

# It takes a village to grammaticalize

Joseph Larson and Patr cia Amaral

Department of Spanish and Portuguese

Indiana University

joelarso@iu.edu

pamaral@iu.edu

## Abstract

This paper investigates the grammaticalization of the noun *caleta* ‘cove, village’ to an intensifier, as part of the system of degree words in Chilean Spanish. We use word embeddings trained on a corpus of tweets to show the ongoing syntactic and semantic change of *caleta*, while also revealing how high degree is expressed in colloquial Chilean Spanish.

## 1 Introduction

Studies of language change using distributional methods have shown the potential of word embeddings (both static and contextualized) to trace syntactic and semantic change over time (Hamilton et al., 2016; Kutuzov et al., 2018; Periti et al., 2024, a.o.).<sup>1</sup> However, such research tends to focus on predicting changes that affect sets of lexical items shifting from one semantic domain to another, which typically reflects cultural and societal changes. Fewer studies have explored both semantic and morphosyntactic change (but see Fonteyn et al. 2022). In this paper, we focus on the semantic and syntactic shift from lexical to grammatical, known as grammaticalization (Meillet, 1912; Hopper and Traugott, 2003), and the stages of this process. Specifically, we study the creation of degree expressions like English *very*, *a lot*.

Traditionally, degree expressions have been associated with adjectives, considered the prototypical gradable category. However, degree modification is also compatible with nouns and verbs, which shows that gradability cuts across syntactic categories (Bolinger, 1972; Neeleman et al., 2004; Doetjes, 2008). As a word becomes a degree expression over time, it typically expands its distribution along different categories: e.g. it first com-

bines with nouns before co-occurring with verbs and adjectives. Hence, the grammaticalization of degree expressions provides insight into the semantics of degree and patterns in the distribution of degree words (Amaral, 2016; Luo et al., 2019). This paper examines an understudied variety, Chilean Spanish, and uses word embeddings to investigate the emerging system of degree words to which one grammaticalized word shifts. We investigate the grammaticalization of *caleta* in Chilean Spanish, from a noun denoting ‘cove, hiding place (where merchandise can be stored)’, ‘village’, as in ex. (1), to a quantifier and degree adverb ‘much, a lot’, as in (2), where *caleta* modifies the verb and denotes high degree.

- (1) Esta experiencia la realizamos  
this experience CL.FEM.SG.ACC do.PST.1PL  
en Zapallar, en la caleta de pescadores  
in Zapallar in the caleta of fishermen  
‘We did this experience in Zapallar, in the fishermen’s cove’
- (2) me gust  caleta  
CL.1SG.DAT like.PST.3SG caleta  
‘I liked it a lot.’

We use word embeddings to examine to what extent the grammaticalization of *caleta* has developed while also shedding light on the system of degree modifiers in Chilean Spanish. We ask, (i) how far along has *caleta* grammaticalized in Chilean Spanish, and (ii) what types of evidence do word embeddings provide of different stages of grammaticalization of degree words?

## 2 Previous Work

Linguists have provided analyses of the gradual process by which lexical items acquire grammatical functions: for example, in this diachronic change, nouns lose their categorial properties like occurring after a determiner or being pluralized.

<sup>1</sup>For a recent state-of-the-art survey comparing different approaches to semantic change using large language models, see Periti and Montanelli 2024.

The grammaticalization of nouns into degree adverbs (e.g. the development from *lot* ‘a set of objects’ to *a lot* ‘much’) is well attested cross-linguistically: other examples are French adverb *beaucoup* from *un beau coup* ‘a good strike’ and English *a bit* from ‘a bite, a portion that fits in the mouth’ (Abeillé et al., 2004; Marchello-Nizia, 2006; Verveckken, 2012; Traugott, 2008; Amaral, 2020).

This research has shown that a typical structure in which nouns occur - modification by a prepositional phrase, as in *a lot* [<sub>PP</sub> *of chairs*], *a mountain* [<sub>PP</sub> *of books*] - provides a starting point for quantity and degree interpretations. This structure undergoes subsequent syntactic reanalysis, where the head noun (e.g. *lot*) loses nominal properties and *a lot of* becomes an adverb modifying the second noun. The development of so-called binominal structures Det  $N_1$  of  $N_2$ , which may or may not further evolve to a fully adverbial category, plays a crucial role in the grammaticalization of degree words. In our study, we also include the structure (*Det*) *caleta* of  $N$ , hence we investigate the distribution of *caleta de*.

As argued by Doetjes, 2008, degree words across languages show a systematic behavior in terms of the words they can modify. These well-attested patterns correspond to types along a continuum of syntactic-semantic word classes, where a degree expression can modify all word classes (like French *trop*, type C) or just a subset of classes, gradable adjectives (like English *very*, type A), see Figure 1. As words develop into one type, they are predicted to modify words in the order along the continuum; for instance, if a word co-occurs with words of category V, it is expected to co-occur with words of category IV before it appears with words of category III.<sup>2</sup> As we investigate whether *caleta* has grammaticalized into a degree word, we will examine its stage of development with respect to Doetjes’ continuum.

While some computational studies of grammaticalization have adopted case-driven approaches similar to ours (Fonteyn and Manjavacas, 2021; Amaral et al., 2023; Nagata et al., 2024), we also

<sup>2</sup>(Doetjes, 2008) differentiates between ‘gradable’ and ‘eventive’ adjectives and verbs by whether or not the modifier is targeting the degree or is quantifying over events. The example she gives is from Dutch: *Jan is veel ziek* ‘Jan is sick a lot’ vs. *Jan is erg ziek* ‘Jan is very sick.’ In the former, *veel* as a quantifier targets eventive adjectives, thus it can only modify the quantity of sick events. In the latter, *erg* expresses the degree of sickness, i.e. the severity of his illness.

Category	Word Class					
I	gradable adjectives	Type A <i>very</i> <sup>E</sup>	Type B <i>erg</i> <sup>D</sup>	Type C <i>trop</i> <sup>F</sup>		
IIa	gradable nominal predicates		<i>očen</i> <sup>R</sup>	<i>muito</i> <sup>P</sup>		
IIb	gradable verbs	Type D <i>beaucoup</i> <sup>F</sup>		<i>molto</i> <sup>I</sup>		
III	eventive verbs		Type E <i>veel</i> <sup>D</sup>			
	eventive adjectives	<i>a lot</i> <sup>E</sup>				
	comparatives		<i>mnogo</i> <sup>R</sup>		Type F <i>a mountain</i> <sup>E</sup>	Type G <i>many</i> <sup>E</sup>
IV	mass nouns					
V	plural nouns					

Figure 1: Typology of degree expressions according to their distribution along a continuum of word classes. Table adapted with modifications from (Doetjes, 2008, 138). Superscripts indicate language: R for Russian, D for Dutch, F for French, E for English, P for Portuguese, and I for Italian.

investigate how a distributional analysis of *caleta* can provide insight on the set of degree expressions currently used in colloquial Chilean Spanish. In other words, we aim to examine not just the grammaticalization of *caleta* but also how this word fits in the system of degree words in Chilean Spanish and types of degree expressions across languages.

### 3 Methodology

#### 3.1 Corpus Creation

To ensure we had a good representation of colloquial Chilean Spanish, we created a subcorpus from an already existing corpus of online data (Ortiz-Fuentes, 2023). The already existing corpus contained roughly 19GB of data, from diverse sources, including news, tweets, online reviews and other miscellaneous web content. We chose to create a subcorpus just from tweets to reduce the computational load for our later experiments and since we only wanted informal instances of language; *caleta* typically only occurs in less formal registers. The resulting subcorpus of 27,306,582 tweets consisted of exactly 342,979,307 tokens. The time span of these tweets is from 2010 to 2020.

### 3.2 Preprocessing

We first normalized the text in the corpus: we removed case, punctuation, diacritics, URLs, hash-tags, and any repeated letters. For this last step, we only allowed double letters where they occur within normative Spanish orthography (i.e.  $\langle r \rangle$ ,  $\langle c \rangle$ ,  $\langle l \rangle$ ), elsewhere only single letters were allowed. Then we input the corpus into a plain text file separated by newlines. The resulting file was then lemmatized using SpaCy’s Spanish lemmatizer (Honnibal et al., 2020).

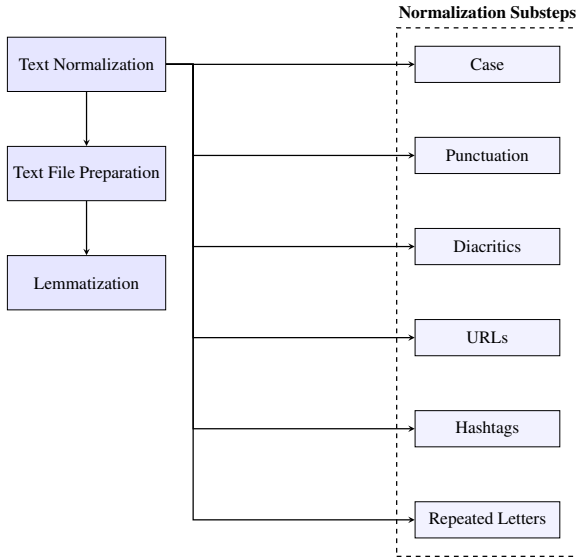


Figure 2: Preprocessing steps.

### 3.3 Model Selection

To represent the distributional patterns of words in our corpus, we decided to use static word embeddings over contextualized word embeddings. Non-contextualized embeddings allow us to compare our target word with other words in Chilean Spanish to examine the current stage of grammaticalization of *caleta* as determined by its closeness to different subsystems in the language.

The algorithm we use is Skip-Gram with Negative Sampling (SGNS) implemented in word2vec (Mikolov et al., 2013) to extract embeddings, based on previous research that showed good results for studies of semantic change (Hu et al., 2022, a.o.). For this reason, we do not consider it necessary to use a more computationally expensive operation (e.g. dynamic word embeddings). We trained each model for five epochs, a minimum token count of 10 and the skip-gram algorithm. Initially, we experimented with several hyperparameters: the window

size, the minimal word count and the vector size. The only hyperparameter that proved to be significant was the window size (see next section for more details). The resulting model used a vector length of 100 and a minimal word count of 10. To verify the validity of the model, we used analogy tests targetting gender-based morphological and semantic relations (see Table 1 for specifics). We found the analogy used always 100% accurate.

Relationship	Word Pair 1		Word Pair 2	
	Word A	Word B	Word A	Word B
Age-based	<i>Hombre</i>	<i>Mujer</i>	<i>Niño</i>	<i>Niña</i>
	'Man'	'Woman'	'Boy'	'Girl'
Familial	<i>Padre</i>	<i>Madre</i>	<i>Hijo</i>	<i>Hija</i>
	'Father'	'Mother'	'Son'	'Daughter'
Feline	<i>Niño</i>	<i>Gato</i>	<i>Niña</i>	<i>Gata</i>
	'Boy'	'Cat (male)'	'Girl'	'Cat (female)'
Canine	<i>Niño</i>	<i>Perro</i>	<i>Niña</i>	<i>Perra</i>
	'Boy'	'Dog (male)'	'Girl'	'Dog (female)'

Table 1: The four analogy tests used to validate Word2Vec model. The equation used was  $WB_2 = WA_1 - WA_2 + WB_1$ .

### 3.4 Window Size

As mentioned in the previous section, the only hyperparameter we adjusted for the model was the window size. We extracted models for  $w = [1, 10]$ . Our hypothesis was that lower window sizes would be more adequate for showing grammaticalization, since the scope of grammatical words like quantifiers lies within its immediate neighbors, whereas higher window sizes show neighbors within the same semantic field (therefore its lexical use). However, since we use a corpus of tweets, window size is fairly limited by the genre itself (a possible limitation we address later).

## 4 Results

### 4.1 Caleta

Here we display only the results of the experiments with a small ( $w = 1$ ) and a large ( $w = 10$ ) window size. This allows us to compare the information obtained by manipulating this parameter. In Figure 3, the word embeddings show both neighbors of the lexical noun and neighbors of the degree word. Nearest neighbors of the noun are toponyms (i.e. names of villages) and other nouns with related meanings (e.g. *playa* ‘beach’ and *balneario* ‘bathhouse’). As for the neighbors of the degree word, we find degree expressions, both adverbs and quantifiers like *mucho* and *ene*, both meaning

The co-occurrence of neighbors of both meanings shows that *caleta* has partially grammaticalized; it still retains its lexical use as a noun. These findings provide evidence for a situation of layering (Hopper, 1991), i.e. the synchronic co-existence of older and more recent functions of a form in a language.

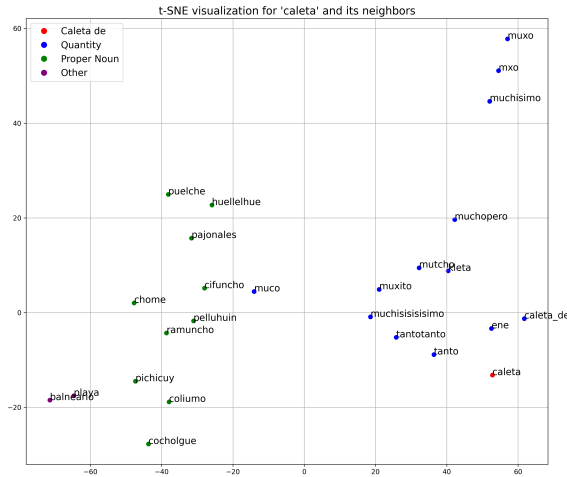


Figure 3: TSNE representation of *caleta* and its top 25 neighbors. Embeddings were created with a window size of 1. Blue corresponds to words that are quantifiers, green corresponds to toponyms (i.e. names of villages), and purple corresponds to semantically related nouns.

If we now use a larger window size, the results are different, with more neighbors associated with the lexical item. In Figure 4 we find the plural noun (*caletas*); as mentioned in historical analyses, the ability to be pluralized is a syntactic property of nouns. This attests to the persistence of some nominal categorial properties of *caleta*. We also find the noun *pescadores* ‘fishermen’, as the noun *caleta* typically refers to a village of fishermen and hence the nouns often co-occur (in *caleta de pescadores*), and related nouns like *muelle* ‘pier’ and *poza* ‘puddle’.

## 4.2 Caleta de

We analyzed the results of *caleta de* separately from those of *caleta* since the former is the vestige of a binominal quantifier preceding the grammaticalization of the latter. Figure 5 and Figure 6 show the TSNE representations of the nearest neighbors of *caleta de*. For the smaller window size, we see other quantifiers like *ene* (more in the next section), *caleta*, etc. The majority of neighbors here are

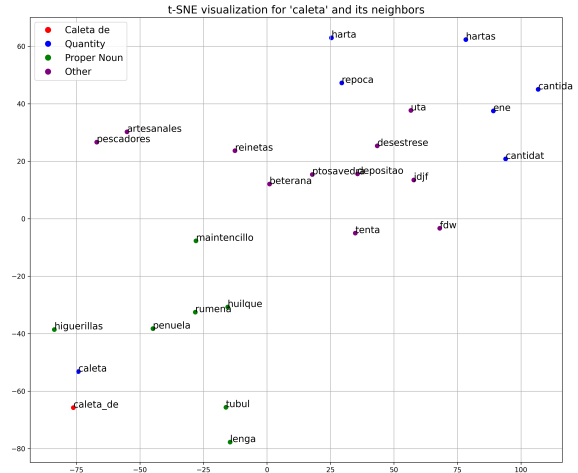


Figure 4: TSNE representation of *caleta* and its top 25 neighbors. Embeddings were created with a window size of 10. Blue corresponds to words that are quantifiers, green corresponds to toponyms (i.e. examples of *caletas*), and purple corresponds to semantically related nouns.

quantifiers in their orthographical variants found in tweets (e.g. *mucho*, *moxo*, *nucho*, etc). Two other words that form part of binominal quantifiers are also present, *monton* and *montones*, both meaning ‘pile’ and ‘piles’, but which have grammaticalized in the same fashion as *caleta* to denote a large quantity (*un montón de N* ‘a lot of N’). In this window size, only one proper noun is present, *Chorromil*, the name of a village. Lastly, we find other quantifiers, like *cualquiers* and *cualesquiers*, both orthographical variations of *cualquier*, ‘whichever’, and *puras*, a determiner in Chilean Spanish.

In the larger window size, we see *caleta* as its nearest neighbor. Other quantifiers like *mucho*, *ene*, *harto*, etc. are present, but they are much further away than semantically related nouns like *pescadores* ‘fishermen’, *artesanales* ‘craftsmen’, *reinetas*, a plural noun denoting a variety of white fish, as well as toponyms that are names of *caletas*. These results show once more how important the hyperparameter of window size is in capturing the grammatical meaning of relatively newly grammaticalized words in a language.

### 4.3 Ene

We decided to display the top 10 neighbors for the word *ene*, since *ene* always appeared as a top neighbor for *caleta* and *caleta de*. *Ene* comes from the Spanish pronunciation of the grapheme  $\langle n \rangle$  and is used in Mathematics to denote an unspeci-





Rank	Word	Score
1	<i>kleta</i> (orthographical variation of <i>caleta</i> )	0.71
2	<i>caleta de</i> 'a lot of'	0.68
3	<i>cantitat</i> ( <i>cantidad</i> , orthographical variation, 'quantity')	0.67
4	<i>graziash</i> ( <i>gracias</i> , orthographical variation, 'thanks')	0.66
5	<i>jsjsjd</i> 'laughter'	0.66
6	<i>harto</i> 'a lot'	0.66
7	<i>puxis</i> (orthographic variation of <i>pucha</i> , 'darn')	0.66
8	<i>autodelicioso</i> (lit. 'self-delicious', term used for masturbation)	0.64
10	<i>muchosaño</i> ( <i>muchos años</i> as one word, 'many years')	0.63

Table 3: Ranked words with their scores (cosine) for *ene* for  $w = 10$ . Bold words correspond to quantifiers.

#### 4.4 Other quantifiers

Lastly, we show word embeddings of other degree words, in this case 'stable' quantifiers in Chilean Spanish: *harto* 'a lot', *mucho* 'a lot', *tanto* 'so many.' It is worth mentioning that unlike *caleta*, *caleta de* and *ene* (which syntactically can be considered degree adverbs), these quantifiers inflect for gender and number when modifying a noun. The purpose of using the lemmatizer was to control for this, but as the results show, some inflected tokens of these quantifiers were not properly lemmatized.

Tables 4, 5, 6, 7, 8 and 9 show the nearest neighbors for *harto*, *mucho* and *tanto* at the two window sizes. For *harto*, we see that the majority of its neighbors are other quantifiers for both window sizes, as well as orthographical variations (e.g. *harrito*, *arto*) and inflected versions of the lexeme, like the feminine form *harta*. Likewise, *tanto* as its neighbors for the smaller window size shows mostly orthographical variations (e.g. *tsnto*, *tabto*), while for the larger window size we can see similar results to *ene*, where nouns like 'laughter' are amongst the neighbors. For *mucho*, we can see mostly orthographical variants for the smaller window size (e.g. *muxo*, *muxho*) and for the larger window size we see less orthographical variations and more of other quantifiers, even its antonym *poco*, which also occurs with intensifying affixes: *re-poco* and *poc-azo* 'very little'.

## 5 Discussion

Our word embedding results for *caleta* show that nowadays the word is used to express high degree. In addition, in our results both the lexical noun and the degree modifier are present. The choice of hyperparameters, specifically window size, has important consequences: a small window size yields nearest neighbors for both forms, while a larger window size results in more neighbors of

Rank	Word (Gloss)	Score
1	<i>arto</i> 'harto' (orthographical variation)	0.94
2	<i>mucho</i> 'a lot'	0.84
3	<i>bastante</i> 'quite'	0.78
4	<i>harrito</i> 'harto' (orthographical variation)	0.74
5	<i>moxo</i> 'mucho' (orthographical variation)	0.72
6	<i>muchísimo</i> 'mucho' (superlative)	0.71
7	<i>muxo</i> 'mucho' (orthographical variation)	0.69
8	<i>mutcho</i> 'mucho' (orthographical variation)	0.68
9	<i>mucho</i> 'mucho' (orthographical variation)	0.67
10	<i>nacho</i> 'mucho' (orthographical variation)	0.66

Table 4: Ranked words with their scores (cosine) for *harto* for  $w = 1$ . Bold words correspond to quantifiers.

Rank	Word (Gloss)	Score
1	<i>arto</i> 'harto' (orthographical variation)	0.81
2	<i>mucho</i> 'a lot'	0.72
3	<i>sosi</i> ( <i>eso sí</i> , abbreviation, 'though')	0.69
4	<i>bastante</i> 'quite'	0.68
5	<i>harta</i> 'a lot'	0.68
6	<i>ene</i> 'a lot'	0.66
7	<i>pucha</i> 'darn'	0.63
8	<i>haarto</i> 'harto' (orthographical variation)	0.63
9	<i>repoco</i> 'poco' (intensifier)	0.63
10	<i>pocazo</i> 'poco' (augmentative)	0.61

Table 5: Ranked words with their scores (cosine) for *harto* for  $w = 10$ . Bold words correspond to quantifiers.

Rank	Word (Gloss)	Score
1	<i>tsnto</i> 'tanto' (orthographical variation)	0.76
2	<i>demasia</i> ( <i>demasiado</i> , phonetic variation, 'too much')	0.70
3	<i>tantotanto</i> 'tanto' (repeated)	0.69
4	<i>mucho</i> 'a lot'	0.69
5	<i>tantoy</i> ( <i>tanto y</i> as one word, 'so much and')	0.69
6	<i>tabto</i> 'tanto' (orthographical variation)	0.68
7	<i>tantísimo</i> 'tanto' (superlative)	0.67
8	<i>tnto</i> 'tanto' (orthographical variation)	0.64
9	<i>tanro</i> 'tanto' (orthographical variation)	0.64
10	<i>mutcho</i> 'mucho' (orthographical variation)	0.64

Table 6: Ranked words with their scores (cosine) for *tanto* for  $w = 1$ . Bold words correspond to quantifiers.

Rank	Word (Gloss)	Score
1	<i>mucho</i> 'a lot'	0.71
2	<i>tsnto</i> 'tanto' (orthographical variation)	0.65
3	<i>tantotanto</i> 'tanto' (repeated)	0.63
4	<i>tantísimo</i> 'tanto' (superlative)	0.60
5	<i>simuchas</i> ( <i>sí muchas</i> as one word, 'yes a lot')	0.60
6	<i>jskdka</i> 'laughter'	0.60
7	<i>jajajajajaun</i> 'laughter'	0.60
8	<i>muchogracias</i> ( <i>muchas gracias</i> as one word, 'thanks a lot')	0.59
9	<i>tisin</i> ( <i>tí sin</i> as one word, 'you (prepositional), without')	0.58
10	<i>pueso</i> (portmanteau of <i>pues eso</i> , 'exactly')	0.58

Table 7: Ranked words with their scores (cosine) for *tanto* for  $w = 10$ . Bold words correspond to quantifiers.

the lexical noun. We hypothesize that this is due to the fact that as a degree word, *caleta* is a modi-

Rank	Word (Gloss)	Score
1	<b><i>muchísimo</i></b> ‘mucho’ (superlative)	0.91
2	<b><i>mxo</i></b> ‘mucho’ (orthographical variation)	0.88
3	<b><i>harto</i></b> ‘a lot’	0.82
4	<b><i>muxo</i></b> ‘mucho’ (orthographical variation)	0.81
5	<b><i>mucjo</i></b> ‘mucho’ (orthographical variation)	0.80
6	<b><i>muchi</i></b> ‘mucho’ (diminutive)	0.77
7	<b><i>muho</i></b> ‘mucho’ (orthographical variation)	0.77
8	<b><i>muxho</i></b> ‘mucho’ (orthographical variation)	0.77
9	<b><i>arto</i></b> ‘harto’ (orthographical variation)	0.76
10	<b><i>nucho</i></b> ‘mucho’ (orthographical variation)	0.75

Table 8: Ranked words with their scores (cosine) for *mucho* for  $w = 1$ . Bold words correspond to quantifiers.

Rank	Word (Gloss)	Score
1	<b><i>muchísimo</i></b> ‘mucho’ (superlative)	0.79
2	<b><i>harto</i></b> ‘a lot’	0.74
3	<b><i>tanto</i></b> ‘so much’	0.71
4	<b><i>poco</i></b> ‘a little’	0.67
5	<b><i>muchoy</i></b> ( <i>mucho</i> y as one word, ‘a lot and’)	0.65
6	<b><i>muccho</i></b> ‘mucho’ (orthographical variation)	0.65
7	<b><i>bastante</i></b> ‘quite’	0.65
8	<b><i>muchopero</i></b> ( <i>mucho pero</i> as one word, ‘a lot but’)	0.64
9	<b><i>aunpero</i></b> ( <i>aún pero</i> as one word, ‘still but’)	0.63
10	<b><i>muchisisismo</i></b> ‘mucho’ (repeated superlative)	0.61

Table 9: Ranked words with their scores (cosine) for *mucho* for  $w = 10$ . Bold words correspond to quantifiers.

fier, and occurs in close adjacency to the modified word. Hence, a small window captures this distribution. On the other hand, as a lexical noun *caleta* is less syntactically constrained, with more positional freedom and semantic content.

While cosine similarity scores give us insight into a changing word’s distribution, they alone do not tell us about its syntactic properties in detail. To better understand *caleta*’s current status as a degree modifier, we performed a *post-hoc* analysis of the top 20 collocates of *caleta* and *caleta de*. We looked specifically at the top tokens that immediately precede and proceed the two strings in our unlemmatized corpus. We were interested in the kinds of words that *caleta* and *caleta de* have come to modify, in accordance to Doetjes’s typology of degree modifiers (see Section 2).

Our analysis shows that *caleta* has evolved extensively beyond its original lexical usage, wherein it was only compatible with count nouns that were semantically related e.g. *pescadores* ‘fishermen’ *camarones* ‘shrimp (plural)’, headed by the preposition *de*. The structure *caleta de* is now compatible with count nouns beyond the semantic domain of a fishing village: *años* ‘years’, *veces* ‘times/instances’ (see (6)), as well as mass nouns e.g. *plata* ‘money (informal)’, *tiempo* ‘time’ (see

(7)). It can also modify comparatives e.g. *mejor* ‘better’, *peor* ‘worse’ (see (9)); eventive verbs e.g. *dormir* ‘to sleep’, *reír* ‘to laugh’ (see (8)); gradable verbs *gustar* ‘to like’, *querer* ‘to want’ (see (2)); and finally gradable nominal predicates<sup>3</sup> e.g. *hambre* ‘hunger’, *pena*, ‘sorrow’, as in (10).

- (6) Hacer caleta de años  
make.PRS.3SG caleta of years  
‘Many years ago’
- (7) es caleta de plata  
be.PRS.3SG caleta of money  
‘it’s a lot of money.’
- (8) Yo igual reí caleta.  
1SG.NOM same laugh.PST.1SG caleta  
‘I laughed a lot, anyway.’
- (9) hay que cuidarse  
be.existential.PRS.3SG that care.INF.REF  
caleta mejor...  
caleta better  
‘one has to take care of themselves much better.’
- (10) Hacer caleta de frío.  
make.PRS.3SG caleta of coldness  
‘It’s really cold.’

There were no cases of *caleta* modifying either eventive adjectives or gradable adjectives within our corpus. This, according to Doetjes’s classification, indicates that *caleta* has evolved into a type D degree modifier. Figure 7 shows *caleta*’s position in this typology, in comparison to the other degree expressions in Chilean Spanish that we have discussed in this paper. Our results align with claims in the literature that Type C and D are the most common in the Romance languages (Doetjes, 2008). Lastly, within our results, *caleta* has no nearest neighbors with Type A modifiers (e.g. *muy* ‘very’), which combine exclusively with gradable adjectives. This is not surprising since Type A modifiers have no overlap in word classes with Type D modifiers; their distributions are disjoint. This highlights how embeddings capture syntactic properties of words, as opposed to just similarity of meaning.

Our study has two main findings, which answer the research questions above. First, we have shown that *caleta* is undergoing grammaticalization: both

<sup>3</sup>Gradable nominal predicates, in Doetjes’s definition, are nouns that are generally the objects of light verb expressions. The examples she gives are from French e.g. *Elle a très soif* ‘She is very thirsty.’ In Spanish, such light verb constructions also exist, so we consider cases like *tener sed* ‘to be thirsty (lit. to have thirst)’ to also be examples of nominal predicates.

Category	Word Class					
I	gradable adjectives	Type A				
IIa	gradable nominal predicates	Type D	Type B	Type C		
IIb	gradable verbs	caleta		harto		
	eventive verbs	ene		bastante		
III	eventive adjectives	mucho		demasiado		
	comparatives	tanto	Type E		Type F un montón cantidad mon- tones	
IV	mass nouns					Type G vario
V	plural nouns					

Figure 7: Degree words found in our results and their corresponding types according to Doetjes’ model; modified table from (Doetjes, 2008, 138)

the older and the new meaning are captured by the word embeddings. Importantly, we see a difference in the results depending on the window size, when compared to other degree words which are grammatical items and not undergoing change, like *mucho* and *harto*. In the latter case, window size does not significantly impact the neighbors. Additionally, our *post-hoc* analysis provided insight on the properties of *caleta* as a degree word.

Second, our word embeddings have allowed us to reveal the inventory of degree words in colloquial Chilean Spanish, including a word that to date had never been investigated, *ene*. These words denote high degree (intensifiers), words that are known to change rapidly due to social and expressive pressure (Ito and Tagliamonte, 2003). Since *caleta* and *ene* are not normative forms, they are left out of traditional studies. This entails that we may miss instances of change possibly of interest to current linguistic theory. Hence, word embeddings can be a tool to study lesser-known subsystems of a language and capture ongoing changes in synchrony.

## 6 Conclusion

Our study contributes to studies of language change by analyzing intensifiers in colloquial Chilean Spanish (an understudied variety) from the past twenty years. We reveal an ongoing change that had not been previously studied. Using spontaneous speech from tweets, we gained access to informal speech where speakers communicate in

an unedited way, which has allowed us to study the use of older and more recent degree expressions. Hence, our study shows how Digital Humanities as an interdisciplinary field can expand our knowledge of low-resource language varieties. In our specific case, the examination of the data through language processing revealed instances of grammaticalization that to the best of our knowledge had not been analyzed before.

We have shown that static word embeddings provide evidence for this change and can reveal meaning relations not previously studied. Moreover, we show that different choices of hyperparameters have an effect on which meaning (the lexical vs. the grammatical) of the word undergoing change, *caleta*, is represented.

Some limitations of our study are due to the genre itself. One such limitation is the difficulty with lemmatization: as we have mentioned, these are tweets, so we find strings that do not conform to normative orthography (for example, typos, abbreviations etc), therefore the lemmatizer has difficulty with detecting words of the same lexeme. In addition, Tweeter users tend to adopt orthographical forms that reflect pronunciation and sometimes are intended to be expressive, like repeating vowels in a word to express a very high degree. Furthermore, using a corpus of tweets means that the character limit has an impact on the possible window sizes. To obviate this problem, further studies on *caleta* could use longer texts that have the same register as tweets, e.g. blog posts.

## References

- Anne Abeillé, Olivier Bonami, Danièle Godard, and Jesse Tseng. 2004. *The Syntax of French de-N’ Phrases*. *Proceedings of the International Conference on Head-Driven Phrase Structure Grammar*, pages 6–26.
- Patrícia Amaral. 2016. *When Something Becomes a Bit*. *Diachronica*, 33:151–186.
- Patrícia Amaral. 2020. Bocado: Scalar Semantics and Polarity Sensitivity. *Zeitschrift für romanische Philologie*, 136(4):1114–1136.
- Patrícia Amaral, Hai Hu, and Sandra Kübler. 2023. Tracing semantic change with distributional methods: The contexts of algo. *Diachronica*, 40(2):153–194.
- Dwight Bolinger. 1972. *Degree Words*. De Gruyter Mouton, Berlin, Boston.
- Jenny Doetjes. 2008. *Adjectives and Degree Modification*. In *Adjectives and Adverbs: Syntax, Semantics,*



- and Discourse, pages 123–155. Oxford University Press.
- Lauren Fonteyn and Enrique Manjavacas. 2021. [Adjusting scope: a computational approach to case-driven research on semantic change](#). In *Proceedings of the Workshop on Computational Humanities Research (CHR 2021)*, volume 2898 of *CEUR Workshop Proceedings*, pages 280–298.
- Lauren Fonteyn, Enrique Manjavacas, and Sara Budts. 2022. Exploring morphosyntactic variation and change with distributional semantic models. *Journal of Historical Syntax*, 6:1–41.
- William L. Hamilton, Jure Leskovec, and Dan Jurafsky. 2016. [Cultural shift or linguistic drift? comparing two computational measures of semantic change](#). In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2116–2121, Austin, Texas. Association for Computational Linguistics.
- Matthew Honnibal, Ines Montani, Sofie Van Landeghem, and Adrienne Boyd. 2020. [spacy: Industrial-strength natural language processing in python](#). *The Journal of Open Source Software*, 5(53):2914.
- Paul Hopper. 1991. On some principles of grammaticalization. In *Approaches to Grammaticalization*, pages 17–35. Benjamins.
- Paul J. Hopper and Elizabeth Closs Traugott. 2003. *Grammaticalization*, 2 edition. Cambridge Textbooks in Linguistics. Cambridge University Press.
- Hai Hu, Patrícia Amaral, and Sandra Kübler. 2022. Word embeddings and semantic shifts in historical spanish: Methodological considerations. *Digital Scholarship in the Humanities*, 37(2):441–461.
- Rika Ito and Sali Tagliamonte. 2003. [Well weird, right dodgy, very strange, really cool: Layering and recycling in english intensifiers](#). *Language in Society*, 32(2):257–279.
- Andrey Kutuzov, Lilja Øvrelid, Terrence Szymanski, and Erik Velldal. 2018. [Diachronic word embeddings and semantic shifts: a survey](#). In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 1384–1397, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- Yiwei Luo, Dan Jurafsky, and Beth Levin. 2019. [From insanely jealous to insanely delicious: Computational models for the semantic bleaching of English intensifiers](#). In *Proceedings of the 1st International Workshop on Computational Approaches to Historical Language Change*, pages 1–13, Florence, Italy. Association for Computational Linguistics.
- Christiane Marchello-Nizia. 2006. *Grammaticalisation et changement linguistique*. De Boeck.
- Antoine Meillet. 1912. L’ évolution des formes grammaticales. *Scientia*, 12:130–148.
- Tomas Mikolov, Kai Chen, G.s Corrado, and Jeffrey Dean. 2013. Efficient Estimation of Word Representations in Vector Space. *Proceedings of Workshop at ICLR*, 2013:1–12.
- Ryo Nagata, Yoshifumi Kawasaki, Naoki Otani, and Hiroya Takamura. 2024. [A Computational Approach to Quantifying Grammaticalization of English Deverbal Prepositions](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 211–220, Torino, Italia. ELRA and ICCL.
- Ad Neeleman, Hans Van de Koot, and Jenny Doetjes. 2004. [Degree expressions](#). *The Linguistic Review*, 21(1):1–66.
- Jorge Ortiz-Fuentes. 2023. [Chilean Spanish Corpus](#).
- Francesco Periti, Pierluigi Cassotti, Haim Dubossarsky, and Nina Tahmasebi. 2024. [Analyzing semantic change through lexical replacements](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4495–4510, Bangkok, Thailand. Association for Computational Linguistics.
- Francesco Periti and Stefano Montanelli. 2024. [Lexical semantic change through large language models: a survey](#). *ACM Comput. Surv.*, 56(11).
- Real Academia Española. 2025. [Diccionario de la lengua española](#).
- Elizabeth Traugott. 2008. Grammaticalization, Constructions and the Incremental Development of Language: Suggestions from the Development of Degree Modifiers in English. *Variation, Selection, Development: Probing the Evolutionary Model of Language Change*, pages 219–250.
- Katrien Verveckken. 2012. [Towards a Constructional Account of High and Low Frequency binominal Quantifiers in Spanish](#). *Cognitive Linguistics*, 23(2).