

GeNRe: a French Gender-Neutral Rewriting System Using Collective Nouns

Anonymous ACL submission

Abstract

A significant portion of the textual data used in the field of Natural Language Processing (NLP) exhibits gender biases, particularly due to the use of masculine generics (masculine words that are supposed to refer to mixed groups of men and women), which can perpetuate and amplify stereotypes. Gender rewriting, a NLP task that involves automatically detecting and replacing gendered forms with neutral or opposite forms (e.g., from masculine to feminine), can be employed to mitigate these biases. Such systems are available for English, Arabic, Portuguese and German, but no French system is available. We create an original French gender-neutral rewriting system using collective nouns, which are gender-fixed in French. This paper presents GeNRe, the very first French gender-neutral rewriting system. We introduce a rule-based system (RBS) tailored for the French language alongside two fine-tuned large language models trained on data generated by our RBS. We also explore the use of instruction models to enhance the performance of our other systems and find that Claude 3 Opus combined with our dictionary achieves results close to our RBS. Through this contribution, we hope to promote the advancement of gender bias mitigation techniques in NLP for French.

1 Introduction

Since the 1970s, a number of psycholinguistics studies have focused on how language influences thoughts (Berlin and Kay, 1969; Kay and McDaniel, 1978). Further studies examining gender in language showed that it could lead to cognitive biases (Jacobson and Insko, 1985; Sczesny et al., 2016), particularly when it comes to the use of masculine generics (masculine words that are supposed to refer to mixed groups of men and women)¹ (Braun et al., 2005; Richy and Burnett,

2021). For example, Stahlberg et al. (2001) showed that when asked to name a celebrity in a certain field in German, respondents were more likely to give the name of a man when a masculine generic was used in the question.

Gender bias in natural language processing (NLP) models is a critical issue that can lead to biased predictions and the amplification of biases present in training data. This problem is particularly relevant for machine translation systems, which are highly susceptible to gender biases when translating between languages with different grammatical gender systems (Savoldi et al., 2021). Data augmentation, which involves balancing the amount of data for all genders in a specific language, has been proposed as a potential solution to debias NLP systems (Zhao et al., 2018). To achieve this goal, current research projects automatically propose alternatives to sentences containing masculine generics, contributing to an NLP task known as “gender rewriting”.

As of yet, gender neutralization techniques have not been developed for French, even though it is a heavily gendered language. Thus, we aim to create a French gender-neutral rewriting system using human collective nouns, defined by Lecolle (2019) as “nouns referring to entities comprised of groups of individuals”². Collective nouns have been widely discussed in the literature, especially when it comes to French (Flaux, 1999; Lammert, 2010; Lammert and Lecolle, 2014; Lecolle, 2019). Since this type of noun has a fixed gender in French³, it is an effective way of achieving gender neutralization. This gender-neutral rewriting system, GeNRe (**G**ender-**N**eutral **R**ewriting System Using French Collective

women pursuing that occupation, instead of using “police officers.”

²In French: « nom désignant une entité composée d’un ensemble d’individus humains. »

³For instance, “la police” (“police”) refers to both policemen and policewomen.

¹There is no strict equivalent in English, but an example could be the use of “policemen” to refer to both men and

Nouns), is the very first gender-neutral rewriting system for French⁴ and could foster the development of other types of gender rewriting systems for that language in the future.

2 The Task of Gender Rewriting

While Alhafni et al. (2022b) were the first to define this task as “gender rewriting,” similar efforts had already been pursued for Arabic (Habash et al., 2019), German (Pomerence, 2022), and English (Sun et al., 2021). Alhafni et al. (2022b) suggest the following definition for this task: “generating alternatives of a given Arabic sentence to match different target user gender contexts.” (2). While this definition works well for Alhafni et al.’s work, as they focus specifically on Arabic and create a system to switch between the masculine gender and the feminine gender, it is not universally applicable. Indeed, among the aforementioned works, several approaches to gender rewriting have been explored: Habash et al. (2019) and Alhafni et al. (2022a) developed a system to transform Arabic sentences with masculine words into sentences with feminine equivalents, and vice versa. Pomerence’s (2022) system seeks to provide inclusive suggestions for input sentences in German; so does Veloso et al.’s (2023) system for Portuguese. Finally, Sun et al. (2021), Vanmassenhove et al. (2021) and He et al. (2021) created systems to neutralize gender in an English input sentence.

We suggest a new, universal definition for the task that works for all languages and all the transformation approaches when it comes to gender based on the latest works mentioned previously:

The use of a gendered input sentence to generate one or several alternative sentences with different gender forms by neutralizing them, choosing inclusive forms or switching to another gender.

3 Gender in French

In French, nouns are classified as either masculine or feminine, and the gender of a noun influences the form of adjectives, pronouns, and verbs that accompany it. The gender of human nouns reflects the sociological gender of the referent (for instance, “danseuse” refers to a female dancer), while gender

of nouns referring to unanimated beings is arbitrary (Watbled, 2012).

The masculine gender is considered to be the “default” gender in French, and can be used in a non-specific context or to refer to groups of people composed of both men and women. The use of masculine as the default gender can however lead to both gender biases and invisibilizing women. As a result, two main writing techniques have been developed to avoid its use: visibilization techniques and neutralization techniques.

Visibilization techniques seek to highlight the feminine ending of words by separating the masculine ending from the feminine one through the use of specific symbols (asterisk, interpunct: *acteur.ice*) or by affecting the feminine ending directly (using capital or bold letters). Neutralization techniques, on the other end, aim to use epicene words, that is words whose form is the same for masculine and feminine (e.g. “spécialiste”, *specialist*), or words that refer to groups of people, such as collective nouns (e.g. “lectorat”, *readership*), these having a fixed gender which is not associated with the genders of the people within that group.

We chose to focus on gender neutralization due to it being a less explored issue in research comparatively to visibilization techniques. Moreover, while collective nouns are a great asset for gender neutralization, their usage is still restricted to a few words and their full potential has not yet been explored.

4 Methodology

We propose three different approaches for the task of gender-rewriting: a rule-based approach, a neural model fine-tuning approach and an instruction model approach. To build the resources used for these systems, we first create a dictionary of French collective nouns and their member noun counterparts, which we describe in Section 4.1. In Section 4.2, we then give details about the datasets that we extracted sentences from for the development of our rule-based system, large language model (LLM) fine-tuning and evaluation. Finally, in Section 4.3, we delve into the specifics of our experimental design with the aforementioned model types.

4.1 Dictionary

First, we manually created a dictionary with French collective nouns and their member noun counter-

⁴Code and data are publicly available on GitHub: <https://github.com/REDACTED>

parts. Two approaches were used to fill this dictionary: literature review, consisting of retrieving collective nouns mentioned in the French literature, and manual collecting, consisting of collecting occurrences on the Internet and in newspaper articles, as well as scraping French Wiktionary pages containing lists of such nouns. We respectively retrieved 210 and 105 nouns using these methods (315 in total). Table 1 contains a few examples of entries in our dictionary.

Collective noun	Member noun (masc. plural)
académie (<i>academy</i>)	académiciens (<i>academicians</i>)
armée (<i>army</i>)	soldats (<i>soldiers</i>)
milice (<i>militia</i>)	miliciens (<i>militiamen/women</i>)
artillerie (<i>artillery</i>)	artilleurs (<i>artillerists</i>)
auditoire (<i>listenership</i>)	auditeurs (<i>listeners</i>)
ballet (<i>ballet</i>)	danseurs (<i>dancers</i>)
police (<i>police</i>)	policiers (<i>police officers</i>)

Table 1: Collective noun-member noun dictionary overview

4.2 Datasets

Using our dictionary, we searched for occurrences of masculine plural member nouns in a French Wikipedia dataset with 1.58 million texts (graelo, 2023)⁵. We extracted 292,076 sentences containing such nouns. In addition, we also extracted French sentences from the Europarl EN-FR corpus (Koehn, 2005), a corpus created from the proceedings of the European Parliament and available in 21 languages, including English and French. This corpus was filtered to include French sentences only, and 106,878 additional sentences were extracted for neural model fine-tuning and evaluation (total 398,954). Both of these corpora are made available for research purposes.

For the rule-based system specifically, tags were automatically added at the beginning and at the end of each member phrase in the extracted sentences, with the ID of the entry in the dictionary. This was done because member nouns may have several

collective noun counterparts, leading to several different sentences being generated in addition to the main one. For instance, the member noun “soldats” (*soldiers*) could well be replaced with collective nouns “armée” (*army*) “bataillon” (*battalion*), “infanterie” (*infantry*) or “régiment” (*regiment*). As we used data generated by our rule-based system for neural model fine-tuning (see Section 4.3.2), this was especially useful to generate all the possible variations of the input sentence, and thus increase the number of examples the models were trained on. Moreover, the use of tags also helps guarantee the member nouns to be replaced in the input sentence, as only those that are between tags will be taken into account. Example 1 shows how these tags are used.

- (1) a. Un historique permet de lister <n-126>les auteurs</n> et de consulter les modifications successives de l’article par <n-68>ses rédacteurs</n>. (A history allows one to list <n-126>the authors</n> and view successive modifications to the article by <n-68>its editors</n>.)

Finally, we created a corpus-specific evaluation dataset comprised of 250 sentences from each corpus (total 500), and we manually gender-neutralized each sentence to have gold sentences.

4.3 Models

In this section, we present three different model types for gender-neutral rewriting: a rule-based model, two neural models, and an instruction model. Each model takes a different approach to the task, allowing us to compare their performance.

4.3.1 Rule-based model

We developed a rule-based system (RBS) to automatically apply the correct syntactic rules when converting a member noun into a collective noun, which leads to number and gender changes in the sentence.

The RBS consists of two main components: a syntactic dependency detection component and a generation component.

The dependency detection component primarily relies on spaCy (Honnibal et al., 2020) with the fr_core_news_sm pipeline as well as a set of rules to detect the words that are syntactically related to the member noun that needs to be replaced.

The generation component replaces each member noun in the sentence with its collective noun

⁵Dataset is available here: <https://huggingface.co/datasets/graelo/wikipedia>. License: CC-BY-SA-3.0

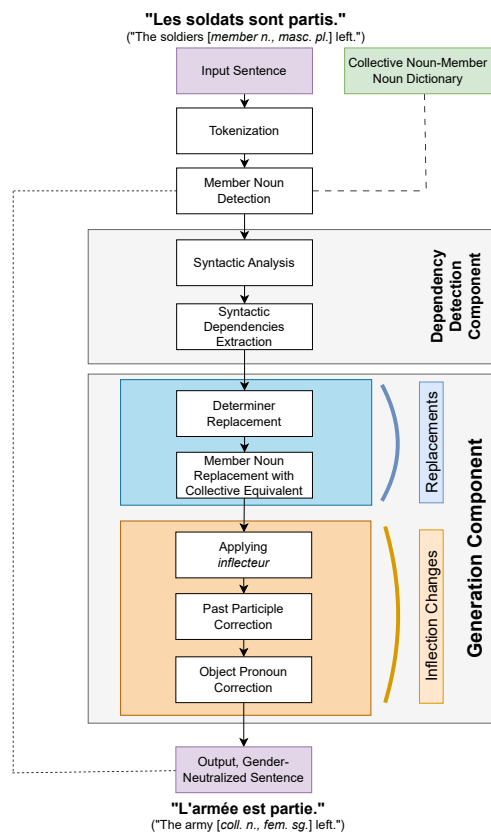


Figure 1: Rule-based model replacement pipeline overview

counterpart found in the dictionary, adjusting the determiner, handling elision, and reinflecting the detected dependencies using *inflecteur* (Chuttarsing, 2021), a Python module leveraging the Delaf French morphological dictionary⁶ and *french-camembert-postag-model*⁷, a CamemBERT-based (Martin et al., 2020) part of speech (POS) tagging model for French. Our RBS also makes additional replacements for past participles and object pronouns as these are not always being well handled by the *inflecteur* Python module. If no member nouns are detected in the sentence, the original sentence will be returned instead as it is already considered gender-neutral. Figure 1 shows an overview of the rule-based model pipeline.

4.3.2 Neural models

Recent research on gender rewriting has focused on training neural models as well as fine-tuning large language models using data generated by RBS to improve task-specific performance. While

⁶<https://uclouvain.be/fr/instituts-recherche/ilc/cental/delaf-2-0.html>

⁷<https://huggingface.co/gilf/french-camembert-postag-model>

some studies (Sun et al., 2021; Veloso et al., 2023) showed a decrease in performance compared to RBS, Vanmassenhove et al. (2021) found a notable improvement of 0.27 in WER. We aim to investigate whether fine-tuning large language models can significantly improve the results of RBS, hypothesizing that the linguistic knowledge acquired by these models during training on large text corpora will help resolve errors in the training corpus and enhance results.

Two Seq2seq large language models (LLMs), t5-small (Raffel et al., 2020) and m2m100_418M (Fan et al., 2020), were selected for the experiments, and were fine-tuned using our two RBS-generated corpora (Wikipedia and Europarl) containing gender-neutralized and non-gender-neutralized sentence pairs. The training dataset for each model consisted of 60,000 sentence pairs per corpus, and the validation dataset had 6,000 (10%). Hyperparameters used for training are available in Appendix A.

4.3.3 Instruction model

The rapid development of LLMs and advances in NLP have demonstrated the ability to manipulate language models’ behavior to predict text continuations and perform specific tasks without explicit training, leading to “instruction models” such as InstructGPT (Ouyang et al., 2022) or, more recently, Mixtral 8x7B Instruct (Jiang et al., 2024). This is primarily achieved through the use of “prompts” or instructions given to the language model (Liu et al., 2021). While some studies have briefly mentioned the potential of instruction models to reduce gender biases in automatically generated texts, and have occasionally experimented with such models⁸, no gender rewriting study has yet conducted a comprehensive analysis of their capabilities for this specific task. As a result, we aimed to leverage this kind of model in order to evaluate its performance for this task. We chose Claude 3 Opus, which is, at the time of writing, considered to be the best model for textual generation according to specific benchmark (Anthropic, 2024).

To comprehensively evaluate the performance of Claude 3 Opus, we designed three distinct types of instructions to test its ability to generate gender-neutral texts. Corresponding prompts are available

⁸For instance, Veloso et al. (2023) tried to make use of ChatGPT to generate gender-inclusive sentences in Portuguese, and suggested that the use of instruction models could prove useful to automatically create gender-inclusive datasets.

in Appendix A.

- The “BASE” instruction provides a basic task description, asking the model to make the sentence inclusive by replacing masculine generics with their collective noun equivalents, without explicitly specifying the replacement word.
- The “DICT” instruction leverages our collective noun dictionary and asks the model to replace masculine generics with their corresponding collective nouns, those being explicitly mentioned. There are two different versions for the “DICT” instruction: “DICT-SG”, used when only one generic masculine noun with a matching collective noun was found in the sentence, and “DICT-PL”, used when several generic masculine nouns with matching collective nouns were found.
- The “CORR” instruction takes sentences generated by our RBS as input and tasks the model with correcting potential errors, such as mismatches between verb and adjective numbers and genders.

5 Results

To evaluate the performance of our different rewriting models, we leverage two evaluation metrics commonly used for the task of gender rewriting: Word Error Rate (WER) and BLEU (Papineni et al., 2002). JiWER 3.0.3⁹ and bleu 0.3¹⁰ Python packages were used with default parameters.

Average results of each model on the two corpora are available in Table 2.

Type	WER (↓)	BLEU (↑)
Baseline (unchanged)	13.35%	80.55
GeNRe-RBS	3.40%	93.43
GeNRe-T5	5.11%	90.68
GeNRe-M2M-100	5.40%	90.17
Claude 3 Opus-BASE	12.16%	82.98
Claude 3 Opus-DICT	3.75%	93.64
Claude 3 Opus-CORR	10.17%	85.13

Table 2: Results by model type. Bold indicates the best results overall.

The RBS and Claude 3 Opus-DICT achieved the best results in our experiments. While the

⁹<https://pypi.org/project/jiwer/>

¹⁰<https://pypi.org/project/bleu/>

RBS model achieved the best WER score, Claude 3 Opus-DICT achieved the highest BLEU score. These results can be explained by the fact that WER and BLEU scores capture distinct aspects of text generation. Due to its reliance on predefined rules, the RBS easily preserves original words and word order, likely leading to a lower WER. On the contrary, instruction models are known to be more prompt to slightly deviate from the original formulation of sentences, which may increase the WER without significantly affecting the BLEU score due to the order of words not being taken into account.

The neural models also showed mostly promising results. Comparing the two of them, they achieved similar results, with the T5 model slightly outperforming M2M-100. However, both models showed a minor decrease in performance compared to the RBS. As a result, similarly to Veloso et al. (2023), we do not find a significant improvement compared to our RBS following fine-tuning.

Moreover, we also provide the distribution of errors made by GeNRe-RBS, GeNRe-T5 and GeNRe-M2M-100 in Figure 2. Error types can be divided into three main categories: POS (ADJ, DET, DET_COREF, PRON_COREF, VERB), text generation (CASE, GEN_FAILURE, SPECIAL_CHAR) and other (ELISION, MISID_NOUN, PUNCT, SEM, UNREPLACED).

Text generation errors, labeled with (N) in Figure 2, are strictly specific to neural models. CASE refers to capitalization errors (missing/extra uppercase or lowercase); GEN_FAILURE refers to token-specific generation errors (for instance, incorrectly replacing a proper name with a non-existent name); SPECIAL_CHAR refers to errors related to special characters (for instance, accents).

When it comes to other errors, ELISION is used when there was an issue with how one or multiple words in the generated sentence were elided¹¹. MISID_NOUN occurs when a word in the automatically annotated corpus was mistaken as a noun. PUNCT refers to errors related to punctuation or typography (double spaces, for example). SEM is used to label automatically generated sentences which cannot be considered semantically correct due to the replacement of the member noun with its collective noun counterpart¹² Finally, UNRE-

¹¹For instance, in French, the masculine determiner “le” and the feminine determiner “la” (*the*) should be elided and written as “l’” when the word that follows begin with a vowel or a mute “h”.

¹²As discussed by Lecolle (2019), collective nouns in

PLACED occurs when a member noun that is in our dictionary was not replaced.

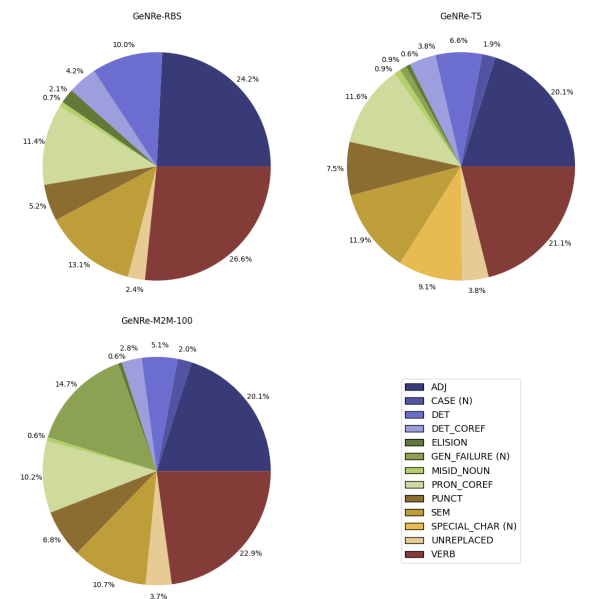


Figure 2: Distribution of errors across GeNRe-RBS, GeNRe-T5 and GeNRe-M2M-100

Across all three models, the most prominent error types are related to verbs and noun cases. Verbs account for 26,6% of errors for GeNRe-RBS, 21,1% for GeNRe-T5, and 22,9% for GeNRe-M2M-100. On the other hand, adjectives account for 24,2% of errors for GeNRe-RBS, and 20,1% for both GeNRe-T5 and GeNRe-M2M-100.

The M2M-100 model is highly prone to making token-specific generation errors (14,7%), this type of error being strictly specific to this model. Similarly, we find that the T5 model also makes specific errors related to the handling of special characters. We discuss these issues more in detail in Section 6.

6 Discussion

A qualitative analysis of the generated sentences revealed that the RBS was making most of its errors when modifying adjectives and verbs. This is not surprising given that these two part-of-speech categories are the ones which require the most complex

French, and more specifically human collective nouns, feature specific semantic characteristics due to how they are used to group human beings under a common denomination, based for example on their profession (« le professorat » [*professorate*]), their social status (« l’aristocratie » [*the aristocracy*]), or their political leaning (« la gauche » [*the left*]). Combining human collective nouns with specific verbs or contexts may thus not be considered semantically correct, and may occur when transforming a sentence. We labeled such transformed sentences with this error.

changes when transitioning from a member noun to a collective noun. Indeed, in French, adjectives undergo a certain number of changes when changing number or gender. Verbs can also have these same changes when used as past participles; otherwise, only number change will affect them. For instance, in Example 2, the verb “seront” (pl., *will be*) should have been changed to “sera” (sg.) to match with the new collective noun “citoyenneté” (*citizenry*).

- (2) a. Cette démarche fera progresser les droits **des citoyens**, car, par l’intermédiaire du Parlement, **les citoyens seront** en contact direct avec la Commission, ce qui lui confèrera une légitimité considérable. [original sent.]
(This approach will increase **citizens’ [masc.]** rights, because, through the Parliament, **citizens will [pl.]** have a direct line to the Commission thereby generating considerable legitimacy.)
- b. Cette démarche fera progresser les droits **de la citoyenneté**, car, par l’intermédiaire du Parlement, **la citoyenneté seront** en contact direct avec la Commission, ce qui lui confèrera une légitimité considérable. [GeNRe-RBS]
(This approach will increase the rights of **the citizenry**, because, through the Parliament, **the citizenry will [pl.]** have a direct line to the Commission thereby generating considerable legitimacy.)
- c. Cette démarche fera progresser les droits **de la citoyenneté**, car, par l’intermédiaire du Parlement, **la citoyenneté sera** en contact direct avec la Commission, ce qui lui confèrera une légitimité considérable. [manual sent.]
(This approach will increase the rights of **the citizenry**, because, through the Parliament, **the citizenry will [sg.]** have a direct line to the Commission thereby generating considerable legitimacy.)

Similarly, in Example 3, the adjective “chargés” (pl., *in charge of*) should match the new singular collective noun “parlement” (*parliament*) and be changed to “chargé”.

- (3) a. Je vous invite à informer **les députés**

467	européens chargés des dossiers agricoles de l'avancement des négociations. [original sent.]	aux sources de financement correspondantes. [GeNRe-FT-M2M-100]	517
468	(I urge you to inform the Members of European Parliament [masc] in charge of [pl.] the agricultural issues about the progress of negotiations.)	(A second factor is the Commission's support for local actors [coll. sg.] who want [sg.] to take part in these programmes, so that they can access the corresponding funding mechanisms.)	518
469			519
470			520
471			521
472			522
473			523
474	b. Je vous invite à informer le parlement européen chargés des dossiers agricoles de l'avancement des négociations. [GeNRe-RBS]	Additionally, the fine-tuned models were capable of utilizing different collective noun equivalences from the dictionary (some collective nouns being associated to the same member noun).	524
475	(I urge you to inform the European parliament in charge of [pl.] the agricultural issues about the progress of negotiations.)	Errors observed in the fine-tuned models and different from the RBS included token generation failures (M2M-100, Example 5, where "Nebski" was generated instead of "Zemski"), and incorrect generation of special characters (T5, as in Example 6 where "main-d'uvre" was generated instead of "main-d'œuvre" [<i>labour</i>]). The first error might come from the multilingual aspect of the model, as it may generate words or mix tokens from other languages, while the second error is probably due to the model being mostly trained on English data. For both models, we also found cases where words were not uppercased correctly, as in Example 7.	525
476			526
477			527
478			528
479	c. Je vous invite à informer le parlement européen chargé des dossiers agricoles de l'avancement des négociations. [manual sent.]		529
480	(I urge you to inform the European parliament in charge of [sg.] the agricultural issues about the progress of negotiations.)		530
481			531
482			532
483			533
484			534
485			535
486			536
487			537
488			538
489			539
490			540
491			541
492			542
493			543
494			544
495			545
496			546
497			547
498			548
499			549
500			550
501			551
502			552
503			553
504			554
505			555
506			556
507			557
508			558
509			559
510			560
511			561
512			562
513			563
514			564
515			565
516			566

- 566 (8) a. Dans une lettre à la **famille** datée 616
567 du 13 juin 1861, Zeng Guofan a 617
568 ordonné à **ses propres navires** de 618
569 surveiller les navires commerciaux bri- 619
570 tanniques après avoir remarqué que des 620
571 marchands étrangers déchargeaient du 621
572 riz à **la rébellion** à Anqing. [GeNRe- 622
573 RBS] 623
574 (In a letter addressed to the **family** and dated 624
575 June 13, 1861, Zeng Guofan ordered **his own** 625
576 **vessels** to monitor British commercial vessels 626
577 after noticing that foreign sellers were giving 627
578 rice to **the rebellion** in Anqing.) 628
- 579 b. Dans une lettre à la **parenté** datée 629
580 du 13 juin 1861, Zeng Guofan 630
581 a ordonné à **sa propre flotte** de 631
582 surveiller les navires commerciaux bri- 632
583 tanniques après avoir remarqué que des 633
584 marchands étrangers déchargeaient du 634
585 riz **aux rebelles** à Anqing. [Claude 3 635
586 Opus-BASE] 636
587 (In a letter addressed to the **kinfolk** and dated 637
588 June 13, 1861, Zeng Guofan ordered **his own** 638
589 **fleet** to monitor British commercial vessels af- 639
590 ter noticing that foreign sellers were giving rice 640
591 to **rebels** in Anqing.) 641

592 Notably, the DICT prompt was observed to gener- 642
593 ate sentences with correct verbs and adjectives, 643
594 indicating its ability to effectively leverage the col- 644
595 lective noun dictionary to produce grammatically 645
596 accurate sentences. We give such an example in 646
597 Example 9. 647

- 598 (9) a. Mais l’armée protestante, toujours 648
599 agressive, **restaient** à la charge des 649
600 habitants et **constituaient** une lourde 650
601 charge. [GeNRe-RBS] 651
602 (But the Protestant army, still aggressive, **re-** 652
603 **mained [pl.]** in the care of the local people and 653
604 **constituted [pl.]** a heavy burden.) 654
- 605 b. Mais l’armée protestante, toujours 655
606 agressive, **restait** à la charge des habi- 656
607 tants et **constituait** une lourde charge. 657
608 [Claude 3 Opus-DICT] 658
609 (But the Protestant army, still aggressive, **re-** 659
610 **mained [sg.]** in the care of the local people and 660
611 **constituted [sg.]** a heavy burden.) 661

612 Nonetheless, among the errors made by Claude 3 661
613 Opus-DICT, we identified instances of unreplaced 662
614 nouns, where the model failed to substitute the mas- 663
615 culine generics with their corresponding collective 664

noun equivalents, such as in Example 10.

- (10) a. Paradoxalement, cette progression en 617
618 voix s’accompagne d’un recul en nom- 619
619 bre d’élus, du fait de la poussée des 620
620 candidats indépendants (pour la plu- 621
621 part de la **représentation** de la commu- 622
622 nauté kurde) et du CHP. [GeNRe-RBS] 623
623 (Paradoxically, this increase in votes paralleled 624
624 a decrease in the number of elected representa- 625
625 tives due to better results for the independent 626
626 candidates (most of them **coming from the** 627
627 **representation** of the Kurdish community) and 628
628 CHP. 629
- b. Paradoxalement, cette progression en 630
631 voix s’accompagne d’un recul en nom- 631
632 bre d’élus, du fait de la poussée des 632
633 candidats indépendants (pour la plu- 633
634 part des **représentants** de la commu- 634
635 nauté kurde) et du CHP. [Claude 3 635
636 Opus-DICT] 636
636 (Paradoxically, this increase in votes paralleled 637
637 a decrease in the number of elected representa- 638
638 tives due to better results for the independent 639
639 candidates (most of them **being representa-** 640
640 **tives** of the Kurdish community) and CHP. 641

7 Conclusion 641

642 Our work represents a step towards addressing 642
643 gender-biased textual data in French. We make 643
644 three key contributions to the task of gender rewrit- 644
645 ing in NLP: 1) a dictionary of French collective 645
646 nouns and their corresponding member nouns, 646
647 which serves as a resource for future research in this 647
648 area; 2) a dataset of gender-neutralized and non- 648
649 gender-neutralized sentences; and 3) a rule-based 649
650 system that effectively gender-neutralizes French 650
651 sentences using collective nouns, laying ground- 651
652 work for further advancements for this task in that 652
653 language. Our experiment combining our manu- 653
654 ally created dictionary with the Claude 3 Opus 654
655 instruction model also shows promise for the use 655
656 of such models for the task of gender rewriting. We 656
657 strongly believe that future research further explor- 657
658 ing the capabilities of these models for that task 658
659 could lead to the development of effective solutions 659
660 for mitigating gender bias in many languages. 660

Ethics Statement 661

662 We did not filter the datasets that were used for the 662
663 development of the RBS and for fine-tuning mod- 663
664 els for harmful, hateful, inappropriate or personal 664

665	content. Considering the sources used to constitute	Brent Berlin and Paul Kay. 1969. <i>Basic Color Terms: Their Universality and Evolution</i> . University of California Press.	714
666	these datasets (Wikipedia and Europarl), we believe		715
667	it very unlikely for those to display such type		716
668	of content. Similarly, when it comes to output	Friederike Braun, Sabine Sczesny, and Dagmar	717
669	sentences generated by the fine-tuned models, since	Stahlberg. 2005. <i>Cognitive Effects of Masculine</i>	718
670	those were trained on replacing specific words or	<i>Generics in German: An Overview of Empirical Findings</i> . <i>Communications</i> , 30(1):1–21.	719
671	part of speech in sentences, the generation of such		720
672	content seems unlikely. As discussed in the paper,	Adrien Chutturarsing. 2021. <i>Inflecteur</i> .	721
673	instruction models are more prone to reformulating		
674	input sentences: while we did not find any inap-	Angela Fan, Shruti Bhosale, Holger Schwenk, Zhiyi	722
675	propriate content in the Claude 3 Opus-generated	Ma, Ahmed El-Kishky, Siddharth Goyal, Man-	723
676	sentences we evaluated, this kind of models may	deep Baines, Onur Celebi, Guillaume Wenzek,	724
677	be trained on such data, which might lead to the	Vishrav Chaudhary, Naman Goyal, Tom Birch, Vi-	725
678	generation of harmful or hateful content.	taliy Liptchinsky, Sergey Edunov, Edouard Grave,	726
		Michael Auli, and Armand Joulin. 2020. <i>Beyond</i>	727
		<i>English-Centric Multilingual Machine Translation</i> .	728
		<i>Preprint</i> , arxiv:2010.11125.	729
679	Limitations		
680	French collective nouns adhere to specific seman-	Nelly Flaux. 1999. <i>À propos des noms collectifs</i> . <i>Revue</i>	730
681	tic rules, which means that their usage may not be	<i>de linguistique romane</i> , (63):471–502.	731
682	universally applicable to all sentences, sometimes		
683	resulting in constructions that appear asemantic.	graelo. 2023. <i>Graelo/wikipedia dataset</i> .	732
684	This limitation is further compounded by the fact	https://huggingface.co/datasets/graelo/wikipedia .	733
685	that only a small subset of these nouns is actively		
686	employed in everyday language by native speakers,	Nizar Habash, Houda Bouamor, and Christine Chung.	734
687	which restricts their versatility and adaptability in	2019. <i>Automatic Gender Identification and Reinflec-</i>	735
688	various linguistic contexts. We however believe	<i>tion in Arabic</i> . In <i>Proceedings of the First Workshop</i>	736
689	that they are good candidates for gender neutral-	<i>on Gender Bias in Natural Language Processing</i> ,	737
690	ization, and the development of our system may	pages 155–165, Florence, Italy. Association for Com-	738
691	help promote a broader use of such nouns. In ad-	putational Linguistics.	739
692	dition, combining our system with a contextual or		
693	semantic analysis framework could help address	Zexue He, Bodhisattwa Prasad Majumder, and Julian	740
694	these issues by ensuring that the collective noun	McAuley. 2021. <i>Detect and Perturb: Neutral Rewrit-</i>	741
695	equivalents are both contextually relevant and se-	<i>ing of Biased and Sensitive Text via Gradient-based</i>	742
696	mantically appropriate.	<i>Decoding</i> . <i>Preprint</i> , arxiv:2109.11708.	743
		Matthew Honnibal, Ines Montani, Sofie Van Lan-	744
		degheem, and Adriane Boyd. 2020. <i>spaCy: Industrial-</i>	745
		<i>strength Natural Language Processing in Python</i> .	746
		Marsha B. Jacobson and William R. Insko. 1985. <i>Use</i>	747
		<i>of nonsexist pronouns as a function of one’s feminist</i>	748
		<i>orientation</i> . <i>Sex Roles</i> , 13(1-2):1–7.	749
697	References		
698	Bashar Alhafni, Nizar Habash, and Houda Bouamor.	Albert Q. Jiang, Alexandre Sablayrolles, Antoine	750
699	2022a. <i>User-Centric Gender Rewriting</i> . In <i>Proceed-</i>	Roux, Arthur Mensch, Blanche Savary, Chris	751
700	<i>ings of the 2022 Conference of the North American</i>	Bamford, Devendra Singh Chaplot, Diego de las	752
701	<i>Chapter of the Association for Computational Lin-</i>	Casas, Emma Bou Hanna, Florian Bressand, Gi-	753
702	<i>guistics: Human Language Technologies</i> , pages 618–	anna Lengyel, Guillaume Bour, Guillaume Lam-	754
703	631, Seattle, United States. Association for Compu-	ple, L�lio Renard Lavaud, Lucile Saulnier, Marie-	755
704	tational Linguistics.	Anne Lachaux, Pierre Stock, Sandeep Subramanian,	756
		Sophia Yang, Szymon Antoniak, Teven Le Scao,	757
		Th�ophile Gervet, Thibaut Lavril, Thomas Wang,	758
		Timothe�e Lacroix, and William El Sayed. 2024. <i>Mix-</i>	759
		<i>tral of Experts</i> . <i>Preprint</i> , arxiv:2401.04088.	760
705	Bashar Alhafni, Nizar Habash, Houda Bouamor,	Paul Kay and Chad K. McDaniel. 1978. <i>The Linguistic</i>	761
706	Ossama Obeid, Sultan Alrowili, Daliyah Alzeer,	<i>Significance of the Meanings of Basic Color Terms</i> .	762
707	Khawlah M. Alshantqiti, Ahmed ElBakry, Muham-	<i>Language</i> , 54(3):610–646.	763
708	ammad ElNokrashy, Mohamed Gabr, Abderrahmane Is-		
709	sam, Abdelrahim Qaddoumi, K. Vijay-Shanker, and	Philipp Koehn. 2005. <i>Europarl: A parallel corpus for</i>	764
710	Mahmoud Zyate. 2022b. <i>The Shared Task on Gender</i>	<i>statistical machine translation</i> .	765
711	<i>Rewriting</i> . <i>Preprint</i> , arxiv:2210.12410.		
		Marie Lammert. 2010. <i>S�manticque et cognition : les</i>	766
712	Anthropic. 2024. <i>The Claude 3 Model Family: Opus,</i>	<i>noms collectifs</i> . Droz, Gen�ve.	767
713	<i>Sonnet, Haiku</i> .		

768	Marie Lammert and Michelle Lecolle. 2014. Les noms collectifs en français, une vue d’ensemble. <i>Cahiers de lexicologie</i> , (105):203–222.	
769		
770		
771	Michelle Lecolle. 2019. <i>Les noms collectifs humains en français. Enjeux sémantiques, lexicaux et discursifs</i> . Lambert-Lucas, Université de Lorraine.	
772		
773		
774	Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2021. Pre-train, Prompt, and Predict: A Systematic Survey of Prompting Methods in Natural Language Processing . <i>Preprint</i> , arxiv:2107.13586.	
775		
776		
777		
778		
779	Louis Martin, Benjamin Muller, Pedro Javier Ortiz Suárez, Yoann Dupont, Laurent Romary, Éric De La Clergerie, Djamel Seddah, and Benoît Sagot. 2020. CamemBERT: A Tasty French Language Model . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7203–7219, Online. Association for Computational Linguistics.	
780		
781		
782		
783		
784		
785		
786		
787	Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback . <i>Preprint</i> , arxiv:2203.02155.	
788		
789		
790		
791		
792		
793		
794		
795	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. BLEU: A method for automatic evaluation of machine translation . In <i>Proceedings of the 40th Annual Meeting on Association for Computational Linguistics - ACL ’02</i> , Philadelphia, Pennsylvania. Association for Computational Linguistics.	
796		
797		
798		
799		
800		
801	David Pomeranke. 2022. INCLUSIFY: A benchmark and a model for gender-inclusive German . <i>Preprint</i> , arxiv:2212.02564.	
802		
803		
804	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer . <i>Preprint</i> , arxiv:1910.10683.	
805		
806		
807		
808		
809	Célia Richy and Heather Burnett. 2021. Démêler les effets des stéréotypes et le genre grammatical dans le biais masculin : une approche expérimentale . <i>GLAD!</i> , (10).	
810		
811		
812		
813	Beatrice Savoldi, Marco Gaido, Luisa Bentivogli, Matteo Negri, and Marco Turchi. 2021. Gender Bias in Machine Translation. <i>Transactions of the Association for Computational Linguistics</i> , 9:845–874.	
814		
815		
816		
817	Sabine Sczesny, Magda Formanowicz, and Franziska Moser. 2016. Can Gender-Fair Language Reduce Gender Stereotyping and Discrimination? <i>Frontiers in Psychology</i> , 7.	
818		
819		
820		
	Dagmar Stahlberg, Sabine Sczesny, and Friederike Braun. 2001. Name Your Favorite Musician: Effects of Masculine Generics and of their Alternatives in German . <i>Journal of Language and Social Psychology</i> , 20(4):464–469.	821
		822
		823
		824
		825
	Tony Sun, Kellie Webster, Apu Shah, William Yang Wang, and Melvin Johnson. 2021. They, Them, Theirs: Rewriting with Gender-Neutral English . <i>Preprint</i> , arxiv:2102.06788.	826
		827
		828
		829
	Eva Vanmassenhove, Chris Emmery, and Dimitar Shterionov. 2021. NeuTral Rewriter: A Rule-Based and Neural Approach to Automatic Rewriting into Gender-Neutral Alternatives . <i>Preprint</i> , arxiv:2109.06105.	830
		831
		832
		833
		834
	Leonor Veloso, Luisa Coheur, and Rui Ribeiro. 2023. A Rewriting Approach for Gender Inclusivity in Portuguese . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 8747–8759, Singapore. Association for Computational Linguistics.	835
		836
		837
		838
		839
		840
	Jean-Philippe Watbled. 2012. Linguistique du genre. <i>L’Harmattan</i> , pages 167–179.	841
		842
	Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2018. Gender Bias in Coreference Resolution: Evaluation and Debiasing Methods . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)</i> , pages 15–20, New Orleans, Louisiana. Association for Computational Linguistics.	843
		844
		845
		846
		847
		848
		849
		850
		851
	A Appendix	852
	A.1 Fine-Tuning Details	853
	Models were trained on a single NVIDIA RTX 4090 GPU. Training time took approximately 3 hours for each model.	854
		855
		856
	A.1.1 GeNRe-T5	857
	BATCH_SIZE = 48	858
	NUM_PROCS = 16	859
	EPOCHS = 5	860
	LEARNING_RATE = 0.0005	861
	WEIGHT_DECAY = 0.02	862
	A.1.2 GeNRe-M2M-100	863
	BATCH_SIZE = 8	864
	NUM_PROCS = 16	865
	EPOCHS = 5	866
	LEARNING_RATE = 0.0005	867
	WEIGHT_DECAY = 0.02	868

869
870
871
872
873
874
875
876
877
878

A.2 Instruction Model Hyperparameters

```
model="claude-3-opus-20240229",
temperature=0,
messages=[
  {"role": "user",
   "content": f"{message}"},
  {"role": "assistant",
   "content": "Here is the
output sentence:"}
]
```

A.3 Instruction Details

Table 3 contains the different types of instructions given to Claude 3 Opus as well as their respective content.

“EXAMPLES” refers to the few-shot sentences given to the instruction model. See Tables 4 and 5 for more information.

“ORIGINAL SENTENCE” is replaced with the sentence containing one or several masculine generic nouns that we want to replace with their collective counterparts. It is part of the prompt in a similar way to the example sentences so that the instruction model is guided towards generating the final, gender-neutralized sentence.

Instruction Type	Content
BASE	Make this French sentence inclusive by replacing generic masculine nouns with their French collective noun equivalents. Generate the final sentence only without any comments nor notes. {EXAMPLES} {ORIGINAL SENTENCE} →
DICT-SG	Make this French sentence inclusive by replacing generic masculine noun {NM} with its respective French collective noun equivalent {NCOLL}. Generate the final sentence only without any comments nor notes. {EXAMPLES} {ORIGINAL SENTENCE} →
DICT-PL	Make this French sentence inclusive by replacing generic masculine nouns {NM1, NM2, ...} with their respective French collective noun equivalents {NCOLL1, NCOLL2, ...}. Generate the final sentence only without any comments nor notes. {EXAMPLES} {ORIGINAL SENTENCE} →
CORR	Correct grammar in this French sentence. Generate the final sentence only without any comments nor notes. {EXAMPLES} {ORIGINAL SENTENCE} →

Table 3: Content of instructions per type given to Claude 3 Opus

A.4 Few-shot sentences given to Claude 3 Opus

Tables 4 and 5 contain the few-shot sentences used respectively for the “BASE” and “DICT” instructions, and the “CORR” instruction. They were formatted as such in the prompt:

[Sentence with masculine generic] → [Gender-neutralized sentence].

893
894
895
896
897
898
899
900

Sentence with masculine generic	Gender-neutralized sentence
Le président de la FIFA Sepp Blatter rejette les accusations des manifestants en les accusant d’opportunisme. (FIFA President Sepp Blatter dismisses the protesters’ accusations as opportunism.)	Le président de la FIFA Sepp Blatter rejette les accusations de la manifestation en l’accusant d’opportunisme. (FIFA President Sepp Blatter dismisses the protest’s accusations as opportunism.)
Les auteurs et les spectateurs ont été satisfaits des réponses des représentants. (Authors and spectators were pleased with the representatives’ responses.)	L’atorat et le public ont été satisfaits des réponses de la représentation . (The authorship and the audience were pleased with the representation’s responses.)
Le vicaire général proposa de disperser les religieux dans d’autres maisons de l’ordre et de procéder à la réfection des bâtiments. (The vicar general suggested to disperse religious people to other houses of the order to repair the buildings.)	Le vicaire général proposa de disperser le couvent dans d’autres maisons de l’ordre et de procéder à la réfection des bâtiments. (The vicar general suggested to disperse the convent to other houses of the order to repair the buildings.)

Table 4: Few-shot sentences for “BASE” and “DICT” instructions. Bold indicates the differences between sentences with masculine generics and gender-neutralized sentences.

RBS-generated sentence with errors	Manual sentence
Le président de la FIFA Sepp Blatter rejette les accusations de la manifestation en les accusant d’opportunisme. L’atorat et le public a été satisfaits des réponses des la représentation.	Le président de la FIFA Sepp Blatter rejette les accusations de la manifestation en l’accusant d’opportunisme. L’atorat et le public ont été satisfaits des réponses de la représentation.
Le vicaire générale proposa de disperser le couvent dans d’autres maisons de l’ordre et de procéder à la réfection des bâtiments.	Le vicaire général proposa de disperser le couvent dans d’autres maisons de l’ordre et de procéder à la réfection des bâtiments.

Table 5: Few-shot sentences for “CORR” instruction. Bold indicates the differences between the RBS-generated sentences with error and the manual, correct sentences.