

Detecting Unassimilated Borrowings in Spanish: An Annotated Corpus and Approaches to Modeling

Anonymous ACL submission

Abstract

This work presents a new resource for borrowing identification and analyzes the performance and errors of several models on this task. We introduce a new annotated corpus of Spanish newswire rich in unassimilated lexical borrowings—words from one language that are introduced into another without orthographic adaptation—and use it to evaluate how several sequence labeling models (CRF, BiLSTM-CRF, and Transformer-based models) perform. The corpus contains 370,000 tokens and is larger, more borrowing-dense, OOV-rich, and topic-varied than previous corpora available for this task. Our results show that a BiLSTM-CRF model fed with either Transformer-based embeddings pretrained on codeswitched data or a combination of contextualized word embeddings (along with character embeddings and Spanish and English subword embeddings) outperforms results obtained by a multilingual BERT-based model.

1 Introduction and related work

Lexical borrowing is the process of bringing words from one language into another (Onysko, 2007; Poplack et al., 1988). Borrowings are a common source of out-of-vocabulary (OOV) words, and the task of detecting borrowings has proven to be useful both for lexicographic purposes and for NLP downstream tasks such as parsing (Alex, 2008a), text-to-speech synthesis (Leidig et al., 2014) and machine translation (Tsvetkov and Dyer, 2016).

Recent work has approached the problem of extracting lexical borrowings in European languages such as German (Alex, 2008b; Garley and Hockenmaier, 2012; Leidig et al., 2014), Italian (Furiassi and Hofland, 2007), French (Alex, 2008a; Chesley, 2010), Finnish (Mansikkaniemi and Kurimo, 2012), Norwegian (Andersen, 2012; Losnegaard and Lyse, 2012), and Spanish (Serigos, 2017), with a particular focus on English lexical borrowings (often called *anglicisms*).

Computational approaches to mixed-language data have traditionally framed the task of identifying the language of a word as a tagging problem, where every word in the sequence receives a language tag (Lignos and Marcus, 2013; Molina et al., 2016; Solorio et al., 2014). As lexical borrowings can be single (e.g. *app*, *online*, *smartphone*) or multi-token (e.g. *machine learning*), they are a natural fit for chunking-style approaches. Álvarez Mellado (2020b) introduced chunking-based models for borrowing detection in Spanish media which were later improved (Álvarez Mellado, 2020a), producing an F1 score of 86.41.

However, both the dataset and modeling approach used by Álvarez Mellado (2020a) had significant limitations. The dataset focused exclusively on a single source of news and consisted only of headlines. The number and variety of borrowings were limited, and there was a significant overlap in borrowings between the training set and the test set, which prevented assessment of whether the modeling approach was actually capable of generalizing to previously unseen borrowings. Additionally, the best results were obtained by a CRF model, and more sophisticated approaches were not explored.

The contributions of this paper are a new corpus of Spanish annotated with unassimilated lexical borrowings and a detailed analysis of the performance of several sequence-labeling models trained on this corpus. The models include a CRF, Transformer-based models, and BiLSTM-CRF with different word, subword and character embeddings (including contextualized embeddings pretrained on codeswitched data). The corpus contains 370,000 tokens and is larger and more topic-varied than previous resources. The test set was designed to be as difficult as possible; it covers sources and dates not seen in the training set, includes a high number of OOV words (92% of the borrowings in the test set are OOV) and is very borrowing-dense (20 borrowings per 1,000 tokens).

042
043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082

Media	Topics	Set(s)
ElDiario.es	General newspaper	Train, Dev.
El orden mundial	Politics	Test
Cuarto poder	Politics	Test
Politikon	Politics	Test
El salto	Politics	Test
La Marea	Politics	Test
Pfkará	Feminism	Test
El blog salmón	Economy	Test
Pop rosa	Gossip	Test
Vida extra	Videogames	Test
Espinof	Cinema & TV	Test
Xataka	Technology	Test
Xataka Ciencia	Technology	Test
Xataka Android	Technology	Test
Genbeta	Technology	Test
Microsiervos	Technology	Test
Agencia Sinc	Science	Test
Diario del viajero	Travel	Test
Bebe y más	Parenthood	Test
Vitónica	Lifestyle & sports	Test
Los otros 18	Sports	Test
Foro atletismo	Sports	Test
Motor pasión	Automobiles	Test

Table 1: Sources included in each dataset split (URLs provided in the appendix)

2 Data collection and annotation

2.1 Contrasting lexical borrowing with codeswitching

Linguistic borrowing can be defined as the transference of linguistic elements between two languages. Borrowing and code-switching have been described as a continuum (Clyne et al., 2003).

Lexical borrowing involves the incorporation of single lexical units from one language into another language and is usually accompanied by morphological and phonological modification to conform with the patterns of the recipient language (Haugen, 1950; Onysko, 2007; Poplack et al., 1988).

On the other hand, code-switches are by definition not integrated into a recipient language, unlike established loanwords (Poplack, 2012). While code-switches require a substantial level of fluency, comply with grammatical restrictions in both languages, and are produced by bilingual speakers in bilingual discourses, lexical borrowings are words used by monolingual individuals that eventually become lexicalized and assimilated as part of the recipient language lexicon until the knowledge of “foreign” disappears (Lipski, 2005).

2.2 Data selection

Our dataset consists of Spanish newswire annotated for unassimilated lexical borrowings. All of the

Set	Tokens	ENG	OTHER	Unique
Training	231,126	1,493	28	380
Development	82,578	306	49	316
Test	58,997	1,239	46	987
Total	372,701	3,038	123	1,683

Table 2: Corpus splits with counts

sources used are European Spanish online publications (newspapers, blogs, and news sites) published in Spain and written in European Spanish.

Data was collected separately for the training, development, and test sets to ensure minimal overlap in borrowings, topics, and time periods. The training set consists of a collection of articles appearing between August and December 2020 in *elDiario.es*, a progressive online newspaper based in Spain. The development set contains sentences in articles from January 2021 from the same source.

The data in the test set consisted of annotated sentences extracted in February and March 2021 from a diverse collection of online Spanish media that covers specialized topics rich in lexical borrowings and usually not covered by *elDiario.es*, such as sports, gossip or videogames (see Table 8).

To focus annotation efforts for the training set on articles likely to contain unassimilated borrowings, the articles to be annotated were selected by first using a baseline model and were then human-annotated. To detect potential borrowings, the CRF model and data from Álvarez Mellado (2020b) was used along with a dictionary look-up pipeline. Articles that contained more than 5 borrowing candidates were selected for annotation.

The main goal of data selection for the development and test sets was to create borrowing-dense, OOV-rich datasets, allowing for better assessment of generalization. To that end, the annotation was based on sentences instead of full articles. If a sentence contained a word either flagged as a borrowing by the CRF model, contained in a wordlist of English, or simply not present in the training set, it was selected for annotation. This data selection approach ensured a high number of borrowings and OOV words, both borrowings and non-borrowings. While the training set contains 6 borrowings per 1,000 tokens, the test set contains 20 borrowings per 1,000 tokens. Additionally, 90% of the unique borrowings in the development set were OOV (not present in training). 92% of the borrowings in the test set did not appear in training (see Table 2).

2.3 Annotation process

The corpus was annotated with BIO encoding using Doccano (Nakayama et al., 2018) by a native speaker of Spanish with a background in linguistic annotation. The annotation guidelines (which we provide in the appendix) were based on those of Álvarez Mellado (2020a) but were expanded to account for a wider diversity of topics. English lexical borrowings were labeled ENG, other borrowings were labeled OTHER. Here is an example from the training set:¹

Benching [ENG], estar en el banquillo de tu
crush [ENG] mientras otro juega de titular.

In order to assess the quality of the guidelines and the annotation, a sample of 9,110 tokens from 450 sentences (60% from the test set, 20% from training, 20% from development) was divided among a group of 9 linguists for double annotation. The mean inter-annotation agreement computed by Cohen’s kappa was 0.91, which is above the 0.8 threshold of reliable annotation (Artstein and Poesio, 2008).

2.4 Limitations

We believe it is best to be upfront about the potential limitations of this resource. The corpus consists exclusively of news published in Spain and written in European Spanish. This fact by no means implies the assumption that European Spanish represents the whole of the Spanish language.

The notion of assimilation is usage-based and community-dependant, and thus the dataset we present and the annotation guidelines that were followed were designed to capture a very specific phenomena at a given time and in a given place: unassimilated borrowings in the Spanish press.

In order to establish whether a given word has been assimilated or not, the annotation guidelines rely on lexicographic sources such as the prescriptivist *Diccionario de la Lengua Española* (Real Academia Española, 2020) by the Royal Spanish Academy, a dictionary that aims to cover world-wide Spanish but whose Spain-centric criteria has been previously pointed out (Blanch, 1995; Fernández Gordillo, 2014). In addition, prior work has suggested that Spanish from Spain may have a higher tendency of anglicism-usage than other Spanish dialects (McClelland, 2021). Consequently, we limit the scope of the dataset to

¹“Benching: being on your crush’s bench while someone else plays in the starting lineup.”

Set	Precision	Recall	F1
Development			
ALL	74.13	59.72	66.15
ENG	74.20	68.63	71.31
OTHER	66.67	4.08	7.69
Test			
ALL	77.89	43.04	55.44
ENG	78.09	44.31	56.54
OTHER	57.14	8.70	15.09

Table 3: CRF results on the development and test set

European Spanish not because we consider that this variety represents the whole of the Spanish-speaking community, but because we consider that the approach we have taken here may not account adequately for the whole diversity in borrowing assimilation within the Spanish-speaking world.

3 Modeling

The corpus was used to evaluate four types of models for borrowing extraction: (1) a CRF model, (2) several Transformer-based models, (3) a BiLSTM-CRF model with different types of unadapted embeddings (word, subword, and character embeddings) and (4) a BiLSTM-CRF model with previously fine-tuned Transformer-based embeddings pretrained on codeswitched data. By *unadapted* embeddings, we mean embeddings that have not been fine-tuned for the task of anglicism detection or a related task (e.g. codeswitching).

Evaluation for all models required extracted spans to match the annotation exactly in span and type to be correct. Evaluation was performed with SeqScore (Palen-Michel et al., 2021), using `conlleval`-style repair for invalid label sequences. All models were trained using an AMD 2990WX CPU and a single RTX 2080 Ti GPU.

3.1 Conditional random field model

As baseline model, we evaluated a CRF model with handcrafted features from Álvarez Mellado (2020b). The model was built using `pycrfsuite` (Korobov and Peng, 2014), a Python wrapper for `crfsuite` (Okazaki, 2007) that implements CRF for labeling sequential data. The model also uses the `Token` and `Span` utilities from `spaCy` library (Honnibal and Montani, 2017). The following handcrafted binary features from Álvarez Mellado (2020b) were used for the model:

- Bias: active on all tokens to set per-class bias
- Token: the string of the token

	Development			Test		
	Precision	Recall	F1	Precision	Recall	F1
BETO						
ALL	73.35	72.11	72.73	86.63	76.65	81.34
ENG	73.28	83.33	77.98	86.73	79.10	82.74
OTHER	100.00	2.04	4.00	71.43	10.87	18.87
mBERT						
ALL	82.73	76.90	79.71	89.98	77.59	83.33
ENG	82.97	87.58	85.21	90.45	80.23	85.03
OTHER	71.43	10.20	17.86	33.33	6.52	10.91

Table 4: Results on the development set and test set for Transformer-based models (BETO, mBERT)

- Uppercase: active if the token is all uppercase
- Titlecase: active if only the first character of the token is capitalized
- Character trigram: an active feature for every trigram contained in the token
- Quotation: active if the token is any type of quotation mark (‘ ’ “ ” « »)
- Suffix: last three characters of the token
- POS tag: part-of-speech tag of the token provided by `spaCy` utilities
- Word shape: shape representation of the token provided by `spaCy` utilities
- Word embeddings: provided by Spanish `word2vec` 300 dimensional embeddings by [Cardellino \(2019\)](#)
- URL: active if the token could be validated as a URL according to `spaCy` utilities
- Email: active if the token could be validated as an email address by `spaCy` utilities
- Twitter: active if the token could be validated as a possible Twitter special token: `#hashtag` or `@username`

A window of two tokens in each direction was used for feature extraction. Optimization was performed using L-BFGS with the following hyperparameter values chosen following the best results from [Álvarez Mellado \(2020b\)](#) were set: $c1 = 0.05$, $c2 = 0.01$. As shown in Table 3, the CRF produced an overall F1 score of 66.15 on the development set (P: 74.13, R: 59.72) and an overall F1 of 55.44 (P: 77.89, R: 43.04) on the test set. The CRF results on our dataset are far below the F1 of 86.41 reported by [Álvarez Mellado \(2020b\)](#), showing the impact that a topically-diverse, OOV-rich dataset can have, especially on test set recall (43.04). These results demonstrate that we have created a more difficult task and motivate using more sophisticated models.

3.2 Transformer-based models

We evaluated two Transformer-based models:

- BETO base cased model: a monolingual BERT

model trained for Spanish ([Cañete et al., 2020](#))

- mBERT: multilingual BERT, trained on Wikipedia in 104 languages ([Devlin et al., 2018](#))

Both models were run using the `Transformers` library by HuggingFace ([Wolf et al., 2020](#)). The same default hyperparameters were used for both models: 3 epochs, batch size 32, and maximum sequence length 256.

As shown in Table 4, the mBERT model performed best. Both models performed better on the test set than on the development set, despite the difference in topics between them, suggesting good generalization. This is a remarkable difference compared to the CRF results, where the CRF performed substantially worse on the test set than the development set.

3.3 BiLSTM-CRF

We explored several possibilities for a BiLSTM-CRF model fed with different types of word and subword embeddings. The purpose was to assess whether the combination of different embeddings that encode different linguistic information could outperform the best results obtained by the Transformer-based models in Section 3.2. All of our BiLSTM-CRF models were built using `Flair` ([Akbik et al., 2018](#)) with default hyperparameters (hidden size = 256, learning rate = 0.1, mini batch size = 32, maximum number of epochs = 150) and used embeddings provided by `Flair`.

3.3.1 Preliminary embedding experiments

We first ran exploratory experiments on the development set with different types of embeddings using `Flair` tuning functionalities. We explored the following embeddings: Transformer embeddings (mBERT and BETO), `fastText` embeddings ([Bojanowski et al., 2017](#)), one-hot embeddings, byte pair embeddings ([Heinzerling and Strube, 2018](#)), and character embeddings ([Lample et al., 2016](#)).

The best results were obtained by a combination of mBERT embeddings and character embeddings (F1: 74.00), followed by a combination of BETO embeddings and character embeddings (F1: 72.09). These results show that using contextualized embeddings unsurprisingly outperform non-contextualized embeddings for this task, and that subword representation is important for the task of extracting borrowings that have not been adapted orthographically. The finding regarding the importance of subwords is consistent with previous work; feature ablation experiments for borrowing detection have shown that character trigram features contributed the most to the results obtained by a CRF model (Álvarez Mellado, 2020b).

The worst result (F1: 39.21) was produced by a model fed with one-hot vectors, and the second-worst result was produced by a model fed exclusively with character embeddings. While it performed poorly (F1: 41.65), this fully unlexicalized model outperformed one-hot embeddings, reinforcing the importance of subword information for the task of unassimilated borrowing extraction.

3.3.2 Optimal embedding combination

In light of the preliminary embedding experiments and our earlier experiments with Transformer-based models, we fed our BiLSTM-CRF model with different combinations of contextualized word embeddings (including BERT embeddings from Devlin et al.), byte-pair embeddings and character embeddings. Table 5 shows development set results from different combinations of embeddings. The best overall F1 on the development set (81.79) was obtained by the combination of BETO embeddings, BERT embeddings and byte-pair embeddings. The model fed with BETO embeddings BERT embeddings, byte-pair embeddings and character embeddings ranked second (F1=81.48).

Several things stand out from the results in Table 5. The BETO+BERT embedding combination consistently works better than mBERT embeddings, and BPE embeddings contribute to better results. Character embeddings, however, seem to produce little effect at first glance. Given the same model, adding character embeddings produced little changes in F1 or even slightly hurt the results (81.48 vs 81.79; 78.96 vs 79.01; 78.96 vs 78.64). Although character embeddings seem to make little difference in overall F1, recall was consistently higher in models that included character embeddings (77.46 vs 77.18; 76.62 vs 76.34;

77.18 vs 74.65), and in fact, the model with BETO+BERT embeddings, BPE embeddings and character embeddings produced the highest recall overall (77.46). This is an interesting finding, as our results from Sections 3.1 and 3.2 as well as prior work (Álvarez Mellado, 2020b) identified recall as weak for borrowing detection models.

The two best-performing models from Table 5 (BETO+BERT embeddings, BPE embeddings and optionally character embeddings) were evaluated on the test set. Table 6 gives results per type on the development and test sets for these two models. For both models, results on the test set were better (F1: 84.02, F1: 84.25) than on the development set (F1: 81.79, F1: 81.48). Although the best F1 score on the development set was obtained with no character embeddings, when run on the test set the model with character embeddings obtained the best score and the highest recall (R: 79.53), which again seems to corroborate the positive impact that character information can have in recall when dealing with previously unseen borrowings.

3.4 Borrowing detection as a transfer learning task from codeswitching

Finally, we decided to explore whether detecting unassimilated lexical borrowings could be framed as transfer learning from language identification in codeswitching. As before, we ran a BiLSTM-CRF model using Flair, but instead of using the unadapted Transformer embeddings, we used codeswitch embeddings², fine-tuned Transformer-based embeddings pretrained for language identification on the Spanish-English section of the codeswitching LinCE dataset (Aguilar et al., 2020).

Table 7 gives results for these models on the development and test sets. The two best-performing models on the development set where the BiLSTM-CRF with codeswitch and BPE embeddings (F1: 80.17) and the BiLSTM-CRF model with codeswitch, BPE and character embeddings (F1: 79.5). None of these outperformed the best development set results obtained by the two best performing BiLSTM models with unadapted embeddings from Section 3.3 (F1: 81.79, F1: 81.48).

However, these two models did outperform the BiLSTM-CRF models with unadapted embeddings when run on the test set; the model with codeswitch, BPE and character embeddings produced the best

²<https://github.com/sagorbrur/codeswitch>

Word embedding	BPE embedding	Char embedding	Precision	Recall	F1
mBERT	-	-	82.27	69.30	75.23
mBERT	-	✓	79.45	72.96	76.06
mBERT	multi	✓	81.37	73.80	77.40
mBERT	es, en	-	83.07	74.65	78.64
mBERT	es, en	✓	80.83	77.18	78.96
BETO, BERT	-	✓	81.44	76.62	78.96
BETO, BERT	-	-	81.87	76.34	79.01
BETO, BERT	es, en	✓	85.94	77.46	81.48
BETO, BERT	es, en	-	86.98	77.18	81.79

Table 5: Overall F1 score obtained on the development set with the BiLSTM-CRF model using different combinations of multilingual word embeddings, subword embeddings and character embeddings

	Development			Test		
	Precision	Recall	F1	Precision	Recall	F1
BETO+BERT and BPE						
ALL	86.98	77.18	81.79	90.56	78.37	84.02
ENG	86.90	88.89	87.88	90.68	80.87	85.49
OTHER	50.00	2.04	3.92	71.43	10.87	18.87
BETO+BERT, BPE and char						
ALL	85.94	77.46	81.48	89.57	79.53	84.25
ENG	86.58	88.56	87.56	89.92	82.08	85.82
OTHER	57.14	8.16	14.29	50.00	10.87	17.86

Table 6: Results on the development set and test set for BiLSTM-CRF model with BETO and BERT embeddings, BPE embeddings and optionally character embeddings

all around results on the test set (F1:85.49), a result that outperformed all results from all previous models both on ENG borrowings and OTHER borrowings. It should be noted that this transfer learning approach is indirectly using more data than just the training data from our initial corpus, as the codeswitch-based BiLSTM-CRF models benefit from the labeled data seen during pretraining for the language-identification task.

4 Error analysis

We compared the different results produced by the best performing model of each type on the test set: (1) the mBERT model, (2) the BiLSTM-CRF with BERT+BETO, BPE and character embeddings and (3) the BiLSTM-CRF model with codeswitch, BPE and character embeddings. We divide the error analysis into two sections. We first analyze errors that were made by all three models, with the aim of discovering which instances of the dataset were challenging for all models. We then analyze unique answers (both correct and incorrect) per model, with the aim of gaining insight on what are the unique characteristics of each model in comparison with other models.

4.1 Errors made by all models

4.1.1 Borrowings labeled as O

There were 137 tokens in the test set that were incorrectly labeled as O by all three models. 103 of these were of type ENG, 34 were of type OTHER. These errors can be classified as follows:

– Borrowings in upper case (12), which tend to be mistaken by models with proper nouns:

Análisis de empresa basados en **Big Data** [ENG].³

– Borrowings in sentence-initial position (9), which were titlecased and therefore consistently mislabeled as O:

Youtuber [ENG], mujer y afroamericana: Candace Owen podría ser la alternativa a Trump.⁴

Sentence-initial borrowings are particularly tricky, as models tend to confuse these with foreign named entities. In fact, prior work on anglicism detection based on dictionary lookup (Serigos, 2017) stated that borrowings in sentence-initial position were rare in Spanish and consequently chose to ignore all foreign words in sentence-initial position

³“Business analytics based on Big Data”

⁴“Youtuber, woman and African-American: Candace Owen could be the alternative to Trump”

	Development			Test		
	Precision	Recall	F1	Precision	Recall	F1
Codeswitch embeddings						
ALL	80.30	74.65	77.37	90.50	77.12	83.28
ENG	80.49	86.27	83.28	90.48	79.74	84.77
OTHER	50.00	2.04	3.92	100.00	6.52	12.24
Codeswitch embeddings + char embeddings						
ALL	77.84	75.21	76.50	89.83	80.39	84.85
ENG	78.01	86.93	82.23	89.80	83.13	86.34
OTHER	50.00	2.04	3.92	100.00	6.52	12.24
Codeswitch embeddings + BPE						
ALL	83.08	77.46	80.17	89.18	81.48	85.16
ENG	82.82	88.24	85.44	89.37	84.18	86.70
OTHER	100.00	10.20	18.52	57.14	8.70	15.09
Codeswitch embeddings + BPE + char embeddings						
ALL	82.48	76.90	79.59	90.30	81.17	85.49
ENG	82.21	87.58	84.81	90.42	83.78	86.97
OTHER	100.00	10.20	18.52	71.43	10.87	18.87

Table 7: Results on the development set and test set for LSTM-CRF model with different combinations of codeswitch embeddings, BPE embeddings and character embeddings

under the assumption that they could be considered named entities. However, these examples (and the difficulty they pose for models) prove that sentence-initial borrowings are not rare and therefore should not be overlooked.

– Borrowings that also happen to be words in Spanish (8), such as the word *primer*, that is a borrowing found in makeup articles (*un primer hidratante*, “an hydrating primer”) but also happens to be a fully Spanish adjective meaning “first” (*primer premio*, “first prize”). Borrowings like these are still treated as fully unassimilated borrowings by speakers, even when the form is exactly the same as an already-existing Spanish word and were a common source of mislabeling, especially partial mismatches in multitoken borrowings: *red* (which exists in Spanish meaning “net”) in *red carpet*, *tractor* in *tractor pulling* or *total* in *total look*.

– Borrowings that could pass as Spanish words (58): most of the mislabeled borrowings were words that do not exist in Spanish but that could orthographically pass for a Spanish word. That is the case of words like *burpees* (hypothetically, a conjugated form of the non-existing verb *burpear*), *gimbal*, *mules*, *bromance* or *nude*.

– Other borrowings (50): a high number of mislabeled borrowings were borrowings that were orthographically implausible in Spanish, such as *trenchs*, *multipads*, *hypes*, *riff*, *scrunchie* or *mint*. The fact that none of our models were able to correctly classify these orthographically implausible examples leaves the door open to further exploration of character-based models and investigating character-

level perplexity as a source of information.

4.1.2 Non-borrowings labeled as borrowings

29 tokens were incorrectly labeled as borrowings by all three models. These errors can be classified in the following groups:

– Metalinguistic usage and reported speech: a foreign word or sentence that appears in the text to refer to something someone said or wrote.

Escribir “**icon pack**” [ENG] en el buscador.⁵

– Lower-cased proper nouns: such as websites.

Hay que acceder a la página **flywithkarolg** [ENG]⁶

– Computer commands: the test set included blog posts about technology, which mentioned computer commands (such as *sudo apt-get update*) that were consistently mistaken by our models as borrowings. These may seem like an extreme case—after all, computer commands do contain English words—but they are a good example of the real data that a borrowing-detection system may encounter.

– Foreign words within proper nouns: lower-cased foreign words that were part of multitoken proper nouns.

La serie “10.000 **ships** [ENG]” cuenta la odisea de la princesa Nymeria.⁷

– Acronyms and acronym expansions:

El entrenamiento HITT (**high intensity interval training** [ENG])⁸

⁵“Type ‘icon pack’ on the search box”

⁶“You need to access the website flywithkarolg”

⁷“The series ‘10,000 ships’ tells the story of princess Nymeria”

⁸“HITT training (High-intensity interval training)”

521	– Assimilated borrowings: certain borrowings	including old borrowings that are considered today	567
522	that are already considered by RAE’s dictionary	as fully assimilated (such as <i>films</i> or <i>sake</i>) or the	568
523	as fully assimilated were labeled by all models as	usage of <i>post</i> as a prefix of latin origin (as in <i>post-</i>	569
524	anglicisms.	<i>produccion</i>), which other models mistook with the	570
525	Labios rojos, a juego con el top [ENG]. ⁹	English word <i>post</i> .	571
526	4.1.3 Type confusion	4.2.3 BiLSTM-CRF with codeswitch	572
527	Three tokens of type OTHER were marked by all	embeddings	573
528	models as ENG. There were no ENG borrowings	The codeswitch-based system incorrectly labeled	574
529	that were labeled as OTHER by all three models.	18 tokens as borrowings, including proper nouns	575
530	Había buffet [ENG] libre. ¹⁰	(7) such as <i>Baby Spice</i> , and fully asimilated bor-	576
531	4.2 Unique answers per model	rowings (5), such as <i>jersey</i> , <i>relax</i> or <i>tutorial</i> .	577
532	We now summarize the unique mistakes and correct	This model correctly labeled 27 tokens that were	578
533	answers made per model, with the aim of under-	mistakenly ignored by other models, including mul-	579
534	standing what data points were handled uniquely	titoken borrowings (<i>dark and gritty</i> , <i>red carpet</i>)	580
535	well or badly by each model.	and other borrowings that were non-compliant with	581
536	4.2.1 mBERT	Spanish orthograpich rules but that were however	582
537	There were 46 tokens that were incorrectly labeled	ignored by other models (<i>messy</i> , <i>athleisure</i> , <i>multi-</i>	583
538	as borrowings only by the mBERT model. These	<i>touch</i> , <i>workaholic</i>).	584
539	include foreign words used in reported speech or	The codeswitch-based model also correctly la-	585
540	acronym expansion (21), proper names (11) and	beled as ○ 16 tokens that the other two models	586
541	already assimilated borrowings (7).	labeled as borrowings, including acronym expan-	587
542	There were 27 tokens that were correctly labeled	sions, lower-cased proper names and orthographi-	588
543	only by the mBERT model. The mBERT model	cally unorthodox Spanish words, such as the ideo-	589
544	was particularly good at detecting the full span	phone <i>tiki-taka</i> or <i>shavales</i> (a non-standard writing	590
545	of multitoken borrowings as in <i>no knead bread</i> ,	form of the word <i>chavales</i> , “guys”).	591
546	<i>total white</i> , <i>wide leg</i> or <i>kettlebell swings</i> (which	5 Conclusion	592
547	were only partially detected by other models) and	We have introduced a new corpus of Spanish	593
548	at detecting borrowings that could pass for Spanish	newswire annotated with unassimilated lexical bor-	594
549	words (such as <i>fashionista</i> , <i>samples</i> , <i>vocoder</i>). In	rowings. The test set has a high number of OOV	595
550	addition, the mBERT model also correctly labeled	borrowings—92% of unique borrowings in the test	596
551	as ○ 12 tokens that the other two models mistook as	set were not seen during training—and is more	597
552	borrowings, including morphologically adapted an-	borrowing-dense and varied than resources previ-	598
553	glicisms, such as <i>craftear</i> (Spanish infinitive of the	ously available. We have used the dataset to explore	599
554	verb <i>to craft</i>) <i>crackear</i> (from <i>to crack</i>) or <i>lookazo</i>	several sequence labeling models (CRF, BiLSTM-	600
555	(augmentative of the noun <i>look</i>).	CRF, and Transformer-based models) for the task	601
556	4.2.2 BiLSTM-CRF with unadapted	of extracting lexical borrowings in an high-OOV	602
557	embeddings	setting. Results show that the BiLSTM-CRF model	603
558	There were 23 tokens that were incorrectly labeled	fed with character embeddings, Spanish and En-	604
559	as borrowings solely by this model, the most com-	glish subword embeddings and either Transformer-	605
560	mon types being assimilated borrowings (such as	-based embeddings pretrained on codeswitched data	606
561	<i>fan</i> , <i>clon</i>) and Spanish words (<i>fiestiones</i>) (9 each).	(F1: 85.49) or a combination of contextualized	607
562	32 tokens were correctly labeled as borrowings	word embeddings (F1: 84.25) produced the best	608
563	only by this model. These include borrowings that	results and outperformed prior models for this task	609
564	could pass for Spanish words (<i>camel</i> , <i>canvas</i>). In	(CRF F1: 55.44) and multilingual Transformer-	610
565	addition, this model also correctly labeled as ○ 6	-based models (mBERT F1: 83.33).	611
566	tokens that the other two mistook as borrowings,	References	612
	⁹ “Red lips, matching top”	Fundéu. Fundación del Español Urgente. http://	613
	¹⁰ “There was a free buffet”	www.fundeu.es .	614

615	Gustavo Aguilar, Sudipta Kar, and Thamar Solorio.	Michael Clyne, Michael G Clyne, and Clyne Michael.	669
616	2020. LinCE: A centralized benchmark for linguistic code-switching evaluation . In <i>Proceedings of the 12th Language Resources and Evaluation Conference</i> , pages 1803–1813, Marseille, France. European Language Resources Association.	2003. <i>Dynamics of language contact: English and immigrant languages</i> . Cambridge University Press.	670
617			671
618		Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. BERT: pre-training of deep bidirectional transformers for language understanding . <i>CoRR</i> , abs/1810.04805.	672
619			673
620			674
621	Alan Akbik, Duncan Blythe, and Roland Vollgraf.	Luz Fernández Gordillo. 2014. La lexicografía del español y el español hispanoamericano. <i>Andamios</i> , 11(26):53–89.	676
622	2018. Contextual string embeddings for sequence labeling. In <i>COLING 2018, 27th International Conference on Computational Linguistics</i> , pages 1638–1649.		677
623			678
624		Cristiano Furiassi and Knut Hofland. 2007. The retrieval of false anglicisms in newspaper texts. In <i>Corpus Linguistics 25 Years On</i> , pages 347–363. Brill Rodopi.	679
625			680
626	Beatrice Alex. 2008a. <i>Automatic detection of English inclusions in mixed-lingual data with an application to parsing</i> . Ph.D. thesis, University of Edinburgh.		681
627			682
628		Matt Garley and Julia Hockenmaier. 2012. Beefmoves: Dissemination, diversity, and dynamics of English borrowings in a German hip hop forum . In <i>Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 135–139, Jeju Island, Korea. Association for Computational Linguistics.	683
629	Beatrice Alex. 2008b. Comparing corpus-based to web-based lookup techniques for automatic English inclusion detection . In <i>Proceedings of the Sixth International Conference on Language Resources and Evaluation (LREC’08)</i> , Marrakech, Morocco. European Language Resources Association (ELRA).		684
630			685
631			686
632			687
633			688
634			689
635	Elena Álvarez Mellado. 2020a. An Annotated Corpus of Emerging Anglicisms in Spanish Newspaper Headlines. In <i>Proceedings of the Fourth Workshop on Computational Approaches to Code Switching</i> , pages 1–8, Marseille, France. European Language Resources Association.	Juan Gómez Capuz. 1997. Towards a typological classification of linguistic borrowing (illustrated with anglicisms in romance languages). <i>Revista alicantina de estudios ingleses</i> , 10:81–94.	690
636			691
637			692
638			693
639			
640			
641	Elena Álvarez Mellado. 2020b. <i>Lázaro: An extractor of emergent anglicisms in Spanish newswire</i> . Master’s thesis, Brandeis University.	Einar Haugen. 1950. The analysis of linguistic borrowing. <i>Language</i> , 26(2):210–231.	694
642			695
643			
644	Gisle Andersen. 2012. Semi-automatic approaches to Anglicism detection in Norwegian corpus data. In Cristiano Furiassi, Virginia Pulcini, and Félix Rodríguez González, editors, <i>The anglicization of European lexis</i> , pages 111–130. John Benjamins.	Benjamin Heinzerling and Michael Strube. 2018. BPEmb: Tokenization-free pre-trained subword embeddings in 275 languages . In <i>Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)</i> , Miyazaki, Japan. European Language Resources Association (ELRA).	696
645			697
646			698
647			699
648			700
649	Ron Artstein and Massimo Poesio. 2008. Inter-coder agreement for computational linguistics. <i>Computational Linguistics</i> , 34(4):555–596.		701
650			702
651		Matthew Honnibal and Ines Montani. 2017. spaCy 2: Natural language understanding with bloom embeddings, convolutional neural networks and incremental parsing. https://spacy.io/ .	703
652	Juan M Lope Blanch. 1995. Americanismo frente a españolismo lingüísticos. <i>Nueva revista de filología hispánica</i> , 43(2):433–440.		704
653			705
654			706
655	Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. Enriching word vectors with subword information . <i>Transactions of the Association for Computational Linguistics</i> , 5:135–146.	Mikhail Korobov and Terry Peng. 2014. Python-crfsuite. https://github.com/scrapinghub/python-crfsuite .	707
656			708
657			709
658		Guillaume Lample, Miguel Ballesteros, Sandeep Subramanian, Kazuya Kawakami, and Chris Dyer. 2016. Neural architectures for named entity recognition . In <i>Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 260–270, San Diego, California. Association for Computational Linguistics.	710
659	Cristian Cardellino. 2019. Spanish Billion Words Corpus and Embeddings. https://crscardellino.github.io/SBWCE/ .		711
660			712
661			713
662	José Cañete, Gabriel Chaperon, Rodrigo Fuentes, Juhui Ho, Hojin Kang, and Jorge Pérez. 2020. Spanish pre-trained bert model and evaluation data. In <i>PML4DC at ICLR 2020</i> .		714
663			715
664			716
665			717
666	Paula Chesley. 2010. Lexical borrowings in French: Anglicisms as a separate phenomenon. <i>Journal of French Language Studies</i> , 20(3):231–251.	Sebastian Leidig, Tim Schlippe, and Tanja Schultz. 2014. Automatic detection of anglicisms for the pronunciation dictionary generation: a case study on our German IT corpus. In <i>Spoken Language Technologies for Under-Resourced Languages</i> .	718
667			719
668			720
			721
			722

723	Constantine Lignos and Mitch Marcus. 2013. Toward	Barbara Plank. 2016. What to do about non-standard	777
724	web-scale analysis of codeswitching. Presented at	(or non-canonical) language in NLP. <i>arXiv preprint</i>	778
725	the 87th Annual Meeting of the Linguistic Society	<i>arXiv:1608.07836</i> .	779
726	of America.		
727	John M Lipski. 2005. Code-switching or borrowing?	Shana Poplack. 2012. What does the nonce borrowing	780
728	No sé so no puedo decir, you know. In <i>Selected</i>	hypothesis hypothesize? <i>Bilingualism: Language</i>	781
729	<i>proceedings of the second workshop on Spanish so-</i>	<i>and Cognition</i> , 15(3):644–648.	782
730	<i>ciolinguistics</i> , pages 1–15. Cascadilla Proceedings		
731	Project Somerville, MA.	Shana Poplack, David Sankoff, and Christopher Miller.	783
732	Gyri Smordal Losnegaard and Gunn Inger Lyse. 2012.	1988. The social correlates and linguistic processes	784
733	A data-driven approach to anglicism identification	of lexical borrowing and assimilation. <i>Linguistics</i> ,	785
734	in Norwegian. In Gisle Andersen, editor, <i>Explor-</i>	26(1):47–104.	786
735	<i>ing Newspaper Language: Using the web to create</i>		
736	<i>and investigate a large corpus of modern Norwegian</i> ,	Chris Pratt. 1980. <i>El anglicismo en el español peninsu-</i>	787
737	pages 131–154. John Benjamins Publishing.	<i>lar contemporáneo</i> , volume 308. Gredos.	788
738	André Mansikkaniemi and Mikko Kurimo. 2012. Un-	Real Academia Española. 2020. Diccionario de la	789
739	supervised vocabulary adaptation for morph-based	lengua española, ed. 23.4. http://dle.rae.	790
740	language models. In <i>Proceedings of the NAACL-</i>	<i>es</i> .	791
741	<i>HLT 2012 Workshop: Will We Ever Really Replace</i>		
742	<i>the N-gram Model? On the Future of Language</i>	Felix Rodríguez González. 1999. Anglicisms in	792
743	<i>Modeling for HLT</i> , pages 37–40. Association for	contemporary Spanish. An overview. <i>Atlantis</i> ,	793
744	Computational Linguistics.	21(1/2):103–139.	794
745	Jacob McClelland. 2021. A brief survey of anglicisms	Jacqueline Rae Larsen Serigos. 2017. <i>Applying corpus</i>	795
746	among spanish dialects. <i>Colorado Research in Lin-</i>	<i>and computational methods to loanword research:</i>	796
747	<i>guistics</i> .	<i>new approaches to Anglicisms in Spanish</i> . Ph.D. the-	797
748	Giovanni Molina, Fahad AlGhamdi, Mahmoud	sis, The University of Texas at Austin.	798
749	Ghoneim, Abdelati Hawwari, Nicolas Rey-	Thamar Solorio, Elizabeth Blair, Suraj Mahar-	799
750	Villamizar, Mona Diab, and Tamar Solorio.	jan, Steven Bethard, Mona Diab, Mahmoud	800
751	2016. Overview for the second shared task on	Ghoneim, Abdelati Hawwari, Fahad AlGhamdi, Ju-	801
752	language identification in code-switched data .	lia Hirschberg, Alison Chang, and Pascale Fung.	802
753	In <i>Proceedings of the Second Workshop on Computa-</i>	2014. Overview for the first shared task on language	803
754	<i>tional Approaches to Code Switching</i> , pages 40–49,	identification in code-switched data . In <i>Proceedings</i>	804
755	Austin, Texas. Association for Computational	<i>of the First Workshop on Computational Approaches</i>	805
756	Linguistics.	<i>to Code Switching</i> , pages 62–72, Doha, Qatar. Asso-	806
757	Hiroki Nakayama, Takahiro Kubo, Junya Kamura,	ciation for Computational Linguistics.	807
758	Yasufumi Taniguchi, and Xu Liang. 2018. doc-	Yulia Tsvetkov and Chris Dyer. 2016. Cross-lingual	808
759	cano: Text annotation tool for human. https:	bridges with models of lexical borrowing. <i>Journal</i>	809
760	://github.com/doccano/doccano .	<i>of Artificial Intelligence Research</i> , 55:63–93.	810
761	Eugenia Esperanza Núñez Nogueroles. 2017. An up-	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien	811
762	to-date review of the literature on anglicisms in	Chaumond, Clement Delangue, Anthony Moi, Pier-	812
763	Spanish. <i>Diálogo de la Lengua</i> , IX, pages 1–54.	ric Cistac, Tim Rault, Remi Louf, Morgan Funtow-	813
764	Naoaki Okazaki. 2007. Crfsuite: a fast implementation	icz, Joe Davison, Sam Shleifer, Patrick von Platen,	814
765	of Conditional Random Fields (CRFs). http://	Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu,	815
766	www.chokkan.org/software/crfsuite/ .	Teven Le Scao, Sylvain Gugger, Mariama Drame,	816
767	Alexander Onysko. 2007. <i>Anglicisms in German: Bor-</i>	Quentin Lhoest, and Alexander Rush. 2020. Trans-	817
768	<i>rowing, lexical productivity, and written codeswitch-</i>	formers: State-of-the-art natural language process-	818
769	<i>ing</i> , volume 23. Walter de Gruyter.	ing . In <i>Proceedings of the 2020 Conference on Em-</i>	819
770	Chester Palen-Michel, Nolan Holley, and Constantine	<i>pirical Methods in Natural Language Processing:</i>	820
771	Lignos. 2021. SeqScore: Addressing barriers to re-	<i>System Demonstrations</i> , pages 38–45, Online. Asso-	821
772	producibile named entity recognition evaluation .	ciation for Computational Linguistics.	822
773	In <i>Proceedings of the 2nd Workshop on Evaluation and</i>		
774	<i>Comparison of NLP Systems</i> , pages 40–50, Punta		
775	Cana, Dominican Republic. Association for Compu-		
776	tational Linguistics.		

823	A Data sources		870
824	See Table 8 for the URLs of the sources used in the		871
825	test set.		872
826	B Annotation guidelines		873
827	B.1 Objective		874
828	This document proposes a set of guidelines for		875
829	annotating emergent unassimilated lexical borrow-		876
830	ings, with a focus on English lexical borrowings		877
831	(or anglicisms). The purpose of these annotation		878
832	guidelines is to assist annotators to annotate unas-		879
833	similated lexical borrowings from English that ap-		880
834	pear in Spanish newswire, i.e. words from English		881
835	origin that are introduced into Spanish without any		882
836	morphological or orthographic adaptation.		883
837	This project approaches the phenomenon of lex-		884
838	ical borrowing from a synchronic point of view,		885
839	which means that we will not be annotating all		886
840	words that have been borrowed at some point of		887
841	the history of the Spanish language (like arabisms),		888
842	but only those that have been recently imported and		889
843	have not been integrated into the recipient language		890
844	(in this case, Spanish).		891
845	B.2 Tagset		892
846	We will consider two possible tags for our anno-		893
847	tation: ENG, for borrowings that come from the		894
848	English language (or anglicisms), and OTHER for		895
849	other borrowings that comply with the following		896
850	guidelines but that come from languages other than		
851	English.		
852	B.3 What an unassimilated lexical borrowing		897
853	is		898
854	In this section we provide an overview of what		899
855	words will be considered as unassimilated lexical		900
856	borrowings for the sake of our annotation project.		901
857	B.3.1 Definition and scope		902
858	The concept of <i>linguistic borrowing</i> covers a wide		903
859	range of linguistic phenomena. We will first pro-		904
860	vide a general overview of what lexical borrowing		905
861	is and what will be understood as an anglicism		906
862	within the scope of this project.		907
863	Lexical borrowing is the incorporation of single		908
864	lexical units from one language (the donor lan-		909
865	guage) into another language (the recipient lan-		910
866	guage) and is usually accompanied by morpho-		911
867	logical and phonological modification to conform		912
868	with the patterns of the recipient language (Haugen,		913
869	1950; Onysko, 2007; Poplack et al., 1988).		
	Anglicisms are lexical borrowings that come		
	from the English language (Gómez Capuz,		
	1997; Pratt, 1980; Rodríguez González, 1999;		
	Núñez Nogueroles, 2017). For our annotation		
	project, we will focus on direct, unassimilated,		
	emerging anglicisms, i.e. lexical borrowings from		
	the English language into Spanish that have re-		
	cently been imported and that have still not been		
	assimilated into Spanish, that is, words like <i>smart-</i>		
	<i>phone</i> , <i>influencer</i> , <i>hype</i> , <i>lawfare</i> or <i>reality show</i> .		
	Although this project focuses on lexical borrow-		
	ings from English, we will also consider borrow-		
	ings from other languages that comply with these		
	guidelines. Borrowings from the English language		
	will be annotated with the tag ENG, while borrow-		
	ings from other languages shall be annotated with		
	the tag OTHER:		
	... financiados a		
	través de la plataforma de		
	[crowdfunding](ENG) del club		
	[gourmet](OTHER) que tengas más		
	cerca ¹¹		
	Other types of borrowings, such as semantic		
	calques, syntactic anglicisms or literal translations		
	will be considered beyond the scope of these an-		
	notation project and will not be covered in these		
	guidelines.		
	B.3.2 Types of lexical borrowing		
	Lexical borrowings can be adapted (the spelling of		
	the word is modified to comply with the phonolog-		
	ical and orthographic patterns of the recipient lan-		
	guage, as in <i>fútbol</i> or <i>tuit</i>) or unadapted (the word		
	preserves its original spelling: <i>millennial</i> , <i>newslet-</i>		
	<i>ter</i> , <i>like</i>). For this annotation project, we will be		
	focusing on unassimilated lexical borrowings: this		
	means that adapted borrowings will be ignored and		
	only unadapted borrowings will be tagged (see sec-		
	tion B.4.2 for a full description on the differences		
	between adapted and unadapted borrowings).		
	B.3.3 Multiword borrowings		
	Lexical borrowings can be both single-token units		
	(<i>online</i> , <i>impeachment</i>), as well as multiword ex-		
	pressions (<i>reality show</i> , <i>best seller</i>). Multitoken		
	borrowings will be labeled as one entity.		
	¹¹ Examples in these guidelines will display the lexical bor-		
	rowing that should be labeled between square brackets, with		
	the the corresponding tag in parentheses. Examples with no		
	words marked with brackets will illustrate cases where no		
	lexical borrowing should be tagged.		

Media	URL
El orden mundial	https://elordenmundial.com/
Cuarto poder	https://www.cuartopoder.es/
Politikon	https://www.politikon.es/
El salto	https://www.elsaltodiario.com/
La Marea	https://www.lamarea.com/
Píkara	https://www.pikaramagazine.com/
El blog salmón	https://www.elblogsalmon.com/
Pop rosa	https://www.poprosa.com/
Vida extra	https://www.vidaextra.com/
Espinof	https://www.espinof.com/
Xataka	https://www.xataka.com/
Xataka Ciencia	https://www.xatakaciencia.com/
Xataka Android	https://www.xatakandroid.com/
Genbeta	https://www.genbeta.com/
Microsiervos	https://www.microsiervos.com/
Agencia Sinc	https://www.agenciasinc.es/
Diario del viajero	https://www.diariodelviajero.com/
Bebe y más	https://www.bebesymas.com/
Vitónica	https://www.vitonica.com/
Los otros 18	https://www.losotros18.com/
Foro atletismo	https://www.foroatletismo.com/
Motor pasión	https://www.motorpasion.com/

Table 8: Media included in the test set

914 imagina ser un '[tech bro]' con
915 millones de dólares (ENG)

916 The annotation should however distinguish be-
917 tween a multitoken borrowing and adjacent borrow-
918 ings. A phrase like *signature look* is a multiword
919 borrowing (the full phrase has been borrowed as a
920 single unit) and should be annotated as such.

921 para recrear su [total look]
922 (ENG)

923 However, a phrase like *look sporty* follows the
924 NAdj order that is typical of Spanish grammar (but
925 impossible in English): these are in fact two sep-
926 arate borrowings (*look* and *sporty*) that have been
927 borrowed independently and happen to be colo-
928 cated in a phrase. The annotation should capture
929 these nuances:

930 un [look] (ENG) [sporty] (ENG)
931 perfecto

932 B.3.4 Origin of the borrowings

933 Establishing the origin of a certain borrowings can
934 sometimes be tricky, as the language of origin can
935 sometimes be disputed. Additionally, certain bor-
936 rowings might have originated in a certain lan-
937 guage, but may have reached the recipient language
938 through another language.

939 In order to establish the origin of borrowings,
940 the origin attributed by reference dictionaries and

institutions (Real Academia Española, 2020; fun)
will be followed.

941 This means that words like *junior* and *senior*
942 (whose frequency and perhaps even their pronun-
943 ciation may have changed due to the influence of
944 English) will still be considered as latinisms, as
945 DLE registers their adapted versions (*júnior* and
946 *sénior*) as such (and mentions no English influ-
947 ence). Similarly, the word *barista* might have en-
948 tered the Spanish language via English, but RAE's
949 Observatorio de Palabras considers it of Italian
950 origin (and should therefore be annotated with
951 OTHER label). 952

953 B.4 What an unasimilated lexical borrowing 954 is not 955

956 In the previous section we provided an overview of
957 what words will be considered as an unassimilated
958 lexical borrowing for the sake of our annotation
959 project. In this section we will cover what an unas-
960 similated lexical borrowing is *not*.

961 There are several phenomena that are close
962 enough to unassimilated borrowing and that can
963 sometimes be mistaken with. In this section we
964 will list what phenomena will not be considered as
965 unassimilated lexical borrowings (and are therefore
966 beyond the scope of our annotation project), as well
967 as provide guidelines in order to distinguish these
968 cases and adjudicate them.

969 We will focus on three main phenomena: assim-
970 ilated borrowings, proper names and code-mixed

971 inclusions.

972 Figure 1 summarizes the decision steps that can
973 be followed when deciding if a certain word should
974 be labeled or not as a lexical borrowing.

975 **B.4.1 Assimilated vs unassimilated** 976 **borrowings**

977 This annotation project aims to capture unassimilated
978 lexical borrowings. As a general rule, all
979 unadapted lexical borrowings should be tagged.
980 This means that direct borrowings that have not
981 gone through any morphological or orthographic
982 modification process should be labeled.

983 Lexical adaptation, however, is a diachronic
984 process and, as a result, what constitutes an unadapted
985 borrowing is not clear-cut. The following guide-
986 lines define what borrowings will be considered
987 as unassimilated (and therefore should be tagged)
988 versus those that have already been integrated into
989 the recipient language (and therefore should not be
990 tagged).

991 **B.4.2 Adapted borrowings**

992 Words that have already gone through orthographi-
993 cal or morphological adaptation (such as *fútbol*,
994 *líder*, *tuit* or *espoiler*) will be considered assimilated
995 and therefore should not be labeled. Partial
996 adaptations (such as *márketing*, where an accent
997 has been added) will also be excluded.

998 Borrowings that have not been adapted but

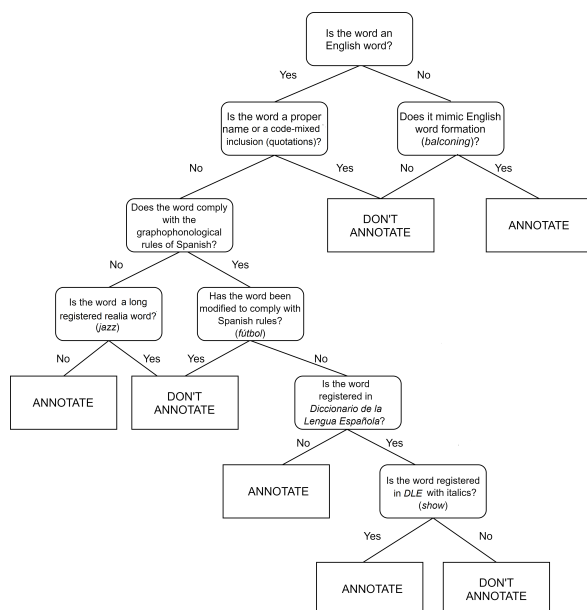


Figure 1: Decision steps to follow during the annotation process to decide whether to annotate a word as an anglicism.

999 whose original spelling complies with grapho-
1000 phonological rules of Spanish (and are therefore
1001 unlikely to be further adapted, such as *bar*, *fan*,
1002 *web*, *internet*, *club*, *set* or *videoclip*) will be tagged
1003 as a borrowing or not depending on how re-
1004 cent or emergent they are. In order to determine
1005 which unadapted, graphophonologically acceptable
1006 borrowings are to be annotated, the latest online
1007 version of the *Diccionario de la lengua española*
1008 (Real Academia Española, 2020) will be consulted
1009 (as of February 2021)¹². If the DLE dictionary
1010 already registers the word with that meaning and
1011 with no italics or quotation marks, then it will be
1012 considered assimilated and therefore should not be
1013 tagged.

1014 This means that a word like *set* (when used to
1015 refer to a collection of things, a television studio
1016 or a part of a tennis match) will be considered as-
1017 similated because it is already registered in DLE
1018 dictionary with no italics, and therefore should not
1019 be labeled as ENG. On the other hand, a word like
1020 *nude*, although its spelling also complies with Span-
1021 ish graphophonological rules, will be considered
1022 an unassimilated borrowing because it has not been
1023 registered yet in the dictionary, and should there-
1024 fore be tagged as such.

1025 ganó el primer set

1026 los tonos '[nude]' (ENG)

1027 It should be noted that this guideline only
1028 applies to lexical borrowings that comply with
1029 graphophonological rules of Spanish. Unadapted
1030 lexical borrowings that do not comply with
1031 graphophonological rules of Spanish (such as *show*,
1032 *look*, etc) will be tagged as borrowing, regardless
1033 of whether the word is included in the dictionary
1034 or not (although see section B.4.6 for exceptions to
1035 this).

1036 It is important to emphasize that, in order for
1037 an unadapted graphophonologically-compliant bor-
1038 rowing to be considered assimilated it should be
1039 registered in the dictionary both without italics and
1040 with the corresponding meaning. For instance, a
1041 word like *top* (that is graphophonologically accept-
1042 able in Spanish) is registered in DLE with no italics,
1043 but it is only registered with the meaning of
1044 a piece of clothing. The word *top* as referring to
1045 the upper part of something (as in *top 5*) is not reg-
1046 istered. Consequently, the borrowing *top* will be
1047 considered assimilated when referring to the piece

¹²<https://dle.rae.es/>

1048	of clothing, but unassimilated when used to talk	B.4.4 Number inflection	1091
1049	about the best elements of a ranking or the upper	Unassimilated borrowings may be incorporated as	1092
1050	part of something.	invariable in number <i>los master</i> , with the same	1093
1051	un top estampado	plural inflection that they had in the donor lan-	1094
1052	el [top] cinco de artistas	guage (<i>los pappardelle</i>) or may form a new plural	1095
1053	(ENG)	that is non-existent in the donor language (<i>los pap-</i>	1096
1054	la [top] desfiló (ENG)	<i>pardelles</i>). For number inflection, we follow the	1097
1055	Similarly, the word <i>post</i> will not be considered	same criteria that DLE (Real Academia Española,	1098
1056	a borrowing when used as a prefix of Latin origin,	2020) follows: a non-Italian plural like <i>pizzas</i> is	1099
1057	but will be labeled with ENG when used to refer	still regarded as unadapted (and therefore should	1100
1058	to something that is published on a social media	be written italicized even when the true Italian plu-	1101
1059	platform.	ral would be <i>pizze</i>). Consequently, non assimilated	1102
1060	el mundo post pandemia	borrowings that have a non-cannonical plural in-	1103
1061	un [post] de Facebook (ENG)	flexion form will still be considered as an unassim-	1104
1062	Additionally, assimilated borrowings can still be	ilated borrowing and labeled as such.	1105
1063	part of new unassimilated borrowings, in which	una serie de animación de	1106
1064	case they will be labeled as such:	[mechas] (OTHER)	1107
1065	un [boys club] (ENG)	B.4.5 Pseudoanglicisms	1108
1066	B.4.3 Words derived from foreign lexemes	Words that do not exist in English (or exist with a	1109
1067	Words derived from foreign lexemes that do not	different meaning) but were coined following En-	1110
1068	comply with Spanish orthotactics but that have	glish morphological paradigm to imitate English	1111
1069	been morphologically derived following the Span-	words (such as <i>footing</i> or <i>balconing</i>) will be anno-	1112
1070	ish paradigm (such as <i>hacktivista</i> , <i>randomizar</i> ,	tated as anglicisms.	1113
1071	<i>shakespeariano</i>) will be considered assimilated and	la imagen del '[balconing]' y	1114
1072	should therefore not be labeled as a borrowing.	las excursiones etílicas (ENG)	1115
1073	Compound names where one of the lexemes is	practicaba [footing] por la	1116
1074	a borrowing will be labeled as a borrowing or not	calle (ENG)	1117
1075	according to the degree of independence among	B.4.6 Realia words	1118
1076	the lexemes. A verb+noun compound (as <i>caza-</i>	Borrowings that refer to culture-specific elements	1119
1077	<i>clicks</i>) will not be labeled as a borrowing, because	(often called <i>realia words</i>) that were imported long	1120
1078	the elements are not independent from one another.	ago but that have remained unadapted will not be	1121
1079	However, noun-noun compounds where each of the	tagged as borrowing. This means that if a borrow-	1122
1080	lexemes work can work independent from one can	ing is not adapted (i.e. its form remained exactly	1123
1081	be labeled as borrowings:	as it came from the donor language) but refers to a	1124
1082	una casa-[loft] (ENG)	particular cultural object that came via the original	1125
1083	Similarly, prefixed borrowings will be labeled	language, that has been registered for a while in	1126
1084	as a borrowing, as long as the borrowing keeps	Spanish dictionaries and is not perceived as new	1127
1085	independence from the prefix:	anymore, then it will not be tagged as a borrowing,	1128
1086	la ex [influencer] (ENG)	even if does not comply with graphophonologic	1129
1087	For prefixed borrowings, it should be checked	rules of Spanish.	1130
1088	whether the prefix can also be considered part of	The purpose of this guideline is to account for	1131
1089	the borrowing:	cultural terms such as <i>pizza</i> , <i>whisky</i> , <i>jazz</i> , <i>blues</i> ,	1132
1090	los [nano influencers] (ENG)	<i>banjo</i> or <i>sheriff</i> . These are all borrowings that are	1133
		reluctant to be adapted or translated, even when	1134
		they have been around in the Spanish language for	1135
		long. The reason is that they refer to cultural inven-	1136
		tions (the name was imported along with the object	1137

1138	it refers to), and, given their cultural significance,		
1139	they never competed with a Spanish equivalent and		
1140	are seen as assimilated.		
1141	Therefore, unadapted borrowings that refer to		
1142	cultural innovations (such as music, cooking, sport		
1143	names etc) and that have been registered for long		
1144	in the Spanish language ¹³ will not be tagged as		
1145	emergent borrowings.		
1146	It should be noted that this only applies to bor-		
1147	rowings that have been around enough time to be		
1148	registered in dictionaries. A word like <i>hip hop</i> is		
1149	a realia word, but it is still recent enough and has		
1150	not been registered in the dictionary. In that case, it		
1151	should be considered as unassimilated and tagged		
1152	as such.		
1153	B.4.7 Latinisms		
1154	Borrowings that were introduced directly from		
1155	Latin language (such as <i>deficit</i> , <i>curriculum</i> , etc)		
1156	will not be considered emergent and therefore will		
1157	not be tagged as a borrowing. However, it should		
1158	be noted that unassimilated borrowings from other		
1159	languages that happen to have a Latin etymology		
1160	and and that are introduced with a distinct meaning		
1161	(such as <i>adlib</i> or <i>premium</i> etc) will still be tagged		
1162	as borrowings.		
1163	B.5 Borrowings vs names		
1164	B.5.1 Proper nouns		
1165	Non-Spanish proper nouns will not be tagged as		
1166	borrowings. These include:		
1167	• person names: <i>Bernie Sanders</i> .		
1168	• organization names: <i>WikiLeaks</i> .		
1169	• product names: <i>Slack</i> .		
1170	• location names: <i>Times Square</i> .		
1171	• dates and celebrations: <i>St. Patrick's Day</i> ,		
1172	<i>Black Friday</i> .		
1173	• event names: <i>Brexit</i> , <i>procés</i> .		
1174	• social and political movements: <i>Black Lives</i>		
1175	<i>Matter</i> , <i>MeToo</i> .		
1176	• treaties and documents: <i>New Deal</i> , <i>Privacy</i>		
1177	<i>Shield</i> , <i>French Tech Visa</i> .		
	¹³ RAE dictionary https://dle.rae.es/ ,		
	Mapa de diccionarios https://webfrrl.rae.es/		
	and CREA http://corpus.rae.es/creanet.html and COR-		
	PES https://webfrrl.rae.es/CORPES/view/		
	inicioExterno.view can be consulted		
	• titles of cultural productions: <i>Stranger Things</i> .		1178
	B.5.2 Borrowings in proper nouns		1179
	Borrowings that appear as part of proper nouns or		1180
	named entities (such as book titles or organization		1181
	names, as in <i>Los Hermanos Podcast</i>) will not be		1182
	labeled as borrowings.		1183
	B.5.3 Proper nouns in borrowings		1184
	Multiword borrowings and expressions can some-		1185
	times include proper nouns. Even when a proper		1186
	noun in isolation cannot be considered a borrowing,		1187
	proper nouns within a borrowed expression will be		1188
	considered part of the borrowing, as long as the		1189
	proper noun is part of the borrowing and is used		1190
	following the grammar rules of the donor language		1191
	(for example, in an English noun noun compound):		1192
	Tecnología [made in Spain]		1193
	(ENG)		1194
	[Google cooking] (ENG)		1195
	B.5.4 Names of institutions and political roles		1196
	Non-Spanish names that refer to political institu-		1197
	tions (such as <i>Parlament</i> or <i>Bundestag</i>) or to politi-		1198
	cal roles and figures (<i>lehendakari</i> , <i>president</i> , <i>con-</i>		1199
	<i>seller</i>) will be excluded and will not be tagged as		1200
	borrowings.		1201
	B.5.5 Words derived from proper nouns		1202
	Words derived from proper nouns (via metonymy		1203
	or eponymy) will not be tagged as a borrowing, as		1204
	long as the relation with the proper noun they come		1205
	from is transparent to the speaker such as:		1206
	• products: <i>un iPhone</i> , <i>un whatsapp</i> , <i>un bizum</i> ,		1207
	<i>un Scalextric</i> , <i>el Satisfyer</i> .		1208
	• works of arts: <i>un monet</i>		1209
	• characters: <i>un frankenstein</i> .		1210
	However, borrowings that originated from a		1211
	proper noun in the donor language but entered the		1212
	Spanish language as common nouns and are cur-		1213
	rently recognized as such, will be labeled as bor-		1214
	rowings. In order to adjudicate which of these		1215
	words are still used in Spanish as proper names and		1216
	which are common nouns, dictionaries and other		1217
	reference works can be consulted.		1218

1219	B.5.6 Names of peoples or languages	acronyms will not be tagged as a borrowing, even if the acronym is of non-Spanish origin	1265
1220	Names of peoples or languages (such as <i>inuit</i>) will not be labeled as borrowings, even if the word is borrowed from another language and is not registered in Spanish dictionaries.		1266
1221		un lector de CD	1267
1222		An acronym however may be tagged as a borrowing if it appears as part of a borrowed multiword expression, as in <i>CD player</i> , <i>peak TV</i> , <i>PC gaming</i> :	1268
1223		un [CD player] (ENG)	1269
1224	B.5.7 Fictitious creatures	Acronym expansions, that is, the expansion of an acronym into the words that form the acronym (that is usually added in between brackets after an acronym has been introduced) will also not be considered a borrowing:	1270
1225	Unadapted names of fictitious creatures (such as <i>hobbit</i> or <i>troll</i>) will be labeled as a borrowing.		1271
1226			1272
1227	En un agujero en el suelo		1273
1228	vivía un [hobbit] (ENG)		1274
1229	B.5.8 Scientific units		1275
1230	Unadapted borrowings that refer to widespread scientific units (such as <i>hertz</i> , <i>newton</i> , <i>byte</i> , etc) will be considered assimilated and should not be tagged as a borrowing	La técnica de PCR (protein chain reaction)	1276
1231			1277
1232		It is important to note that for a sequence to be considered as an acronym expansion it must appear after the acronym has been introduced and serve as a gloss to it (so that it expands what the letters in the acronym stand for). Usages where the full sequence is introduced in the text and later on acronymized for the sake of brevity can still be considered as borrowings.	1278
1233			1279
1234	B.5.9 Names of species		1280
1235	Scientific names of a species (such as Latin names) will not be tagged as a lexical borrowing (<i>anisakis</i>). Names of fruit, vegetable and plant varieties (such as <i>manzana golden</i> , <i>patatas Kennebec</i> or <i>aguacate Hass</i>) will also be excluded.	Utilizaron técnicas de [Machine Learning] (también conocido como ML) (ENG)	1281
1236			1282
1237			1283
1238			1284
1239			1285
1240	B.6 Borrowings vs other code-mixed inclusions		1286
1241			1287
1242	Borrowing (using units from one language in another language) and code-switching (intertwining segments of different languages in the same discourse) have frequently been described as a continuum (Clyne et al., 2003), with a fuzzy frontier between the two. As a result, it can be difficult to tell the difference between borrowing and other code-mixed inclusions. The following guidelines can assist annotators adjudicate edge cases.		1288
1243			1289
1244		B.6.2 Digits	1290
1245		Similarly to proper nouns, digits in isolation cannot be considered borrowings. As a result, we cannot take for granted that digits within the surroundings of a borrowing will automatically be part of the borrowing.	1291
1246			1292
1247		[top ten] (ENG)	1293
1248		[top] 10 (ENG)	1294
1249		However, if the word order of the tokens makes it clear that the digit is part of a multitoken borrowing (because the order complies with the grammatical structure of an English noun-noun compound), we can label it as part of the borrowing:	1295
1250			1296
1251	When in doubt while dealing with code-mixed inclusion, the annotator may find it helpful to ask the following question as a rule of thumb: would it make sense to have this non-Spanish word registered in a dictionary of Spanish? If the answer is no (for instance, because the word reflects the literal quotation of what someone said or because the inclusions is metalinguistic usage rather than borrowing), then we are probably not in front of a borrowing but of another type of code-mixed inclusion (and should not be tagged as a borrowing).		1297
1252			1298
1253			1299
1254			1300
1255			1301
1256		los [10% banks] (ENG)	1302
1257			1303
1258		B.6.3 Metalinguistic usage	1304
1259		Non-Spanish words that appear to refer to the word itself in linguistic discourse and do not cover a lexical gap will not be tagged as a borrowing:	1305
1260			1306
1261		El término viene de la palabra 'ghost', 'fantasma' en inglés	1307
1262	B.6.1 Acronyms and acronym expansions		1308
1263	We consider acronyms to be a different phenomenon from borrowings. Consequently,		1309
1264			

1310	It should be noted that the newer, less adapted,	may not be suitable if applied to a project with a dif-	1357
1311	less transparent a new word is, the more likely that	ferent scope. These are some of the shortcomings	1358
1312	the speaker will be aware of the decoding difficulty	and limitations that these guidelines may have.	1359
1313	it may pose to the reader and will decide to add		
1314	some sort of metalinguistic strategy or awareness	B.7.1 Text genre	1360
1315	around it, in the form of metacomments, word-	These guidelines were designed to specifically cap-	1361
1316	pointers, meaning explanations, etc (<i>known as, so</i>	ture borrowings in a corpus of Spanish newswire.	1362
1317	<i>called</i>). Borrowings with these types of signals	Newswire is a very specific genre of text that by	1363
1318	will still be considered borrowings, as long as they	no means represent the whole of a language (Plank,	1364
1319	are covering a lexical gap.	2016).	1365
1320	True metalinguistic usage where the foreign	B.7.2 Donor language	1366
1321	word covers no lexical gap but exclusively pro-	These guidelines were created with English lexical	1367
1322	vides linguistic information (such as etymological	borrowings in mind, which are the most frequent	1368
1323	information) will not be considered a borrowing.	source of borrowing today in the Spanish press.	1369
1324	B.6.4 Literal quotations	Although the criteria can be applied to other lan-	1370
1325	Words or sequences in languages other than Span-	guages as well (and in fact the annotation tagset we	1371
1326	ish that are reflecting literally what someone said	propose includes the tag OTHER to account for bor-	1372
1327	or wrote (as in a quotation, a statement or a slogan)	rowings from other languages other than English),	1373
1328	will not be considered a borrowing.	a more fine-grained approach would require further	1374
1329	El eslogan 'Make America Great	guidelines.	1375
1330	Again'	B.7.3 Synchronic approach to borrowing	1376
1331	Es uno de los primeros	This project approaches emergent, unassimilated	1377
1332	resultados de Google cuando	lexical borrowing in a synchronic fashion. The pro-	1378
1333	alguien busca "remote work in	cess of borrowing and the notion of assimilation is,	1379
1334	Spain" (trabajo en remoto en	however, time-dependent. A diachronic approach	1380
1335	España).	to lexical borrowing would require a wider scope,	1381
1336	B.6.5 Expressions	a different theoretical framework and an expanded	1382
1337	In general terms, multiword borrowings will be	set of criteria.	1383
1338	tagged as borrowings. However, phrases and ex-	B.7.4 Geographic variety	1384
1339	pressions that are not integrated into the sentence	The guidelines in this document were designed to	1385
1340	will be excluded. This means that autonomous ex-	capture borrowings used in Spanish newspapers,	1386
1341	pressions that are rather code switched sentences	that is, written in the variety of Spanish that is spo-	1387
1342	(rather than real borrowings) that work as a unit	ken in Spain and may not be suitable to account	1388
1343	totally independently of the rest of the linguistic	for other dialects. For instance, according to the	1389
1344	context (and that we would not expect to be reg-	guidelines we have just introduced, a word like <i>liv-</i>	1390
1345	istered in a dictionary) will not be considered or	<i>ing</i> (that is used heavily in some Latin American	1391
1346	tagged as a borrowing.	varieties to refer to the living room) would be con-	1392
1347	La innovación y la competencia	sidered unassimilated. It is arguable whether these	1393
1348	tan escasas en la radiotelevisión	criteria would be suitable for a project that tried to	1394
1349	o peor aún en Internet ("the	capture emergent lexical borrowings in Argentinian	1395
1350	winners takes all" o "most").	text, for example.	1396
1351	B.7 Limitations of these guidelines		
1352	These guidelines are intended to assist annotators		
1353	when labeling lexical borrowings. These guide-		
1354	lines, however, were created with a specific goal in		
1355	mind (to capture unassimilated English lexical bor-		
1356	rowings from a corpus of Spanish newswire) and		