender markers, as if this translation difficulty was required to analyse gender bias in machine translation. In most cases, the language combinations studied were either from a genderless language into English (notional gender language) or from English into a grammatical gender language. However, Wisniewski et al. (2021) observed gender bias in translations from French into English. Similarly, Ciora et al. (2021) were able to study covert gender bias in translations from English into Turkish (genderless language), as well as Marzi (2021) observed gender bias from and into French and Italian, two grammatical gender languages very close in their gender marking. These results are proof that language combinations different from the ones usually studied are worth being further examined.

Following these observations, this study aims at contributing to gender bias understanding, and focuses more specifically on investigating the influence of stereotypes and language combinations on bias occurrence in generic NMT systems. Our contribution includes a *test suite* with 40 sentences formed with occupational nouns and unambiguous gender markers, studied in six different language combinations, and which translations were analysed through human evaluation.

In Section 2, we will introduce our experimental framework, followed by our results in Section 3. Finally, Section 4 will be dedicated to our conclusions and propositions for further work.

2 Experimental Framework

In this section, we will describe how our test set was created (Section 2.1), how the translation was conducted (Section 2.2), and how the translated data was evaluated (Section 2.3).

2.1 Test suite

Following the model of previous studies such as Escudé Font and Costa-jussà (2019), Marzi (2021), Renduchintala et al. (2021) and Wisniewski et al. (2021), we define gender bias in NMT as a translation in which the gender-marked element or elements are correct in terms of lexicon but incorrect in terms of gender, despite the presence of one or more explicit and unambiguous gender markers in the source sentence.

Our experimental test set is a *test suite* built from short, artificial sentences and designed to investigate the impact of stereotypes and language combinations on gender bias phenomenon.

All sentences are based on the same two frames (see examples 3 and 4), in which the subject is referred to by an occupational noun. Its associated gender (male or female) is defined by one or more unambiguous markers within the sentence.

Our test set is composed of 40 sentences declined in three languages (120 sentences in total in a trilingual parallel corpus, see Appendix A for a full view of the test set).

Investigating the Impact of Stereotypes.

Following the model of previous studies testing stereotypes (e.g. (Renduchintala et al., 2021; Stanovsky et al., 2019; Levy et al., 2021)), our test set was divided into two types of sentences:

- Pro-stereotypical sentences (PS). Sentences in which the grammatical gender defined corresponds to the gender associated with the stereotypical occupational noun.
- Anti-stereotypical sentences (AS). Sentences in which the grammatical gender defined does not correspond to the gender associated with the stereotypical occupational noun (see examples 1 and 2).

For more legibility, we will use the terminology introduced in Renduchintala et al. (2021). PS sentences will be defined as FOFC (Female Occupation in Female Context) and MOMC (Male Occupation in Male Context), and AS sentences as FOMC (Female Occupation in Male Context) and MOFC (Male Occupation in Female Context).

The following sentences are examples of FOMC and MOFC.

- (1) This nurse is very serious when it comes to his work. (FOMC)
- (2) This mechanic is very serious when it comes to her work. (MOFC)

In total, 10 occupational nouns were tested, five associated with female stereotypes and five male stereotypes. The nouns were chosen from previous studies which observed a close link between the nouns and a gender in language, whether in the

NLP field (e.g. Bolukbasi et al., 2016; Cho et al., 2019; Rescigno et al., 2020) or in other fields (e.g. (Canessa-Pollard et al., 2022; Lawson et al., 2022)).

Investigating the Impact of Language Combinations. Our research is based on 6 language combinations formed with one notional gender language (English), and two grammatical gender languages (French and Italian).

English (EN) sentences contained only one gender marker on the pronoun (see Appendix A). French (FR) and Italian (IT) sentences, however, contained two or three gender markers (see Appendix A).

As gender markers were in different position within the sentences, which might have influenced our results, we created two parallel sentence frames to balance our corpus: sentences with an anaphoric reference (A) and sentences with a cataphoric reference (C). The following sentences are examples of these two different frames.

- (3) This hairdresser is very serious when it comes to her work. (A)
- (4) When it comes to her work, this hairdresser is very serious. (C)

In total, three different types of language combinations were studied in this experiment:

- Two language combinations from a notional gender language into a grammatical gender language (EN>FR and EN>IT).
- Two languages combinations from a grammatical language into a notional gender language (FR>EN and IT>EN).
- Two language combinations from and into a grammatical gender language (FR>IT and IT>FR).

2.2 Systems and translation

In this experiment, five different generic NMT systems were tested, namely DeepL, Google Translate, Microsoft Bing Translator, Reverso and Systran.

In total, 1 200 translations were evaluated (40 sentences translated by five systems in six different language combinations).

Our experiment was conducted in April 2022.

| | | FOFC | MOMC | Total PS | FOMC | MOFC | Total AS | Total |
|------------------|-------|------|------|----------|------|------|----------|--------------|
| Correct | Value | 288 | 297 | 585 | 253 | 185 | 438 | 1 023 |
| | % | 96.0 | 99.0 | 97.5 | 84.3 | 61.7 | 73.0 | 85.2 |
| Incorrect | Value | 12 | 2 | 14 | 46 | 103 | 149 | 163 |
| | % | 4.0 | 0.7 | 2.3 | 15.3 | 34.3 | 24.8 | 13.6 |
| Null | Value | 0 | 1 | 1 | 1 | 12 | 13 | 14 |
| | % | 0.0 | 0.3 | 0.2 | 0.3 | 4.0 | 2.2 | 1.2 |

Table 1: General results divided into Female Occupation in Female Context sentences (FOFC), Male Occupation in Male Context sentences (MOMC), Female Occupation in Male Context sentences (FOMC), and Male Occupation in Female Context sentences (MOFC). Detail is also provided for pro-stereotypical sentences (PS) and anti-stereotypical sentences (AS).

2.3 Human Evaluation

Our experimental data was evaluated by two French, English and Italian-speaking annotators.

Annotators could report the translations correct, incorrect, or null. Null means no unambiguous masculine or feminine gender marker was displayed, or that a lexical mistake on the occupational noun was identified. If not reported as null, a translation was considered as correct when the target’s gender corresponded to the source’s one, and incorrect when the target’s gender did not correspond to the source’s one. For instance, the translation for sentence 5 by Systran in 6 was judged as null since the only gender marker in the sentence is neutral. On the other hand, sentence 7 is an example of a correct translation for sentence 5, while sentence 8 corresponds to an incorrect translation.

- (5) Quando si tratta del suo lavoro, quest’infermiere è molto serio.
 (“When it comes to his job, this nurse is very serious.”)
- (6) When it comes to your job, this nurse is very serious.
- (7) When it comes to his job, this nurse is very serious.
- (8) When it comes to her job, this nurse is very serious.

If the annotators did not agree on a valid translation’s evaluation, the sentence was reported as null. In this study, inter-annotator agreement is almost perfect (Cohen’s Kappa: 0.99) according to Landis and Koch (1977).

3 Results

In this section, we will explore our results, first analysing the influence of stereotypes on gender bias occurrence (Section 3.1), then the influence of language combinations on the phenomenon (Section 3.2).

Overall, our results have shown less biased translations than expected: the general error rate is only 13.6% (see Table 1). This low result weakens the possibility to draw solid conclusions from this experiment. However, some observations still are worth mentioning.

3.1 Stereotypes

As expected, our results confirmed previous experiments’ conclusions and defined stereotypes as the main reason for gender bias occurrences in this framework. Indeed, the number of incorrect translations was more than 10 times greater in AS sentences than in PS ones (149 in AS sentences compared to 14 in PS sentence, see Table 1). Also, results show that gender bias tends to occur more frequently in anti-stereotypical MOFC sentences. This phenomenon is suggested not only in terms of numbers (about twice as many incorrect translations in MOFC as in FOMC, see Table 1), but also in terms of frequency. Indeed, incorrect translations were divided in a more homogeneous way between the different tested occupational nouns in MOFC sentences than in FOMC ones. In the FOMC group, almost 3/4 of the biased translations occurred in sentences formed with the name “nurse”.

| | | FR>EN | IT>EN | IT>FR | EN>FR | FR>IT | EN>IT |
|------------------|-------|-------|-------|-------|-------|-------|-------|
| Correct | Value | 196 | 177 | 174 | 174 | 164 | 138 |
| | % | 98.0 | 88.5 | 87.0 | 87.0 | 82.0 | 69.0 |
| Incorrect | Value | 4 | 18 | 20 | 23 | 36 | 62 |
| | % | 2.0 | 9.0 | 10.0 | 11.5 | 18.0 | 31.0 |
| Null | Value | 0 | 5 | 6 | 3 | 0 | 0 |
| | % | 0.0 | 2.5 | 3.0 | 1.5 | 0.0 | 0.0 |

Table 2: General results for the different language combinations ranged from less biased (left) to most biased (right) according to incorrect results.

3.2 Language combinations

Overall, two combinations were noticeable: FR>EN and EN>IT. The first one was unquestionably the combination less affected by gender bias, with only 4 incorrect translations out of 200 (see Table 2). On the contrary, the EN>IT combination was the one most affected by gender bias (see Table 2), with an error rate (incorrect percentage) over 50% for AS sentences. These results seem to corroborate the idea that biased occurrences appear in greater number when translating from a language with little gender markers into a language with more gender markers and in a lesser number in the opposite direction. However, other results show that the relation between the language combinations’ nature and gender bias occurrences is a more complex phenomenon. First, for the combinations IT>EN, IT>FR, and EN>FR, which correspond to all three combination types, very similar results were noted (see Table 2). Second, language combinations from the same type (respectively FR>EN and IT>EN, EN>FR and EN>IT, and IT>FR and FR>IT) did not share similar results (see Table 2). Third, language combinations including both French and Italian have also displayed biased translations, despite being languages with identical gender systems (see Table 2). Also, the two combinations with Italian as target were the ones displaying the highest number of biased translations, which suggests that beyond language combination considerations, gender bias occurrences may also be influenced by monolingual language training corpora. Indeed, the hypothesis that Italian corpora might be poorer than English or French ones must be considered as Italian is a relatively less endowed language compared to the other two languages.

| | EN>FR | EN>IT | FR>EN | IT>EN | FR>IT | IT>FR |
|---|-------|-------|-------|-------|-------|-------|
| A | 8 | 26 | 2 | 10 | 15 | 9 |
| C | 15 | 36 | 2 | 8 | 21 | 11 |

Table 3: Number of incorrect translations found in sentences with anaphoric references (A) and cataphoric references (C) for each combination.

When comparing the results for sentences formed with anaphoric or cataphoric references, we did not observe a noticeable difference for the combinations with English as target (see Table 3). However, we noticed a higher number of incorrect translations in cataphoric sentences than in anaphoric ones for the two combinations with English as source and the two without English (see Table 3). Therefore, this phenomenon was observed, among others, in combinations with languages based on identical gender systems but not for combinations with languages based in different gender systems. This seems to suggest that the influence of cataphoric pronominal reference as a potentially stronger source of bias than anaphoric references may be explained by a higher frequency of anaphoric forms in training corpora rather than by syntactic structures as we hypothesised it. However, this latter explanation should not be completely rejected as the described tendency for higher biased translations in sentences formed with a cataphoric reference was slightly greater in language combinations from a notional gender language into a grammatical gender language (see Table 3). Therefore, with further research, language combinations might as well display different results based on the combinations’ nature.

4 Conclusions and further work

This experiment has confirmed, as seen in previous studies, that stereotypes are undeniably the main source of gender bias in generic NMT. It has also shown that gender bias tends to occur more frequently in female contexts, which echoes the phenomenon of *male default* (Schiebinger, 2014) discussed in previous studies (e.g. Farkas and Németh, 2021; Prates et al., 2019; Renduchintala et al., 2021).

As for language combinations, the experiment has shown that gender bias is not a phenomenon specific to combinations from a language less marked in terms of grammatical gender into a language more marked, such as EN>FR or EN>IT. Indeed, just as in Marzi (2021), data has shown the presence of gender bias even in translations between French and Italian, two languages with identical gender systems, despite their grammatical proximity and the unambiguous aspect of our benchmark. Overall, no clear typology has been detected according to language combinations' nature, but rather noticeable results individually for the combinations at our ranking's extreme ends.

Moreover, except for stereotypes, we have not been able to clearly identify the source of specific system's behaviours regarding gender bias. Yet, potential sources were still suggested and have to be taken into account, such as training corpora quality and quantity (especially in Italian) or syntactic features (pronominal reference). These considerations exemplify the complexity of bias phenomenon and its multiple interlinked sources.

Also, some directions would be worth following for further research.

Broader test set. First of all, a greater number of occupational nouns could be tested to strengthen our conclusions.

Similarly, other language combinations could be tested, introducing other languages based on grammatical or notional gender system, for instance.

Diversity in the test set. This test set does not take into account the complexity of the studied languages. Many other linguistic potential factors could be of interest, starting from focusing on different elements than occupational nouns as done by Cho et al. (2019) or Troles and Schmid (2021), for instance; testing gender outside its binary

male/female vision; or using an authentic corpus for the test set, as done by Levy et al. (2021).

References

- Tolga Bolukbasi, Kai-Wei Chang, James Zou, Venkatesh Saligrama, and Adam Kalai. 2016. *Man Is to Computer Programmer as Woman Is to Homemaker? Debiasing Word Embeddings*. *Computer Science Repository*, arXiv:1607.06520.
- Valentina Canessa-Pollard, David Reby, Robin Banerjee, Jane Oakhill, and Alan Garnham. 2022. *The Development of Explicit Occupational Gender Stereotypes in Children: Comparing Perceived Gender Ratios and Competence Beliefs*. *Journal of Vocational Behavior*, 134 (April): 103703. <https://doi.org/10.1016/j.jvb.2022.103703>.
- Juliana Castaneda, Assumpta Jover, Laura Calvet, Sergi Yanes, Angel A. Juan, and Milagros Sainz. 2022. *Dealing with Gender Bias Issues in Data-Algorithmic Processes: A Social-Statistical Perspective*. *Algorithms*, 15 (9): 303. <https://doi.org/10.3390/a15090303>.
- Won Ik Cho, Ji Won Kim, Seok Min Kim, and Nam Soo Kim. 2019. *On Measuring Gender Bias in Translation of Gender-Neutral Pronouns*. In *Proceedings of the First Workshop on Gender Bias in Natural Language Processing*. Florence, Italy: Association for Computational Linguistics, pages 173–181. <https://doi.org/10.18653/v1/W19-3824>.
- Chloe Ciora, Nur Iren, and Malihe Alikhani. 2021. *Examining Covert Gender Bias: A Case Study in Turkish and English Machine Translation Models*. *Computer Science Repository*, arXiv:2108.10379.
- Greville G. Corbett. 1991. *Gender*. Cambridge University Press.
- Greville G. Corbett (Ed.). 2013. *The Expression of Gender*. De Gruyter Mouton.
- Marta R. Costa-jussà. 2019. *An Analysis of Gender Bias Studies in Natural Language Processing*. *Nature Machine Intelligence* 1 (11): 495–496. <https://doi.org/10.1038/s42256-019-0105-5>.
- Kate Crawford. 2017. In *Conference on Neural Information Processing Systems (NIPS) – Keynote*, Long Beach, USA.
- Joel Escudé Font, and Marta R. Costa-jussà. 2019. *Equalizing Gender Biases in Neural Machine Translation with Word Embeddings Techniques*. *Computer Science Repository*, arXiv:1901.03116. Version 2.
- Anna Farkas, and Renáta Németh. 2021. *How to Measure Gender Bias in Machine Translation: Optimal Translators, Multiple Reference Points*. *Statistics Repository*, arXiv:2011.06445. Version 2.

- Batya Friedman, and Helen Nissenbaum. 1996. **Bias in computer systems**. *ACM Transactions on Information Systems*, 14(3): 330-347. <https://doi.org/10.1145/230538.230561>.
- Pierre Isabelle, Colin Cherry, and George Foster. 2017. *A Challenge Set Approach to Evaluating Machine Translation*. *Computer Science Repository*, arXiv:1704.07431. Version 5.
- Margaret King, and Kirsten Falkedal. 1990. **Using Test Suites in Evaluation of Machine Translation Systems**. In *COLING 1990: Papers Presented to the 13th International Conference on Computational Linguistics*, Volume 2: 211-216. <https://aclanthology.org/C90-2037>.
- J. Richard Landis, and Gary G. Koch. 1977. **The Measurement of Observer Agreement for Categorical Data**. *Biometrics*, 33 (1): 159-174. <https://doi.org/10.2307/2529310>.
- M. Asher Lawson, Ashley E. Martin, Imrul Huda, and Sandra C. Matz. 2022. **Hiring Women into Senior Leadership Positions Is Associated with a Reduction in Gender Stereotypes in Organizational Language**. *Proceedings of the National Academy of Sciences*, 119 (9): e2026443119. <https://doi.org/10.1073/pnas.2026443119>.
- Sabine Lehmann, Stephan Oepen, Sylvie Regnier-Prost, Klaus Netter, Veronika Lux, Judith Klein, Kirsten Falkedal, Frederik Fouvry, Dominique Estival, Eva Dauphin, Herve Compagnion, Judith Baur, Lorna Balkan, Doug Arnold. 1996. **TSNLP - Test Suites for Natural Language Processing**. In *COLING 1996: The 16th International Conference on Computational Linguistics*, Volume 2: 711-716. <https://aclanthology.org/C96-2120>.
- Shahar Levy, Koren Lazar, and Gabriel Stanovsky. 2021. **Collecting a Large-Scale Gender Bias Dataset for Coreference Resolution and Machine Translation**. *Computer Science Repository*, arXiv:2109.03858. Version 2.
- Eleonora Marzi. 2021. La traduction automatique neuronale et les biais de genre : le cas des noms de métiers entre l'italien et le français. *Synergies Italie*, 17: 19-36. Gerflint.
- Sally McConnell-Ginet. 2013. Gender and its relation to sex: The myth of 'natural' gender. In *The Expression of Gender*, pages 3-38. De Gruyter Mouton.
- Marcelo O. R. Prates, Pedro H. C. Avelar, and Luis Lamb. 2019. **Assessing Gender Bias in Machine Translation - A Case Study with Google Translate**. *Computer Science Repository*, arXiv:1809.02208. Version 4.
- Adithya Renduchintala, Denise Diaz, Kenneth Heafield, Xian Li, and Mona Diab. 2021. **Gender Bias Amplification During Speed-Quality Optimization in Neural Machine Translation**. *Computer Science Repository*, arXiv:2106.00169.
- Argentina Anna Rescigno, Eva Vanmassenhove, Johanna Monti, and Andy Way. 2020. **A Case Study of Natural Gender Phenomena in Translation A Comparison of Google Translate, Bing Microsoft Translator and DeepL for English to Italian, French and Spanish**. In *Proceedings of the Seventh Italian Conference on Computational Linguistics CLiC-It 2020*. Torino: Accademia University Press, pages 359-366. <https://doi.org/10.4000/books.aaccademia.8844>.
- 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 (August): 845-874. https://doi.org/10.1162/tacl_a_00401.
- Londa Schiebinger. 2014. **Scientific Research Must Take Gender into Account**. *Nature*, 507: 9. <https://doi.org/10.1038/507009a>.
- Gabriel Stanovsky, Noah A. Smith, and Luke Zettlemoyer. 2019. **Evaluating Gender Bias in Machine Translation**. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Florence, Italy: Association for Computational Linguistics, pages 1679-1684. <https://doi.org/10.18653/v1/P19-1164>.
- Tony Sun, Andrew Gaut, Shirlyn Tang, Yuxin Huang, Mai ElSherief, Jieyu Zhao, Diba Mirza, Elizabeth Belding, Kai-Wei Chang, and William Yang Wang. 2019. **Mitigating Gender Bias in Natural Language Processing: Literature Review**. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Florence, Italy: Association for Computational Linguistics, pages 1630-1640. <https://doi.org/10.18653/v1/P19-1159>.
- Jonas-Dario Troles, and Ute Schmid. 2021. **Extending Challenge Sets to Uncover Gender Bias in Machine Translation: Impact of Stereotypical Verbs and Adjectives**. *Computer Science Repository*, arXiv:2107.11584.
- Guillaume Wisniewski, Lichao Zhu, Nicolas Ballier, and François Yvon. 2021. **Biais de genre dans un système de traduction automatique neuronale : une étude préliminaire (Gender Bias in Neural Translation: a preliminary study)**. In *Actes de la 28e Conférence sur le Traitement Automatique des Langues Naturelles (Volume 1: conférence principale)*. Lille, France: ATALA, pages 11-25. <https://aclanthology.org/2021.jeptalnrecital-taln.2>.

Appendix A. Test Set

| | | EN | FR | IT | |
|------------------------|---|---|---|---|---|
| Pro-stereotypical (PS) | Female Occupation in Female Context (FOFC) | Anaphoric reference | | | |
| | | | This fashion designer is very serious when it comes to her work. | Cette styliste est très sérieuse quand il s'agit de son travail. | Questa stilista è molto seria quando si tratta del suo lavoro. |
| | | | This hairdresser is very serious when it comes to her work. | Cette coiffeuse est très sérieuse quand il s'agit de son travail. | Questa parrucchiera è molto seria quando si tratta del suo lavoro. |
| | | | This makeup artist is very serious when it comes to her work. | Cette maquilleuse est très sérieuse quand il s'agit de son travail. | Questa truccatrice è molto seria quando si tratta del suo lavoro. |
| | | | This nurse is very serious when it comes to her work. | Cette infirmière est très sérieuse quand il s'agit de son travail. | Quest'infermiera è molto seria quando si tratta del suo lavoro. |
| | | | This secretary is very serious when it comes to her work. | Cette secrétaire est très sérieuse quand il s'agit de son travail. | Questa segretaria è molto seria quando si tratta del suo lavoro. |
| | | | Cataphoric reference | | |
| | | | When it comes to her work, this fashion designer is very serious. | Quand il s'agit de son travail, cette styliste est très sérieuse. | Quando si tratta del suo lavoro, questa stilista è molto seria. |
| | | | When it comes to her work, this hairdresser is very serious. | Quand il s'agit de son travail, cette coiffeuse est très sérieuse. | Quando si tratta del suo lavoro, questa parrucchiera è molto seria. |
| | | | When it comes to her work, this makeup artist is very serious. | Quand il s'agit de son travail, cette maquilleuse est très sérieuse. | Quando si tratta del suo lavoro, questa truccatrice è molto seria. |
| | | When it comes to her work, this nurse is very serious. | Quand il s'agit de son travail, cette infirmière est très sérieuse. | Quando si tratta del suo lavoro, quest'infermiera è molto seria. | |
| | | When it comes to her work, this secretary is very serious | Quand il s'agit de son travail, cette secrétaire est très sérieuse. | Quando si tratta del suo lavoro, questa segretaria è molto seria. | |
| | | Anaphoric reference | | | |
| | | This CEO is very serious when it comes to his work. | Ce PDG est très sérieux quand il s'agit de son travail. | Quest'amministratore delegato è molto serio quando si tratta del suo lavoro. | |
| | | This engineer is very serious when it comes to his work. | Cet ingénieur est très sérieux quand il s'agit de son travail. | Quest'ingegnere è molto serio quando si tratta del suo lavoro. | |
| | | This mechanic is very serious when it comes to his work. | Ce mécanicien est très sérieux quand il s'agit de son travail. | Questo meccanico è molto serio quando si tratta del suo lavoro. | |
| | | This pilot is very serious when it comes to his work. | Ce pilote est très sérieux quand il s'agit de son travail. | Questo pilota è molto serio quando si tratta del suo lavoro. | |
| | | This police officer is very serious when it comes to his work. | Ce policier est très sérieux quand il s'agit de son travail. | Questo poliziotto è molto serio quando si tratta del suo lavoro. | |
| | | Cataphoric reference | | | |
| | | When it comes to his work, this CEO is very serious. | Quand il s'agit de son travail, ce PDG est très sérieux. | Quando si tratta del suo lavoro, quest'amministratore delegato è molto serio. | |
| | When it comes to his work, this engineer is very serious. | Quand il s'agit de son travail, cet ingénieur est très sérieux. | Quando si tratta del suo lavoro, quest'ingegnere è molto serio. | | |
| | When it comes to his work, this mechanic is very serious. | Quand il s'agit de son travail, ce mécanicien est très sérieux. | Quando si tratta del suo lavoro, questo meccanico è molto serio. | | |
| | When it comes to his work, this pilot is very serious. | Quand il s'agit de son travail, ce pilote est très sérieux. | Quando si tratta del suo lavoro, questo pilota è molto serio. | | |
| | When it comes to his work, this police officer is very serious. | Quand il s'agit de son travail, ce policier est très sérieux. | Quando si tratta del suo lavoro, questo poliziotto è molto serio. | | |

| | | | | | |
|---|--|----------------------|--|---|--|
| Anti-stereotypical (AS) | Female occupation in Male Context (FOMC) | Anaphoric reference | This fashion designer is very serious when it comes to his work. | Ce styliste est très sérieux quand il s'agit de son travail. | Questo stilista è molto serio quando si tratta del suo lavoro. |
| | | | This hairdresser is very serious when it comes to his work. | Ce coiffeur est très sérieux quand il s'agit de son travail. | Questo parrucchiere è molto serio quando si tratta del suo lavoro. |
| | | | This makeup artist is very serious when it comes to his work. | Ce maquilleur est très sérieux quand il s'agit de son travail. | Questo truccatore è molto serio quando si tratta del suo lavoro. |
| | | | This nurse is very serious when it comes to his work. | Cet infirmier est très sérieux quand il s'agit de son travail. | Quest'infermiere è molto serio quando si tratta del suo lavoro. |
| | | | This secretary is very serious when it comes to his work. | Ce secrétaire est très sérieux quand il s'agit de son travail. | Questo segretario è molto serio quando si tratta del suo lavoro. |
| | | Cataphoric reference | When it comes to his work, this fashion designer is very serious. | Quand il s'agit de son travail, ce styliste est très sérieux. | Quando si tratta del suo lavoro, questo stilista è molto serio. |
| | | | When it comes to his work, this hairdresser is very serious. | Quand il s'agit de son travail, ce coiffeur est très sérieux. | Quando si tratta del suo lavoro, questo parrucchiere è molto serio. |
| | | | When it comes to his work, this makeup artist is very serious. | Quand il s'agit de son travail, ce maquilleur est très sérieux. | Quando si tratta del suo lavoro, questo truccatore è molto serio. |
| | | | When it comes to his work, this nurse is very serious. | Quand il s'agit de son travail, cet infirmier est très sérieux. | Quando si tratta del suo lavoro, quest'infermiere è molto serio. |
| | | | When it comes to his work, this secretary is very serious. | Quand il s'agit de son travail, ce secrétaire est très sérieux. | Quando si tratta del suo lavoro, questo segretario è molto serio. |
| | Male Occupation in Female Context (MOFC) | Anaphoric reference | This CEO is very serious when it comes to her work. | Cette PDG est très sérieuse quand il s'agit de son travail. | Quest'amministratrice delegata è molto seria quando si tratta del suo lavoro. |
| | | | This engineer is very serious when it comes to her work. | Cette ingénieure est très sérieuse quand il s'agit de son travail. | Quest'ingegnera è molto seria quando si tratta del suo lavoro. |
| | | | This mechanic is very serious when it comes to her work. | Cette mécanicienne est très sérieuse quand il s'agit de son travail. | Questa meccanica è molto seria quando si tratta del suo lavoro. |
| | | | This pilot is very serious when it comes to her work. | Cette pilote est très sérieuse quand il s'agit de son travail. | Questa pilota è molto seria quando si tratta del suo lavoro. |
| | | | This police officer is very serious when it comes to her work. | Cette policière est très sérieuse quand il s'agit de son travail. | Questa poliziotta è molto seria quando si tratta del suo lavoro. |
| | | Cataphoric reference | When it comes to her work, this CEO is very serious. | Quand il s'agit de son travail, cette PDG est très sérieuse. | Quando si tratta del suo lavoro, quest'amministratrice delegata è molto seria. |
| | | | When it comes to her work, this engineer is very serious. | Quand il s'agit de son travail, cette ingénieure est très sérieuse. | Quando si tratta del suo lavoro, quest'ingegnera è molto seria. |
| | | | When it comes to her work, this mechanic is very serious. | Quand il s'agit de son travail, cette mécanicienne est très sérieuse. | Quando si tratta del suo lavoro, questa meccanica è molto seria. |
| | | | When it comes to her work, this pilot is very serious. | Quand il s'agit de son travail, cette pilote est très sérieuse. | Quando si tratta del suo lavoro, questa pilota è molto seria. |
| When it comes to her work, this police officer is very serious. | | | Quand il s'agit de son travail, cette policière est très sérieuse. | Quando si tratta del suo lavoro, questa poliziotta è molto seria. | |