

# Say “No” to Missing Polarity: A Negation Enrichment of the Porttinari UD Treebank

Isaac Souza de Miranda Junior<sup>\*◇</sup>, Oto Araújo Vale<sup>\*</sup>, Marie-Catherine de Marneffe<sup>◇</sup>

Universidade Federal de São Carlos<sup>\*</sup>, Université catholique de Louvain<sup>◇</sup>  
isc\_jr@live.com, otovale@ufscar.br, marie-catherine.demarneffe@uclouvain.be

## Abstract

Negation is a central phenomenon in linguistics: every language has some way of expressing the difference between an affirmative sentence and a negative one (Horn and Wansing, 2025). However, the treatment of negation remains uneven in Natural Language Processing (Jimenez-Zafra et al., 2017; Jiménez-Zafra et al., 2020). This paper presents the enrichment of a Brazilian Portuguese corpus with negation-related morphological information within the Universal Dependencies (UD) framework (Nivre et al., 2020; de Marneffe et al., 2021). We enrich the Porttinari-base corpus (Duran et al., 2023) by systematically adding the UD morphological features `Polarity=Neg` and `PronType=Neg` for 18 negation-related lexical items. The enrichment only modifies the morphological features, leaving tokenization and dependency structure unchanged. To evaluate the computational results of this enrichment, we present an experiment using the Brazilian Portuguese parser PortParser (Lopes and Pardo, 2024), which we trained both on the original Porttinari-base data (Duran et al., 2023) and on our enriched version. Our results show that after enrichment, the parser’s performance remains stable, and the newly introduced features are being learned.

**Keywords:** linguistic annotation, Universal Dependencies, corpus enrichment, Brazilian Portuguese, negation

## 1. Introduction

The ability to express opposition between utterances with respect to their truth conditions is widely regarded as a defining property of human language (Horn, 2001). Although some animal species can be trained to respond to notions such as rejection or absence, they do not exhibit propositional negation, i.e., the capacity to produce utterances whose truth conditions systematically contradict those of another utterance (Heine and Kuteva, 2007). Negation is therefore an intrinsic component of human linguistic competence, closely linked to logical evaluation, inference, and discourse/pragmatic level (Grice, 1989; Schwenter, 2016).

From a linguistic perspective, negation raises well-known challenges related to scope and quantification, which have been extensively discussed in semantics (Russell, 1905; Wittgenstein, 1922; Frege, 1948; Grice, 1989). Comparatively, less work on negation has been done in Natural Language Processing (NLP) (Jimenez-Zafra et al., 2017; Jiménez-Zafra et al., 2020). In terms of negation and scope detection, work mainly happened in the context of biomedical research (Chapman et al., 2001; Mutalik et al., 2001; Goldin and Chapman, 2003).

Universal Dependencies (UD) is an open community project dedicated to developing treebank annotations (collections of syntactically annotated texts) that are consistent across multiple languages (Nivre et al., 2020). It has features for the annotation of negation-related morphology, but such features remain unevenly used across the UD treebanks (Findlay and Haug, 2025), including in resources for Portuguese. Following the UD guidelines for negation-related morphology, this work

presents an enrichment of the Porttinari-base corpus (Duran et al., 2023), adding negation information at the UD morphological level, following explicit linguistic criteria grounded in descriptive studies of Brazilian Portuguese on negation (Mito, 1992, 1998; Neves, 2000; Miranda Junior, 2022).

We evaluate the impact of our morphological enrichment of the corpus in a parsing setup to assess whether the new features can be learned and what their potential impact on parsing performance is. We train a Brazilian Portuguese parser, PortParser (Lopes and Pardo, 2024), on the original and enriched versions of the Porttinari-base corpus and compare the performance. This experiment aims to investigate whether negation-related morphology can be integrated into the corpus while retaining parsing performance. Our results show that the parser’s performance remains stable, and that the newly introduced features were learned with high reliability.

The paper is structured as follows. In Section 2, we survey related work. In Section 3, we present the data and enrichment procedure. In Section 4, we describe the parsing evaluation setup. We report results in Section 5 and conclude in Section 6.

## 2. Related Work

In this section, we first briefly survey NLP work focusing on negation. We then discuss the UD principles for annotating negation as well as the current state of the Portuguese corpora available in the UD version 2.17 (Nov. 2025).

## 2.1. Negation in Natural Language Processing and Annotated Corpora

Negation has long been recognized as a challenging phenomenon in NLP due to its interactions with scope, polarity, and focus. Early computational work focused on identifying negation cues and their scope, particularly in biomedical and clinical domains (Chapman et al., 2001; Mutalik et al., 2001; Goldin and Chapman, 2003). This line of research was later extended through annotated corpora and shared tasks targeting negation and speculation, notably with BioScope (Szarvas et al., 2008), an open source biomedical corpus annotated to identify negation and uncertainty terms, along with their syntactic scopes. It contains over 20,000 sentences extracted from medical reports, articles, and scientific abstracts for training computational models. Several BioScope-based scope-resolution systems for biomedical text incorporate syntactic parse structures in their modeling, aiming to detect negation cues and negation scope, see Morante and Daelemans (2009) and Zou et al. (2013) as examples.

More recently, large-scale surveys have highlighted the persistent difficulty NLP faces in handling negation, especially in tasks that involve scope, focus, and semantic interpretation (Jiménez-Zafra et al., 2020). Kletz (2025) conducted a large survey on how BERT-based and generative AI models encode and process negation. His results show that small models have some difficulty processing negation, but this difficulty decreases as the amount and quality of data increase. Thus, most of the work concentrates on semantic or discourse-level problems (cue detection, scope resolution, textual inference, and sentiment-related phenomena). At the same time, comparatively less attention is paid to making negation explicit in low-level linguistic strata, such as morphology.

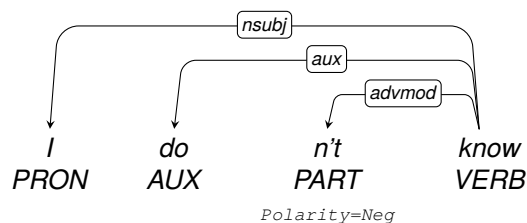
The present study is situated in the need for marking negation-related lexical items in the morphological layer of a UD dependency treebank.

## 2.2. Negation in Universal Dependencies

Within Universal Dependencies (UD) (Nivre et al., 2020; de Marneffe et al., 2021), negation is treated as a cross-linguistically relevant phenomenon and may be encoded through morphological features and lexical distinctions. Rather than being represented by a single word class, negative meaning may surface through sentential negators, negative pronouns, determiners, adverbs, or other lexical categories, depending on the language. The UD guidelines capture this phenomenon with two ded-

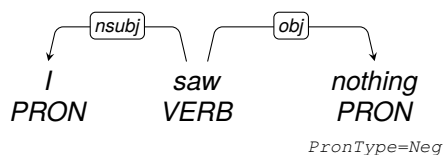
icated morphological features: `Polarity=Neg`<sup>1</sup> and `PronType=Neg`.<sup>2</sup> The first feature is the typical marking for negative polarity, while the second is specifically dedicated to pronouns that express negative meaning. We can see how both features are used in examples 1 and 2.

**Example 1.** ‘I don’t know.’



In UD, the negator *n't* (lemma *not*) is attached as `advmod` to the main verb and carries the feature `Polarity=Neg`.

**Example 2.** ‘I saw nothing.’



Here, the negative item *nothing* functions as the object of the verb and is morphologically specified as `PronType=Neg`.

Recent empirical work has shown that negation-related annotation is applied unevenly across UD treebanks. Based on a survey of the Universal Dependencies v2.15 (Nov. 2024), Findlay and Haug (2025) report that 224 out of 296 treebanks (76%) use `Polarity=Neg`, while only 99 (33%) use `PronType=Neg`. Moreover, 62 treebanks (21%) do not use any of the features, indicating that a substantial portion of the UD resources do not annotate negation at the morphological level. Findlay and Haug (2025) also discuss the limitations of the UD framework to account for semantic effects of negation. Using morphological features on negative items cannot account for phenomena such as double negation or negative concord, and negation scope. We agree with this statement; however, marking negative items is the first step towards representing more complex phenomena linked to negation.

We replicate (Findlay and Haug, 2025)'s survey on the most recent UD release (v2.17, Nov. 2025). The results show a comparable distribution despite the repository's growth. Out of 339 treebanks, 248

<sup>1</sup><https://universaldependencies.org/u/feat/Polarity.html>

<sup>2</sup><https://universaldependencies.org/u/feat/PronType.html>

(73%) use `Polarity=Neg`, 120 (35%) use `Pron-Type=Neg`, and 74 (21%) do not use any of the features. In absolute terms, the number of treebanks employing negation-related features has increased; proportionally, however, the distribution remains remarkably stable. Approximately one quarter of the UD treebanks still omit explicit morphological marking of negation, and the adoption of `Pron-Type=Neg` continues to lag substantially behind `Polarity=Neg`.

### 2.2.1. Negation-Related Annotation in Portuguese UD Treebanks

Since our study focuses on the systematic enrichment of negation-related features, it is important to situate Porttinari-base within the broader landscape of available Portuguese resources in UD. We therefore provide a comparative overview of existing treebanks in the v2.17 (Nov. 2025) release, considering their size, domain, annotation history, and alignment with UD standards.

**UD\_Portuguese Bosque** (Rademaker et al., 2017) is the UD conversion of the Bosque corpus. Bosque is a subcorpus of Floresta Sintática (Afonso et al., 2002), a large Brazilian Portuguese treebank annotated morphologically and syntactically with the PALAVRAS parser (Bick, 2014). The Floresta Sintática comprises four subcorpora: Bosque (fully revised), Selva (partially revised), and Floresta Virgem and Amazônia (unrevised). The UD Portuguese Bosque corpus contains 9,357 sentences, totaling 227,827 tokens. Rather than being annotated from scratch under UD, it was produced via a context-sensitive Constraint Grammar-based conversion with the PALAVRAS parser (Bick, 2014), followed by targeted manual corrections to ensure UD compliance. The resource includes text from the news domain of both European Portuguese material from CETEMPúblico (Rocha and Santos, 2000) and Brazilian Portuguese material from CETEN-Folha (Afonso et al., 2000), preserving part of the rich linguistic information of the original annotation while providing a UD v2 representation.

**CINTIL-UDep** is a Portuguese dependency resource derived from the manually validated CINTIL treebank collection (Branco et al., 2022). The corpus contains 38,400 sentences, totaling 475,860 tokens. The data came largely from the news domain. Its dependency layer was originally produced by expert linguists, supported by deep grammatical processing with LXGram (Costa and Branco, 2010) and double-blind annotation with adjudication. UD annotations were obtained via deterministic rule-based conversion from the original dependency scheme, resulting in a UD compliant treebank that substantially expands the volume of manually validated Portuguese data (Branco et al., 2022).

**DANTEStocks** is a UD annotated treebank of

Brazilian Portuguese tweets from the stock market domain, developed within the Porttinari initiative (Felippo et al., 2024). The corpus contains 4,042 tweets (80,998 tokens) and treats each tweet as a single syntactic unit, preserving fragmentation and platform-specific conventions typical of user-generated content, such as hashtags, mentions, URLs, and non-standard orthography. Annotation followed a semi-automatic workflow, alternating automatic processing and expert revision, guided by UD specific recommendations for Brazilian Portuguese and Twitter phenomena (Felippo et al., 2024).

**Portuguese-GSD** is the converted part of the Google Universal Dependency treebanks for Portuguese, originally introduced in a cross-linguistic effort to harmonize dependency annotation across six languages (German, English, Swedish, Spanish, French, and Korean) (McDonald et al., 2013). The corpus has 12,020 sentences (318,666 tokens) from news and blogs.<sup>3</sup> The original annotation was produced through automatic preprocessing and manual revision, followed by a harmonization stage to ensure cross-linguistic consistency (McDonald et al., 2013). There is not much information on how the Portuguese version was adapted from the original data, but the GitHub changelogs indicate that it was incorporated into the UD framework and converted to UD v2. The changelogs also show that there was some iterative validation and revision, including corrections to tokenization, lemma completion, and dependency labeling.<sup>4</sup>

**Portuguese-PUD** corresponds to the Portuguese portion of the parallel dataset released for the CoNLL 2017 shared task on multilingual parsing (Zeman et al., 2017). It comprises 1,000 sentences, totaling 23,407 tokens, aligned one-to-one across 15 languages. The corpus domain consists of news and Wikipedia texts. According to the corpus GitHub page, the 1,000 English sentences were professionally translated into Portuguese, after which morphological and syntactic annotation was applied and later converted to UD v2. The entire treebank is labeled as a test set and was designed primarily for cross-linguistic evaluation rather than supervised training.<sup>5</sup>

**PetroGold** (Souza and Freitas, 2022) is a UD treebank of Brazilian Portuguese from the academic domain, specifically in the oil & gas sector. The corpus has 8,946 sentences, totaling 250,605 tokens. It comprises complete theses and dissertations. Annotation was generated automat-

<sup>3</sup><https://universaldependencies.org/>

<sup>4</sup>[https://github.com/UniversalDependencies/UD\\_Portuguese-GSD/tree/master?tab=readme-ov-file](https://github.com/UniversalDependencies/UD_Portuguese-GSD/tree/master?tab=readme-ov-file)

<sup>5</sup>[https://github.com/UniversalDependencies/UD\\_Portuguese-PUD/tree/master](https://github.com/UniversalDependencies/UD_Portuguese-PUD/tree/master)

Corpus	# Sent.	# Tokens	% Neg Ann.
Bosque	9,357	227,827	13.62
GSD	12,020	318,666	11.41
PUD	1,000	23,407	9.50
PetroGold	8,946	250,605	5.87
CINTIL	38,400	475,860	0.00
DANTEStocks	4,042	80,998	0.00
Porttinari-base	8,418	168,080	0.00

Table 1: Numbers (#) of sentences and tokens in the Portuguese treebanks of UD v2.17 (Nov. 2025) as well as percentage of sentences that contain at least one token annotated with `Polarity=Neg` or `PronType=Neg`.

	Bosque	GSD	PUD	PetroGold
não	□	□	□	□
nada	■	□		
nenhum	■	■		
nenhuma	■	■		
nem			□	
nunca			□	

Table 2: Lexical items across Portuguese UD treebanks with morphological annotation for negation. □ indicates presence of the `Polarity=Neg` feature; ■ presence of the `PronType=Neg` feature. (*nada* has both in Bosque as it can be a pronoun or an adverb.)

ically and subsequently subjected to systematic manual revision within the Petrolês Project (Applied Computational Intelligence Laboratory (ICA), PUC-Rio, 2025). PetroGold served as an empirical testbed for evaluating revision methodologies in dependency treebanks, including inter-annotator disagreement, inconsistent n-grams, and rule-based verification, which have been shown to yield measurable gains in intrinsic parsing metrics (Souza and Freitas, 2022).

**Porttinari-base** (Duran et al., 2023) (original vs. enriched) is our target dataset and is described in detail in Section 3.

A summary of corpora statistics is given in Table 1. The inventory of lexical items related to negation annotated in each treebank is shown in Table 2.

The treebanks also differ in how explicitly they document negation annotation. Bosque-UD provides a well-motivated account of guideline changes across UD releases, including the shift from the dependency relation `neg` to the morphological feature `Polarity=Neg`. It also states that the sentential negator *não* should be tagged strictly as an adverb (ADV) and offers guidance for negative concord, in which items such as *nada* receive negative polarity features independently of their syntactic function (e.g., adverbial vs. object).

DANTEStocks reports decisions tailored to user-generated content, including how frequent orthographic variants in social media — such as *nao* (without diacritics) — are integrated into clause structure, typically attaching to the predicate via `advmod`. By contrast, CINTIL-UDep and PetroGold provide no technical discussion of negation in their corpus description papers. This uneven documentation mirrors the quantitative picture: across Portuguese UD treebanks, explicit polarity marking is attested but typically limited to a narrow set of frequent negative lexical items, often centered on *não* and indefinite pronouns (*nada* ‘nothing’, *nenhum* ‘none’ and *ninguém* ‘no one’), as shown in Table 2.

Overall, the Portuguese UD treebanks exhibit substantial heterogeneity both in coverage and in the documentation of negation-related annotation, which makes the present work relevant in both aspects: the enrichment and its documentation.

### 3. Data

This study focuses on Porttinari-base (Duran et al., 2023), a gold-standard dependency treebank for Brazilian Portuguese annotated under UD. Porttinari-base contains 8,418 sentences (168,080 tokens) drawn from journalistic texts and constitutes the manually revised core of the Porttinari project. Sentences were selected to maximize syntactic diversity while avoiding instances that are overly short or excessively long, with a preference for 10–40 tokens (Duran et al., 2023).

Porttinari-base was initially annotated automatically using UDPipe 2 (Straka, 2018) and then manually revised for lemmas, UPOS, morphological features, and dependency relations. The enrichment proposed here only modifies the morphological layer (UFeats) in a deterministic, lexeme-driven way, leaving sentence boundaries, tokenization, and dependency structure unchanged.

Our enrichment decisions are informed by descriptive studies of Brazilian Portuguese negation (Mito, 1992, 1998; Neves, 2000) and by the taxonomy of items related to negation systematized in Miranda Junior (2022). In Miranda Junior (2022), lexical items related to negation are organized into a set of descriptive-functional categories that show the functions these items can perform, including sentential negation, non-sentential negation, double negation, negative concord, negative pronouns, negative conjunction, negative preposition, and negation reinforcement.

Miranda Junior (2022)’s taxonomy contains 19 lexical items related to negation. All items, except one (*salvo* ‘except/unless’), appear in the Porttinari-base data. These 18 items (24 lemma–UPOS combinations) are given in Tables 3 and 4. They cover canonical sentential negation, negative adverbs,

negative pronouns/determiners, conjunctions, and a restricted class of prepositional exclusion operators. Our proposal is therefore intentionally broad in scope: it targets *negation-related* lexical items rather than only prototypical sentential negation.

This broader scope is compatible with the UD treatment of polarity. UD guidelines state that `Polarity` is typically used for overt grammatical negation, including function words in languages that negate through separate particles, whereas pronominal forms are handled through `PronType=Neg`. In our annotation, this distinction is preserved: pronouns and determiners, such as *nada*, *ninguém*, and *nenhum*, are annotated with `PronType=Neg`, whereas non-pronominal items that overtly contribute negative meaning are annotated with `Polarity=Neg`. This includes not only canonical negators such as *não*, *nunca*, and *jamais*, but also a restricted class of exclusion operators when they occur with the lemma–UPOS combinations defined in our inventory.

This is particularly relevant for the class of *prepositions of exclusion* described by Neves (2000). Although these forms do not usually negate the main clause in the same way as the canonical negation form *não*, they delimit a domain by excluding one of its possible members, thereby contributing a negative interpretation in context. When these elements occur in a sentence as a preposition, the sentence can often be paraphrased by changing the preposition to *não*, without substantial change in the interpretation, as can be seen in Example 3. We therefore treat those cases as negation-related operators for our annotation purposes.

**Example 3.** Exclusion preposition with negation-related function:

Todos foram à festa, exceto João.  
everyone go.PST.3PL to.the party except João

‘Everyone went to the party, except João.’

Paraphrase:

Todos foram à festa, não o João.  
everyone go.PST.3PL to.the party, not the João

‘Everyone went to the party, but not João.’

We annotated all tokens whose lemma matched one of the items in Tables 3 and 4, adding morphological information depending on their UPOS: non-pronominal items received a `Polarity=Neg` feature; pronouns and determiners were updated by replacing `PronType=Ind` with `PronType=Neg`, while preserving indefiniteness (`Definite=Ind`), as shown in Example 4.

**Example 4.** Previously underspecified negative pronoun:

Nada foi feito.  
nothing be.PST.3SG do.PTCP

‘Nothing was done.’

Negative items		Exclusion operators	
Lemma	UPOS	Lemma	UPOS
nunca never	ADV	exclusive excluded	ADP
jamais never	ADV	excluso excluded	ADP
sequer even	ADV	exceto except	ADP
tampouco neither	ADV	fora out	ADP
tampouco nor	CCONJ	afora apart	ADP
nada nothing	ADV	menos minus	ADP
nem not even	ADV	senão except / but	ADP
nem nor	CCONJ	tirante excluding	ADP
nem nor	SCONJ	sem without	ADP
não not	ADV		
não no	INTJ		

Table 3: Lemma–UPOS combinations enriched with `Polarity=Neg`, grouped into traditional negation items and exclusion operators.

Lemma	UPOS
nada nothing	PRON
ninguém nobody / no one	PRON
nenhum no / none	DET
nenhum none	PRON

Table 4: Lemma–UPOS combinations receiving the `PronType=Neg` and `Definite=Ind` features.

Before:

nada nada PRON Gender=Masc  
Number=Sing  
PronType=Ind

After:

nada nada PRON Gender=Masc  
Number=Sing  
PronType=Neg  
Definite=Ind

**Example 5.** Previously unmarked negative adverb:

Ele nunca respondeu.  
he never answer.PST.3SG

‘He never answered.’

Before:

nunca nunca ADV -

After:

nunca nunca ADV Polarity=Neg

When a token had no negation-related morphological feature (as in Example 5), enrichment consisted solely of adding the appropriate UD feature conditioned on the lemma–UPOS pair. No structural changes were introduced, and dependency relations were left untouched; the intervention was strictly morphological.

Table 5 quantifies the presence of items related to negation in the corpus. The enrichment affected 1,890 tokens across 1,563 sentences, corresponding to 1.12% of all tokens and 18.57% of all sentences in Porttinari-base. Of these, 1,711 tokens receive `Polarity=Neg`, while 179 tokens are revised from `PronType=Ind` to `PronType=Neg` and are additionally specified for `Definite=Ind`.

Measure	Count
Total tokens in Porttinari-base	168,080
Enriched tokens	1,890
<code>Polarity=Neg</code> added	1,711
<code>PronType=Neg</code> revised	179
Total sentences in Porttinari-base	8,418
Sentences affected	1,563

Table 5: Extent of the negation-related enrichment in Porttinari-base.

#### 4. Parsing-based Evaluation Setup

We evaluate the impact of the enrichment in a UD parsing pipeline using PortParser (Lopes and Pardo, 2024), a neural dependency parser designed for Brazilian Portuguese. PortParser performs UPOS tagging, morphological feature prediction (UFeats), and dependency parsing, which makes it suitable for assessing both the parser’s ability to learn the enriched morphological representation and the extent to which this enrichment affects syntactic structure prediction.

Following the best parser configuration, according to Lopes and Pardo (2024), we generated five randomized train/dev/test sets of Porttinari-base using a 70/10/20 proportion (5,893/842/1,683 sentences). In our experiment, only the train and test sets were used, for training and evaluation respectively, but we followed the set distribution of Lopes and Pardo (2024).

For each split, we trained two models: (i) a baseline model trained on the original corpus and (ii) a model trained on the negation-enriched version of the same sentences. The only experimental variable is the presence or absence of the negation-related UFeats introduced in the enriched corpus. Training follows the two-stage process reported as best-performing in the original PortParser study: 80 epochs with the encoder frozen, followed by 20 fine-

Run	LAS			UAS		
	Base	Enr.	$\Delta$	Base	Enr.	$\Delta$
1	92.51	92.31	-0.20	94.33	94.13	-0.20
2	91.89	92.13	+0.24	93.78	94.04	+0.26
3	91.74	91.67	-0.07	93.59	93.53	-0.06
4	92.07	91.90	-0.17	94.11	93.85	-0.26
5	91.52	91.73	+0.21	93.47	93.72	+0.25
<b>Mean</b>	91.95	91.95	0.00	93.86	93.85	-0.01

Table 6: LAS and UAS accuracy across five randomized splits in percentage.

tuning epochs. Tokenization is controlled across conditions by using the UDPipe 2 UD tokenizer (Straka, 2018) for all splits.

We report standard parsing metrics for syntactic dependencies (UAS and LAS). Since the intervention modifies only a restricted portion of the morphological layer, the evaluation must distinguish between overall parser compatibility and its ability to learn the newly introduced features. We therefore also evaluate the prediction of negation-related features, as well as those already present in the original annotation.

## 5. Results

Importantly, the modifications introduced in the treebank are deterministic and at the level of the morphological features: no syntactic structures were altered, and no additional tokens or dependency relations were introduced. For this reason, substantial shifts in parsing performance are not expected. The question is not whether a morphology-only intervention would improve attachment accuracy, but whether our negation-related enrichment can be introduced while keeping the same quality of the parser’s predictions and whether the enriched features are learned reliably.

Table 6 gives the LAS and UAS accuracy for each run. The difference in the UAS mean between training on the base corpus and on the enriched one is negligible ( $-0.01\%$ ), indicating that enrichment does not affect unlabeled structural decisions. The mean LAS is identical across configurations. Differences fluctuate across runs but do not exhibit systematic improvement or degradation.

We further conducted an error analysis focusing on dependency relations, with at least 100 LAS errors in the test set for both training versions. Table 7 gives the mean variation in the error rates for these most frequent dependency relations. It also indicates how many of the 1,890 enriched tokens appear in these dependency relations, as well as how often dependencies point to an enriched token.

Dep.	Token	% aff.	Base		Enr.		$\Delta$
			%	(n)	%	(n)	
nsubj	62	0.65	7.87	(151)	7.82	(150)	-0.04
obj	41	0.56	6.37	(93)	6.38	(93)	+0.01
obl	10	0.11	16.35	(294)	15.76	(283)	-0.60
advmod	1,434	23.30	10.64	(129)	10.23	(124)	-0.42
amod	0	0.00	7.43	(99)	7.47	(99)	+0.04

Table 7: Mean LAS error by dependency relation on the test set across the five runs, considering only relations with at least 100 LAS errors in both training conditions. The table reports how many of the 1,890 enriched tokens occur with each relation (Token), as well as how often dependencies point to an enriched token (% aff.).

As shown, the percentage differences between conditions are low overall, with a maximum variation of 0.6%. For `advmod`, the dependency relation most affected by our enrichment (where 75% of our changes lie (1,434 tokens out of 1,890)), we see a slight upgrade of 0.42%.

Run	Base %	Enr. %	$\Delta$
1	98.29	98.27	-0.02
2	98.25	98.22	-0.03
3	98.23	98.20	-0.03
4	98.27	98.28	+0.01
5	98.28	98.25	-0.03
<b>Mean</b>	98.26	98.24	-0.02

Table 8: Accuracy of UFeats without negation-related features.

We also evaluate the accuracy of the morphological features. Table 8 gives the mean accuracy over all morphological features (UFeats), excluding the two we worked on; the negligible difference ( $-0.02$  percentage points) indicates that enrichment does not degrade pre-existing morphological categories.

We also perform a targeted evaluation of the enriched features: we measure mean precision, recall, and F1 for the negation-related features introduced by the enrichment across the five randomized runs.

Feature	Precision	Recall	F1
<code>Polarity=Neg</code>	99.87	99.16	99.51
<code>PronType=Neg</code> <code>+Definite=Ind</code>	97.75	99.39	98.84

Table 9: Mean Precision, Recall, and F1-score (%) on the test set across five runs for negation-related features introduced in the enriched annotation.

Table 9 directly evaluates the parser’s ability to recover the newly enriched negation-related features. This step is crucial because stable

aggregate parsing metrics alone do not demonstrate that the additional representation has been learned. By measuring prediction quality specifically for `Polarity=Neg` and for the enriched `PronType=Neg|Definite=Ind` features, we can assess whether the enrichment is not only compatible with the parser but also recoverable from contextual evidence at inference time.

The results show that the new features are learned with high reliability. For `Polarity=Neg`, the parser achieves 99.87% precision, 99.87% recall, and 99.51% F1. For negative pronouns, results are likewise strong, with 97.75% precision, 99.39% recall, and 98.84% F1.

The parser outputs point to a consistent pattern. Enrichment increases morphological granularity without destabilizing syntactic predictions: UAS remains effectively unchanged, mean LAS is identical across conditions (Table 6), and LAS error by high-frequency relations shows only small, non-systematic fluctuations (Table 7). In the shared morphological space, Base-UFeats accuracy is preserved (Table 8). The targeted results in Table 9 complement this picture by showing that the newly introduced negation-related features are also learned with high reliability, demonstrating near-ceiling performance for `Polarity=Neg` and strong recovery of the enriched negative-pronoun features. Together, these findings support the view that the enrichment makes negation-related morphology more explicit without disrupting syntactic parsing.

## 6. Conclusions and Future Work

In this work, we described the process of enriching the Brazilian Portuguese UD treebank Porttinari-base (Duran et al., 2023) with negation-related UD features, following linguistically motivated criteria grounded in descriptive studies of Brazilian Portuguese.

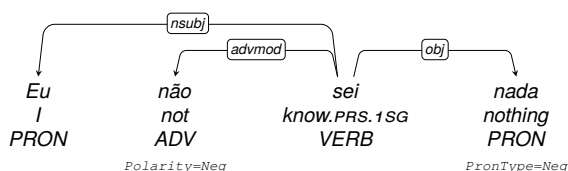
We also evaluated the impact of our intervention on parsing performance. We trained the PortParser on the original version of Porttinari-base and on our enriched version, to compare performance. Both versions of the PortParser perform equally well. Aggregate syntactic accuracy (UAS/LAS) remains stable, base morphological accuracy is preserved, and the newly introduced features (`Polarity=Neg` and `PronType=Neg`) are learned with high reliability.

The enrichment, as expected, does not yield measurable gains in intrinsic parsing metrics (UAS or LAS), but its relevance extends beyond parsing accuracy itself. Negation plays a central role in inference, polarity sensitivity, and semantic interpretation. An inadequate representation has been shown to negatively affect downstream tasks,

such as sentiment analysis and information extraction (Jiménez-Zafra et al., 2020; Findlay and Haug, 2025). From this perspective, making negation morphologically explicit constitutes a relevant contribution to the NLP field: it provides a reliable structural signal that can be exploited by higher-level models built on top of dependency parsing, including semantic parsing, polarity-aware modeling, and discourse-level analysis.

Our enrichment does not directly capture relations between negative elements, including those involved in negative concord or double negation. These phenomena are particularly relevant in Brazilian Portuguese, though. Negative concord is the main form in Portuguese, as in Example 6, where the standard negation *não* co-occurs with the negative pronoun *nada*, but the sentence expresses a single semantic negation ('I know nothing').

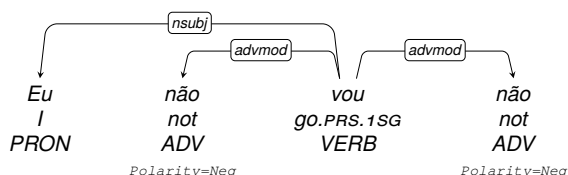
**Example 6.** *Eu não sei nada.*



'I know nothing.'

In Brazilian Portuguese, we can have occurrences of double negation, as in Example 7:

**Example 7.** *Eu não vou não.*



'I am not going.'

After our enrichment in Porttinari-base, we thus have cases similar to Examples 6 and 7, where two items have features indicating negative polarity (not necessarily the same feature).

As emphasized by Findlay and Haug (2025), marking negation on two elements may negatively affect sentiment analysis. An extraction tool may apply polarity-altering rules twice, resulting in positive polarity, as reported by Kanayama and Iwamoto (2020). Future work will aim to extend the representation of negation within the UD framework to clarify the relationship between elements that are related through negation, as in Examples 6 and 7.

## 7. Acknowledgements

Isaac Souza de Miranda Jr. is a PhD. student supported by grants #2023/01892-4 and #2025/20010-8, São Paulo Research Foundation (FAPESP).

Marie-Catherine de Marneffe is a research associate of the Fonds de la Recherche Scientifique – FNRS.

This work was carried out in part at the Center for Artificial Intelligence of the University of São Paulo (C4AI, <http://c4ai.inova.usp.br/>), with support from grant #2019/07665-4, São Paulo Research Foundation (FAPESP), and IBM Corporation. The project was also supported by the Ministry of Science, Technology and Innovation, with resources from Law No. 8.248, of October 23, 1991, within the scope of PPI-SOFTEX, coordinated by Softex and published as Residence in TIC 13, DOU 01245.010222/2022-44.

## 8. Bibliographical References

Susana Afonso, Eckhard Bick, Renato Haber, and Diana Santos. 2002. *Floresta Sintáctica: A treebank for Portuguese*. In *Proceedings of the Third International Conference on Language Resources and Evaluation (LREC 2002)*, Las Palmas, Canary Islands, Spain. European Language Resources Association (ELRA).

Applied Computational Intelligence Laboratory (ICA), PUC-Rio. 2025. *Petrolês Project*. Accessed: 2026-02-15.

Eckhard Bick. 2014. *Palavras: A Constraint Grammar-Based parsing system for Portuguese*. In Tony Berber Sardinha and Thelma de Lourdes São Bento Ferreira, editors, *Working with Portuguese Corpora*, pages 279–302. Bloomsbury Academic.

António Branco, João Ricardo Silva, Luís Gomes, and João António Rodrigues. 2022. *Universal Grammatical Dependencies for Portuguese with CINTIL data, LX Processing and CLARIN Support*. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference (LREC 2022)*, pages 5617–5626, Marseille, France. European Language Resources Association.

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*, 34:301–310.

Francisco Costa and António Branco. 2010. *LX-Gram: A deep linguistic processing grammar for Portuguese*. In *Computational Processing of the Portuguese Language*, pages 86–89, Berlin, Heidelberg. Springer Berlin Heidelberg.

Marie-Catherine de Marneffe, Christopher D. Manning, Joakim Nivre, and Daniel Zeman. 2021.

- Universal dependencies. *Computational Linguistics*, 47(2):255–308.
- Magali Duran, Lucelene Lopes, Maria das Graças Nunes, and Thiago Pardo. 2023. [The dawn of the Porttinari multigenre treebank: Introducing its journalistic portion](#). In *Anais do XIV Simpósio Brasileiro de Tecnologia da Informação e da Linguagem Humana*, pages 115–124, Belo Horizonte, MG, Brazil. SBC.
- Ariani Di Felippo, Maria das Graças Volpe Nunes, and Bryan Khelven Barbosa. 2024. [A dependency treebank of tweets in Brazilian Portuguese: Syntactic annotation issues and approach](#). In *Proceedings of the 15th Brazilian Symposium in Information and Human Language Technology*, pages 86–94, Belém do Pará, Brazil. Association for Computational Linguistics.
- Jamie Yates Findlay and Dag Trygve Truslew Haug. 2025. [Negation in Universal dependencies](#). In *Proceedings of the Eighth Workshop on Universal Dependencies (UDW, SyntaxFest 2025)*, pages 70–79, Ljubljana, Slovenia. Association for Computational Linguistics.
- Gottlob Frege. 1948. Sense and reference. *The Philosophical Review*, 57(3):209–230.
- Ilya Goldin and Wendy W. Chapman. 2003. Learning to detect negation with ‘not’ in medical texts. In *Proceedings of Workshop on Text Analysis and Search for Bioinformatics, ACM SIGIR*.
- Paul Grice. 1989. *Studies in the way of words*. Harvard University Press.
- Bernd Heine and Tania Kuteva. 2007. *The genesis of grammar: A reconstruction*. OUP Oxford.
- Laurence R. Horn. 2001. *A natural history of negation*, 2nd edition. CSLI Publications.
- Laurence R. Horn and Heinrich Wansing. 2025. [Negation](#). In Edward N. Zalta and Uri Nodelman, editors, *The Stanford Encyclopedia of Philosophy*, Spring 2025 edition. Metaphysics Research Lab, Stanford University.
- Salud M. Jimenez-Zafra, María T. Martín-Valdivia, Eugenio Martínez-Cámara, and Luis A. Ureña-Lopez. 2017. Studying the scope of negation for Spanish sentiment analysis on Twitter. *IEEE Transactions on Affective Computing*, 10(1):129–141.
- Salud M. Jiménez-Zafra, Rose Morante, María T. Martín-Valdivia, and Luis A. Ureña-Lopez. 2020. [Corpora annotated with negation: An overview](#). *Computational Linguistics*, 46(1):1–52.
- Hiroshi Kanayama and Ran Iwamoto. 2020. [How universal are Universal dependencies? Exploiting syntax for multilingual clause-level sentiment detection](#). In *Proceedings of the 12th Language Resources and Evaluation Conference (LREC 2020)*, Marseille, France. European Language Resources Association (ELRA). Licensed under CC-BY-NC.
- David Kletz. 2025. [Négation et modèles de langue pré-entraînés](#). Ph.D. thesis, Université Sorbonne Nouvelle - Paris 3. Thèse de doctorat dirigée par Amsili, Pascal et Candito, Marie-Hélène Sciences du langage Paris 3 2025.
- Lucelene Lopes and Thiago Pardo. 2024. [Towards Portparser-a highly accurate parsing system for Brazilian Portuguese following the Universal dependencies framework](#). In *Proceedings of the 16th International Conference on Computational Processing of Portuguese-Vol. 1*, pages 401–410.
- Ryan McDonald, Joakim Nivre, Yvonne Quirnbach-Brundage, Yoav Goldberg, Dipanjan Das, Kuzman Ganchev, Keith Hall, Slav Petrov, Hao Zhang, Oscar Täckström, Claudia Bedini, Núria Bertomeu Castelló, and Jungmee Lee. 2013. [Universal dependency annotation for multilingual parsing](#). In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 92–97, Sofia, Bulgaria. Association for Computational Linguistics.
- Carlos Miotto. 1992. [Negação sentencial no português brasileiro e teoria da gramática](#). Ph.D. thesis, Universidade Estadual de Campinas.
- Carlos Miotto. 1998. [Tipos de negação](#). *Cadernos de Estudos Linguísticos*, 34.
- Isaac Souza de Miranda Junior. 2022. [Não é nada não: uma análise das lexias negativas do português para anotação nas Universal Dependencies](#). Ph.D. thesis, Universidade Federal de São Carlos.
- Roser Morante and Walter Daelemans. 2009. [A metalearning approach to processing the scope of negation](#). In *Proceedings of the Thirteenth Conference on Computational Natural Language Learning (CoNLL-2009)*, pages 21–29, Boulder, Colorado. Association for Computational Linguistics.
- Pradeep G. Mutalik, Aniruddha Deshpande, and Prakash M. Nadkarni. 2001. [Use of general-purpose negation detection to augment concept indexing of medical documents: A quantitative study using the UMLS](#). *Journal of the American Medical Informatics Association*, 8(6):598–609.

- Maria Helena de Moura Neves. 2000. *Gramática de usos do português*. Editora UNESP, São Paulo.
- Joakim Nivre, Marie-Catherine de Marneffe, Filip Ginter, Jan Hajič, Christopher D. Manning, Sampo Pyysalo, Sebastian Schuster, Francis Tyers, and Daniel Zeman. 2020. [Universal dependencies v2: An evergrowing multilingual treebank collection](#). In *Proceedings of the Twelfth Language Resources and Evaluation Conference (LREC 2020)*, pages 4034–4043.
- Alexandre Rademaker, Fabricio Chalub, Livy Real, Claudia Freitas, Eckhard Bick, and Valeria Paiva. 2017. [Universal dependencies for Portuguese](#). In *Proceedings of the Fourth International Conference on Dependency Linguistics (Depling 2017)*, pages 197–206.
- Paulo Alexandre Rocha and Diana Santos. 2000. [CETEMPúblico: Um corpus de grandes dimensões de linguagem jornalística portuguesa](#). In *V Encontro para o Processamento Computacional da Língua Portuguesa Escrita e Falada (PROPOR 2000)*, pages 131–140, Atibaia, SP, Brazil. ICMC/USP.
- Bertrand Russell. 1905. On denoting. *Mind*, 14(56):479–493.
- Scott A. Schwenter. 2016. Some issues in negation in Portuguese. In *The Handbook of Portuguese Linguistics*, pages 425–440.
- Elvis Souza and Claudia Freitas. 2022. [Polishing the gold – how much revision do we need in treebanks?](#) In *Proceedings of the Universal Dependencies Brazilian Festival*, pages 1–11, Fortaleza, Brazil. Association for Computational Linguistics.
- M. Straka. 2018. [UDPipe 2.0 prototype at CoNLL 2018 UD shared task](#). In *Proceedings of the CoNLL 2018 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies*, pages 197–207.
- 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.
- Ludwig Wittgenstein. 1922. *Tractatus logico-philosophicus*. Kegan Paul, Trench, Trubner & Co., London. With an introduction by Bertrand Russell.
- Daniel Zeman, Martin Popel, Milan Straka, Jan Hajič, Joakim Nivre, Filip Ginter, Juhani Luotolahti, Sampo Pyysalo, Slav Petrov, Martin Potthast, Francis Tyers, Elena Badmaeva, Memduh Gokirmak, Anna Nedoluzhko, Silvie Cinková, Jan Hajič jr., Jaroslava Hlaváčová, Václava Kettnerová, Zdeňka Urešová, Jenna Kanerva, Stina Ojala, Anna Missilä, Christopher D. Manning, Sebastian Schuster, Siva Reddy, Dima Taji, Nizar Habash, Herman Leung, Marie-Catherine de Marneffe, Manuela Sanguinetti, Maria Simi, Hiroshi Kanayama, Valeria de Paiva, Kira Droganova, Héctor Martínez Alonso, Çağrı Çöltekin, Umut Sulubacak, Hans Uszkoreit, Vivien Macketanz, Aljoscha Burchardt, Kim Harris, Katrin Marheinecke, Georg Rehm, Tolga Kayadelen, Mohammed Attia, Ali Elkahky, Zhuoran Yu, Emily Pitler, Saran Lertpradit, Michael Mandl, Jesse Kirchner, Hector Fernandez Alcalde, Jana Strnadová, Esha Banerjee, Ruli Manurung, Antonio Stella, Atsuko Shimada, Sookyung Kwak, Gustavo Mendonça, Tatiana Lando, Rattima Nitisaroj, and Josie Li. 2017. [CoNLL 2017 shared task: Multilingual parsing from raw text to Universal dependencies](#). In *Proceedings of the CoNLL 2017 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies*, pages 1–19, Vancouver, Canada. Association for Computational Linguistics.
- Bowei Zou, Guodong Zhou, and Qiaoming Zhu. 2013. [Tree Kernel-based negation and speculation scope detection with structured syntactic parse features](#). In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 968–976. Association for Computational Linguistics.

## 9. Language Resource References

- Afonso, S. and Simões, A. and Frankenberg-Garcia, A. and Pinto, A. and Barreiro, A. and Maia, B. and Mota, C. and Oliveira, D. and Bick, E. and Ranchhod, E. 2000. [CETEN-Folha Corpus: Brazilian Portuguese newspaper texts](#). European Language Resources Association (ELRA). Catalog ELRA-U-W 0010.