La traduction dans une perspective de genre Enjeux politiques, éditoriaux et professionnels

Édité par Sara Amadori, Cécile Desoutter, Chiara Elefante et Roberta Pederzoli



S. Amadori - C. Desoutter - C. Elefante - R. Pederzoli (cur.) - La traduction dans une perspective de genre: enjeux politiques, éditoriaux et professionnels - Milano, LED, 2022 - ISSN 2283-5628 - ISBN 978-88-7916-997-4 https://www.ledonline.it/index.php/LCM-journal/pages/view/LCM-series



http://www.ledonline.it/LCM-Journal

La Collana / The Series

Dipartimento di Lingue, Letterature, Culture e Mediazioni Università degli Studi di Milano

> DIREZIONE / EDITOR-IN-CHIEF Marie-Christine Jullion

Comitato di direzione / Editors

Marina Brambilla - Maria Vittoria Calvi - Lidia Anna De Michelis Giovanni Garofalo - Dino Gavinelli - Antonella Ghersetti - Maria Grazia Guido Elena Liverani - Stefania Maci - Andrea Maurizi - Chiara Molinari Stefano Ondelli - Davide Papotti - Francesca Santulli - Girolamo Tessuto Giovanni Turchetta - Stefano Vicari

Comitato di redazione / Sub-Editors

Maria Matilde Benzoni - Paola Cotta Ramusino Mario de Benedittis - Kim Grego - Giovanna Mapelli - Bettina Mottura Mauro Giacomo Novelli - Letizia Osti Maria Cristina Paganoni - Giuseppe Sergio - Virginia Sica

Comitato scientifico internazionale / International Scientific Committee

James Archibald - Natalija G. Bragina - Kristen Brustad - Giuditta Caliendo Giorgio Fabio Colombo - Luciano Curreri - Hugo de Burgh - Anna De Fina Daniel Dejica - Claudio Di Meola - Denis Ferraris - Lawrence Grossberg Stephen Gundle - Décio de Alencar Guzmán - Matthias Heinz Rosina Márquez-Reiter - Samir Marzouki - John McLeod Estrella Montolío Durán - M'bare N'gom - Christiane Nord Daragh O'Connell - Roberto Perin - Giovanni Rovere Lara Ryazanova-Clarke - Françoise Sabban - Paul Sambre Srikant Sarangi - Kirk St. Amant - Junji Tsuchiya - Xu Shi

All works published in this series have undergone external peer review. Tutti i lavori publicati nella presente Collana sono stati sottoposti a peer review da parte di revisori esterni. ISSN 2283-5628 ISBN 978-88-7916-997-4

Copyright © 2022

IED Edizioni Universitarie di Lettere Economia Diritto Via Cervignano 4 - 20137 Milano www.lededizioni.com - www.ledonline.it - E-mail: led@lededizioni.com

I diritti di riproduzione, memorizzazione e archiviazione elettronica, pubblicazione con qualsiasi mezzo analogico o digitale (comprese le copie fotostatiche, i supporti digitali e l'inserimento in banche dati) e i diritti di traduzione e di adattamento totale o parziale sono riservati per tutti i paesi.

Le fotocopie per uso personale del lettore possono essere effettuate nei limiti del 15% di ciascun volume/fascicolo di periodico dietro pagamento alla SIAE del compenso previsto dall'art. 68, commi 4 e 5, della legge 22 aprile 1941 n. 633.

Le riproduzioni effettuate per finalità di carattere professionale, economico o commerciale o comunque per uso diverso da quello personale possono essere effettuate a seguito di specifica autorizzazione rilasciata da: AIDRO, Corso di Porta Romana n. 108 - 20122 Milano E-mail segreteria@aidro.org <mailto:segreteria@aidro.org> sito web www.aidro.org <http://www.aidro.org/>

Volume pubblicato con il contributo del Dipartimento di Lingue, Letterature e Culture Straniere Università degli Studi di Bergamo

In copertina: Elio Pastore, Umanità in cammino (Moving people #36) Tecnica mista su carta acquerello Canson Infinity Etching, cm 33 × 33 - 2021 www.eliopastore.it

Videoimpaginazione: Paola Mignanego *Stampa:* Logo

Table de matières

Traduction et genre: engagement éthique et défis professionnels Sara Amadori - Cécile Desoutter - Chiara Elefante - Roberta Pederzoli	7
"Thematic Adaptation": On Localizing the Language of "Global Feminism" and Gender Politics in Transnational Feminist Translation Practice and Studies <i>Luise von Flotow</i>	17
Barbara Bray (1924-2010) comme médiatrice interculturelle à la BBC de 1953 à 1972 <i>Pascale Sardin</i>	33
"Le professeur est très intelligent / La prof est très attirante": Recognizing and Reducing Gender Bias in Neural Machine Translation <i>Giuseppe Sofo</i>	49
Queering the Gender Binary American Trans-Themed YA Literature and Its Translation into Italian <i>Beatrice Spallaccia</i>	69
Section thématique L'engagement politique, intellectuel et traductif de l'édition jeunesse indépendante entre la France et l'Italie	
Édition jeunesse généraliste, traduction et questions de genre: analyse comparée du geste éditorial de Babalibri et de L'école des loisirs <i>Sara Amadori</i>	97
Édition pour la jeunesse indépendante entre engagement éthique, traduction et questions de genre: le geste éditorial de Camelozampa et Settenove <i>Roberta Pederzoli</i>	113
Une maison d'édition pour la jeunesse indépendante et militante: engagement, traduction et questions de genre chez Lo Stampatello <i>Valeria Illuminati</i>	131
Les Auteur.es	149

S. Amadori - C. Desoutter - C. Elefante - R. Pederzoli (cur.) - La traduction dans une perspective de genre: enjeux politiques, éditoriaux et professionnels - Milano, LED, 2022 - ISSN 2283-5628 - ISBN 978-88-7916-997-4 https://www.ledonline.it/index.php/LCM-journal/pages/view/LCM-series

"Le professeur est très intelligent / La prof est très attirante"

Recognizing and Reducing Gender Bias in Neural Machine Translation

Giuseppe Sofo

DOI: https://dx.doi.org/10.7359/997-2022-gsof

Abstract

The encounter between Translation Studies and Gender Studies has proven extremely important for the investigation of the issues related to gender representation that are implicit in every natural language, and that processes of translation help to futher unveil. The rapid evolution of Neural Machine Translation over the last years, and the overwhelming presence of machine translation in the contemporary society, call for an investigation of how these issues of gender representation have evolved in the context of machine translation. After outlining the encounter between Translation Studies and Gender Studies, as well as the evolution of machine translation over the last decades, this article will focus on the definition of the issue of gender bias in Neural Machine Translation, on its origins and consequences, and on the attempts that have been made at eliminating or at least reducing gender bias in machine translation.

Keywords: translation; gender; gender bias; Neural Machine Translation; Translation Studies.

Mots-clés: traduction; genre; biais de genre; traduction automatique neuronale; traductologie.

The cultural turn in Translation Studies has opened new paths of investigation into the practice and process of translation, which have revealed how translation can become the ideal space to understand the inner nature of languages, and the cultures they represent. In this sense, the encounter between Translation Studies and Gender Studies in particular, has proven extremely important for the investigation of the issues related to gender representation that are implicit in every natural language, and that the process of translation helps to further unveil.

The rapid evolution of Neural Machine Translation over the last years, and the overwhelming presence of machine translation in the contemporary society, call for an investigation of how these issues of gender representation have evolved in the context of machine translation. After outlining the encounter between Translation Studies and Gender Studies, and the evolution of machine translation, this article will focus on the definition of the issue of gender bias in Neural Machine Translation, on its origins and consequences, and on the attempts that have been made at eliminating or at least reducing gender bias in machine translation.

1. The encounter between Translation Studies and Gender Studies

Reflection on translation and gender originated at the threshold between disciplines and between languages, intended both as natural languages and as disciplinary languages. During the cultural turn in Translation Studies in the 1990s, the décentrement of the discipline occurred through a shift of its focus from the centrality of European literatures, and in particular sacred and ancient literatures (especially Greek and Latin), to other contexts in which the encounter between languages is primarily the result of a conflict. Initially, the focus on colonial and imperialist violence proposed a decentralization of the field of translation. Then, particularly in Canada, the advent of "feminist translation" emphasized the patriarchal nature of language and culture, due to its underrepresentation of the feminine presence in language. Translation thus became "a veritable political tool" (de Lotbinière-Harwood 1991, 27) to counter the domination of the masculine over the feminine that is common to many languages, and which the famous French formula "le masculine l'emporte sur le féminin" expresses very eloquently, beyond its purely linguistic value.

Translation scholars such as Sherry Simon and Luise von Flotow have extensively investigated how translations try to favour the emersion of the feminine presence in language as evidence of an "anti-traditional, aggressive and creative approach to translation" (von Flotow 1991, 70). This is what de Lotbinière-Harwood refers to when she talks about the relevance of these practices to the feminist struggle: "far from being neutral, the act of translating constitutes a speech full of consequences. In addition to being a way of passing from one language to another, translation is also a place of power. For feminist translators, it represents a space to invest, a power to exercise" (de Lotbinière-Harwood 1991, 12). In this sense, translation becomes a practice of resistance to the silent power of language, a space to be conquered in order to open up other spaces, in which the translator, instead of hiding, openly manifests his or her presence by inscribing it in the text. These practices of subversive translation have often been dismissed from the outside as purely ideological ploys. However, their space of action is not limited to the purely political, or politico-cultural, sphere; the motivations and consequences of these practices are in fact purely linguistic and, as such, they can tell us a great deal about the history and future of languages, and specialty languages.

Sherry Simon was right, then, when she wrote in 1996 that "taken together, translation and gender seem to offer a particularly attractive matrix through which to investigate issues of identity in language" (Simon 1996, X). Pascale Sardin's words can help to understand the reason of this "attractive matrix", when she highlights that "la traduction, en tant que transfert culturel où se cristallisent de nombreux enjeux doxiques, constitue [...] un espace privilégié de manifestation de la question du genre" (Sardin 2009, 10). Much has been done in the last thirty years, and yet Gender Studies and Translation Studies continue to meet in fruitful ways, and gender issues are increasingly decisive in the evolution of natural languages.

Gender Studies have in fact not only contributed to the formation and transformation of sectorial languages, but also of "natural" languages. The increasing use of so-called "inclusive writing" – forms of language used to avoid perpetuating gender bias or prejudices, even unintentionally – raises questions of translation, both in the strict and in the broader senses (Sofo 2019a). This practice leads in fact to a transformation of the norms of usage as much as to a transformation of the language itself. At the same time, it gives rise to new translation problems because of the very different forms that these codes have taken in each language. In fact these codes depend on the nature of each single language and on the different degrees of activism in each culture. The field of inclusive writing is a further demonstration of the possibility of change and evolution in any language, as proven by the attempts to undo the prevalence of the masculine as a symbolic construct in society and language – whether this happens through a "demasculinization", a "feminization" or a "neutralization" of gender. Even more interesting, perhaps, are the attempts to find solutions that challenge the binary conception of languages based on a distinction between masculine and feminine, in an effort to shape an inclusive language for nonbinary gender identities, while acting on the characteristics of the languages themselves. This is the case of the use of the asterisk, "-u", or "a" as gender-neutral suffixes in Italian, of the increasingly common use in activist English-speaking circles of "they" as a singular subject, or of the alphabet created by the Swiss designer Tristan Bartolini to express the masculine and feminine forms at the same time (Bartolini 2020). These solutions challenge the rules of the languages involved, but at the same time they are based on a deep understanding of their characteristics and on a creative intervention to expand their possibilities.

The consequences of the encounter between gender and translation are then manifold and can reveal a lot about the future of translation, both in the field of "biotranslation" (Froeliger 2013, 20), as human translation is increasingly defined in the French-speaking context, and in the field of machine translation. An in-depth understanding of the role played by translation in our daily lives cannot in fact be limited to discovering what this practice has been able to transmit from one culture to another, and from one system of thought to another. It must above all seek to understand how this has been done and what "resistances" (Berger 2016, 8-9) translation has faced in this process. If this is crucial to every act of translation, this is even more true for any product of machine translation, and of Neural Machine Translation, for two reasons that are less obvious than they should probably be. On the one hand, the overwhelming presence of these products in our reality shaped by interactions with digital devices, make Neural Machine Translation very influential in the way we conduct our lives; on the other hand, the lack of transparency of the algorithms at the basis of these state-of-theart translation systems.

2. The advent of Neural Machine Translation

The history of Machine Translation spans over more than seven decades, and it is one of successes and failures, which have led to the development of different tools used in the process of translating from one language into another, from systems of Machine Translation in the proper sense, to Computer-Assisted Translation (CAT) tools. The first attempts to create Machine Translation software can be dated between the second half of the 1940s and the first half of the 1950s, and until the 1980s, Machine Translation systems were largely based on a combination of dictionaries and grammar rules. Since 2016, Machine Translation systems have moved to a drastically different paradigm, known as Neural Machine Translation, based on a neural network that is responsible for the encoding and the decoding of the source text into a target text, without a predefined set of rules. The deep-learning system of Neural Machine Translation has very quickly conquered the market and has made the results of Machine Translation much more trustworthy.

This has fostered an unprecedented presence of Machine Translation in everyday life, especially in the field of localization, now largely based on unedited or post-edited Machine Translation. Human translators have in fact often become "post-editors", revising texts translated by a machine, or even "pre-editors", modifying the original text so that it can be processed more easily by the machine. This means that most of the texts that millions of people interact with every day in the digital world, for very different and at times crucial purposes connected to their daily life (from receiving information to communicating with other people or institutions, from reading an instruction manual to understanding the rules of admission in a country), have been produced by Machine Translation. Most readers are not aware of this, or at least not enough. The issue of non-transparency of Neural Machine Translation, in the sense that readers tend to ignore who has produced the translation they are reading (and, at times, even the fact that it is a translation of a text originally written in another language), could be solved through forms of labelling, as Melby (2022) has suggested, to help users distinguish between human and machine translation¹.

This is why we can no longer discuss the way we communicate without discussing the way we translate, and the way we interact with digital tools in translation, as well as why we need to make the presence of Machine Translation transparent for its users. Instead of wondering

¹ Other products of artificial intelligence that are subject to bias specify this clearly for their users through disclaimers. For example, the image generator DALL·E mini, includes a disclaimer on its "bias and limitations" that reads: "While the capabilities of image generation models are impressive, they may also reinforce or exacerbate societal biases. While the extent and nature of the biases of the DALL·E mini model have yet to be fully documented, given the fact that the model was trained on unfiltered data from the Internet, it may generate images that contain stereotypes against minority groups. Work to analyze the nature and extent of these limitations is ongoing, and will be documented in more detail in the DALL·E mini model card" (DALL·E mini 2022).

what translators have gained and/or lost with the advent of machine translation, we should ask what people in general, and not only those actively involved in the translation process, can gain from their interaction with these tools. Van de Meer wrote (2010): "Translation in the 21st century will be a basic utility for everyone on the planet – a human right that everyone can demand and expect. Electricity, water, roads, the Internet, and language translation are all part of the basic services that help drive civilization as we know it". The world we live in shows that people can benefit greatly from fast, free, and sufficiently reliable machine translation systems. In particular, Neural Machine Translation has made possible to translate huge amounts of texts that no human translator could have handled, thus opening the door to a "democratization" of the process of translation.

Could the interaction between Translation Studies and Gender Studies also benefit from this evolution? Natural languages have been accused of carrying a gender bias that is considered dangerous for the people who speak those languages and shape their identity through them (Sofo 2019b), while human translation has often been accused of reinforcing this bias. Given all this could machine translation, and Neural Machine Translation in particular, with the "neutrality" that is usually attributed to machines, help us to reduce this bias instead?

3. Gender bias in Neural Machine Translation

Gender bias is the tendency to favour one gender over the other, based on gender stereotypes (the generalization of characteristics and attributes of a certain group according to their gender) rather than actual differences between individuals, and has mostly taken the form of a favouritism of men in many fields. The issue of gender bias in natural languages is usually attributed to a specific human intention of erasing, hiding, or reducing the role of women in society. As Rachele Raus highlights for French, with words that could apply to many other languages, "la difference de genre grammatical semble se rattacher à la presence d'une forte idéologie andro-centrique sous-jacente" (Raus 2004, 2). Michard and Viollet have explored about the relationship between "marquage du genre et biais androcentrique dans la structure et le contenu de la langue" (Michard et Viollet 1991, 102). Richy and Burnett have gone perhaps even further for French: they demonstrate that "le genre grammatical masculin en français crée un biais masculin dans l'interprétation de syntagmes nominaux référant aux humains allant au-delà des stéréotypes, qui pourtant jouent aussi un rôle" (Richy et Burnett 2021, 2).

Most researchers who have worked on making natural languages a more inclusive tool of communication have in fact highlighted very strongly and very convincingly the role of men and of chauvinist societies in the "masculinisation" of languages, against the idea that languages were "born this way" (see for French: Viennot 2015, 2017; for Italian: Robustelli 2014). This would suggest that the alleged neutrality of machine translation could thus serve as a tool to re-establish a wider "neutrality" (in the sense of an absence of bias) in the translation process between two natural languages, avoiding human intervention, since the discussion about gender issues in language, "including its human and emotional components, should be basically indifferent to a machine" (One Word 2021).

However, matters are much more complicated than this, and machine translation, like many other products of artificial intelligence, is rather characterized by what is being increasingly defined as "machine bias", the phenomenon according to which "machine learning algorithms are prone to reinforce or amplify human biases" (Farkas and Németh 2022, 1). Several examples could be given, involving different languages and situations. To give just a few examples, if we translate from English into French the sentence: "The surgeon spoke to the nurse", using Google Translate, we will get: "Le chirurgien a parlé à l'infirmière", although the sentence could be correctly translated in three other ways, with different genders attributed to both persons involved². Google Translate's choice of the first version depends on the fact that statistically, according to the data it relies on ³, there are more male than female surgeons, and more female than male nurses.

Even more interesting, and surely more problematic, are other cases in which Machine Translation does not only rely on statistical data on the actual presence of men and women in one profession or the other, but also on gender stereotypes. To give a simple and very eloquent ex-

² "Le chirurgien a parlé à l'infirmier"; "La chirurgienne a parlé à l'infirmière"; "La chirurgienne a parlé à l'infirmier".

³ It is also important to highlight that corpora do not reflect present-day society, but rather a much more layered history of texts. As Jonathan Davis writes: "inherent social biases, such as gender or ethnicity bias, manifest themselves in datasets which are essentially historical records of the given society", which means that "if more men have *bistorically* been doctors than women, a machine learning model trained on historical data will learn that doctors are more likely to be male than female, irrelevant of the *current* gender split amongst doctors" (Davis 2020).

ample, if we translate from Italian into French using Google Translate the sentence: "L'insegnante è molto intelligente", we get "Le professeur est très intelligent"; if we simply replace the adjective "intelligent" with "attractive", we get: "La prof est très attirante", which does not only switch to the feminine, but also shortens "professeur" (or "professeure") into "prof", a more colloquial version of the term.

The surgeon spoke to the nurse	\rightarrow	Le chirurgien a parlé à l'infirmière
L'insegnante è molto intelligente	\rightarrow	Le professeur est très intelligent
L'insegnante è molto attraente	\rightarrow	La prof est très attirante
(Google Translate, 28.6.2022)		

Of course, bias does not only exist in machine translation, but in all products of artificial intelligence, and gender bias is not the only kind of bias produced by this phenomenon. In May 2016, the independent non-profit investigative journalism platform ProPublica was able to show that COMPAS (Correctional Offender Management Profiling for Alternative Sanctions), the AI tool used by U.S. courts to predict future criminals, had a bias against black people, i.e., it was more likely to identify innocent black people as guilty, and more likely to identify guilty white people as innocent (Angwin *et al.* 2016).

The definition of bias is also rather ambiguous. According to the Merriam-Webster Online Dictionary, it is a "systematic error introduced into sampling or testing by selecting or encouraging one outcome or answer over others" (Merriam-Webster). Michal Měchura provides two distinct definitions that speak volumes about this concept: "the technical definition of bias is that it is the tendency of an automated system to make the same kind of assumptions again and again. Then there is the popular definition of bias, which is basically the same but with an implication of offense, harm, and injustice" (Měchura 2022).

This double definition is very important in the context of machine translation. Of course, "machine bias" is not produced "intentionally" by the machine, but is mostly due to its learning through corpora (although not only, as we will discuss later on). Just like a miseducated child will not become better than the education he or she was given, if we train machine translation with biased data, the machine will produce biased results. However, the effects produced by "machine bias" are just as likely to cause "offense, harm, and injustice" as any other kind of bias, and, as we saw, rather than simply reproducing bias, machines may amplify it.

Hence, focusing on these issues is extremely relevant, in order to act on the products of Neural Machine Translation systems which are not only non-transparent for the user, but quite often also for their creators. As Wisniewski, Zhu, Ballier and Yvon and Zhu argue, "contrary to previous generations of MT engines where transfer rules were quite transparent, understanding this *information flow* [from the encoder to the decoder] within state-of-the-art neural MT systems is a challenging task, and a key step for their interpretability" (Wisniewski *et al.* 2021b).

Gender bias refers to "gender choices made in a translation that are not present in the source text" (Pérez Piñeiro 2021), and more precisely to all the issues resulting from an underrepresentation of one gender over the other when these choices are made. Wisniewski, Zhu, Ballier and Yvon have identified the following kinds of gender bias that could occur in machine translation:

(a) le fait que des erreurs de traduction sont plus fréquentes pour des énoncés qui mettent en scène des participantes de genre féminin; (b) le fait que des traductions rendent linguistiquement explicite le genre des actants évoqués, alors que l'intention du locuteur peut être de le laisser ambigu; (c) le fait que ces explicitations privilégient des assignations stéréotypiques, confortant, voire renforçant des préjugés sexistes dans les textes traduits. (Wisniewski *et al.* 2021a, 12)

This situation produces two different kinds of problems, ranging from "fausser la manière dont certains groupes sont représentés dans les textes (*representational harm*)" to "conduire à un service de moindre qualité pour les femmes (*allocational harm*)" (*ibidem*).

The main issue of gender bias in machine translation is then linked to a problem that exists also in human translation, i.e., the diversity between languages that have different ways of expressing gender, which in turn leads to ambiguities in the translation, at times unresolvable. While some languages have only masculine and feminine forms, some also have neutral forms, while others are gender-neutral. Of course, this creates problems to human translators as well, who must decide between different available options when translating from a gender-neutral language (such as English, Turkish or Finnish) into a language that must express a gender (such as French, Italian, Spanish, German or Arabic). While a human translator might be able to infer the correct choice from other parts of the text, from the context, or from his or her own experience, obviously this is much more complicated for machine translation, which mostly works at the level of the sentence, and might not be able to retrieve all the information it needs to make a proper choice and ultimately has to "guess" the correct choice.

Two kinds of gender biases result from this process. First, all relevant research shows that, when it has to guess, machine translation very often chooses the masculine over the feminine; second, machine translation is more able to guess correctly the feminine, when the pronouns or nouns are associated with adjectives, or names of occupations that are stereotypically attributed to women. Names of occupations are an ideal field of action for this kind of research, not only because they help to understand all the layers of bias that are included in the process, but also because this is historically the most active lexical field in the process of "demasculinization" that languages have undergone over the last decades (for French: Yaguello 1989; Cerquiglini 1999; for Italian: Sabatini 1987; Robustelli 2014), the first one in which women claimed a "name" of their own. Studies on machine translation of sentences including male and female co-referents, connected to names of occupations, have been carried out for several different language pairs and cultural contexts (Stanovsky et al. 2019; Luccioli et al. 2020; Richy et Burnett 2021: Farkas and Németh 2022).

Stanovsky, Smith, and Zettlemoyer presented in 2019 "the first large-scale multilingual evaluation of gender-bias in machine translation", translating two datasets (Winogender and WinoBias) of sentences from English into eight different target languages (French, Spanish, Italian, Russian, Ukranian, Hebrew, Arabic, and German). Sentences describe "a scenario with human entities, who are identified by their role [...], and a pronoun [...], which needs to be correctly resolved to one of the entities", such as "the doctor asked the nurse to help her in the procedure" (Stanovsky et al. 2019, 1679). The results showed that most tested systems perform very poorly, and that "the best performing model on each language often does not do much better than a random guess for the correct inflection" (ibid., 1681). Even more interestingly, "all tested systems have a significant and consistently better performance when presented with pro-stereotypical assignments (e.g., a female nurse), while their performance deteriorates when translating antistereotypical roles (e.g., a male receptionist)" (ibid., 1682). This obviously leads to the conclusion that "MT models are significantly prone to translate based on gender stereotypes rather than more meaningful context" (ibid., 1683).

In a case study conducted between English and Italian, Luccioli, Dolei and Xausa were able to show that MT systems "exhibit a statistical bias towards male defaults, as well as a tendency to reproduce gender stereotypes" (2020, 44). Even more interestingly, by using three adjectives with different stereotypical distribution as modifiers of names of occupations - "wise", stereotypically attributed to men; "beautiful", stereotypically attributed to women; "strong", equally distributed -, they proved that "'stereotypical' adjectives can affect the MT output", because "adding the 'stereotypical' female adjective *beautiful*, the systems' performance improves with the female co-referent and worsens with the male one. On the contrary, when adding the 'stereotypical' male modifier wise, the systems' accuracy is maintained consistently high with the male co-referent and decreases even further when the co-referent is her", while "the baseline set and the set with modifier strong show that the best performance is achieved with the male co-referent; conversely, the accuracy is rather low when the co-referent is her" (ibidem). This confirms that MT systems tend to largely prefer the masculine when they have no additional information or undecisive additional information, but increase their correctness for both male and female co-referents (and especially for the latter) when stereotypical information is added to the sentence.

In a series of experiments conducted between French and English on sentences such as "L'actrice a terminé son travail", Wisniewski, Zhu, Ballier, and Yvon have been able to identify not only machine mistakes, but also something much more important than that, i.e., the process through which Neural Machine Translation systems produce these mistakes. They began by constructing a set of 388 sentences on the model "The [noun] completed [his/her] work", where the noun is the name of a profession, with the respective French model "[det] [nom] a fini son travail", where [det] is the determiner preceding the noun, and [nom] always the name of a profession. This once again proved that Neural Machine Translation systems have huge difficulties in predicting gender information during the translation (Wisniewski et al. 2021a, 16). However, results became even more interesting once they replaced the article with the epicene determiner "chaque", thus leaving the mark of gender only to the profession. What happens is that the tested system JoeyNMT "ne commet aucune erreur en transférant en anglais le genre du nom de métier masculin, alors qu'il se trompe presque systématiquement (94.55% d'erreurs) pour les féminins" (ibid., 17).

During these experiments, they discovered that manipulating the representation of the possessive pronoun "son" did not have any impact on the choice of "his/her" by the system, which preferred to rely on the name of profession. This shows that "ce n'est pas parce qu'une information 'linguistique' est encodée dans les representations neuronales qu'elle est exploitée par le réseau" (*ibid.*, 19). This introduces a further complexity to the process, because it reveals that some linguistic infor-

mation might be completely ignored by the system, which prefers to rely on other (perhaps less relevant) information for its process, thus producing a less successful result.

In two forthcoming articles, they were able to further investigate the causes of these deviations, succeeding in revealing some of the previously unknown details of the process of these Machine Translation systems and giving us the opportunity to enter the "black box" of Neural Machine Translation. In the first article, they show not only that the system makes more mistakes for the feminine than for the masculine, but also that "la différence est plus marquée pour l'encodeur que pour le décodeur. Il apparait donc que, pour 'corriger' les prédictions erronées des exemples féminins, le réseau de neurones cherche principalement à mieux extraire des informations de la source d'où peut être extraite l'information de genre" (Wisniewski *et al.*, forthcoming-a).

In the second article, they point out that the prediction of the masculine form "his" and of the feminine form "her" happens in completely different ways. In fact, their most recent contribution illustrates that "counter-intuitively, the TM does not need to use information from the source sentence to predict *his* while this is necessary for the prediction of *her*" (Wisniewski *et al.*, forthcoming-b, 1). As the authors explain more in detail:

Above all, these results show that the TM does not need to use the same information to decide between *her* and *his*: for the latter, it can rely on the target context only; but, for the former, it must learn to transfer information from the source, which it only does imperfectly (at least for our system). This suggests that gender bias could partly result from using cross-entropy as a loss function: indeed, the model appears to be quite capable of learning to predict *his* for the wrong reason and without necessarily taking into account the gender information present in the source, which weakens the estimation of the parameters necessary to take into account gender information (e.g. cross attention). (Wisniewski *et al.*, forthcoming-b, 4)

This ground-breaking discovery opens a completely different perspective, which partially goes against one of the most common assumptions about gender bias in machine translation. It is usually assumed that machine bias is due to a reflection of human bias in corpora, this points to a possible different perception of the concept of "machine bias" itself, since this kind of bias might derive from the machine itself, which treats the masculine form as the expected form and the feminine as a possible "deviation" from this norm. This also has some roots in human data, since the machine is trained with corpora where the masculine is largely more present than the feminine ⁴, and it is therefore much easier for the machine to treat feminine forms as possible "deviations" from the norm of the masculine. Ironically, though, the pathbreaking work carried out by researchers in linguistics to reduce the presence of the masculine and highlight the equal status of masculine and feminine forms, thus promoting non-binary systems for languages, is drastically overturned by artificial intelligence. For the machine, in fact, "binarism" is not a philosophical choice, but rather a construction paradigm.

4. Possible solutions for "bias reduction"

Any method to attenuate or mitigate the effects of gender bias in machine translation must be preceded by a careful study of the biases themselves, and of the reasons behind their formation. It is therefore crucial to:

étudier les mécanismes internes des systèmes de traduction neuronaux et notamment les flux d'informations entre l'encodeur et le décodeur. En effet, une meilleure compréhension de ces flux permettrait d'une part d'expliquer les biais de genre qui émaillent les traductions de ces derniers, étape nécessaire pour la correction de ceux-ci, et d'autre part de progresser vers une meilleure compréhension de la capacité de généralisation de ces architectures, premier pas pour développer des systèmes capables d'apprendre à partir de moins de données. (Wisniewski *et al.*, forthcoming-a)

Over the last years, several possible solutions have been tested to reduce this bias, both by academic researchers and by companies. We could call this an effort of "bias reduction", to make use of a concept proposed by Google AI, which defines it as a "new method of evaluation [...] which measures the relative reduction of bias" (Johnson 2020) from one system to another.

While it might seem that the most obvious solution would be to train machines with corpora that are not gender biased, this is easier said than done. In fact, "bias is inherent in human texts, and it is not necessarily a source of unfairness" (Farkas and Németh 2021, 1), and producing completely unbiased corpora seems to be on the one hand

⁴ This is due to different reasons, including a larger presence of texts written by men as well as a larger presence of texts about men, which does not only characterize corpora chosen for training, but also the texts available on the Internet as a whole.

almost impossible, and on the other hand questionable because it would produce an artificial corpus detached from reality, ultimately producing results that would therefore be less trustworthy.

Wisniewski, Zhu, Bellier, and Yvon have discussed three different approaches used for bias reduction:

Mesurer les biais permet aussi d'évaluer l'impact de travaux visant à les atténuer dans des traductions automatiques. Ces travaux mobilisent principalement trois types de techniques [...]. Une première consiste à manipuler les représentations lexicales. [...] Les techniques de pré-annotation [...] insèrent dans le texte source des marques explicites de genre, qui vont servir à orienter le système vers des traductions correctes. [...] Une troisième famille d'approches manipule les distributions des données d'apprentissage en s'appuyant sur des méthodes d'augmentation de données (*counterfactual data augmentation (CDA*)). (Wisniewski *et al.* 2021a, 22)

In December 2018, Google Translate began to offer "gender-specific translations" for certain language pairs ⁵ in order to allow the user to choose between the masculine and the feminine form when both were possible. As Kuczmarski and Johnson point out, this is a technique "to generate both a masculine and a feminine translation for gender-neutral text, thereby reducing or eliminating gender bias in machine translation", based on three main components: "detection, generation of alternatives, and validation" (Kuczmarski and Johnson 2018). To be more specific, this happens through a "three-step approach, which involved detecting gender-neutral queries, generating gender-specific translations and checking for accuracy" (Johnson 2020).

Less than two years later, in April 2020, after trying to expand the same system to more language pairs and realizing the shortcomings of this approach, Google has produced an entirely different process, based on a "rewriting" approach. Instead of identifying gender-neutral queries before translating, this new approach first "generate[s] the initial translation", and "the translation is then reviewed to identify instances where a gender-neutral source phrase yielded a gender-specific translation" (*ibidem*). If this is the case, the system applies "a sentence-level rewriter to generate an alternative gendered translation. Finally, both the initial and the rewritten translations are reviewed to ensure that the only difference is the gender" (*ibidem*). The result of this change in paradigm is an improved "bias reduction" from 60% to 95%.

⁵ Initially, for translations from Turkish into English and then for translations from English into Spanish, which is "the most popular language-pair in Google Translate" (Johnson 2020).

Another issue stems from the fact that "by assuming a femalemale dichotomy and by emphasizing language which reflects the two categories, linguists may be reinforcing biological essentialism, even if they emphasize that language, like gender, is learned behavior" (Bergvall and Bing 1998, 505). While claiming that grammatical gender and social gender are not related, studies focusing only on the masculine-feminine dichotomy tend in fact to erase the presence of neutral and non-binary forms, which are increasingly common in natural languages and which pose a whole different set of translation problems, as I suggested before. If the translation of inclusive language is already an issue for human translators, machine translation adds the problem of identifying new forms of writing that challenge traditional grammar, that is not to be given for granted. Machine Translation does not always respond very well to the new forms of inclusion that are being experimented in recent years in different languages. A study by the German translation company One Word focusing on this issue has shown for example that different NMT systems react very differently to signs like the *point milieu* used in French (One Word 2021), and the same could be said for other forms which subvert the rules of natural languages.

The study by One Word suggests to carefully select the MT system and "to be aware of potential stumbling blocks". However, it also mentions the fact that "the gender specifications must already be included in the training material and implemented consistently", since by doing so "the machine could also be trained to translate correctly into the desired gender from the target language" (*ibidem*). The same was already proven in 2019 in a study which proposed "the new task of re-writing gendered sentences to be gender-neutral with the underlying goal of developing gender inclusive natural language systems" (Sun *et al.* 2019, 4). This study was also able to show "how it is possible to train a model to do this without any human-labeled data, by automatically generating large amounts of data by rule" (*ibid.*, 4-5), thus removing the problem of human labelling, which would make the task very difficult to sustain over large amounts of data.

Researchers are thus beginning to focus on these issues as well, to "extend previous machine translation coreference to the translation of gender-neutral language, which may be used by non-binary individuals or to avoid the social impact of using gendered language" (Saunders, Sallis, and Byrne 2020; Bradley *et al.* 2019). We will see what results these new studies will produce, and how this evolution of natural languages will affect the usability of Machine Translation.

5. CONCLUSION

As we have seen, gender bias in machine translation results from three components of the process. The first is human bias, because behind any product of artificial intelligence, there are real people with their own biases, who choose the data that constitute machine translation models. Second, biased data affect machine translation because the corpora we use to train machine translation are largely biased, since they reproduce the biases and "history" of our societies, which are only very lately working on reducing bias. Last but not least, there seems to be an embedded bias in Neural Machine Translation, whose reasons are still not entirely transparent, but whose effects on the prevalence of masculine over feminine forms are very evident in the translations produced by these systems.

Current research on this topic demonstrates that the more we understand the nature of this bias, the more we might find ways to reduce it. Post-editors, whose work is to revise the output of machine translation, must pay particular attention to this gender bias and try to restore the presence of feminine or neutral and non-binary forms that machine translation tends to erase. More than anything else, however, the real challenge is to make readers or users of these translations aware of the fact that they are reading a product of machine translation, and to make them aware of the possible shortcomings of these translations. And this is perhaps also the best possible solution to better understand the advantages and disadvantages of the large-scale use of Neural Machine Translation.

References

- Angwin, Julia, Jeff Larson, Surya Mattu, and Lauren Kirchner. 2016. "Machine Bias: There's Software Used across the Country to Predict Future Criminals. And It's Biased against Blacks". *ProPublica*, May 23. [28/06/2022]. https://www.propublica.org/article/machine-bias-risk-assessments-incriminal-sentencing
- Bartolini, Tristan. 2020. "Alphabet Inclusif". [28/06/2022]. https://www.profa.ch/wp-content/uploads/2020/10/alphabet_inclusif.png
- Berger, Anne Emmannuelle. 2016. "Gender Springtime in Paris: A Twenty-First-Century Tale of Seasons". *differences: A Journal of Feminist Cultural Studies* 27 (2): 1-26.

- Bergvall, Victoria L., and Janet M. Bing. 1998. "The Question of Questions: Beyond Binary Thinking". In *Language and Gender: A Reader*, edited by Jennifer Coates, 495-510. Oxford: Blackwell.
- Bradley, Evan D., Julia Salkind, Ally Moore, and Sofi Teitsort. 2019. "Singular "They' and Novel Pronouns: Gender-Neutral, Nonbinary, or Both?". *Proceedings of the Linguistic Society of America* 4 (36): 1-7.
- Cerquiglini, Bernard, éd. 1999. Femme, j'écris ton nom. Paris: La documentation française.
- DALL·E mini. 2022. https://dallemini.com/
- de Lotbinière-Harwood, Susanne. 1991. *Re-Belle et Infidèle. La traduction comme* pratique de réécriture au féminin / The Body Bilingual: Translation as a *Rewriting in the Feminine*. Montréal - Toronto: Les Éditions du Remue-Ménage - The Women's Press.
- Farkas, Anna, and Renáta Németh. 2021. "How to Mesure Gender Bias in Machine Translation: Real-World Oriented Machine Translators, Multiple Reference Points". Social Sciences & Humanities Open 5: 1-11.
- Flotow, Luise von. 1991. "Feminist Translation: Contexts, Practices and Theories". *Traduire la théorie* 4 (2, II semestre): 69-84.
- Frœliger, Nicolas. 2013. *Les Noces de l'analogique et du numérique. De la traduction pragmatique*. Paris: Les Belles Lettres.
- Google Translate. [28/06/2022]. https://translate.google.com/
- Johnson, Melvin. 2020. "A Scalable Approach to Reducing Gender Bias in Google Translate". Google AI Blog: The Latest from Google Research. Last modified April 22. [28/06/2022]. https://ai.googleblog.com/2020/04/a-scalable-approach-to-reducinggender.html
- Kuczmarski, James, and Melvin Johnson. 2018. "Gender-Aware Natural Language Translation". *Technical Disclosure Commons*, October 8. [28/06/2022]. https://www.tdcommons.org/dpubs_series/1577/
- Luccioli, Alessandra, Ester Dolei, and Chiara Xausa. 2020. "Investigating Gender Bias in Machine Translation: A Case Study between English and Italian". *MediAzioni* 29 (*Metodi e ambiti nella ricerca sulla traduzione, l'interpretazione e l'interculturalità*): B29-B49.
- Měchura, Michal. 2022. "What You Need to Know About Bias in Machine Translation". *Slator: Language Industry Intelligence*. Last modified May 27. [28/06/2022].

https://slator.com/what-you-need-to-know-about-bias-in-machine-translation/

Melby, Alan. 2022. "Confining MT; Reclaiming Localization". Paper given at the *Tralogy III Conference* (Paris, April 7-8, 2022).

Merriam-Webster . [28/06/2022]. https://www.merriam-webster.com

Michard, Claire, et Catherine Viollet. 1991. "Sexe et genre en linguistique. Quinze ans de recherches féministes aux États-Unis et en R.F.A.". Unité/Diversité 4 (2): 97-128. One Word. 2021. "Gendering in Machine Translation – How It Works and What to Watch out For: Can Gendering Be Done by Machines?". Last modified June 17. [28/06/2022].

https://www.oneword.de/en/gendering-in-machine-translation/

- Pérez Piñeiro, Pablo. 2021. "Managing Gender Bias in Machine Translation". RWS. Last modified August 27. [28/06/2022]. https://www.rws.com/blog/gender-bias-machine-translation-language-weaver/
- Raus, Rachele. 2004. "La Linguistique française et les études de genre. Le 'discours polémique' des femmes et l'imaginaire linguistique". Lecture online course Introduzione agli studi di genere. [28/06/2022].
 https://www.cirsde.unito.it/sites/c555/files/allegatiparagrafo/04-05-2016/3._ la_linguistique_francaise_et_les_etudes_de_genre_fr.pdf
- Richy, Célia, et Heather Burnett. 2021. "Démêler les effets des stéréotypes et le genre grammatical dans le biais masculin. Une approche expérimentale". *GLAD!* 10: 1-31.
- Robustelli, Cecilia. 2014. Donne, grammatica e media. Suggerimenti per l'uso dell'italiano. Roma: GiULiA.
- Sabatini, Alma. 1987. *Il sessismo nella lingua italiana*. Roma: Presidenza del Consiglio dei Ministri.
- Sardin, Pascale. 2009. "Colloque sentimental". Palimpsestes 22: 9-21.
- Saunders, Danielle, Rosie Sallis, and Bill Byrne. 2020. "Neural Machine Translation Doesn't Translate Gender Coreference Right Unless You Make It". In *Proceedings of the Second Workshop on Gender Bias in Natural Language Processing*, edited by Marta R. Costa-Jussà, Christian Hardmeier, Will Radford, and Kellie Webster, 35-43. Stroudsburg (PA): Association for Computational Linguistics.
- Simon, Sherry. 1996. Gender in Translation: Cultural Identity and the Politics of Transmission. London New York: Routledge.
- Sofo, Giuseppe. 2019a. "Traduction du langage inclusif et échanges entre le français et l'italien". *Savoirs en prisme* 10. https://savoirsenprisme.com/numeros/10-2019-les-nouvelles-formesdecriture/traduction-du-langage-inclusif-et-echanges-entre-le-francais-etlitalien/
- Sofo, Giuseppe. 2019b. "Il genere della traduzione. Per una traductologie d'intervention". de genere – Rivista di studi letterari, postcoloniali e di genere 5 (Il genere della traduzione / Le genre de la traduction / The Gender and Genre of Translation, a cura di Giuseppe Sofo e Anne Emmanuelle Berger): XIII-XXX. [28/06/2022].

http://www.degenere-journal.it/?journal=degenere&page=article&op=view &path%5B%5D=119

Stanovsky, Gabriel, Noah A. Smith, and Like Zettlemoyer. 2019. "Evaluating Gender Bias in Machine Translation". In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics Association, 1679-1684. Firenze: Association for Computational Linguistics.

S. Amadori - C. Desoutter - C. Elefante - R. Pederzoli (cur.) - La traduction dans une perspective de genre: enjeux politiques, éditoriaux et professionnels - Milano, LED, 2022 - ISSN 2283-5628 - ISBN 978-88-7916-997-4 https://www.ledonline.it/index.php/LCM-journal/pages/view/LCM-series

Sun, Tony, Kellie Webster, Apu Shah, William Yang Wang, and Melvin Johnson. 2021. "They, Them, Theirs: Rewriting with Gender-Neutral English". [28/06/2022]. https://dordine.org/pdf/2102.06788.pdf

https://arxiv.org/pdf/2102.06788.pdf

- van de Meer, Jaap. 2010. "Where Are Facebook, Google, IBM and Microsoft Taking Us?". *TAUS Articles*. Last modified August 2. [28/06/2022]. https://www.taus.net/think-tank/articles/translate-articles/where-arefacebook-google-ibm-and-microsoft-taking-us
- Viennot, Éliane, éd. 2015. *L'Académie contre la langue française. Le dossier "féminisation*". Donnemarie-Dontilly: Éditions iXe.
- Viennot, Éliane. 2017. Non, le masculin ne l'emporte pas sur le féminin. Petite bistoire des résistances de la langue française. Édition augmentée. Donnemarie-Dontilly: Éditions iXe.
- Wisniewski, Guillaume, Lichao Zhu, Nicolas Ballier, et François Yvon. 2021a. "Biais de genre dans un système de traduction automatique neuronale. Une étude préliminaire". Dans *Traitement Automatique des Langues Naturelles*, 11-25. Lille. hal-03265895.
- Wisniewski, Guillaume, Lichao Zhu, Nicolas Ballier, and François Yvon. 2021b. "Screening Gender Transfer in Neural Machine Translation". In Proceedings of the Fourth BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP, edited by Jasmijn Bastings, Yonatan Belinkov, Emmanuel Dupoux, Mario Giulianelli, Dieuwke Hupkes, Yuval Pinter, and Hassan Sajjad, 311-321. Stroudsburg (PA): Association for Computational Linguistics.
- Wisniewski, Guillaume, Lichao Zhu, Nicolas Ballier, et François Yvon. Forthcoming-a. "Flux d'informations dans les systèmes encodeur-décodeur. Application à l'explication des biais de genre dans les systèmes de traduction automatique". Dans Actes de la 29^e Conférence sur le Traitement Automatique des Langues Naturelles (Avignon, 27 juin - 1 juillet 2022).
- Wisniewski, Guillaume, Lichao Zhu, Nicolas Ballier, and François Yvon. Forthcoming-b. "Investigating the Roots of Gender Bias in Machine Translation: Observations on Gender Transfer between French and English". NAACL (ACL Rolling Review).

Yaguello, Marina. 1989. Le sexe des mots. Paris: Belfond.

Résumé

La rencontre entre la traductologie et les études de genre s'est avérée extrêmement importante pour l'investigation des questions liées à la représentation du genre qui sont implicites dans toute langue naturelle, et que les processus de traduction aident à dévoiler davantage. L'évolution rapide de la traduction automatique neuronale (NMT) au cours des dernières années et la présence croissante de la traduction automatique dans la sphère contemporaine nous obligent à nous intéresser à la manière dont ces questions de représentation du genre ont évolué dans le contexte de la traduction automatique. Après une introduction à la rencontre entre la traductologie et les études de genre, et à l'évolution de la traduction automatique au cours des dernières décennies, cet article se concentre sur la définition de la question du biais de genre dans la traduction automatique neuronale, sur ses raisons et ses conséquences, et sur les tentatives qui ont été faites pour éliminer ou du moins réduire le biais de genre dans la traduction automatique.