

FReND: A French Resource of Negation Data

Hafida Le Cloirec - Ait Yahya*, Olga Seminck*, Pascal Amsili

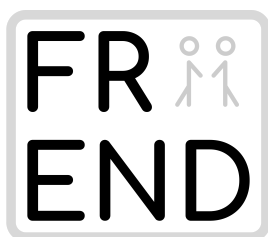
Lattice UMR 8094: CNRS, ENS-PSL Université, Université Sorbonne Nouvelle
1 rue Maurice Arnoux, 92120 Montrouge, France

hafida.le-cloirec@sorbonne-nouvelle.fr, olga.seminck@cnrs.fr, pascal.amsili@ens.fr

Abstract

FReND is a freely available corpus of French language in which negations are hand-annotated. Negations are annotated by their *cues* and *scopes*. Comprising 590 K tokens and over 8.9 K negations, it is the largest dataset available for French. A variety of types of textual genres are covered: literature, blog posts, Wikipedia articles, political debates, clinical reports and newspaper articles. As the understanding of negation is not yet mastered by current state of the art AI-models, FReND is not only a valuable resource for linguistic research into negation, but also as training data for AI tasks such as negation detection.

Keywords: Negation, Corpus, French



French Resource of Negation Data

1. Introduction

Negation is an important component of language and a crucial aspect when extracting information from text as it inverses the truth value of propositions. The need to detect negations and the exact propositions that fall under their scope was first underlined in the domain of biomedical texts (Chapman et al., 2001). Indeed, understanding negation is very important in the automatic selection of patients for clinical trials and in making automatically correct decisions about clinical reports. With this aim in view, fifteen years ago the Bioscope Corpus (Szarvas et al., 2008) was developed and in 2012, the task of “detecting the scope of negation” was formalized during the 2012 *SEM Shared Task (Morante and Blanco, 2012). In order to have more annotated data — and also because the community recognized the importance of negation scope detection beyond the domain of bio-medical texts — the Conan Doyle Corpus (Morante and Daelemans, 2012) was developed and used in this Shared Task.

Negation data continue to play a crucial role in the field of NLP. Not only was it found that the ‘understanding’ of negation is limited by the language models of the BERT-generation (Ettinger, 2020; Kassner and Schütze, 2020; Hossain et al., 2020;

Kletz et al., 2023), the latest generation of very large language models still also largely underperforms on tasks and benchmarks featuring negation (Ye et al., 2023; Jang et al., 2023; García-Ferrero et al., 2023; Chen et al., 2023). It is therefore necessary to compile and use negation data in order to assess and improve language models’ performance on common sense reasoning and natural language inference tasks (Hosseini et al., 2021; Hossain et al., 2022; Kletz et al., 2023; García-Ferrero et al., 2023; Truong et al., 2022).

Whereas most corpora have been developed for the English language, Jiménez-Zafra et al. (2020) provide a list of negation corpora developed until 2020 that also includes resources for Spanish, Swedish, Chinese, Dutch, German, and Italian. French remains an under-resourced language in the matter of negation. The only resources covering this language (and also Brazilian Portuguese) are the ESSAI and CAS corpora created by Daloux et al. (2021) who annotated clinical reports for a total of 238 K tokens, 10.4 K sentences and 1829 ‘negative sentences’¹. With the goal of enabling more in-depth research into negation in the French language, we decided to build the FReND corpus, covering text genres, annotated for negation, that are not yet available for French: newspaper texts, Wikipedia articles, political debate, literary texts and blog post texts. This resource, which is two- to threefold the size of the ESSAI and CAS corpora covering these new genres, will enable researchers to better understand negation in French and will provide training data for machine learning tasks. Moreover, it is important to note that our resource only contains natural data, in contrast with most of the datasets used to measure and

¹We find the term *negative sentences* somewhat confusing as negations scope over propositions of which a sentence can contain more than one.

*These authors contributed equally to this work.

enhance the performance of language models on NLI that contain artificially augmented data following specific patterns (Hosseini et al., 2021; Kletz et al., 2023; Garcia-Ferrero et al., 2023; Truong et al., 2022).

2. Methods

2.1. Choice of Texts

We aimed for diverse genres of texts in our corpus. Furthermore, we wished to include only texts that were eligible for distribution with an open licence. We chose to include the following three corpora: the sequoia corpus (Candito and Seddah, 2012) featuring news articles, Wikipedia articles, clinical reports and political debates; the fr-litbank corpus (Lattice, 2023) that contains literary texts; and payetoncorpus (CLLE, 2021), a corpus of blog posts with *meToo* testimonies.

2.2. Annotation Guidelines

The annotation of each text consists of:

1. The identification and delimitation of the word(s) that are negation cues.
2. The identification and delimitation of the part of the sentence affected by the negation cue, i.e. the scope.
3. The association of the scope with its cue.

In the examples given in this guide, cues will be marked in bold and scopes in brackets.

2.2.1. Cue

A negation cue is a linguistic form that has the role of the logical operator \neg that inverses the truth value of a proposition. In French, negation cues can be expressed by:

- one word (*sans*; without)
- several words (*ne ... pas*; not, *plus ... jamais*; never)
- a morpheme (*il-*, *in-*; e.g. *illégal*; illegal)

In Appendix A, we present the list of the most frequent cues in which all these options can be observed.

We have defined one test to determine whether a morpheme is a negation cue or not, and three tests on negation words in general.

MORPHEME TEST To determine if a morpheme has a negative meaning, we verify:

- if it is transparent; i.e. if when we remove the morpheme, we get a word that makes sense.
- if the meaning is compositional, i.e. if the meaning of morpheme+word is equivalent to the meaning of not+word.

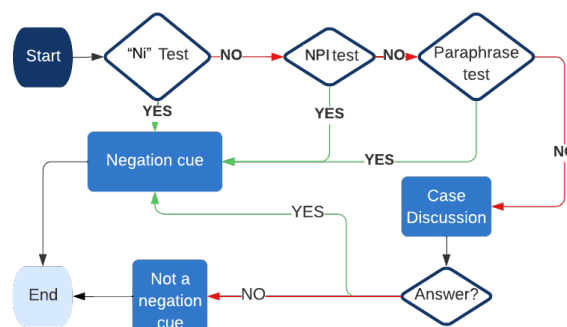


Figure 1: Cue Tests

If and only if these two tests are conclusive, then the morpheme can be annotated as a negation cue. In example (1), *imprudent* in French is fully compositional, *prudent* being an existing word and the meaning of *imprudent* being *not prudent*.

- (1) [Watson a été] **im**[prudent].
Watson was imprudent.

NI TEST: If *ni* (nor) can be added as a conjunction to a word or a phrase, that word or phrase can be considered as a negation cue, as in example (2). If it cannot be added, that should not be considered as indicating that it is not a negation cue: further analysis is needed as in (3) where we can see that it still passes the NPI Test (*qui que ce soit*; ‘anyone’ being an NPI, see below) and that we can consider it a negation.

- (2) Il a laissé tomber ses amis **sans** les aider **ni** les écouter.
He let his friends down without helping nor listening to them.
- (3) a. ?Il a laissé tomber ses amis au lieu de les aider ni les écouter.
He let his friends down instead of helping nor listening to them.
 b. Il a laissé tomber ses amis au lieu d’aider qui que ce soit.
He let his friends down instead of helping anyone.

NPI TEST: Negative Polarity Items are expressions that can only appear in a negative context, or in semantic contexts that have affinities with negation (*du tout* in example (4)).

- (4) a. [Je] n’[ai] **pas** [aimé le film du tout].
I did not like the movie at all.
 b. *J’ai aimé le film du tout.
I liked the movie at all.

If a cue passes the NPI test (i.e. can be combined with an NPI and form a semantically correct sentence, see example (4)), we can consider it to be a

negation cue. But if it fails the NPI test, that does not necessarily mean it is not a negation cue.

PARAPHRASE TEST: We can check whether the paraphrase *Ce n'est pas le cas que...* 'It is not the case that...' can create a sentence with the same semantic content if we delete/replace the cue. If this is the case, we consider that the proposition contains a negation, see example (5).

- (5) a. [Jean] n'[a] **pas** [poursuivi ses rêves].
Jean didn't chase his dreams.
b. Ce n'est pas le cas que Jean a poursuivi ses rêves.
It is not the case that Jean chased his dreams.

Figure 1 describes the testing process for cues.

2.2.2. Scope

The scope is the part of the proposition affected by the negation. We consider that negation is an operation that affects a complete proposition and not only the Verb Phrase where it is syntactically realized most often in French. Therefore, in most cases the scope comprises the whole proposition even when negation appears only within the VP. The scope can be situated to the right or left of the cue, extend on both sides (6), or be discontinuous (7). In case of several cues, scopes may overlap, for example in example (8) where there is a double negation (which leads to a positive sense, just as in English).

- (6) [c'est] **pas** [grave]
it doesn't matter
- (7) Evidemment, [les candidats], eux, n'[ont fait l'objet d']**aucune** [remarque sur leur situation personnelle]...
Of course, the male candidates, on the other hand, did not receive any comments on their personal situation.
- (8) a. [Pécuchet] **ne** [put s'empêcher de dire]:
Pécuchet could not help saying:
b. Pécuchet ne put [s']**empêcher de** [dire]:
Pécuchet could not help saying:

Scope annotation depends on the cue, the context and the syntactic structure of the sentence. We have designed some tests in order to determine the extent of the scope.

PARAPHRASE TEST: To find out which part of the sentence is affected by a negation cue, we can use the paraphrase test with the expression *Il n'est pas le cas que...* (It is not the case that...) or *Il est faux que...* (It is false that...). We can see from the result of the tests (9-b) and (9-c) that the scope is as in example (9-a).

- (9) a. [Ce schtroumpf] n'[est] **pas** [jaune].
This smurf is not yellow.
b. *Ce n'est pas le cas que* [ce schtroumpf est jaune].
It is not the case that this smurf is yellow.
c. *Il est faux que* [ce schtroumpf est jaune].
It is false that this smurf is yellow.

QUESTION-ANSWER TEST: The question-answer test allows us to determine which elements of the sentence should be included in the scope. For example, for (10), we could wonder whether 'to the bridge' should be included or excluded from the scope.

- (10) Nous n'avons pas conduit jusqu'au pont.
We did not drive to the bridge.

The question-answer test helps us to decide on this problem. The test involves asking whether the element in question can be inferred by the negated sentence. If this is not the case, it means that negation scopes over this element. For example (10), the question is:

- Can we infer that they drove to the bridge? → NO

The negation therefore scopes over the "drive to the bridge" section of the sentence and includes the whole sentence (except the cue) as in (11):

- (11) [Nous] n' [avons] **pas** [conduit jusqu'au pont].
We did not drive to the bridge.

2.3. Team and Training

Annotation guidelines were developed by three linguists: a PhD-student, her PI and a research engineer. For the annotation, 5 undergraduate students from Sorbonne Nouvelle's Linguistics Department were recruited and instructed by the linguists. They were given a copy of the annotation guidelines and annotated separately the *fr-wiki* subcorpus of the sequoia corpus (comprising around 19 K words). The PhD student and the engineer also annotated this same section. All 7 annotators discussed their annotations together during three sessions and adjudged on conflicts (so that *fr-wiki* could be integrated in the final corpus). The annotation guidelines were completed by adding the decisions on cases not foreseen in advance. This was the end of the training, but undergraduate students were instructed to contact the linguists for help in case of conflicts they were not able to solve. All 7 annotators annotated the same amount of texts.

2.4. Annotation Tools

We used the brat annotation tool (Stenetorp et al., 2012). This tool allows annotators to select spans

of text by highlighting them with a mouse and then selecting from a scroll-down menu whether it is a cue or a scope. Moreover, it is possible to select a discontinuous span. This is necessary in French as cues often come in multiple parts (e.g. *ne ... pas*; not). Annotators can establish the negation relationship between cues and scopes by dragging their mouse between two spans of text.

2.5. Corpus Format

The corpus is distributed in an XML-format inspired by the Bioscope Corpus (Szarvas et al., 2008). We transformed the brat annotations to this format to make it more human-readable. The XML format also allowed us to integrate useful information in attributes, such as the name of the annotator(s), the source of the text and whether the text is in the training, development or test split when the corpus is used for machine learning. We predefined the splits to ensure that researchers will produce comparable results when using the corpus. The machine learning split has the following distribution: 4/6 training, 1/6 development, 1/6 test. Each sub-corpus of FReND can be found under a tag `<DocumentSet>` which contains the original documents from the sequoia, fr-litbank and payetoncorpus. These documents are labeled `<Document>`. Inside these tags, we can find the tags `<DocumentPart>`. Documents have been truncated to fit easily on one computer screen to make the annotation in brat easier. For each `<DocumentPart>` one can find the annotator(s) and whether the texts have been adjudged.

2.6. Adjudication

Besides the *fr-wiki* subcorpus that was annotated by all 7 annotators and adjudged, 40% of the texts of the corpus underwent adjudication and were annotated by two people. Once they had saved their individual annotations in an online drive, conflicts were detected after running the two annotations through the brat library, developed to calculate inter-annotator agreement for brat annotations (Kolditz et al., 2019). The two annotators then discussed the conflicts and conducted another round of annotation to resolve them using the annotation guidelines. The resolved conflicts were then saved in a gold version of the annotation.

3. Results

3.1. Corpus Statistics

We counted the number of tokens and sentences using the Spacy toolkit (Honnibal et al., 2020). The number of negations and scopes and cues were counted using X-path queries. The corpus statistics for the sequoia corpus, fr-litbank and payetoncorpus can be found in Table 1 and detailed statis-

tics of the sequoia corpus that is made up of a diversity of text genres in Appendix B.

| | # K tokens | # K sentences | # negations | # negs/# sent. | # negs/K toks | # negs no scp. |
|-------------|------------|---------------|-------------|----------------|---------------|----------------|
| seq. | 69 | 3.1 | 673 | 0.22 | 9.8 | 4 |
| fr-lb. | 224 | 11.7 | 3238 | 0.28 | 14.4 | 91 |
| ptc. | 297 | 18.0 | 5047 | 0.28 | 17.0 | 314 |
| tot. | 590 | 32.8 | 8958 | 0.27 | 15.2 | 409 |

Table 1: Corpus Statistics. Number of: tokens, sentences estimated by Spacy, negations, negations per sentence, negations per 1000 tokens, and negations without a scope for the sequoia corpus (seq.), fr-litbank (fr-lb.) and payetoncorpus (ptc.) that make up FReND.

3.2. Inter-annotator Agreement

Inter-annotator agreement was calculated on 40% of the texts, except the fr-wiki subcorpus of sequoia that served as a training set for all annotators. The remaining 60% was annotated by only one person. To give an estimation of the quality of the annotation, we measured for each annotator the F1-score for cues and scopes by comparing their annotations to the gold version (established during the adjudication phase) using the brat library (Kolditz et al., 2019). Table 2 gives the scores for so-called ‘exact matches’ of spans (also the ‘instance’ score) and Table 3 presents the F1-score when we do not use exact matches of spans, but when each token composing the span is counted separately. It can be seen that the quality of the annotation is higher for cues than for scopes (see Table 2 where the macro-average for cues is 0.91 against 0.81 for scopes when measured on instances). However, the token-scores indicate that this difference is probably due to scopes being a lot longer than cues, with the result that many annotation errors on scope are not the result of a significant disagreement between annotators but of forgetting a (small) part of the scope when annotating (see Table 3 where the macro-average for cues is 0.96 and for scopes 0.91).

3.3. Distribution and Licence

FReND, its DTD and complete guidelines in French are distributed under the Creative Commons License: Attribution-NonCommercial-ShareAlike 4.0 International (CC BY-NC-SA 4.0) and can be downloaded here:

<https://github.com/lattice-8094/FReND>.

| | Cue | Scope | Total |
|------------------|-------------|-------------|-------------|
| A | 0.92 | 0.82 | 0.87 |
| B | 0.96 | 0.90 | 0.93 |
| C | 0.92 | 0.80 | 0.86 |
| D | 0.89 | 0.78 | 0.84 |
| E | 0.88 | 0.69 | 0.79 |
| F | 0.94 | 0.86 | 0.90 |
| G | 0.89 | 0.76 | 0.82 |
| Macr.Avg. | 0.91 | 0.81 | 0.86 |

Table 2: F1-score per instance for the annotators A to G and macro-average.

| | Cue | Scope | Total |
|------------------|-------------|-------------|-------------|
| A | 0.96 | 0.92 | 0.93 |
| B | 0.98 | 0.96 | 0.96 |
| C | 0.95 | 0.90 | 0.91 |
| D | 0.94 | 0.88 | 0.90 |
| E | 0.94 | 0.86 | 0.88 |
| F | 0.97 | 0.94 | 0.94 |
| G | 0.95 | 0.92 | 0.93 |
| Macr.Avg. | 0.96 | 0.91 | 0.92 |

Table 3: F1-score per token for the annotators A to G and macro-average.

4. Contributions of FReND

The availability of corpora annotated with negation is essential when training automatic negation processing systems. They should also prove useful in view of the general trend in the NLP field towards the increasing use of Large Language Models (LLMs). As hallucination is an important issue with LLMs, having an annotated corpus can help researchers who develop LLMs to better capture which parts of the texts contain facts and which do not by using the negation data as special features in a fine tuning phase.

The fact that negation can be expressed through different morpho-syntactic mechanisms depending on the language under study, highlights the importance of developing data in various languages. With the development of the FReND corpus, French is no longer an under-resourced language. Therefore, the phenomenon of negation can now be compared more easily between French and other languages. The corpus statistics reveal that negation has a high frequency in the blog corpus (payetoncorpus), literature (fr-litbank) and transcribed debates whereas it is scarcer in biomedical and newspaper texts (see Table 1 and Appendix B). The decision to include a wide variety of text genres in FReND allows us to study this type of differences. And last but not least, we developed new linguistic tests to annotate negation, for example the NPI and the NI tests. As the inter-annotator agreement was rather high, we can conclude that these tests are effective and could

be used for the study of negation and the development of new corpora.

5. Acknowledgements

Our thanks go to the interns who worked with us on this project for their meticulous annotation, and their commitment and enthusiasm during the discussions we had about the annotation guidelines and borderline cases. Thank you Erika Bisbau, Cédric Dumont, Natacha Miniconi, Faustine Monthubert and Rossana Verdier.

6. Bibliographical References

Marie Candito and Djamé Seddah. 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 *Proceedings of the Joint Conference JEP-TALN-RECITAL 2012, volume 2: TALN*, pages 321–334, Grenoble, France. ATALA/AFCP.

Wendy W Chapman, Will Bridewell, Paul Hanbury, Gregory F Cooper, and Bruce G Buchanan. 2001. A simple algorithm for identifying negated findings and diseases in discharge summaries. *Journal of biomedical informatics*, 34(5):301–310.

Jiangjie Chen, Wei Shi, Ziquan Fu, Sijie Cheng, Lei Li, and Yanghua Xiao. 2023. *Say what you mean! large language models speak too positively about negative commonsense knowledge*. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9890–9908, Toronto, Canada. Association for Computational Linguistics.

CLLE. 2021. *Payetoncorpus*. ORTOLANG (Open Resources and TOols for LANGUAGE) –www.ortolang.fr.

Clément Dalloux, Vincent Claveau, Natalia Grabar, Lucas Emanuel Silva Oliveira, Claudia Maria Cabral Moro, Yohan Bonescki Gumiel, and Deborah Ribeiro Carvalho. 2021. Supervised learning for the detection of negation and of its scope in french and brazilian portuguese biomedical corpora. *Natural Language Engineering*, 27:181–201.

Allyson Ettinger. 2020. *What BERT is not: Lessons from a new suite of psycholinguistic diagnostics for language models*. *Transactions of the Association for Computational Linguistics*, 8:34–48.

- Iker García-Ferrero, Begoña Altuna, Javier Alvez, Itziar Gonzalez-Dios, and German Rigau. 2023. [This is not a dataset: A large negation benchmark to challenge large language models](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 8596–8615, Singapore. Association for Computational Linguistics.
- Matthew Honnibal, Ines Montani, Sofie Van Landeghem, and Adriane Boyd. 2020. ["spacy: Industrial-strength natural language processing in python"](#).
- Md Mosharaf Hossain, Dhivya Chinnappa, and Eduardo Blanco. 2022. [An analysis of negation in natural language understanding corpora](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 716–723, Dublin, Ireland. Association for Computational Linguistics.
- Md Mosharaf Hossain, Venelin Kovatchev, Pranoy Dutta, Tiffany Kao, Elizabeth Wei, and Eduardo Blanco. 2020. [An analysis of natural language inference benchmarks through the lens of negation](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9106–9118, Online. Association for Computational Linguistics.
- Arian Hosseini, Siva Reddy, Dzmitry Bahdanau, R Devon Hjelm, Alessandro Sordani, and Aaron Courville. 2021. [Understanding by understanding not: Modeling negation in language models](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1301–1312, Online. Association for Computational Linguistics.
- Joel Jang, Seonghyeon Ye, and Minjoon Seo. 2023. [Can large language models truly understand prompts? a case study with negated prompts](#). In *Proceedings of The 1st Transfer Learning for Natural Language Processing Workshop*, volume 203 of *Proceedings of Machine Learning Research*, pages 52–62. PMLR.
- Salud María Jiménez-Zafra, Roser Morante, María Teresa Martín-Valdivia, and L. Alfonso Ureña-López. 2020. [Corpora annotated with negation: An overview](#). *Computational Linguistics*, 46(1):1–52.
- Nora Kassner and Hinrich Schütze. 2020. [Negated and misprimed probes for pretrained language models: Birds can talk, but cannot fly](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7811–7818, Online. Association for Computational Linguistics.
- David Kletz, Pascal Amsili, and Marie Candito. 2023. [The self-contained negation test set](#). In *Proceedings of the 6th BlackboxNLP Workshop: Analyzing and Interpreting Neural Networks for NLP*, pages 212–221.
- Tobias Kolditz, Christina Lohr, Johannes Hellrich, Luise Modersohn, Boris Betz, Michael Kiehn-topf, and Udo Hahn. 2019. [Annotating german clinical documents for de-identification](#). In *Med-Info*, pages 203–207.
- Lattice. 2023. fr-litbank. <https://github.com/lattice-8094/fr-litbank>.
- Roser Morante and Eduardo Blanco. 2012. [*SEM 2012 shared task: Resolving the scope and focus of negation](#). In **SEM 2012: The First Joint Conference on Lexical and Computational Semantics – Volume 1: Proceedings of the main conference and the shared task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation (SemEval 2012)*, pages 265–274, Montréal, Canada. Association for Computational Linguistics.
- Roser Morante and Walter Daelemans. 2012. [ConanDoyle-neg: Annotation of negation in Conan Doyle stories](#). In *Proceedings of the Eighth International Conference on Language Resources and Evaluation, Istanbul*, pages 1563–1568.
- Pontus Stenetorp, Sampo Pyysalo, Goran Topić, Tomoko Ohta, Sophia Ananiadou, and Jun'ichi Tsujii. 2012. [brat: a web-based tool for NLP-assisted text annotation](#). In *Proceedings of the Demonstrations at the 13th Conference of the European Chapter of the Association for Computational Linguistics*, pages 102–107, Avignon, France. Association for Computational Linguistics.
- György Szarvas, Veronika Vincze, Richárd Farkas, and János Csirik. 2008. [The BioScope corpus: annotation for negation, uncertainty and their scope in biomedical texts](#). In *Proceedings of the Workshop on Current Trends in Biomedical Natural Language Processing*, pages 38–45, Columbus, Ohio. Association for Computational Linguistics.
- Thinh Hung Truong, Yulia Otmakhova, Timothy Baldwin, Trevor Cohn, Jey Han Lau, and Karin Verspoor. 2022. [Not another negation benchmark: The NaN-NLI test suite for sub-clausal negation](#). In *Proceedings of the 2nd Conference*

of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 883–894, Online only. Association for Computational Linguistics.

Mengyu Ye, Tatsuki Kuribayashi, Jun Suzuki, Goro Kobayashi, and Hiroaki Funayama. 2023. [Assessing step-by-step reasoning against lexical negation: A case study on syllogism](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 14753–14773, Singapore. Association for Computational Linguistics.

A. Appendix: Frequency of Cues

| Cue | Type | Frequency |
|-----------------|-------------|-----------|
| ne pas | multi-token | 3506 |
| pas | ADV | 784 |
| sans | ADV | 616 |
| in | morpheme | 561 |
| non | ADV | 463 |
| ne plus | multi-token | 342 |
| ne rien | multi-token | 341 |
| ne jamais | multi-token | 270 |
| ne | ADV | 250 |
| im | morpheme | 127 |
| ne point | multi-token | 88 |
| rien | ADV | 70 |
| ne aucune | multi-token | 63 |
| dés | morpheme | 48 |
| ne aucun | multi-token | 46 |
| refuse | V | 44 |
| rien ne | multi-token | 43 |
| mal | ADV | 43 |
| personne ne | multi-token | 42 |
| jamais ne | multi-token | 34 |
| dé | morpheme | 32 |
| jamais | ADV | 31 |
| ne guère | multi-token | 28 |
| personne n' | multi-token | 25 |
| a refusé | multi-token | 25 |
| aucun ne | multi-token | 24 |
| ir | morpheme | 24 |
| ne ni ni | multi-token | 23 |
| aucune ne | multi-token | 22 |
| sauf | ADV | 22 |
| ne personne | multi-token | 21 |
| aucune | DET | 21 |
| ne pas ni | multi-token | 20 |
| il | morpheme | 20 |
| sans aucune | multi-token | 19 |
| rien n' | multi-token | 18 |
| ne plus rien | multi-token | 17 |
| refuser | V | 17 |
| refusé | V | 17 |
| plus | ADV | 16 |
| ne plus jamais | multi-token | 15 |
| sans ni | multi-token | 14 |
| ne pas du tout | multi-token | 13 |
| ne pas non plus | multi-token | 12 |
| aucun | DET | 12 |
| refusent | V | 12 |
| ne toujours pas | multi-token | 12 |
| non pas | multi-token | 11 |
| refusait | V | 11 |
| ne pu | multi-token | 11 |
| mé | morpheme | 11 |

Table 4: The most frequent negation cues of the FReND corpus (frequency more than 10). Cues with *n'* have been fused with cues with *ne*.

B. Appendix: Detailed Corpus Statistics of the Sequoia Corpus

| | # K tokens | # K sentences | # negations | # negs/# sent. | # negs/K toks | # negs no scp. |
|------|------------|---------------|-------------|----------------|---------------|----------------|
| wiki | 23 | 1.00 | 160 | 0.22 | 7.1 | 0 |
| emt. | 10 | 0.44 | 112 | 0.25 | 11.1 | 0 |
| emd. | 10 | 0.56 | 146 | 0.26 | 15.2 | 1 |
| ann. | 11 | 0.53 | 75 | 0.14 | 6.6 | 2 |
| Eupr | 15 | 0.56 | 180 | 0.32 | 17.0 | 3 |

Table 5: Number of: tokens, sentences estimated by Spacy, negations, negations per sentence, negations per 1000 tokens, and negations without a scope for all the subcorpora of the sequoia corpus (i.e. fr-wiki (wiki) containing Wikipedia articles in French, emea-fr-test (emt.) and emea-fr-dev (emd.) containing biomedical reports, annodis (ann.) containing newspaper articles and Europar (Eupr) containing transcribed debates from the European parliament).