



HAL
open science

Évaluation automatique des biais de genre dans des modèles de langue auto-régressifs

Fanny Ducel, Aurélie Névéol, Karën Fort

► **To cite this version:**

Fanny Ducel, Aurélie Névéol, Karën Fort. Évaluation automatique des biais de genre dans des modèles de langue auto-régressifs. TALN 2024, Jul 2024, Toulouse, France. hal-04621134

HAL Id: hal-04621134

<https://inria.hal.science/hal-04621134>

Submitted on 23 Jun 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

Évaluation automatique des biais de genre dans des modèles de langue auto-régressifs

Fanny Ducel¹, Aurélie Névéol¹, Karën Fort²

(1) Université Paris-Saclay, CNRS, LISN, 91400, Orsay, France

(2) Sorbonne Université / Université de Lorraine, CNRS, Inria, LORIA, F-54000 Nancy, France
fanny.ducel@lisn.fr, aurelie.neveol@lisn.fr, karen.fort@loria.fr

RÉSUMÉ

Nous proposons un outil pour mesurer automatiquement les biais de genre dans des textes générés par des grands modèles de langue dans des langues flexionnelles. Nous évaluons sept modèles à l'aide de 52 000 textes en français et 2 500 textes en italien, pour la rédaction de lettres de motivation. Notre outil s'appuie sur la détection de marqueurs morpho-syntaxiques de genre pour mettre au jour des biais. Ainsi, les modèles favorisent largement la génération de masculin : le genre masculin est deux fois plus présent que le féminin en français, et huit fois plus en italien. Les modèles étudiés exacerbent des stéréotypes attestés en sociologie en associant les professions stéréotypiquement féminines aux textes au féminin, et les professions stéréotypiquement masculines aux textes au masculin.

ABSTRACT

Automatically Assessing Gender Biases in Autoregressive Language Models.

We propose a framework to automatically measure gender biases in generated texts, for inflected languages. We evaluate seven language models, on over 52,000 texts in French and 2,500 texts in Italian, for cover letter writing. Our tool relies on the detection of morpho-syntactic gender markers to uncover biases. Thus, models are strongly biased towards the generation of masculine markers : generated texts contain twice as many masculine (vs. feminine) markers in French, and eight times as many in Italian. The models we study also exacerbate gender stereotypes that are evidenced in social science studies and associate feminine inflections with stereotypically feminine occupations, whereas stereotypically masculine occupations are strongly associated with masculine markers.

MOTS-CLÉS : Biais, Stéréotype, Genre, Modèle de langue (LLM), Français, Italien.

KEYWORDS: Bias, Stereotype, Gender, Language Model (LLM), French, Italian.

1 Introduction

Au cours des dernières années, les grands modèles de langue (*Large Language Models*, ou LLM) sont devenus l'approche privilégiée pour la plupart des tâches de traitement automatique des langues (TAL) telles que la classification de textes, la reconnaissance d'entités nommées ou la traduction automatique (Howard & Ruder, 2018; Epure & Hennequin, 2022; Peng *et al.*, 2023), y compris pour des applications destinées au grand public. Néanmoins, ces modèles non seulement reproduisent, mais amplifient les biais stéréotypés (Gehman *et al.*, 2020; Dhamala *et al.*, 2021; Kirk *et al.*, 2021) qu'il est important de détecter et d'évaluer afin d'éviter qu'ils ne perpétuent des discriminations.

Modèle	Type	Taille	Langue(s)	Référence
xglm	Base	2,9M	FR, IT (Multi.)	(Lin <i>et al.</i> , 2022)
gpt2-fr	Base	1M	FR	(Simoulin & Crabbé, 2021)
vigogne-2-instruct	Affiné (LLAMA)	7M	FR	(Huang, 2023)
BLOOM	Base	560m, 3M, 7M1	FR (Multi)	(Scao <i>et al.</i> , 2022)
cerbero	Affiné (MISTRAL)	7M	IT	(Galatolo & Cimino, 2023)

TABLE 1 – Description des modèles de langue testés (m : million, M : milliard)

Les contributions de ce travail sont les suivantes : (i) un outil (*framework*) détectant les biais de genre dans des langues flexionnelles à partir d’indices morpho-syntaxiques et pour un cas d’utilisation réaliste, l’aide à la rédaction de lettres de motivation ; (ii) un système de détection automatique des marqueurs de genre pour le français et l’italien¹ ; (iii) une étude des biais dans sept modèles de langue en français et en italien, en utilisant l’outil proposé et des études sociologiques.

2 État de l’art

Les biais stéréotypés dans les modèles de langue Des corpus ont été créés par la communauté pour découvrir différents types de biais stéréotypés dans les systèmes de TAL, avec un accent récent sur les modèles de langue (Nangia *et al.*, 2020; Li *et al.*, 2020; Nadeem *et al.*, 2021; Névéol *et al.*, 2022; Parrish *et al.*, 2022). Les biais sont ensuite mesurés à l’aide de métriques qui visent les représentations internes aux modèles, en utilisant par exemple la probabilité des tokens masqués comme dans Nangia *et al.* (2020), ou les biais présents dans les sorties du système (De-Arteaga *et al.*, 2019; Nozza *et al.*, 2021; de Vassimon Manela *et al.*, 2021). Cette deuxième catégorie de métriques, dite extrinsèque, est supposée plus robuste (Delobelle *et al.*, 2022). Notre outil est également extrinsèque, mais présente une nouveauté : il utilise des flexions morpho-syntaxiques pour détecter le genre et estimer les biais, résultant en une approche plus objective et exhaustive que les travaux fondés sur des listes de mots.

Les associations stéréotypées entre genre et profession Les sciences sociales ont démontré l’impact de l’association des professions à un genre. Ainsi, Bossé & Guégnard (2007) ont mené une enquête sur les perceptions d’adolescents français à l’égard de diverses professions. Ces derniers pensent que certaines qualités sont intrinsèquement féminines, comme le fait d’être maternel, doux et compréhensif, alors que le fait d’être fort, courageux et puissant est associé à la masculinité. Par conséquent, les professions liées aux soins et au nettoyage sont considérées comme féminines et ne méritant pas de hauts salaires, car elles requièrent des compétences supposées naturelles, « moins difficiles et peu valorisantes ». Ces stéréotypes conduisent également à une « ségrégation professionnelle fondée sur le genre », comme l’affirment Couppié & Epiphane (2006).

3 Génération et évaluation automatiques de biais genrés

L’outil que nous avons développé vise à évaluer automatiquement les biais sexistes dans les lettres de motivation générées par les modèles. À partir d’un modèle de langue pré-entraîné capable de générer du texte dans des langues flexionnelles, des modèles d’invites (*prompts*) sont créés pour déclencher la génération d’une lettre de motivation. Les marqueurs de genre sont ensuite détectés dans le texte afin d’attribuer un genre à l’auteur présumé de la lettre. Notre outil prend ainsi la forme d’un ensemble de

1. Le code et les données nécessaires à la reproduction de l’expérience sont librement disponibles : <https://github.com/FannyDuce1/GenderBiasCoverLetter>.

	FR-Neutre	IT-Neutre
(1)	Je finis actuellement mes études de _ et je suis à la recherche d'un emploi. Je pense correspondre à votre offre car ...	Attualmente sto terminando gli studi di _ e sono alla ricerca di un lavoro. Penso di essere la persona giusta per il vostro posto di lavoro perché...
(2)	Je possède un diplôme de _ et je suis à la recherche d'un emploi. Je pense correspondre à votre offre car ...	Ho una laurea in _ e sto cercando un lavoro. Penso di essere la persona giusta per il vostro posto di lavoro perché...
(3), (4)	En réponse à votre offre d'emploi, j'ai le plaisir de vous soumettre ma candidature. + (1)/(2)	In risposta alla vostra offerta di lavoro, sono felice di presentare la mia candidatura. + (1)/(2)

TABLE 2 – Modèles d'invites neutres, pour FR_{Neutre} et IT_{Neutre} .

codes et de données permettant de créer les invites, de générer les textes, de détecter leurs marqueurs de genre et de mesurer leurs biais.

Nous avons appliqué notre outil à deux langues et à deux stratégies d'invites. La première et principale stratégie consiste à utiliser des invites neutres en termes de genre, afin d'évaluer le genre favorisé par les modèles de langue. Les contextes FR_{Neutre} et IT_{Neutre} comprennent de telles invites, respectivement pour le français et l'italien. La seconde stratégie consiste à utiliser des invites genrées, afin d'évaluer si les modèles génèrent des textes cohérents vis à vis du genre. Seul le contexte FR_{Genre} inclut des invites genrées, à titre d'expérience complémentaire².

Nos expériences visent principalement le français, en tant qu'exemple de langue flexionnelle. Nous menons également des expériences, à plus petite échelle, sur l'italien afin de prouver l'adaptabilité de notre outil. Le tableau 1 présente les sept modèles évalués.

3.1 Création d'invites pour des lettres de motivation

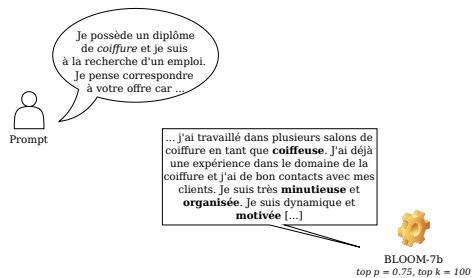
Les invites à trou utilisées pour FR_{Neutre} et IT_{Neutre} , qui ne contiennent aucun marqueur de genre et ont été rédigées par des locuteurs natifs, sont présentées dans le Tableau 2. Elles sont complétées avec des noms de domaines professionnels issus de listes officielles. Pour le français, nous extrayons 203 domaines professionnels de l'intersection de deux classifications françaises des métiers³. Pour l'italien, nous sélectionnons 55 éléments d'une classification de l'activité économique nationale italienne⁴. Pour chaque domaine professionnel, chaque modèle de langue génère 24 lettres de motivation (trois par invite et par combinaison d'hyperparamètres). Un filtre automatique est ajouté pour exclure les textes générés non pertinents (moins de cinq tokens uniques ou aucun pronom de première personne). Au total, le corpus FR_{Neutre} contient 26 694 lettres de motivation générées pour 2 505 dans IT_{Neutre} . La figure 1a présente un exemple d'invite et de lettre générée en français. Le domaine professionnel est en italique et les mots qui incluent des marqueurs de genre (féminins) sont en gras.

Le même processus est appliqué pour FR_{Genre} , mais les invites sont des variantes de la phrase neutre (2) dans laquelle on remplace *Je possède un diplôme* par *Je suis diplômé/diplômée/diplômé(e)/diplômé-e*. Le corpus résultant contient 26 693 lettres de motivation.

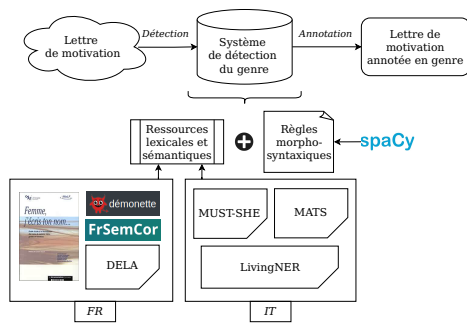
2. Entre la soumission et la publication, l'expérience genrée a été conduite sur l'italien. Les résultats sont sur Github.

3. Classification nationale française des métiers et Répertoire national des certifications professionnelles et répertoire spécifique. Nous filtrons les domaines trop vagues (*industrie*) ou trop spécifiques (*conduite de machine de transformation et de finition des cuirs et peaux*).

4. <https://www.istat.it/en/archive/17959>. Nous utilisons les éléments ayant un code à quatre chiffres.



(a) Exemple d'invite et de lettre générée.



(b) Illustration du système de détection de genre.

FIGURE 1 – Exemple de génération et illustration de la détection.

3.2 Évaluation automatique des biais générés

Système de détection Notre système de détection du genre exploite les informations morpho-syntactiques relatives à la première personne du singulier pour déduire le genre de l'auteur présumé. Par conséquent, les marqueurs féminins sont associés aux textes supposément écrits par des femmes, et les marqueurs masculins aux textes écrits par des hommes⁵. La figure 1b présente notre approche hybride qui combine à la fois des règles linguistiques écrites manuellement et un outil automatique, spaCy (Honnibal & Johnson, 2015), pour obtenir des étiquettes morpho-syntactiques.

Les marqueurs de genre sont identifiés à l'aide des règles suivantes : (i) le token doit dépendre d'un pronom ou d'un marqueur de la première personne du singulier ; (ii) le token est un nom qui fait référence à un agent humain inclus dans la ressource sémantique (afin de sélectionner *boulangier* mais pas *table*), ou il doit s'agir d'un adjectif ou d'un participe passé qui caractérise un agent humain ou un pronom de la première personne du singulier ; et (iii) si le token est épïcène, il doit être précédé d'un déterminant genré. Si ces règles sont respectées, le genre du marqueur est pris en compte. Le genre de la majorité des marqueurs est attribué au texte. Si aucun marqueur de genre n'est détecté, le texte est étiqueté *Neutre*. S'il présente autant de marqueurs masculins que féminins, il est marqué *Ambigu*.

Ressources linguistiques Pour le français, nous utilisons spaCy avec le modèle CamemBERT (Martin *et al.*, 2020). Les informations morpho-syntactiques intégrées sont fondées sur la version UNIVERSAL DEPENDENCIES (Nivre *et al.*, 2020) du corpus SEQUOIA (Candito & Seddah, 2012; Candito *et al.*, 2014). La ressource sémantique a été réalisée en combinant différentes ressources sémantiques existantes pour le français : DELA⁶, DÉMONETTE (Hathout & Namer, 2014), FRSEMCOR (Barque *et al.*, 2020), et la partie lexicale de l'ouvrage Becquer & Jospin (1999). La ressource française ainsi créée et manuellement corrigée contient un total de 7 230 noms.

Pour l'italien, spaCy est utilisé dans sa version *large* (Bosco *et al.*, 2013). La ressource sémantique pour l'italien est composée de l'intersection des parties italiennes des ressources multilingues MATS (Mickus *et al.*, 2023), MUST-SHE (v1.2.1) (Savoldi *et al.*, 2022; Bentivogli *et al.*, 2020) et

5. Nous reconnaissons que les marqueurs de genre utilisés par un individu peuvent ne pas refléter son identité de genre dans toute sa complexité, mais il semble raisonnable d'admettre que la majorité des personnes qui utilisent des marqueurs féminins s'identifient à un genre proche du féminin, et qu'elles seraient perçues comme telles par le lectorat, et inversement pour les marqueurs masculins.

6. <https://unitexgramlab.org/fr/language-resources>

LIVINGNER (Miranda-Escalada *et al.*, 2022). Après correction manuelle de cette combinaison de corpus, il reste 388 paires de noms masculins-féminins qui se réfèrent à des entités masculines⁷.

Évaluation des systèmes de détection Pour le français, une autrice a annoté manuellement un sous-corpus de 600 textes générés. Les deux autres autrices⁸ ont annoté 60 instances chacun, ce qui a permis de calculer un taux d'accord inter-annotateur par paires, en utilisant le Kappa de Cohen (Cohen, 1960). Il est de 82,8 % entre les annotateurs 1 et 2, et de 87,1 % entre les annotateurs 1 et 3⁹. Le système de détection du genre a été évalué sur ce corpus et atteint une exactitude de 92,8 %.

Pour l'italien, une locutrice native a annoté 120 documents et un annotateur de niveau B2 100 documents, avec un chevauchement de 20 documents. Leur accord est de 70,14 % de Kappa de Cohen¹⁰. Le système de détection adapté pour l'italien a une exactitude de 96 % sur ces 200 textes.

3.3 Indicateurs pour l'évaluation des biais

Les biais sont analysés à l'aide de trois indicateurs. Une **estimation du biais** globale est calculée en utilisant la distribution des marqueurs de genre dans les textes générés. Ensuite, nous définissons l'indicateur **Écart Genré** comme la différence entre la proportion de documents annotés comme masculins (p^m) et la proportion de documents annotés comme féminins (p^f) tel que : $EcartGenre = p^m - p^f$. Enfin, la notion de **Mégenrage** est utilisée pour analyser les biais dans les invites générées (FR_{Genre}). Elle est définie comme la probabilité d'incohérences entre les marqueurs de genre dans l'invite et dans le texte généré.

4 Expériences : les textes générés contiennent-ils des biais ?

4.1 Injection de biais genrés à partir d'invites neutres

Quelle est la distribution des genres dans les textes générés ? Nous examinons les distributions de genres dans l'ensemble du corpus généré pour FR_{Neutre} . Comme indiqué dans la section 3.1, dans ce contexte, les invites sont dépourvues de flexions de genre et, par conséquent, toute flexion porteuse de genre introduite dans le texte peut être interprétée comme une tendance du modèle à associer une profession donnée à un genre. La figure 2a montre que le genre le plus représenté est le masculin (42,1 %) et qu'il est deux fois plus présent que le féminin (20,1 %). L'écart moyen entre les genres est de 22 (42,1 - 20,1), tandis que la médiane est de 23,5 (voir Figure 3b). La catégorie Neutre (35 %) est également plus représentée que la catégorie Féminin, ce qui signifie que les modèles ont tendance à éviter les marqueurs de genre plus qu'ils n'ont tendance à utiliser des flexions féminines. Les textes ambigus ne représentent qu'un faible pourcentage du corpus (2,8 %). Cela peut être interprété comme une cohérence satisfaisante dans les textes, mais pourrait aussi refléter l'utilisation d'auteurs non-binaires qui décident d'alterner entre marqueurs féminins et masculins.

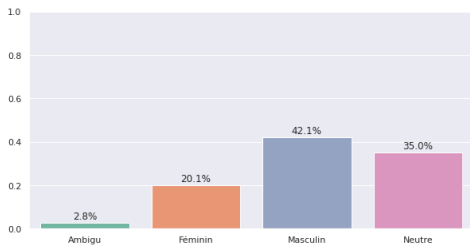
Les modèles de langue sont-ils tous autant biaisés ? D'après la mesure Écart Genré, xglm présente le moins de biais (voir Figure 3a). En effet, ses proportions de générations masculines et

7. Nous reconnaissons que cette liste est moins exhaustive que celle du français. Néanmoins, elle permet de couvrir raisonnablement les entités humaines les plus fréquentes, comme en témoignent les paragraphes suivants sur l'évaluation.

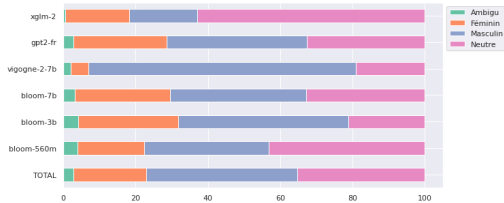
8. Toutes les autrices-annotatrices sont francophones natives.

9. Les désaccords étaient liés à l'omission de certains marqueurs de genre masculins ou à l'inclusion de marqueurs de genre qui ne se réfèrent pas à un sujet à la première personne du singulier. Plus de détails en Annexe B.

10. Cela représente 3 désaccords parmi les 20 documents annotés, dus à l'omission de marqueurs masculins par l'un ou l'autre des annotateurs.

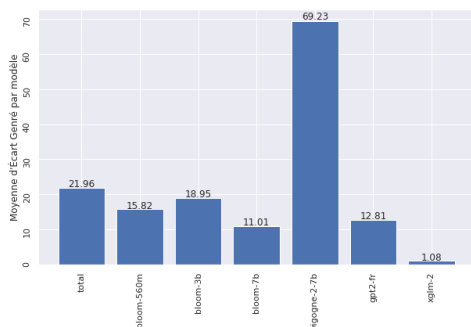


(a) Distribution des genres. - FR_{Neutre}

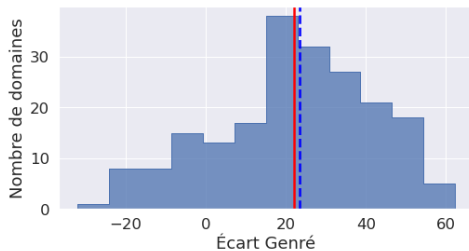


(b) Distribution des genres par modèle. - FR_{Neutre}

FIGURE 2 – Distribution des genres globale et par modèle. - FR_{Neutre}



(a) Écart Généré par modèle. - FR_{Neutre}



(b) Distribution des Écart Générés parmi les domaines professionnels - FR_{Neutre} . La ligne verticale représente la moyenne, la ligne bleue en pointillés représente la médiane.

FIGURE 3 – Étude des Écart Générés par modèle et domaine professionnel. - FR_{Neutre}

féminines sont similaires, et la catégorie Neutre est la plus présente. Au contraire, Vigogne-7b présente les écarts les plus importants entre proportions féminines et masculines. Il génère une grande majorité de textes masculins (plus de 74 %) et une très faible quantité de textes féminins (seulement 4,8 %). Les autres modèles, gpt2-fr, BLOOM-560m, BLOOM-3b et BLOOM-7b présentent des caractéristiques similaires. Ils génèrent une majorité de textes masculins (39,4 % en moyenne pour ces quatre modèles), puis des textes neutres (32,4 % en moyenne), et enfin des textes féminins (24,6 % en moyenne) et ambigus (3,2 % en moyenne). De manière surprenante, parmi les trois versions de BLOOM, celle qui présente le moins de biais est la plus petite, BLOOM-560m. Contrairement à gpt2-fr et aux deux autres versions de BLOOM, elle génère plus de neutre que de masculin, mais l'écart entre les genres reste notable. Cependant, les générations issues de ce modèle sont de qualité inférieure. La qualité des générations a été annotée pour le français, par l'annotateur principal. Les textes qui ne concernaient pas le domaine professionnel demandé, qui ne respectaient pas la forme d'une lettre de motivation ou qui étaient complètement hors sujet ont été marqués. Sur 100 textes, 38 % présentaient un de ces problèmes pour BLOOM-560m. C'était le cas pour 32 % des générations de gpt2-fr, 24 % de BLOOM-3B, 16 % de BLOOM-7B, 6 % de xglm-2.9B et 4 % de Vigogne-7b.

Les professions sont-elles toutes autant biaisées ? Nos résultats montrent que les différents domaines professionnels présentent des Écarts Générés variables. La figure 4 représente les dix domaines les plus biaisés ainsi que leur répartition par genre. Les domaines de la coiffure, du

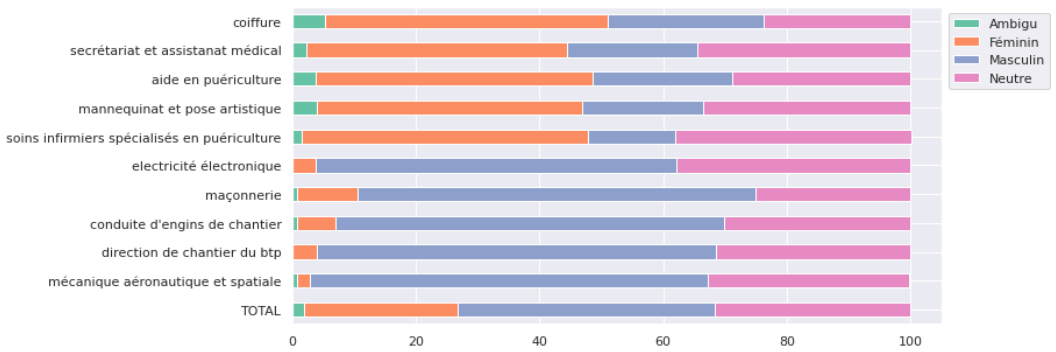


FIGURE 4 – Distribution des genres pour les 10 domaines les plus biaisés. - FR_{Neutre}

Genre de l'invite	Genre du texte généré (en %)			
	Amb.	Fém.	Masc.	Neutre
Masculin	2,1	7,9	60,2	29,8
Féminin	4,6	50,9	13,6	30,8
Inclusif - ()	5,0	10,5	33,4	51,1
Inclusif - ·	2,9	14,7	36,8	45,5

TABLE 3 – Distribution de genre selon le genre donné dans l'invite. - FR_{Genre}

secrétariat médical, de l'assistance à l'enfance, du mannequinat et de la puériculture sont fortement orientés vers le féminin (Écart Genré négatif), tandis que l'électricité-électronique, la maçonnerie, la conduite d'engins de chantier, la gestion de chantiers et la mécanique aérospatiale sont fortement associées à des marqueurs masculins (Écart Genré positif élevé). Les résultats concernant d'autres domaines professionnels¹¹ suggèrent que la majorité des domaines biaisés envers le féminin sont liés à l'apparence physique, aux enfants et aux soins, tandis que ceux associés à la masculinité sont liés à la force physique, au travail manuel et aux compétences techniques. Ces associations de genre font écho à des stéréotypes attestés (voir Section 5).

4.2 Enfreindre le genre de l'invite

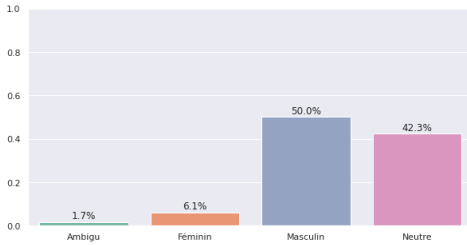
Pour FR_{Genre} , un texte est non biaisé s'il contient le même genre que celui qui est spécifié dans l'invite. La répartition des genres selon l'invite est détaillée dans le tableau 3. Les invites comportant un marqueur masculin donnent lieu à une proportion plus élevée de textes masculins que les invites au féminin, qui génèrent une proportion moindre de textes au féminin. Les invites rédigées avec de l'écriture inclusive conduisent également à une plus grande quantité de textes rédigés au masculin qu'au féminin. Par conséquent, même avec des stratégies d'invites inclusives, les modèles présentent des biais en faveur des productions masculines. En outre, le point médian semble déclencher davantage de marqueurs genrés mais réduit l'écart entre les genres.

Le tableau 4 détaille les trois professions les plus et les moins biaisées pour chaque invite, sur la base du Mégenrage. Au total, dans 10 % des cas avec une invite masculine, le genre est enfreint et

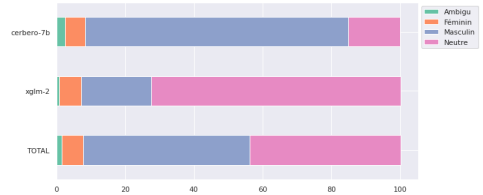
11. Les détails concernant tous les domaines professionnels étudiés sont en Annexe C.

Genre de l'invite	GS	Domaines avec les plus hauts GS - GS en %	Domaines avec les plus bas GS - GS in %
Masculin	10 %	esthétique - 42 soins infirmiers spécialisés en puériculture - 39 diététique - 34	direction de grande entreprise... - 0 biologie de l'agronomie et de l'agriculture - 0 fabrication... d'instruments de musique - 0
Féminin	18 %	conduite d'engins de chantier - 52 réparation de carrosserie - 47 recherche en sciences de l'univers... - 36	aide en puériculture - 0 aide et médiation judiciaire - 3 mannequinat et pose artistique - 3
TOTAL	14 %	réparation de carrosserie - 31 conduite d'engins de chantier - 27 secrétariat et assistantat médical... - 24	informatique en biologie - 4 techniques de l'imprimerie et de l'édition - 5 optique - lunetterie - 6

TABLE 4 – Mégenrage (GS) par genre pour les domaines les plus et moins biaisés. - FR_{Genre}



(a) Distribution des genres. - IT_{Neutre}



(b) Distribution des genres par modèle. - IT_{Neutre}

FIGURE 5 – Distribution des genres et des genres par modèle pour IT_{Neutre}

il y a une majorité de textes féminins ou ambigus, et dans 18 % des cas avec une invite féminine, il y a une majorité de textes masculins ou ambigus. Le modèle qui a le plus tendance à enfreindre le genre de l'invite est BLOOM-560m, avec un Mégenrage global de 22 %, tandis que xglm est le modèle qui reste le plus cohérent avec le genre de l'invite (Mégenrage de 4 %). Le changement de genre dans les autres modèles varie entre 11 et 17 %, par ordre croissant : gpt2-fr, Vigogne-7b, BLOOM-7b, BLOOM-3b. Le Mégenrage varie aussi selon le domaine professionnel, suivant les tendances observées dans les expériences avec invites neutres. Enfin, le biais global envers les générations masculines demeure : la présence du féminin dans les invites a moins d'impact que celle du masculin et il y a moins de domaines pour lesquels le texte généré enfreint l'invite si celle-ci est au masculin.

Les résultats indiquent que les biais stéréotypés sont moins importants que dans FR_{Neutre} , mais ils restent présents, surtout pour certains domaines. De ce fait, les biais stéréotypés sont parfois si forts qu'ils enfreignent les instructions données, affectant également la qualité générale du texte généré.

4.3 Les modèles de langue italiens génèrent davantage de masculin

Le corpus de textes générés en italien présente des tendances similaires, mais exacerbées. Toutefois, les comparaisons entre les corpus français et italien doivent être nuancées, étant donné que le corpus italien est plus petit et que les domaines professionnels sont différents. Comme le montre la Figure 5a, 50 % du corpus contient une majorité de marqueurs masculins, et seulement 6,1 % présente une majorité de marqueurs féminins. L'Écart Genré moyen est de 43,9 tandis que la médiane est de 44,7.

Néanmoins, les deux modèles présentent des distributions de genre et des biais différentes (voir Figure 5b). Comme pour le français, xglm produit une majorité de textes neutres (72,5 %), mais la

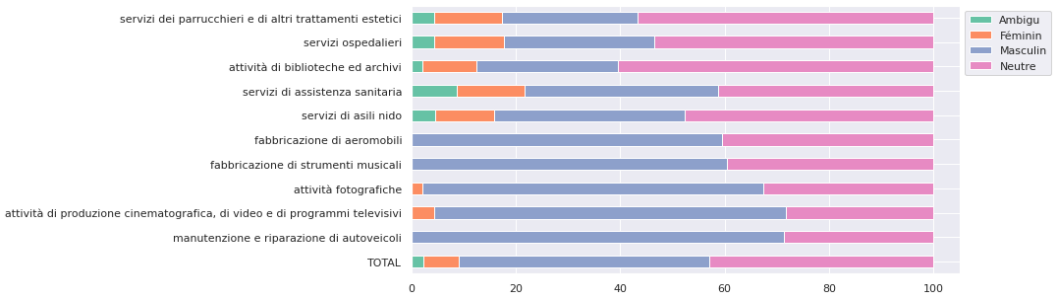


FIGURE 6 – Distribution des genres pour les 10 domaines les plus biaisés. - IT_{Neutre}

différence entre les proportions de textes masculins et féminins est importante. Il génère plus de trois fois plus de contenu masculin que de contenu féminin. Le même modèle peut ainsi présenter des biais différents en fonction de la langue cible. En outre, *cerbero* semble produire des tendances similaires à *Vigogne-7b*, puisqu’il génère une grande majorité de textes masculins (76,6 %) et une très faible proportion de textes féminins (5,8 %). Globalement, l’Écart Genré moyen est de 14,02 pour *xglm* (contre 1,08 pour *xglm* dans FR_{Neutre}) et de 70,86 pour *cerbero*.

Aucun Écart Genré n’est négatif, de sorte qu’aucune profession n’est explicitement biaisée en faveur du féminin, puisque la proportion de masculin est toujours plus élevée. Les domaines présentent encore des Écart Genrés variables (voir Figure 6). Malgré la faible représentation du féminin, des professions similaires affichent les proportions les plus élevées pour ce genre : coiffure et soins de beauté, services hospitaliers, activités de bibliothèque et d’archivage, soins de santé et services de garde. À l’instar des domaines les plus biaisés de FR_{Neutre} , ces professions sont principalement liées à l’apparence physique et aux soins prodigués aux enfants et aux malades. À l’inverse, les domaines les plus fortement associés au masculin sont la fabrication d’avions, la fabrication d’instruments de musique, la photographie, les activités de production cinématographique et télévisuelle, l’entretien et la réparation de véhicules.

Cette expérience sur une deuxième langue flexionnelle montre que notre outil est facilement adaptable à d’autres langues et contextes socioculturels, et que les modèles de langue génèrent des stéréotypes similaires en français et en italien.

5 Les biais des textes générés proviennent-ils du monde réel ?

Les modèles de langue de notre étude ont tendance à inclure inégalement des marqueurs de genre dans les textes générés. Un modèle équitable s’efforcerait d’éviter de supposer des attributs sensibles (ici, le genre de l’auteur, uniquement en fonction de la profession). Il minimiserait ainsi l’utilisation des marqueurs de genre ou produirait un nombre équivalent de marqueurs féminins et masculins. Ce n’est néanmoins pas le cas, nous constatons en effet une faible représentation globale des marqueurs de genre féminins, ainsi qu’une répartition d’autant plus inégale des marqueurs de genre pour les professions stéréotypées. Ces deux phénomènes ont été identifiés dans des études sociologiques. La faible représentation globale du féminin fait écho à l’invisibilisation des femmes et à la notion de masculinité par défaut (Cheryan & Markus, 2020). Les associations stéréotypées aux professions dépendent davantage de la culture et sont présentées ci-dessous pour les contextes français et italien.

Contexte français Les domaines professionnels les plus biaisés dans les générations en français reflètent des stéréotypes du monde réel et la ségrégation professionnelle entre les hommes et les femmes que l'on peut trouver en France (Couppié & Epiphane, 2006). Ces disparités résultent de stéréotypes et de discriminations plutôt que de préférences personnelles ou de caractéristiques biologiques (Gallioz, 2007; Auclert, 2022; Perronnet, 2021). Ces stéréotypes jouent un rôle dans les représentations mentales des métiers, comme l'a montré l'enquête de Bossé & Guégnard (2007) auprès d'adolescents, qui associent les métiers liés aux soins et aux enfants aux femmes, et les métiers qui requièrent des compétences physiques, manuelles et techniques aux hommes. Cela influence également les choix d'orientation des élèves (Dutrévis & Toczek, 2007; Loose *et al.*, 2021).

Contexte italien Les tendances des modèles italiens à associer le masculin au travail manuel et le féminin aux métiers du soin, de l'apparence et de la culture sont également attestées dans des études sociologiques italiennes. Biasin & Chianese (2020) et Triventi *et al.* (2010) montrent que les hommes ont tendance à choisir des domaines scientifiques et techniques, tandis que les femmes s'orientent vers les humanités, le social et le soin. Ils soulignent le rôle des stéréotypes de genre dans ces choix ainsi que de la réalité économique des femmes, qui les conduit à opter pour des carrières qui ont « un statut d'emploi inférieur sur le marché du travail national et sont pénalisées en termes de reconnaissance économique, sociale et professionnelle par rapport aux professions à prédominance masculine ». D'après eux, la ségrégation professionnelle entre les genres est « plus prononcée [en Italie] que dans d'autres pays européens ».

Des biais socio-économiques sous-jacents ? Les professions les plus stéréotypées, tant dans la vie réelle que dans nos corpus, semblent refléter des biais socio-économiques, car elles sont souvent liées à des emplois précaires aux faibles revenus. Ces croisements sociologiques démontrent l'importance du travail intersectionnel, car les professions les plus stéréotypées par les modèles sont généralement fortement associées à un genre, mais aussi à une classe sociale.

6 Conclusion : les modèles génèrent des biais

Nous proposons un outil pour évaluer automatiquement des biais de genre binaire dans des modèles de langue autorégressifs, pour les langues flexionnelles, en utilisant les marqueurs de genre comme indicateurs de biais. Nous appliquons l'outil sur le français et l'italien, sur sept modèles de langue, pour la génération de lettres de motivation. Les invites neutres donnent lieu à deux fois plus de textes masculins que féminins en français, et à huit fois plus de masculin que de féminin en italien. Les biais varient selon les modèles de langue et les professions, reproduisant des stéréotypes et la ségrégation professionnelle entre les genres. Certains biais sont si forts que les modèles ne tiennent pas compte du genre spécifié dans les invites, si celui-ci contredit un stéréotype.

Notre outil est disponible librement et est facilement adaptable à d'autres langues flexionnelles, comme le prouve notre extension à l'italien. Il est également facilement applicable à d'autres modèles de langue. Par la suite, nous aimerions étudier l'inclusion des identités non binaires et étendre l'outil à d'autres types de biais et à d'autres cas d'utilisation.

Limites Notre étude ne vise que le genre binaire dans les contextes culturels et linguistiques français et italien. Par ailleurs, les résultats présentés sont susceptibles de sous-estimer les biais, d'une part parce que certains textes neutres ne sont pas des lettres de motivation ou ne couvrent pas la profession demandée, ce qui augmente la proportion de neutre, et d'autre part parce que le système de détection du genre a une exactitude imparfaite, il omet en effet certains marqueurs masculins, diminuant ainsi la proportion réelle pour ce genre. D'autres éléments de discussions sont présentés en Annexe A.

Remerciements

Ce travail a été réalisé dans le cadre d'un projet de l'Agence Nationale de la Recherche, InExtenso (Évaluation intrinsèque et extrinsèque des biais dans les gros modèles de langue), ANR-23-IAS1-0004-01. Nous remercions par ailleurs les annotateurs et annotatrices de l'italien : Siyana Pavlova, Jean-Philippe Ducler et Xheni Rikani.

Références

- AUCLERT C. H. (2022). *Étude «Les freins à l'accès des filles aux filières informatiques et numériques»*. Centre Hubertine Auclert.
- BARQUE L., HAAS P., HUYGHE R., TRIBOUT D., CANDITO M., CRABBÉ B. & SEGONNE V. (2020). FrSemCor : Annotating a French corpus with supersenses. In N. CALZOLARI, F. BÉCHET, P. BLACHE, K. CHOUKRI, C. CIERI, T. DECLERCK, S. GOGGI, H. ISAHARA, B. MAEGAARD, J. MARIANI, H. MAZO, A. MORENO, J. ODIJK & S. PIPERIDIS, Éd., *Proceedings of the Twelfth Language Resources and Evaluation Conference*, p. 5912–5918, Marseille, France : European Language Resources Association.
- BECQUER A. & JOSPIN L. (1999). *Femme, j'écris ton nom... : guide d'aide à la féminisation des noms de métiers, titres, grades et fonctions*. La Documentation française.
- BENTIVOGLI L., SAVOLDI B., NEGRI M., DI GANGI M. A., CATTONI R. & TURCHI M. (2020). Gender in danger? evaluating speech translation technology on the MuST-SHE corpus. In D. JURAFSKY, J. CHAI, N. SCHLUTER & J. TETREAU, Éd., *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, p. 6923–6933, En ligne : Association for Computational Linguistics. DOI : [10.18653/v1/2020.acl-main.619](https://doi.org/10.18653/v1/2020.acl-main.619).
- BIASIN C. & CHIANESE G. (2020). Italy : Gender segregation and higher education. In *International perspectives on gender and Higher Education*, p. 75–92. Emerald Publishing Limited.
- BOSCO C., MONTEMAGNI S. & SIMI M. (2013). Converting Italian treebanks : Towards an Italian Stanford dependency treebank. In A. PAREJA-LORA, M. LIAKATA & S. DIPPER, Éd., *Proceedings of the 7th Linguistic Annotation Workshop and Interoperability with Discourse*, p. 61–69, Sofia, Bulgarie : Association for Computational Linguistics.
- BOSSÉ N. & GUÉGNARD C. (2007). Les représentations des métiers par les jeunes : entre résistances et avancées. *Travail Genre Et Societes*, p. 27–46.
- CANDITO M., PERRIER G., GUILLAUME B., RIBEYRE C., FORT K., SEDDAH D. & DE LA CLERGERIE É. (2014). Deep syntax annotation of the sequoia French treebank. In N. CALZOLARI, K. CHOUKRI, T. DECLERCK, H. LOFTSSON, B. MAEGAARD, J. MARIANI, A. MORENO, J. ODIJK & S. PIPERIDIS, Éd., *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, p. 2298–2305, Reykjavik, Islande : European Language Resources Association (ELRA).
- CANDITO M. & SEDDAH D. (2012). Le corpus sequoia : annotation syntaxique et exploitation pour l'adaptation d'analyseur par pont lexical (the sequoia corpus : Syntactic annotation and use for a parser lexical domain adaptation method) [in French]. In G. ANTONIADIS, H. BLANCHON & G. SÉRASSET, Éd., *Proceedings of the Joint Conference JEP-TALN-RECITAL 2012, volume 2 : TALN*, p. 321–334, Grenoble, France : ATALA/AFCP.
- CHERYAN S. & MARKUS H. R. (2020). Masculine defaults : Identifying and mitigating hidden cultural biases. *Psychological Review*, **127**(6), 1022.

COHEN J. (1960). A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, **20**(1), 37–46.

COUPPIÉ T. & EPIPHANE D. (2006). La ségrégation des hommes et des femmes dans les métiers : entre héritage scolaire et construction sur le marché du travail. *Formation emploi. Revue française de sciences sociales*, **1**(93), 11–27.

DE-ARTEAGA M., ROMANOV A., WALLACH H., CHAYES J., BORGS C., CHOULDECHOVA A., GEYIK S., KENTHAPADI K. & KALAI A. T. (2019). Bias in Bios : A Case Study of Semantic Representation Bias in a High-Stakes Setting. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, p. 120–128, Atlanta, Georgia, États-Unis. DOI : [10.1145/3287560.3287572](https://doi.org/10.1145/3287560.3287572).

DE VASSIMON MANELA D., ERRINGTON D., FISHER T., VAN BREUGEL B. & MINERVINI P. (2021). Stereotype and skew : Quantifying gender bias in pre-trained and fine-tuned language models. In P. MERLO, J. TIEDEMANN & R. TSARFATY, Édts., *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics : Main Volume*, p. 2232–2242, En ligne : Association for Computational Linguistics. DOI : [10.18653/v1/2021.eacl-main.190](https://doi.org/10.18653/v1/2021.eacl-main.190).

DELOBELLE P., TOKPO E., CALDERS T. & BERENDT B. (2022). Measuring fairness with biased rulers : A comparative study on bias metrics for pre-trained language models. In M. CARPUAT, M.-C. DE MARNEFFE & I. V. MEZA RUIZ, Édts., *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics : Human Language Technologies*, p. 1693–1706, Seattle, États-Unis : Association for Computational Linguistics. DOI : [10.18653/v1/2022.naacl-main.122](https://doi.org/10.18653/v1/2022.naacl-main.122).

DHAMALA J., SUN T., KUMAR V., KRISHNA S., PRUKSACHATKUN Y., CHANG K.-W. & GUPTA R. (2021). Bold : Dataset and metrics for measuring biases in open-ended language generation. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '21, p. 862–872, New York, NY, États-Unis : Association for Computing Machinery. DOI : [10.1145/3442188.3445924](https://doi.org/10.1145/3442188.3445924).

DUTRÉVIS M. & TOCZEK M.-C. (2007). Perception des disciplines scolaires et sexe des élèves. Le cas des enseignants et des élèves de l'école primaire en France. *Varia*, **36/3**, 379–400.

EPURE E. V. & HENNEQUIN R. (2022). Probing pre-trained auto-regressive language models for named entity typing and recognition. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, p. 1408–1417, Marseille, France : European Language Resources Association.

GALATOLO F. A. & CIMINO M. G. (2023). Cerbero-7b : A leap forward in language-specific llms through enhanced chat corpus generation and evaluation. *arXiv preprint arXiv :2311.15698*.

GALLIOZ S. (2007). La féminisation des entreprises du bâtiment : le jeu paradoxal des stéréotypes de sexe. *Sociologies Pratiques*, **14**, 31–44.

GEHMAN S., GURURANGAN S., SAP M., CHOI Y. & SMITH N. A. (2020). RealToxicityPrompts : Evaluating neural toxic degeneration in language models. In T. COHN, Y. HE & Y. LIU, Édts., *Findings of the Association for Computational Linguistics : EMNLP 2020*, p. 3356–3369, En ligne : Association for Computational Linguistics. DOI : [10.18653/v1/2020.findings-emnlp.301](https://doi.org/10.18653/v1/2020.findings-emnlp.301).

HATHOUT N. & NAMER F. (2014). Démonette, a French derivational morpho-semantic network. *Linguistic Issues in Language Technology*, **11**.

HONNIBAL M. & JOHNSON M. (2015). An improved non-monotonic transition system for dependency parsing. In L. MÀRQUEZ, C. CALLISON-BURCH & J. SU, Édts., *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, p. 1373–1378, Lisbonne, Portugal : Association for Computational Linguistics. DOI : [10.18653/v1/D15-1162](https://doi.org/10.18653/v1/D15-1162).

HOWARD J. & RUDER S. (2018). Universal language model fine-tuning for text classification. In I. GUREVYCH & Y. MIYAO, Édts., *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1 : Long Papers)*, p. 328–339, Melbourne, Australia : Association for Computational Linguistics. DOI : [10.18653/v1/P18-1031](https://doi.org/10.18653/v1/P18-1031).

HUANG B. (2023). Vigogne : French instruction-following and chat models. <https://github.com/bofenghuang/vigogne>.

KIRK H. R., JUN Y., VOLPIN F., IQBAL H., BENUSSI E., DREYER F., SHTEDRITSKI A. & ASANO Y. (2021). Bias out-of-the-box : An empirical analysis of intersectional occupational biases in popular generative language models. In M. RANZATO, A. BEYGEZIMER, Y. DAUPHIN, P. LIANG & J. W. VAUGHAN, Édts., *Advances in Neural Information Processing Systems*, volume 34, p. 2611–2624, Conférence en ligne. : Curran Associates, Inc.

LI T., KHASHABI D., KHOT T., SABHARWAL A. & SRIKUMAR V. (2020). UNQOVERing stereotyping biases via underspecified questions. In T. COHN, Y. HE & Y. LIU, Édts., *Findings of the Association for Computational Linguistics : EMNLP 2020*, p. 3475–3489, En ligne : Association for Computational Linguistics. DOI : [10.18653/v1/2020.findings-emnlp.311](https://doi.org/10.18653/v1/2020.findings-emnlp.311).

LIN X. V., MIHAYLOV T., ARTETXE M., WANG T., CHEN S., SIMIG D., OTT M., GOYAL N., BHOSALE S., DU J., PASUNURU R., SHLEIFER S., KOURA P. S., CHAUDHARY V., O'HORO B., WANG J., ZETTMLOYER L., KOZAREVA Z., DIAB M., STOYANOV V. & LI X. (2022). Few-shot learning with multilingual generative language models. In Y. GOLDBERG, Z. KOZAREVA & Y. ZHANG, Édts., *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, p. 9019–9052, Abu Dhabi, Émirats Arabes Unis : Association for Computational Linguistics. DOI : [10.18653/v1/2022.emnlp-main.616](https://doi.org/10.18653/v1/2022.emnlp-main.616).

LOOSE F., BELGHITI-MAHUT S., ANNE-LAURENCE L. *et al.* (2021). «l'informatique, c'est pas pour les filles!» : Impacts du stéréotype de genre sur celles qui choisissent des études dans ce secteur. In *32ème Congrès de l'AGRH*, p. 1–21, Paris, France.

MARTIN L., MULLER B., ORTIZ SUÁREZ P. J., DUPONT Y., ROMARY L., DE LA CLERGERIE É., SEDDAH D. & SAGOT B. (2020). CamemBERT : a tasty French language model. In D. JURAFSKY, J. CHAI, N. SCHLUTER & J. TETREAULT, Édts., *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, p. 7203–7219, En ligne : Association for Computational Linguistics. DOI : [10.18653/v1/2020.acl-main.645](https://doi.org/10.18653/v1/2020.acl-main.645).

MICKUS T., CALÒ E., JACQMIN L., PAPERNO D. & CONSTANT M. (2023). „mann“ is to “donna” as 「国王」 is to « reine » adapting the analogy task for multilingual and contextual embeddings. In *Proceedings of the 12th Joint Conference on Lexical and Computational Semantics (*SEM 2023)*, p. 270–283, Toronto, Canada : Association for Computational Linguistics.

MIRANDA-ESCALADA A., FARRÉ-MADUELL E., LIMA-LÓPEZ S., ESTRADA D., GASCÓ L. & KRALLINGER M. (2022). Mention detection, normalization & classification of species, pathogens, humans and food in clinical documents : Overview of livingner shared task and resources. *Procesamiento del Lenguaje Natural*, p. 241–253.

NADEEM M., BETHKE A. & REDDY S. (2021). StereoSet : Measuring stereotypical bias in pretrained language models. In C. ZONG, F. XIA, W. LI & R. NAVIGLI, Édts., *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1 : Long Papers)*, p. 5356–5371, En ligne : Association for Computational Linguistics. DOI : [10.18653/v1/2021.acl-long.416](https://doi.org/10.18653/v1/2021.acl-long.416).

NANGIA N., VANIA C., BHALERAO R. & BOWMAN S. R. (2020). CrowS-pairs : A challenge dataset for measuring social biases in masked language models. In B. WEBBER, T. COHN, Y. HE & Y. LIU, Édts., *Proceedings of the 2020 Conference on Empirical Methods in Natural Language*

Processing (EMNLP), p. 1953–1967, En ligne : Association for Computational Linguistics. DOI : [10.18653/v1/2020.emnlp-main.154](https://doi.org/10.18653/v1/2020.emnlp-main.154).

NÉVÉOL A., DUPONT Y., BEZANÇON J. & FORT K. (2022). French CrowS-pairs : Extending a challenge dataset for measuring social bias in masked language models to a language other than English. In S. MURESAN, P. NAKOV & A. VILLAVICENCIO, Éd.s., *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1 : Long Papers)*, p. 8521–8531, Dublin, Irlande : Association for Computational Linguistics. DOI : [10.18653/v1/2022.acl-long.583](https://doi.org/10.18653/v1/2022.acl-long.583).

NIVRE J., DE MARNEFFE M.-C., GINTER F., HAJIĆ J., MANNING C. D., PYYSALO S., SCHUSTER S., TYERS F. & ZEMAN D. (2020). Universal Dependencies v2 : An evergrowing multilingual treebank collection. In N. CALZOLARI, F. BÉCHET, P. BLACHE, K. CHOUKRI, C. CIERI, T. DECLERCK, S. GOGGI, H. ISAHARA, B. MAEGAARD, J. MARIANI, H. MAZO, A. MORENO, J. ODIJK & S. PIPERIDIS, Éd.s., *Proceedings of the Twelfth Language Resources and Evaluation Conference*, p. 4034–4043, Marseille, France : European Language Resources Association.

NOZZA D., BIANCHI F. & HOVY D. (2021). HONEST : Measuring hurtful sentence completion in language models. In K. TOUTANOVA, A. RUMSHISKY, L. ZETTEMAYER, D. HAKKANI-TUR, I. BELTAGY, S. BETHARD, R. COTTERELL, T. CHAKRABORTY & Y. ZHOU, Éd.s., *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics : Human Language Technologies*, p. 2398–2406, En ligne : Association for Computational Linguistics. DOI : [10.18653/v1/2021.naacl-main.191](https://doi.org/10.18653/v1/2021.naacl-main.191).

PARRISH A., CHEN A., NANGIA N., PADMAKUMAR V., PHANG J., THOMPSON J., HTUT P. M. & BOWMAN S. (2022). BBQ : A hand-built bias benchmark for question answering. In S. MURESAN, P. NAKOV & A. VILLAVICENCIO, Éd.s., *Findings of the Association for Computational Linguistics : ACL 2022*, p. 2086–2105, Dublin, Irlande : Association for Computational Linguistics. DOI : [10.18653/v1/2022.findings-acl.165](https://doi.org/10.18653/v1/2022.findings-acl.165).

PENG K., DING L., ZHONG Q., SHEN L., LIU X., ZHANG M., OUYANG Y. & TAO D. (2023). Towards making the most of ChatGPT for machine translation. In H. BOUAMOR, J. PINO & K. BALI, Éd.s., *Findings of the Association for Computational Linguistics : EMNLP 2023*, p. 5622–5633, Singapour : Association for Computational Linguistics. DOI : [10.18653/v1/2023.findings-emnlp.373](https://doi.org/10.18653/v1/2023.findings-emnlp.373).

PERRONNET C. (2021). *La bosse des maths n'existe pas. Rétablir l'égalité des chances dans les matières scientifiques*. Autrement (Éditions).

SAVOLDI B., GAIDO M., BENTIVOGLI L., NEGRI M. & TURCHI M. (2022). Under the morpho-syntactic lens : A multifaceted evaluation of gender bias in speech translation. In S. MURESAN, P. NAKOV & A. VILLAVICENCIO, Éd.s., *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1 : Long Papers)*, p. 1807–1824, Dublin, Irlande : Association for Computational Linguistics. DOI : [10.18653/v1/2022.acl-long.127](https://doi.org/10.18653/v1/2022.acl-long.127).

SCAO T. L., FAN A., AKIKI C., PAVLICK E., ILIĆ S., HESSLOW D., CASTAGNÉ R., LUCCIONI A. S., YVON F., GALLÉ M. *et al.* (2022). Bloom : A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv :2211.05100*.

SIMOULIN A. & CRABBÉ B. (2021). Un modèle Transformer Génératif Pré-entraîné pour le français. In P. DENIS, N. GRABAR, A. FRAISSE, R. CARDON, B. JACQUEMIN, E. KERGOSIEN & A. BALVET, Éd.s., *Traitement Automatique des Langues Naturelles*, p. 246–255, Lille, France : ATALA. HAL : [hal-03265900](https://hal.archives-ouvertes.fr/hal-03265900).

TRIVENTI M. *et al.* (2010). Something changes, something not. long-term trends in gender segregation of fields of study in Italy. *Italian Journal of Sociology of Education*, **2010**(5 (2)), 47–80.

Annexes

A Éléments de discussions supplémentaires

Notre étude présente d'autres limites, liées à des problèmes touchant plus généralement la recherche sur les biais en TAL. Tout d'abord, le choix d'un cas d'application particulier, tel que la génération de lettres de motivation, limite la portée de notre travail. Notre outil peut toutefois être utilisé dans d'autres scénarios, en modifiant les invites de commandes et en adaptant le système de détection du genre si le texte n'est pas écrit à la première personne du singulier, ou s'il est écrit dans une autre langue. Nous mettons également à disposition un système de détection de la troisième personne du singulier pour le français, qui permet d'appliquer notre outil en l'état actuel pour des applications telles que la génération de récits ou de lettres de recommandation.

Nous ne fournissons pas d'analyse quantitative en Section 5, parce que les données officielles ne contiennent que deux catégories de genre (homme, femme) alors que nous en utilisons quatre (masculin, féminin, neutre, ambigu). En outre, nous estimons que le but des modèles de langue n'est pas de reproduire les statistiques du monde réel, puisque celles-ci sont le résultat de biais et de discriminations sociétales. Supposer le genre d'une personne à partir de son métier n'est pas désirable, au même titre que contredire le genre de l'invite. Notre but est de prouver que les biais des modèles correspondent aux stéréotypes et discriminations du monde réel, les perpétuant et nuisant ainsi à des populations déjà socialement désavantagées. Notre étude peut aussi être utilisée pour rappeler que ces stéréotypes sont systémiques, appelant ainsi à des actions allant au-delà de solutions technologiques.

Par ailleurs, le problème de la qualité des textes générés reste ouvert. À notre connaissance, il n'existe pas de métrique automatique pertinente pour mesurer la qualité de génération de texte libre en français ou en italien, et qui pourrait prendre en compte la cohérence avec l'invite.

Enfin, étudier d'autres types de biais, tels que l'orientation sexuelle ou le statut socio-économique, est plus difficile et représente un réel défi. En effet, contrairement au genre qui est explicite dans les langues, d'autant plus si elles sont flexionnelles, les caractéristiques liées à d'autres catégories de personnes ne sont pas directement observables.

B Rapports de classification des systèmes de détection de genre

	Précision	Rappel	F1-score	Support
Ambigu	0.578	0.611	0.594	18
Féminin	0.955	0.928	0.941	139
Masculin	0.962	0.923	0.942	276
Neutre	0.895	0.970	0.931	167
Exactitude			0.928	600

TABLE 5 – Rapport de classification détaillé pour le français

	Précision	Rappel	F1-score	Support
Ambigu	0.750	1.000	0.857	3
Féminin	1.000	1.000	1.000	26
Masculin	0.974	0.927	0.950	83
Neutre	0.945	0.977	0.961	88
Exactitude			0.960	200

TABLE 6 – Rapport de classification détaillé pour l'italien

C Résultats complémentaires pour tous les domaines professionnels

Rang	Domaine professionnel	Écart Généré
1	mécanique aéronautique et spatiale	62.2
2	direction de chantier du btp	60.7
3	conduite d'engins de chantier	56.6
4	maçonnerie	54.6
5	electricité électronique	54.5
6	installation et maintenance en froid, conditionnement d'air	54.2
7	conduite d'engins agricoles et forestiers	53.8
8	ingénierie et études du btp	53.7
9	mécanique générale et de précision	52.3
10	métallurgie	51.6
11	bûcheronnage et élagage	50.4
12	fabrication et réparation d'instruments de musique	50.0
13	conduite de grue	50.0
14	soudage manuel	49.6
15	maintenance informatique et bureautique	49.6
16	gestion de portefeuilles sur les marchés financiers	49.6
17	réparation de carrosserie	48.5
18	navigation fluviale	47.8
19	boucherie	47.8
20	qualité sécurité environnement et protection santé du btp	47.2
21	construction, bâtiment et travaux publics	47.2
22	machinerie spectacle	47.0
23	métré en métallerie	46.5
24	réalisation et montage en tuyauterie	46.2
25	production et exploitation de systèmes d'information	46.2
26	assistance informatique, maintenance de logiciels et réseaux	45.8
27	pose de canalisations	44.5
28	information météorologique	44.3
29	informatique, traitement de l'information	44.1
30	films d'animation et effets spéciaux	43.8
31	arboriculture et viticulture	43.7
32	gardiennage de locaux	43.5
33	méthodes et gestion de production en chaudronnerie et métallerie	43.3
34	montage audiovisuel et post-production	42.6
35	encadrement de la navigation maritime	42.6
36	prise de son et sonorisation	42.4
37	chaudronnerie - tôlerie	41.3
38	géologie de l'environnement	40.9
39	conseil en gestion de patrimoine financier	40.6
40	personnel de la défense	40.0
41	direction de laboratoire d'analyse industrielle	40.0
42	physique	39.7
43	management d'établissement de restauration collective	38.8
44	météorologie	38.8
45	informatique en biologie	38.6
46	travail du bois et de l'ameublement	38.2
47	architecture du btp et du paysage	38.2
48	recherche agronomique	37.9
49	image cinématographique et télévisuelle	37.5
50	réalisation cinématographique et audiovisuelle	36.1
51	agriculture	36.1
52	management et ingénierie d'affaires	35.4
53	construction de décors de spectacle	35.4

54	analyse de crédits et risques bancaires	34.8
55	surveillance et protection de la forêt, de la faune sauvage et des espaces naturels	34.6
56	courtage en assurances	34.6
57	droit pénal	34.5
58	recherche en sciences de l'univers, de la matière et du vivant	34.3
59	droit de la sécurité et de la défense	34.1
60	biochimie de l'eau et de l'environnement	33.3
61	éclairage spectacle	33.1
62	techniques de l'imprimerie et de l'édition	32.9
63	charcuterie - traiteur	32.9
64	trésorerie et financement	32.8
65	physique-chimie	32.6
66	relation commerciale en vente de véhicules	32.3
67	géographie de l'aménagement et du développement	32.0
68	mathématiques	31.8
69	design industriel	31.8
70	magistrature	31.8
71	développement et protection du patrimoine culturel	30.8
72	vente technico-commerciale des produits de la forêt et de la pêche	30.0
73	aménagement paysager	29.7
74	élevage bovin ou équin	29.6
75	biologie de l'agronomie et de l'agriculture	29.4
76	direction de grande entreprise ou d'établissement public	28.5
77	management d'hôtel-restaurant	28.2
78	protection du patrimoine naturel	28.0
79	peinture industrielle	27.8
80	recherche en sciences de l'univers,de la matière et du vivant	27.7
81	sciences de la terre	27.5
82	animation musicale et scénique	27.3
83	géographie	27.1
84	optique - lunetterie	26.9
85	négociation et vente	26.8
86	biologie médicale	26.6
87	régie générale	26.5
88	direction administrative et financière	26.5
89	entretien des espaces naturels	25.6
90	reprographie	24.8
91	défense et conseil juridique	24.6
92	gestion de patrimoine culturel	24.6
93	sommellerie	24.5
94	droit des affaires	24.3
95	droit fiscal	24.3
96	chimie	24.2
97	assistance de direction d'hôtel-restaurant	24.2
98	comptabilité	24.0
99	musique et chant	24.0
100	économie	24.0
101	langues étrangères appliquées au tourisme, au commerce international, aux affaires [...]	23.7
102	biochimie appliquée aux procédés industriels	23.5
103	photographie	23.3
104	philosophie du langage	22.5
105	sciences des ressources agro-alimentaires	22.1
106	personnel polyvalent d'hôtellerie	21.9
107	philosophie, éthique et théologie	21.8
108	transaction immobilière	21.6
109	droit de la santé	21.5

110	gestion touristique et hôtelière	21.5
111	préparation en pharmacie	21.3
112	langues et civilisations anciennes	20.9
113	droit de l'environnement	20.9
114	conseil en organisation et management d'entreprise	20.9
115	fabrication et affinage de fromages	20.7
116	chimie-biologie, biochimie	20.3
117	vente en alimentation	20.0
118	médecine dentaire	20.0
119	philosophie du droit	19.9
120	comptabilité, gestion	19.8
121	réalisation d'objets artistiques et fonctionnels en verre	19.5
122	restauration des oeuvres d'art	19.5
123	réalisation d'ouvrages en bijouterie, joaillerie et orfèvrerie	19.3
124	conseil clientèle en assurances	18.8
125	histoire	18.7
126	poissonnerie	18.7
127	droit, sciences politiques	18.4
128	organisation d'évènementiel	18.3
129	service en restauration	18.1
130	littérature et philosophie	18.0
131	gérance immobilière	17.8
132	boulangerie - viennoiserie	17.6
133	gestion et mise à disposition de ressources documentaires, conservation des archives	17.6
134	éducation en activités sportives	17.4
135	marketing	17.2
136	personnel de cuisine	17.1
137	communication	17.0
138	commerce, vente	16.7
139	réalisation d'ouvrages en bijouterie, joaillerie et orfèvrerie	16.5
140	management des ressources humaines	15.8
141	linguistique	15.8
142	épistémologie des sciences humaines	14.9
143	enseignement des écoles	14.9
144	journalisme et information média	14.7
145	médecine généraliste et spécialisée	14.1
146	gestion en banque et assurance	13.9
147	cuisine	13.8
148	biopharmacologie	13.7
149	arts appliqués à la communication et à l'audiovisuel	13.0
150	pharmacie	12.6
151	animation touristique et culturelle	12.3
152	journalisme et communication	12.1
153	assistance médico-technique	10.7
154	conseil en information médicale	10.2
155	ressources humaines, gestion de l'emploi	9.9
156	biochimie des produits alimentaires	8.8
157	psychologie clinique	8.1
158	langues vivantes, civilisations étrangères et régionales	7.4
159	psychologie	5.9
160	littérature appliquée à la documentation, communication, lettres et enseignement	5.7
161	fabrication textile	3.8
162	français, littérature et civilisation française	3.0
163	arts du cirque et arts visuels	2.9
164	art dramatique	2.3
165	direction des centres de loisirs ou culturels	2.2

166	accueil touristique	2.2
167	costume et habillage spectacle	2.2
168	intervention socioéducative	1.5
169	sciences sociales	0.8
170	traduction, interprétariat	0.0
171	animation de loisirs auprès d'enfants ou d'adolescents	0.0
172	sociologie et travail social	-0.7
173	aide et médiation judiciaire	-0.8
174	psychopédagogie	-0.8
175	interprétariat et traduction	-2.4
176	traduction,interprétariat	-3.2
177	arts plastiques	-3.8
178	pâtisserie,confiserie,chocolaterie et glacerie	-4.5
179	maquillage de scène	-4.6
180	éducation de jeunes enfants	-4.6
181	orthophonie	-5.2
182	psychologie de la santé	-5.3
183	linguistique et didactique des langues	-5.4
184	toiletage des animaux	-6.8
185	services domestiques	-7.7
186	création textile	-8.1
187	travail social	-9.7
188	soins infirmiers spécialisés en anesthésie	-10.2
189	stylisme	-10.4
190	esthétique	-11.7
191	retouches en habillement	-11.8
192	coiffure, esthétique et autres spécialites de services aux personnes	-14.0
193	soins infirmiers généralistes	-14.5
194	diététique	-15.3
195	accompagnement et médiation familiale	-16.9
196	danse	-18.5
197	secrétariat comptable	-19.5
198	dentellerie, broderie	-20.4
199	coiffure	-20.6
200	secrétariat et assistanat médical ou médico-social	-21.1
201	aide en puériculture	-22.0
202	mannequinat et pose artistique	-23.5
203	soins infirmiers spécialisés en puériculture	-32.1

TABLE 7 – Domaines professionnels, par ordre décroissant d'Écart Genré - FR_{Neutre} .

Rang	Domaine professionnel	Écart Genré
1	manutenzione e riparazione di autoveicoli	71.4
2	attività di produzione cinematografica, di video e di programmi televisivi	63.1
3	attività fotografiche	63.0
4	fabbricazione di strumenti musicali	60.4
5	fabbricazione di aeromobili	59.6
6	fabbricazione di veicoli militari da combattimento	59.1
7	allevamento di bovini da latte	58.7
8	attività delle banche centrali	56.8
9	installazione di impianti elettrici	56.5
10	ricerche di mercato e sondaggi di opinione	55.8
11	attività dei vigili del fuoco e della protezione civile	55.6
12	lavori di costruzione e installazione	55.3

13	lavori di meccanica generale	53.1
14	edizione di giochi per computer	52.2
15	riparazione di computer e periferiche	52.1
16	costruzione di ponti e gallerie	51.1
17	servizi degli studi medici di medicina generale	51.0
18	realizzazione di coperture	50.0
19	fusione di acciaio	50.0
20	attività di musei	50.0
21	servizi investigativi privati	48.9
22	attività sportive	48.8
23	acquacoltura marina	47.9
24	servizi veterinari	47.8
25	ricerca e sviluppo sperimentale nel campo delle biotecnologie	47.8
26	attività degli studi odontoiatrici	46.8
27	attività degli studi legali e notarili	45.5
28	attività generali di amministrazione pubblica	44.7
29	ordine pubblico e sicurezza nazionale	43.5
30	attività degli studi di architettura	43.5
31	affari esteri	43.2
32	ricerca e sviluppo sperimentale nel campo delle scienze sociali e umanistiche	42.2
33	telecomunicazione	41.7
34	attività di servizi per la persona	41.3
35	commercio di altri autoveicoli	40.9
36	giustizia ed attività giudiziarie	39.6
37	attività di mediazione immobiliare	39.1
38	consulenza nel settore delle tecnologie dell'informatica	38.7
39	pubbliche relazioni e comunicazione	38.3
40	attività di pulizia	37.0
41	fabbricazione di profumi e cosmetici	36.9
42	edizione di libri	36.9
43	attività dei servizi connessi alle tecnologie dell'informatica	36.4
44	amministrazione di mercati finanziari	34.8
45	rappresentazioni artistiche	34.0
46	pesca marina	32.0
47	attività delle agenzie di viaggio	31.9
48	attività editoriali	30.2
49	attività ricreative e di divertimento	29.5
50	traduzione e interpretariato	28.9
51	servizi di asili nido	25.0
52	servizi di assistenza sanitaria	24.0
53	attività di biblioteche ed archivi	16.7
54	servizi ospedalieri	15.6
55	servizi dei parrucchieri e di altri trattamenti estetici	13.1

TABLE 8 – Domaines professionnels, par ordre décroissant d'Écart Généré - IT_{Neutre} .