

# Dialectic Bias in Toxicity Detection of Google’s Perspective API

## A study with five parallel corpora

Hanxin Xia | 3417418

[hanxin.xia@uni-duesseldorf.de](mailto:hanxin.xia@uni-duesseldorf.de)

### Abstract

Hate Speech detection has been a heated topic in NLP/NLU since the spread of social media. Despite the existence of various automatic toxicity scoring models, their biases against language varieties that align with minority identities have been discussed by researchers (Sap et al., 2019). In this work, we focus on Google’s Perspective API and analyze its toxicity scores on four synthetic corpora with different dialect features compared to an original master corpus. We discovered that Perspective scored all parallel corpora significantly differently than the master corpus, which further confirms the existence of the biases. Moreover, we discovered the capping phenomenon with the toxicity scores of the gold toxic sentences: Gold non-toxic sentences are more likely to suffer from toxicity score increase, when converted to a dialect, compared to gold toxic sentences.

This research aims to test the performance of Google’s Perspective API on five parallel corpora. Despite being a continuously improving model on toxicity scoring (Google, 2024), the bias against non-standard English is still present in the results. The model scored (state: March 2024) texts from the four English dialects: African American Vernacular English (AAVE), Nigerian and Indian dialects and Singaporean English (Singlish) significantly differently than the original texts of standard English. Thus, we posit that there is a continued need for optimization in the hate speech detection of online contents for better mitigating biases that may infringe upon the online representation of individuals from marginalized communities.

### 2 Background

### 1 Introduction

With the increasing popularity of social media and the daily widening user bases, the moderation of the contents on online communication platforms have been a crucial part of maintaining the friendliness for the users and building healthy communities. In pursuit of this objective, the platforms have implemented various means including reactive methods such as community guidelines, user report and post-removal of harmful contents and proactive methods either with sensitive words and language pattern rules (Gitari et al., 2015) or machine learning algorithms for automated hate speech detection (Alrehili, 2019).

However current detoxifying approaches are proven to be not perfect. While the reactive methods are commonly subject to deficiency and passivity, the automatic proactive methods often discriminate against minority aligned language varieties other than Standard American English (SAE) (Sap et al., 2019; Zhou et al., 2021).

### 3 Introduction

The term “named entity” was first defined in the Sixth Message Understanding Conference (MUC-6) (Grishman and Sundheim, 1996). In the publicly available description, at least three types of language expressions were categorized as such entities: “unique identifiers” of entities (organizations, persons, locations), times (dates, times), and quantities (monetary values, percentages). The extraction of such entities from a sentence is, however, not a trivial job. Outside the most obvious dictionary matching (Higashinaka et al., 2012; Shang et al., 2018), several other hybrid approaches have been proposed over the years, especially for bio-medical and other domain-specific terms (Rocktäschel et al., 2012; Lou et al., 2020). Nowadays, with pre-trained language models spaCy, StanfordNLP as such, named entity recognition tasks can be easily carried out as part of the standard pipeline during text annotation.

This research aims to investigate the multi-purposed language model spaCy’s performance on named entity recognition tasks when it comes to multilingual data, above all where two languages are used interchangeably in one sentence, the so-called insertional code-switching, as follows:

- (1) @esangiecarajo you asked .. *ya lo borre* so chill *jajajajaja*

The sentence starts with a matrix sentence (L1) in English. An embedded clause in Spanish is then inserted (L2, italicized). The language is switched back to English. Followed by modal particles again in Spanish. We want to know how difficult it is for a monolingual spaCy model to deal with such inserted tokens as in this sentence.

### 4 Research Questions

The research questions can be summarized into three major points:

Firstly, how many inserted L2 tokens that are not named entities themselves are falsely annotated as named entities by spaCy? Named entities in nature are unique to the language’s vocabulary. Since the inserted L2 tokens are “foreign” to a language-specific model, we would expect that some of the unique words would be annotated as NEs.

Secondly, among the tokens that are falsely tagged as named entities, how many are just inserted non-NE L2 tokens? Namely, how many

errors are caused by code-switching? By analyzing this aspect, we are able to estimate how troublesome code-switching actually is for a spaCy model. Furthermore, we would also be able to determine whether it is reasonable to use spaCy on the named entity recognition tasks of code-mixed data.

Finally, we also calculate how many inserted L2 tokens that are actually named entities are correctly identified as NEs by the L1 model. This also helps us to understand how effective spaCy is, if we are only interested in extracting NEs tokens from the embedded language.

### 5 Data

The data comes from Computational Approaches to Linguistic Code-Switching (CALCS) (Aguilar et al., 2018), which are accessible through the LinCE Benchmark website (<https://ritual.uh.edu/lince/datasets>). The specific subset used in this research is the train set in Spanish-English (SPA-ENG) from CALCS Shared Task 2018.

The original data is stored in minimal CoNLL-U format. Sentences are separated by empty lines. The tokens in the sentence are all on individual lines. The tokens are marked with language tags (lang1 for English, lang2 for Spanish). Whether the tokens are named entities and to what type of named entities they belong to, is also marked overtly.

Overall, there are 33,611 sentences, among which 8,692 contain both English and Spanish tokens (code-mixed). By comparing the absolute language tag amount in each sentence, we get 1,478 instances with English as the matrix language, and 7,214 are Spanish dominant.

### 6 Methodology

Since spaCy models are usually built on monolingual data, the choice of which language-specific model should be used to annotate the current sentence needs to be made based on individual cases. The general pipeline goes as follows: 1) Before applying the NLP model, make sure both languages are present in the current sentence; 2) Determine the matrix language (L1) of the sentence; 3) Choose the spaCy model for L1 to annotate the whole sentence, regardless of the inserted L2 tokens; 4) Retrieve named entity recognition results from linguistic features (<https://spacy.io/usage/linguistic-features>) built in spaCy’s standard pipeline. For

efficiency reasons, here we only use small variants of models for both languages: `en_core_web_sm` and `es_core_news_sm`.

Since spaCy only works on strings, the raw token chain of each sentence will be concatenated in the first step. The new text string will be passed to the language model. This, however, causes alignment issues when spaCy's tokenization results differ from the gold tokens in the original sentence dataframe. The solution is to store the named entities extracted by spaCy simply as list items. Then we go through the original token column. If the token is found in the spaCy's NE list, we mark it as Yes, if not, as 0, following the original label scheme.

The focus of this experiment lies in whether the model could sufficiently distinguish between inserted code-switched normal words and named entities. The types of NEs the results tokens are assigned to are less relevant. Hence, we replace all specific named entity types in both gold tags and spaCy results with a unified value Yes.

For the results of each sentence, four types of information are saved: the matrix language, the gold language tags of all tokens in the sentence, the gold named entity property, and the NER results returned by spaCy model. The latter three are stored in a list aligned to the gold token list of the original sentence, on which the error analyses can be conducted.

## 7 Error Analyses

For the first research question, we extracted all indices of all code-switched tokens in each sentence. That is, if the matrix language of the sentence is identified as English, we extract all Spanish tokens, and vice versa. We also want to exclude cases where the inserted tokens are named entities themselves. With the list of target words' indices, the corresponding NER results from the automatic annotation will be compared. The target tokens falsely tagged as Yes by spaCy will be filtered out. As for results, we have 38.38% of normal L2 tokens tagged as named entities.

To investigate the relation between the performance drop and the insertions of L2 tokens, we first collect all named entity tags that are falsely given by spaCy. Then we map these errors to the language tags. If the token on which the error occurs is code-switched, it will be saved to the error list. Dividing NER errors on code-switching points

by total NER errors, we get results of 27.19%. That is, over a quarter of named entity recognition errors are simply due to L2 token insertions.

Lastly, we turn to the few cases where the inserted L2 tokens are actually named entities themselves. We collect all the L2 tokens in a sentence. Then we save those that are tagged as Yes in the normalized gold labels. Our goal is to find out the cases where spaCy also gives Yes tags to these target words (L2 NEs). From the results, we see that almost 70% of these true L2 named entities are correctly tagged as NEs. The error rate is relatively low compared to earlier analyses. However, it is unclear whether this is because spaCy could overcome the "language barrier" while dealing with L2 named entities or simply by chance. Given the fact that by assigning NE tags to all indiscriminately would even achieve 100% accuracy, the actual performances of the models are still questionable.

## 8 Conclusion

In this experiment, we looked into monolingual spaCy language models' performances on named entity recognition tasks with code-mixed data. We see that about 40% of inserted L2 tokens are falsely recognized as named entities, which indicates that code-switching does pose a challenge to monolingual language models in NER. However, under 30% of the errors are directly caused by code-switching. To what degree CS affects performance is therefore unclear.

## References

- Gustavo Aguilar, Fahad AlGhamdi, Victor Soto, Mona Diab, Julia Hirschberg, and Tamar Solorio. 2018. [Named Entity Recognition on Code-Switched Data: Overview of the CALCS 2018 Shared Task](#). In *Proceedings of the Third Workshop on Computational Approaches to Linguistic Code-Switching*, pages 138–147, Melbourne, Australia. Association for Computational Linguistics.
- Ahlan Alrehili. 2019. [Automatic hate speech detection on social media: A brief survey](#). In *2019 IEEE/ACS 16th International Conference on Computer Systems and Applications (AICCSA)*, pages 1–6.
- Njagi Dennis Gitari, Zhang Zuping, Hanyurwimfura Damien, and Jun Long. 2015. A lexicon-based approach for hate speech detection. *International Journal of Multimedia and Ubiquitous Engineering*, 10(4):215–230.
- Google. 2024. [Perspective api](#). Accessed: 2024-06-19.

- Ralph Grishman and Beth M Sundheim. 1996. Message understanding conference-6: A brief history. In *COLING 1996 Volume 1: The 16th International Conference on Computational Linguistics*.
- Ryuichiro Higashinaka, Kugatsu Sadamitsu, Kuniko Saito, Toshiro Makino, and Yoshihiro Matsuo. 2012. Creating an extended named entity dictionary from wikipedia. In *Proceedings of COLING 2012*, pages 1163–1178.
- Yinxia Lou, Tao Qian, Fei Li, and Donghong Ji. 2020. [A graph attention model for dictionary-guided named entity recognition](#). *IEEE Access*, 8:71584–71592.
- Tim Rocktäschel, Michael Weidlich, and Ulf Leser. 2012. [ChemSpot: a hybrid system for chemical named entity recognition](#). *Bioinformatics*, 28(12):1633–1640.
- Maarten Sap, Dallas Card, Saadia Gabriel, Yejin Choi, and Noah A. Smith. 2019. [The risk of racial bias in hate speech detection](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1668–1678, Florence, Italy. Association for Computational Linguistics.
- Jingbo Shang, Liyuan Liu, Xiang Ren, Xiaotao Gu, Teng Ren, and Jiawei Han. 2018. [Learning named entity tagger using domain-specific dictionary](#).
- Xuhui Zhou, Maarten Sap, Swabha Swayamdipta, Yejin Choi, and Noah Smith. 2021. [Challenges in automated debiasing for toxic language detection](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 3143–3155, Online. Association for Computational Linguistics.

## A Data Access

Project work and detailed instructions on how to retrieve the data needed for the analyses are accessible under: [DOI 10.17605/OSF.IO/ZSM43](https://doi.org/10.17605/OSF.IO/ZSM43).