

False Friends or Cognates? A Cross-lingual Semantic Ambiguity Evaluation for Galician, Portuguese and Spanish

Anonymous ACL submission

Abstract

The linguistic proximity between Galician, Portuguese, and Spanish results in a lexical overlap that often conceals semantic interference. This is particularly evident in false friends, posing a challenge for NLP systems. In this work, we assess whether state-of-the-art language models can identify and process false friends among these languages. We introduce six cross-lingual datasets –created using semi-automatic methods or manual construction and all carefully verified– covering cognates and false friends. We evaluate a broad range of encoder and decoder models of varying sizes via zero-shot and few-shot settings. Our results show that closed-weight models achieve the highest accuracy and medium-weight models demonstrate a strong balance of efficiency, while smaller models struggle with systematic issues and biases. Notably, we find that linguistic proximity itself introduces errors: closely related language pairs tend to perform worse, reflecting the challenge of semantic discrimination due to lexical overlap.

1 Introduction

Closely related languages exhibit a high degree of lexical and orthographic similarity, which can facilitate cross-lingual understanding but also give rise to systematic semantic ambiguity (Kallini et al., 2025). This issue is especially salient for closely related Romance languages such as Spanish, Portuguese, and Galician, which share a substantial portion of their vocabulary (Garcia et al., 2018), increasing the risk of semantic interference.

In bilingual and multilingual contexts, formally similar words are categorized as cognates: words across different languages that share a common etymological origin and maintain the same meaning, for example the Galician and Portuguese form *ponte* with the meaning of ‘bridge’. However, cognates may diverge semantically over time;

such cases are commonly referred to as false cognates or false friends. These are pairs that, despite their shared origin, have developed distinct meanings (Dominguez and Nerlich, 2002; Allan, 2009; Chamizo-Domínguez, 2012), often leading to cross-lingual interference. For instance, the Portuguese adjective *esquisito* means ‘strange’ or ‘odd’, whereas the Spanish adjective *exquisito* denotes something ‘refined’ or ‘excellent’. As similar words are often assumed to share meaning, these phenomena are problematic both for second-language learners and bilingual speakers (Durán Escribano, 2004; Brenders et al., 2011), and also for computational language models (Limisiewicz et al., 2023; Cahyawijaya et al., 2025).

The proliferation of multilingual Large Language Models (LLMs) raises the question of whether these models can handle cross-lingual semantic ambiguity. In this context, false friends and cognates constitute a valuable diagnostic of lexical processing. Analyzing whether language models can identify false friends, and determining which contexts and cases are more challenging, is crucial for assessing their cross-lingual capabilities.

Our study investigates the ability of state-of-the-art models to detect and process false friends and cognates in closely related varieties (Galician, Portuguese, and Spanish), which coexist in regions such as the Iberian Peninsula and Latin America, and are frequently considered jointly in evaluation efforts and multilingual modeling (Gonzalez-Agirre et al., 2025; Ángel González et al., 2026). To this end, we introduce new cross-lingual datasets of cognates and false friends in context for these three languages, combining automatically extracted sentences with manually created examples. We evaluate these cross-lingual phenomena using encoder and decoder models. First, we establish an unsupervised baseline and train a logistic regression classifier to assess whether cognates and false friends can be distinguished

083 from contextualized embeddings. We then evaluate
084 current LLMs on a contextual semantic judgment
085 task, where models must decide whether two target
086 words are semantically equivalent in context, evalu-
087 ating both zero-shot and few-shot settings. Our
088 results, complemented by qualitative analysis and
089 a preliminary evaluation in machine translation,
090 show that although closed-weight LLMs demon-
091 strate strong disambiguation abilities, cross-lingual
092 semantic ambiguity remains challenging for most
093 models, particularly in closely related language
094 pairs. All new resources will be freely released.

095 2 Related work

096 Beyond traditional Word Sense Disambiguation,
097 which identifies the intended meaning of a word in
098 context typically by linking it to predefined sense
099 inventories (Navigli, 2009; Bevilacqua et al., 2021),
100 Word-in-Context (WiC) has become increasingly
101 popular for evaluating contextualized word mean-
102 ing in an inventory-free manner. WiC (Pilehvar and
103 Camacho-Collados, 2019) has reformulated lexi-
104 cal semantic evaluation as a binary classification
105 task, determining whether a target word is used
106 with the same meaning in two different contexts.
107 This approach avoids reliance on explicit sense in-
108 ventories and has been extended to multilingual
109 and cross-lingual settings, including XL-WiC (Ra-
110 ganato et al., 2020) and MCL-WiC (Martelli et al.,
111 2021), making it particularly suitable for analyzing
112 multilingual LLMs.

113 In this regard, while multilingual LLMs pro-
114 vide coverage across dozens of languages, their
115 performance often degrades in non-English con-
116 texts (Zhang et al., 2023), which can be attributed
117 in large part to the predominance of English data
118 in their training corpora (Xu et al., 2025). In cross-
119 lingual scenarios, lexical overlap often introduces
120 semantic interference in closely related languages,
121 where orthographic similarity (categorized as cog-
122 nates or false friends) may mask divergent senses.

123 Prior works in this area have focused on the au-
124 tomatic identification of cognates and false friends
125 at the word level. Early approaches addressed the
126 identification of cognates and false friends using
127 bilingual corpora (Mitkov et al., 2007). Subsequent
128 studies have relied on word embeddings and vector-
129 space models for the automatic detection of false
130 friends (Torres and Aluísio, 2011; Ljubešić and
131 Fišer, 2013; Castro et al., 2018; Perrián-Pascual
132 and Fernández Martínez, 2025), as well as for cog-

133 nate detection (Batsuren et al., 2019; Akavarapu
134 and Bhattacharya, 2024). Alternative approaches
135 have treated false friend identification as a classifi-
136 cation task (Nikov et al., 2024). Other studies have
137 proposed unsupervised methods for constructing
138 multilingual lexicons of false friends across multi-
139 ple languages Uban and Dinu (2020), or conduct
140 linguistic analyses of semantic false friends and
141 borrowings accompanied by automatic correction
142 methods (Uban et al., 2025). In addition, it has
143 been also examined the effect of vocabulary over-
144 lap in multilingual LLMs (Kallini et al., 2025). De-
145 spite these advancements, recent studies have high-
146 lighted limitations in the contextual disambigua-
147 tion capabilities of multilingual LLMs. Specifi-
148 cally, Cahyawijaya et al. (2025) demonstrates that
149 state-of-the-art models continue to struggle with
150 false friends, where shared orthography across lan-
151 guages triggers incorrect semantic transfer despite
152 conflicting contextual cues.

153 Some studies on false friends and cognates
154 have focused on descriptive analyses or vocabu-
155 lary lists between Galician, Spanish, and Por-
156 tuguese (Bragado Trigo, 2006; Figueroa, 1995),
157 with small initial WiC datasets available for Gali-
158 cian and Spanish (Abuín and Garcia, 2025). For
159 the Spanish–Portuguese language pair, recent com-
160 putational approaches have been proposed (Castro
161 et al., 2018; Uban et al., 2025). Including Gali-
162 cian constitutes a compelling challenge, as it is
163 closely related to Portuguese while strongly influ-
164 enced by Spanish. In this work, we introduce new
165 WiC-like multilingual datasets of false friends and
166 cognates and present a systematic evaluation of
167 current LLMs using a range of methods.

168 3 Dataset construction

169 To evaluate the cross-lingual abilities of LLMs on
170 lexical semantics, we build a multilingual WiC-
171 style dataset including cognates and false friends.
172 Each instance consists of a word pair in two lan-
173 guages and two contextualized instances, illustrat-
174 ing the sense of the target words (see Table 1).

175 **Languages:** Galician, Portuguese and Span-
176 ish constitute a closely related group of Ibero-
177 Romance languages. While Galician and Por-
178 tuguese have historically been considered varieties
179 of the same language, Galician has been influenced
180 by Spanish due to its status as a lower-prestige
181 language in Spain and an official orthography mod-
182 eled after Spanish conventions (Samartim, 2012).

Dataset	W1	S1	W2	S2	Label	POS
ES-PT	TALHER	Deve ser comido apenas com as mãos, jamais com talheres . <i>It should only be eaten with your hands, never with cutlery.</i>	TALLER	También tiene un moderno taller de restauración. <i>It also has a modern restoration workshop.</i>	False Friend	N
GL-PT	VOTO	Só se permite un voto por familia. <i>Only one vote is allowed per family.</i>	VOTO	Havia apenas dois votos a favor da proposta. <i>There were only two votes in favor of the proposal.</i>	Cognate	N

Table 1: Example instances from our dataset. Columns report the language pair (Dataset), the target words (W_1/W_2) together with their contextual sentences (S_1/S_2), the semantic relation label (False Friend or Cognate), and the Part-of-Speech (POS).

183 Although Portuguese and Spanish also share simi- 221
184 larities, they exhibit more pronounced orthographic 222
185 differences. These asymmetric statuses, combined 223
186 with the fact that in some cases the languages co- 224
187 exist in overlapping territories and are included in 225
188 multilingual LLMs and related initiatives, make 226
189 this triad a compelling testbed for assessing cross- 227
190 lingual semantic disambiguation. 228

191 **Semantic phenomena:** The datasets incorporate 230
192 both *cognates* and *false friends*, further categorized 231
193 into *total*, whose meanings have diverged across 232
194 languages (e.g., *carpeta*, ES ‘folder’ vs. *carpete*, 233
195 PT ‘carpet’), and *partial* false friends, where at 234
196 least one sense differs while others may still over- 235
197 lap (e.g., *balón*, GL ‘ball’ vs. *balão*, PT ‘balloon’ 236
198 or ‘ball’). The datasets include strict homographs 237
199 (e.g., *propina*, meaning *tip* in ES vs. *fee* in PT) 238
200 as well as similar forms resulting from language- 239
201 specific evolution. Specifically, we account for sys- 240
202 tematic morphological variations (e.g., PT *trabalho* 241
203 vs. ES *trabajo*) and phonological divergences, such 242
204 as the preservation of diphthongs in Galician and 243
205 Portuguese in contrast to the Spanish simplifica- 244
206 tion (e.g., GL/PT *touro* vs. ES *toro*). These forms 245
207 capture the subtle lexical nuances and inter-lingual 246
208 overlaps that pose significant challenges for word 247
209 sense disambiguation in multilingual models. 248

210 **Data compilation:** The datasets were con- 249
211 structed using two complementary resources, one 250
212 targeting false friends and the other cognates. False 251
213 friends pairs were manually compiled from lan- 252
214 guage learning materials and online pedagogical re- 253
215 sources for the relevant language pairs (a complete 254
216 list of resources is provided in Appendix B). These 255
217 items correspond to words that are known to be 256
218 problematic for second-language learners or bilin- 257
219 guals due to cross-linguistic similarity. Each pair 258
220 was then annotated as either a total or partial false

friend depending on the degree of semantic overlap. 221
Regarding cognate pairs, they were automatically 222
extracted from WordNet (Miller et al., 1990) by 223
identifying synsets shared between the two lan- 224
guages using the Multilingual Central Repository 225
(MCR) 3.0 (Gonzalez-Agirre et al., 2012). In order 226
to maintain balance between categories, an equiva- 227
lent number of cognate pairs was selected in each 228
respective language pair. To obtain example sen- 229
tences, we first obtained 2024 Wikipedia dumps 230
for each language, which were tokenized and lem- 231
matized using FreeLing (Padró and Stanilovsky, 232
2012). We then extracted up to ten contextualized 233
sentences per word for each language, and a na- 234
tive or near-native speaker of the three languages 235
with a background in Linguistics selected the most 236
representative examples. 237

238 **Annotation task:** The resulting sentences were 239
239 validated by two annotators with a linguistic back- 240
ground. The task consisted of verifying (i) whether 241
the target word in each sentence correctly reflected 242
the assigned sense, and (ii) whether the seman- 243
tic relation remained valid in context. Disagree- 244
ments were discussed, and when no consensus was 245
reached, the pairs were discarded to maintain high- 246
quality and balanced datasets.

247 **Human-authored subset:** In addition to the 248
248 semi-automatically constructed data, we derived 249
a human-authored subset from the validated in- 250
stances. For each language pair, we randomly 251
sampled between 50 and 60 word pairs and asked 252
annotators to produce new sentences containing 253
the same target words and preserving the original 254
senses. These newly created sentences underwent 255
the same validation procedure as described above. 256

257 **Dataset statistics:** Table 2 presents an overview 258
of the dataset per language pair. In total, our 259
data contain 704 validated cross-lingual word pairs,

each one associated with two sentences, thus totaling 1,408 examples. Each set is balanced in number of cognates and false friends, the latter divided into total and partial, as detailed in Appendix 8¹.

Pair	Cogn.	FFs	Total	Sents.
ES-PT	182	182	364	728
GL-ES	75	75	150	300
GL-PT	95	95	190	380
Total	357	367	704	1,408

Table 2: Distribution of cognates (Cogn.) and false friends (FFs) in the dataset (semi-automatic and human-created subsets), and number of sentences (Sents).

4 Experimental Setup

Our experiments follow two complementary paradigms: (i) WiC-based standard methods on encoder models, and (ii) prompt-based approaches using LLMs, in zero-shot and few-shot settings.

4.1 Encoder experiments

Unsupervised Baseline: We implement a baseline based on the cosine similarity between the contextualized embeddings of the target word pair. Following a binary classification approach, a pair is labeled as *same sense* if its similarity exceeds a given threshold, and as *different sense* otherwise. We explore thresholds from 0.00 to 1.00 in steps of 0.02 across all layers, and select the layer-threshold combination that yields the highest accuracy.

Logistic Regression: In line with the supervised framework described by Wang et al. (2019), we train logistic regression classifiers using the concatenated contextualized embeddings of the target word pair. To maintain consistency across experiments, we select the layer with the highest performance on the baseline experiment. We employ the minicons library (Misra, 2022)² to extract contextualized word representations from HuggingFace’s Transformers (Wolf et al., 2020). For evaluation on a specific language pair, the model is trained using the combined data from the remaining two pairs.

Models: We use the following multilingual encoders: mBERT (cased and uncased) (Devlin et al., 2019), XLM-RoBERTa (base and large) (Conneau

et al., 2020), and XL-Lexeme, an XLM-RoBERTa-large trained with a combination of various WiC-like datasets (Cassotti et al., 2023).

4.2 Prompt-based LLMs

We evaluate LLMs in zero-shot and few-shot ($k = 2$) configurations. Given target words W_1 and W_2 , and their corresponding sentences S_1 and S_2 (each from a different language), the model must perform a binary decision: YES if the target words are semantically equivalent, or NO if they refer to distinct senses. To enable automated parsing, we constrain the output format and require a brief justification for each decision (see Appendix E and F for the prompts used in both settings). The justification is not used for scoring, and it is included only to support qualitative error analyses. For the few-shot setting, we provide two balanced examples specific to each language pair: one representing a cognate (YES), and one a false friend (NO). These examples are excluded from the test set. Prompts remain fixed across all models to ensure a fair comparison.

Models: For the prompt-based evaluation, we consider a broad range of open and closed-weight multilingual LLMs. The closed-weight group includes Gemini 2.0 and 2.5 (Comanici et al., 2025), GPT-4o-mini and GPT-4.1 (OpenAI et al., 2024), and Grok-3 and Grok-4. Among the open-weight models, we selected Gemma-2 (Team et al., 2024) and Gemma-3 (Team et al., 2025), Qwen-3 (Yang et al., 2025), Phi-3 (Abdin et al., 2024), Phi-4 (Abdin et al., 2025), Deepseek-V3 (DeepSeek-AI et al., 2025), Llama 3 series (Grattafiori et al., 2024) and Mistral (Jiang et al., 2023, 2024). These models span various scales, from 3B to over 200B parameters. For clarity, Table 3 summarizes the evaluated models, grouped by availability, size, and family.

5 Results

This section presents the experimental results, using the average of each language pair based on both subsets of the data (semi-automatic and manual).³

Encoder models: Among the encoder models, XL-Lexeme—an XLM-RoBERTa-large adapted to WiC-like tasks—achieved the best performance in the baseline method (see Table 9 in Appendix C for the complete baseline results). XL-Lexeme

¹For these pairs, 534 contextual sentence pairs were obtained semi-automatically, while 170 were human-produced (see Table 8 for a more detailed explanation).

²<https://github.com/kanishkamisra/minicons>

³Overall, results tend to be slightly higher on the manual subset. Full results for both subsets are provided in the appendices.

Type	Size	Family	Models
Open-weight	Small [$<10B$]	Google Meta Microsoft Mistral AI	gemma-3-4b-it, gemma-2-9b-it llama-3.1-8b-instruct, llama-3.2-3b-instruct phi-3-mini-128k-instruct ministral-3b, mixtral-8x7b-instruct
	Medium [10B - 50B]	Google Microsoft Mistral AI Qwen	gemma-3-12b-it, gemma-3-27b-it phi-3-medium-128k-instruct, phi-4-reasoning-plus [†] , phi-4-multimodal-instruct mistral-small-24b-instruct-2501, mistral-small-3.2-24b-instruct qwen3-30b-a3b-instruct-2507
	Large [$>50B$]	DeepSeek Meta Qwen	deepseek-chat-v3.1 [†] llama-3.1-70b-instruct, llama-3.2-70b-instruct qwen3-next-80b-a3b-instruct, qwen3-235b-a22b-2507
Closed-weight		Google OpenAI xAI	gemini-2.5-flash, gemini-2.5-flash-lite, gemini-2.0-flash-001 gpt-4o-mini, gpt-4.1, gpt-4.1-mini grok-3-mini, grok-4-fast, grok-4

Table 3: Categorization of the evaluated LLMs by access type (Open vs. Closed), parameter scale (Small, $<10B$); Medium, 10B–50B; and Large, $>50B$), and architectural family. Parameter counts for closed-weight models are omitted due to proprietary restrictions. While the majority of our evaluated models follow the standard instruction-tuning paradigm, we include several reasoning models marked with [†].

embeddings were subsequently used to train logistic regression classifiers, whose best results are reported in the first rows of Table 5. Despite optimizing both the layer and decision threshold for each language pair, this method underperformed the baseline in all cases (except GL–PT, where it was one accuracy point higher), yielding an average accuracy of 70.16 versus the baseline’s 77.80.

Zero vs. few-shot: For the LLMs, overall, few-shot prompting yields a consistent, but modest, improvement over zero-shot performance across the majority of the datasets (see Table 4). Consequently, the remainder of our analysis focuses on the few-shot paradigm, as it represents the upper bound of the models’ capabilities in this task.

Dataset	Zero-shot	Few-shot
ES–PT	82.80	84.15
GL–ES	80.65	82.35
GL–PT	74.45	78.50
Overall	79.30	81.60

Table 4: Comparison of zero-shot and few-shot mean performance (%) across datasets.

Open vs. closed-weight performance: While open-weight models like Gemma-3-27B, Qwen3-235B, and Llama-3.3-70B show competitive results, closed-weight models obtain the best performance, in particular, the Grok family achieving over 90% of accuracy (Table 5).

Model scale performance: Most models under 10B parameters struggle significantly with the task, particularly the Llama-3.2-3b (52.20%). Nevertheless, it should be noted that baseline experiments outperform those obtained by models of this size. In the medium-sized category (10B–50B), performance stabilizes: Mistral-small-3.2-24b and Gemma-3-27b-it demonstrate high efficiency, with both exceeding the 80% accuracy. Notably, Gemma-3-27B surpasses much larger models, such as Llama-3.3-70B (85.70%) and DeepSeek-Chat-v3.1 (85.13%).

Language pair: Overall, linguistic distance tends to be inversely related to model performance: pairs with greater distance (e.g., ES–PT) are generally easier to resolve, while closely related pairs (e.g., GL–PT) tend to be more challenging, with ES–GL exhibiting intermediate performance. For instance, Gemini-2.0-flash exhibits a performance drop from 90.80% for the GL–ES pair to 81.75% in GL–PT. This degradation is expected, given the linguistic proximity between Galician and Portuguese, which increases the difficulty of identifying semantic divergences. Additionally, the substantially larger presence of Portuguese data in training corpora may introduce a bias toward Portuguese interpretations. Only the highest-performing closed models, such as Grok-4-fast, have demonstrated consistent efficiency on GL–PT (94.05%).

Analysis of classification bias: Beyond performance metrics, we analyze model behavior by mea-

Encoder models						
Method	Family	Model	ES-PT	GL-ES	GL-PT	AVG
Baseline	XLM (large)	XL-Lexeme	81.50	81.50	70.50	77.80
Log. Regression	XLM (large)	XL-Lexeme	75.50	63.50	71.50	70.16
Open-weight LLMs						
Size	Family	Model	ES-PT	GL-ES	GL-PT	AVG
<i>Small</i> [<10B]	Google	gemma-2-9b-it	72.35	69.00	68.85	70.00
		gemma-3-4b-it	73.55	78.00	71.35	74.30
	Meta	llama-3.1-8b-instruct	66.70	65.00	65.65	65.78
		llama-3.2-3b-instruct	52.10	50.00	54.50	52.20
	Microsoft	phi-3-mini-128k	68.50	66.00	72.90	70.15
	Mistral AI	ministral-3b	76.00	72.50	61.65	70.05
mixtral-8x7b		68.85	76.00	62.10	68.98	
<i>Medium</i> [10B–50B]	Google	gemma-3-12b-it	85.80	82.50	78.95	82.41
		gemma-3-27b-it	87.95	87.50	83.95	86.46
	Microsoft	phi-3-medium-128k	71.55	66.00	72.90	70.15
		phi-4-multimodal	86.80	86.00	82.50	85.10
		phi-4-reasoning-plus	84.30	83.50	79.50	82.43
	Qwen	qwen3-30b	89.40	81.50	78.95	83.28
		Mistral AI	mistral-small-24b	87.40	86.65	77.60
	mistral-small-3.2-24b		91.40	85.50	81.15	86.01
	<i>Large</i> [>50B]	Qwen	deepseek-chat-v3.1	91.05	83.50	80.55
qwen3-235b			90.40	90.00	89.05	89.81
qwen3-next-80b			87.60	86.50	86.60	86.90
Meta		llama3-3.3-70b-instruct	91.75	84.50	80.85	85.70
		Closed-weight LLMs				
	Google	gemini-2.0-flash	90.80	87.50	81.75	86.68
		gemini-2.5-flash	92.55	91.50	84.30	89.61
		gemini-2.5-flash-lite	82.90	91.50	78.55	84.28
	OpenAI	gpt-4o-mini	85.95	79.50	77.40	80.95
		gpt-4.1	88.40	93.00	87.87	89.55
		gpt-4.1-mini	91.25	89.00	81.75	87.33
	xAI	grok-3-mini	93.20	95.50	88.75	92.48
		grok-4	96.05	91.50	89.60	92.38
		grok-4-fast	95.40	96.00	94.05	95.15

Table 5: Accuracy (%) results across models and language pairs. LLM performance is based on few-shot prompting. The final column (AVG) denotes the macro-average accuracy across all datasets.

390 suring the deviation, defined as the absolute imbalance
391 between YES and NO in predictions. Results
392 in Table 12 (Appendix H, which reports models
393 with the highest deviation) show that performance
394 drops correlate with systematic bias ($r = -0.73$,
395 $\rho = -0.66$), rather than with random errors. This
396 pattern is particularly prevalent in models under
397 10B parameters. Specifically, Gemma-2-9B is
398 the only model exhibiting a consistent negative
399 bias, with a tendency to over-label pairs as false
400 friends. In contrast, models such as Llama-3.2-
401 3b and Mixtral-8x7b show a strong positive bias,
402 classifying the majority of instances as cognates.

403 6 Error Analysis

404 Following the general results from the previous
405 section, we perform a more detailed error analysis.

406 6.1 Error types

407 While the previous results included all 29 mod-
408 els, we now focus on a representative subset of
409 five architectures for a more detailed analysis:
410 Grok-4 (as state-of-the-art model), Qwen-3-235B
411 (high-capacity open-weight), Gemma-3-27B (high-
412 efficiency), Llama-3.1-8B and Gemma-2-9B (rep-
413 resentative of smaller open models). This selection
414 allows us to determine how different models handle
415 the linguistic phenomenon in our datasets.

416 **Cognates and False Friends:** As shown in the
417 left section of Table 6, Gemma-2-9B emerges as
418 an outlier, exhibiting a disproportionately high er-
419 ror rate for cognates. This pattern confirms a sys-
420 tematic bias toward over-predicting false friends,
421 causing the model to fail at recognizing semantic
422 equivalence.

423 Overall, the results indicate that performance

Pair	Model	Error type (%)			POS errors (%)		
		COGN	PFFs	TFFs	A	N	V
ES-PT	llama-3.1-8b	10.99	65.71	46.75	23.33	43.22	13.95
	gemma-2-9b-it	45.60	8.57	1.30	50.00	16.58	37.21
	gemma-3-27b-it	6.59	25.71	16.88	12.22	12.56	25.58
	qwen3-235b	8.79	22.86	2.60	10.00	11.56	16.28
	grok-4	6.04	8.57	0.00	4.44	6.53	6.98
	AVERAGE	15.60	26.28	13.50	19.99	18.09	20.00
GL-ES	llama-3.1-8b	18.95	52.38	53.13	28.57	33.71	43.90
	gemma-2-9b-it	43.16	3.17	12.50	50.00	34.83	21.95
	gemma-3-27b-it	11.58	7.94	9.38	21.43	8.99	17.07
	qwen3-235b	2.11	11.11	21.88	7.14	10.11	14.63
	grok-4	3.16	3.17	3.13	7.14	3.37	4.88
	AVERAGE	15.79	15.55	20.00	22.85	18.20	20.48
GL-PT	llama-3.1-8b	18.95	47.62	56.25	12.20	38.61	42.31
	gemma-2-9b-it	54.74	4.76	15.63	48.78	30.69	34.62
	gemma-3-27b-it	13.68	14.29	25.00	19.51	17.82	15.38
	qwen3-235b	9.47	9.52	31.25	12.20	17.82	7.69
	grok-4	5.26	6.35	9.38	12.20	6.93	0.00
	AVERAGE	20.42	16.50	27.50	20.97	22.37	20.00

Table 6: Error analysis by error type (cognates, partial false friends, total false friends) and Part-of-Speech (adjectives, nouns, verbs).

varies across language pairs. For the ES-PT pair, partial false friends constitute the most challenging category. By contrast, for the pairs including Galician (GL-ES and GL-PT), total false friends dominate the error profile. This suggest that, in the presence of strong similarity, models have greater difficulty detecting diverging meanings. The full distribution per word type is detailed in Table 8 in Appendix A.

Part-of-Speech: A POS analysis was conducted to identify grammatical categories more prone to misclassifications, evaluated using category-specific error rates. Due to their dominance (91.47% of all instances, full distribution in Table 13, Appendix I), we focus on nouns (N), adjectives (A), and verbs (V). As shown in Table 6, no overall trend is observed across these categories, suggesting that performance is driven by model capacity and language pair, rather than by the specific POS.

6.2 Qualitative analysis

We complement the quantitative results with a qualitative analysis of individual model errors.

6.2.1 Most frequent misclassifications

To further investigate the qualitative nature of errors, we focus on word pairs most frequently misclassified across models, selecting those incorrectly classified by at least half of them. Table 7 reports

the most problematic cases for each language pair. Overall, the analysis reveals that both cognates and false friends are major sources of error. For instance, models often treat *marmelada* (PT ‘marmalade’) and *mermelada* (ES ‘jam’) as equivalent, overlooking subtle lexical differences. Similarly, in the GL-ES pair, *almorzar* causes confusion due to mismatched meanings (GL ‘breakfast’ vs. ES ‘lunch’). In GL-PT, *bacharelato* is misinterpreted, as it denotes different educational levels in each language. Many cases involve partial false friends, further complicating the task for the models.

Word-Pairs	Pair	Type	Err.
marmelada–marmelada	ES-PT	FF	27/29
individualizar–individualizar	ES-PT	COG	27/29
grasa–graxa	ES-PT	FF	24/29
presentar–apresentar	ES-PT	COG	23/29
batata–batata	ES-PT	FF	21/29
almorzar–almorzar	GL-ES	FF	25/29
almorzo–almuerzo	GL-ES	FF	23/29
costume–costumbre	GL-ES	COG	21/29
reserva–reserva	GL-ES	COG	20/29
bacharelato–bacharelato	GL-PT	FF	29/29
motorista–motorista	GL-PT	FF	27/29
presente–presente	GL-PT	COG	27/29
copo–copo	GL-PT	COG	27/29
sobrenome–sobrenome	GL-PT	FF	27/29

Table 7: Word pairs with the highest consensus misclassifications across models. Err. column are proportions of models (out of 29) that failed on each item.

463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511

6.2.2 Model reasoning and justification

To provide a more comprehensive analysis, we further examine models’ justification and reasoning. Thus, we manually analyzed the rationales from the two smaller and generally underperforming models: Llama-3.1-8B and Gemma-2-9b.

Cross-lingual interference: Our analysis shows that a recurrent error arises from the models insufficient understanding of Galician, often confusing it with Spanish. For instance, on the mentioned *almorzar*, models justify their decision by stating that both sentences refer to lunch, projecting the Spanish meaning onto Galician term. Similar cross-lingual interference is observed in pairs such as *rato* (GL ‘mouse’ vs. ES ‘a while’) or *ano* (GL ‘year’ vs. ES ‘anus’), where explanations consistently reflect the Spanish interpretation also in Galician.

Hallucinations: In addition, both models exhibit hallucinations when lexical knowledge is lacking. For example, *toro* (GL ‘slice’) is incorrectly described as a type of fish, suggesting that models do not know the actual meaning in Galician and they attempt to infer it from context.

Overgeneralization: We also observe cases where models try to justify similarity. For *oficina* (PT ‘workshop’ vs. ES ‘office’), Llama-3.1-8B overgeneralized by claiming that both refer to ‘a workplace where administrative tasks are carried out’, ignoring the Portuguese sense. We also observed cases of semantic over-interpretation in Gemma-2-9B: For *húmido* (GL/PT ‘humid’), the model incorrectly classified it as a false friend, arguing that ‘the first refers to humid weather, while the second refers to a humid climate’. This finding suggests that, despite the model’s ability to capture the overarching concept, its inherent bias may lead to the differentiation of identical concepts.

6.3 Translation analysis

In order to further investigate the potential implications of mishandling false friends, we perform an experiment on a machine translation task. We evaluate false friends retention by translating sentences using a subset of LLMs (see Appendix K). To automate interference detection, we lemmatize the output with FreeLing and identify cases where the model incorrectly preserves the false friend lemma in the target language.

Our analysis shows that, while models generally handle frequent terms well, there are persistent

cases of false friend interference. For instance, in Galician–Spanish translation, the Galician verb *chocar* (GL ‘to brood’) is incorrectly translated following its Spanish false friend (ES ‘to collide’). Similarly, in PT–ES pairs, the term *salsa* (PT ‘parsley’) is mistranslated as *salsa* (ES ‘sauce’) instead of *perejil*. These errors, while infrequent, can hamper the reliability of machine translation systems when dealing with closely-related languages.

7 Conclusion and Future Work

This paper presented a comprehensive evaluation of how current LLMs handle cross-lingual semantic ambiguity in related linguistic varieties: Galician, Portuguese and Spanish. For this purpose, we introduced a novel, WiC-style manually validated dataset of cognates and false friends. Our experimental results led to several findings. First, while specialized multilingual encoders like XL-Lexeme provide a strong baseline, they are often surpassed by medium and large generative models. Specifically, the few-shot setting provides a consistent improvement over zero-shot and embeddings-based configurations. Closed-weight models (e.g., Grok-4-fast) achieve superior performance, but large and medium-sized open-weight models (e.g., Gemma-3-27B or Qwen-235B) show competitive results. However, we observe a significant performance drop in smaller models (<10B), suggesting that a minimum scale is required to resolve semantic interference in related varieties. The outcomes of a qualitative analysis show that greater linguistic distance between languages reduces semantic ambiguity, thereby facilitating disambiguation. Consequently, the GL–PT pair is the most challenging configuration across all experiments, indicating that languages with high lexical overlap pose difficulties for models both in distinguishing senses and in treating them as separate languages. Finally, false friends—both total and partial—remain a persistent challenge for current LLMs.

In future work, we plan to extend this study by expanding the datasets with more examples, richer contexts, and additional languages. Future analyses may also consider the role of lexical frequency and degrees of ambiguity (e.g., as defined in WordNet) in model performance. Finally, extending the assessment to additional downstream tasks (expanding the analysis on machine translation) would help assess the impact in real-world scenarios.

512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560

561 Limitations

562 **Data:** In our study there are some limitations
563 regarding the dataset. First, the use of sentences
564 from Wikipedia may not fully reflect language in
565 other domains. Second, while the datasets are bal-
566 anced, the overall number of instances remains
567 relatively small. Similarly, the fact that we deal
568 with three specific languages pairs may prevent us
569 from drawing more general conclusions about the
570 cross-lingual capabilities of the evaluated models.
571 Furthermore, the datasets are focused on nouns, ad-
572 jectives, and verbs, with a limited representation of
573 other lexical categories. Finally, there is an inher-
574 ent risk of data contamination, as the false-friends
575 sources and the contextual sentences might have
576 been included in the model’s pre-training corpora.

577 **Models:** In our evaluation we include closed-
578 weight models for which information regarding
579 the architecture and training data is not publicly
580 available. Additionally, our prompt-based exper-
581 iments are limited to specific zero and few-shot
582 configurations; alternative strategies could yield
583 different performance patterns.

584 **Justification analysis:** In this study, we ask mod-
585 els to provide a short justification for their predic-
586 tions. However, we do not conduct a quantitative
587 analysis of the logical consistency of these justi-
588 fications. For that reason, there is the possibility
589 of models hallucinating while still producing the
590 correct answer. A rigorous human evaluation is re-
591 quired to confirm whether models understand these
592 semantic distinctions.

593 References

594 Marah Abdin, Sahaj Agarwal, Ahmed Awadallah, Vid-
595 hisha Balachandran, Harkirat Behl, Lingjiao Chen,
596 Gustavo de Rosa, Suriya Gunasekar, Mojan Java-
597 heripi, Neel Joshi, Piero Kauffmann, Yash Lara,
598 Caio César Teodoro Mendes, Arindam Mitra, Be-
599 smira Nushi, Dimitris Papailiopoulos, Olli Saarikivi,
600 Shital Shah, Vaishnavi Shrivastava, and 4 others.
601 2025. [Phi-4-reasoning technical report](#). *Preprint*,
602 arXiv:2504.21318.

603 Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed
604 Awadallah, Ammar Ahmad Awan, Nguyen Bach,
605 Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat
606 Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck,
607 Sébastien Bubeck, Martin Cai, Qin Cai, Vishrav
608 Chaudhary, Dong Chen, Dongdong Chen, and 110
609 others. 2024. [Phi-3 technical report: A highly capa-
610 ble language model locally on your phone](#). *Preprint*,
611 arXiv:2404.14219.

Marta Vázquez Abuín and Marcos Garcia. 2025. [WiC
evaluation in Galician and Spanish: Effects of dataset
quality and composition](#). In *Proceedings of the
14th Joint Conference on Lexical and Computational
Semantics (*SEM 2025)*, pages 172–178, Suzhou,
China. Association for Computational Linguistics. 612
613
614
615
616
617

VSDS Mahesh Akavarapu and Arnab Bhattacharya. 618
2024. Automated cognate detection as a supervised 619
link prediction task with cognate transformer. In 620
*Proceedings of the 18th Conference of the European
Chapter of the Association for Computational Lin-
guistics (Volume 1: Long Papers)*, pages 965–975. 621
622
623

Keith Allan. 2009. *Concise encyclopedia of semantics*. 624
Elsevier. 625

Khuyagbaatar Batsuren, Gabor Bella, and Fausto 626
Giunchiglia. 2019. [CogNet: A large-scale cognate
database](#). In *Proceedings of the 57th Annual Meet-
ing of the Association for Computational Linguistics*,
pages 3136–3145, Florence, Italy. Association for 627
Computational Linguistics. 628
629
630

Michele Bevilacqua, Tommaso Pasini, Alessandro Ra- 631
ganato, and Roberto Navigli. 2021. [Recent trends
in word sense disambiguation: A survey](#). In *Pro-
ceedings of the Thirtieth International Joint Confer-
ence on Artificial Intelligence, IJCAI-21*, pages 4330–
4338. International Joint Conferences on Artificial 632
Intelligence Organization. Survey Track. 633
634
635
636
637
638

Iago Bragado Trigo. 2006. [Sobre a amizade \(léxica\)
galiza-portugal: os falsos amigos galego-portugués-
español](#). *Madrygal. Revista de Estudios Gallegos*,
9:33–42. 639
640
641
642

Pascal Brenders, Janet G. van Hell, and Ton Dijk- 643
stra. 2011. [Word recognition in child second
language learners: Evidence from cognates and false
friends](#). *Journal of Experimental Child Psychology*,
109(4):383–396. 644
645
646
647

Samuel Cahyawijaya, Ruochen Zhang, Jan Chris- 648
tian Blaise Cruz, Holy Lovenia, Elisa Gilbert, Hi- 649
roki Nomoto, and Alham Fikri Aji. 2025. [Thank
you, stingray: Multilingual large language models
can not \(yet\) disambiguate cross-lingual word senses](#).
In *Findings of the Association for Computational
Linguistics: NAACL 2025*, pages 3228–3250, Al- 650
buquerque, New Mexico. Association for Computa- 651
tional Linguistics. 652
653
654
655
656

Pierluigi Cassotti, Lucia Siciliani, Marco DeGemmis, 657
Giovanni Semeraro, and Pierpaolo Basile. 2023. [XL-
LEXEME: WiC pretrained model for cross-lingual
LEXical sEMantic change](#). In *Proceedings of the
61st Annual Meeting of the Association for Compu-
tational Linguistics (Volume 2: Short Papers)*, pages
1577–1585, Toronto, Canada. Association for Com- 658
putational Linguistics. 659
660
661
662
663
664

Santiago Castro, Jairo Bonanata, and Aiala Rosá. 2018. 665
[A high coverage method for automatic false Friends](#) 666

667	detection for Spanish and Portuguese. In <i>Proceedings of the Fifth Workshop on NLP for Similar Languages, Varieties and Dialects (VarDial 2018)</i> , pages 29–36, Santa Fe, New Mexico, USA. Association for Computational Linguistics.	Aitor Gonzalez-Agirre, Marc Pàmies, Joan Llop, Irene Baucells, Severino Da Dalt, Daniel Tamayo, José Javier Saiz, Ferran Espuña, Jaume Prats, Javier Aula-Blasco, Mario Mina, Adrián Rubio, Alexander Shvets, Anna Sallés, Iñaki Lacunza, Iñigo Pikabea, Jorge Palomar, Júlia Falcão, Lucía Tormo, and 4 others. 2025. <i>Salamandra technical report</i> . Preprint, arXiv:2502.08489.	720
668			721
669			722
670			723
671			724
672	Pedro J Chamizo-Domínguez. 2012. <i>Semantics and pragmatics of false friends</i> . Routledge.		725
673			726
674			727
675	Gheorghe Comanici, Eric Bieber, Mike Schaeckermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, and 1 others. 2025. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. <i>arXiv preprint arXiv:2507.06261</i> .	Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, and 542 others. 2024. <i>The llama 3 herd of models</i> . Preprint, arXiv:2407.21783.	728
676			729
677			730
678			731
679			732
680			733
681	Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. <i>Unsupervised cross-lingual representation learning at scale</i> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 8440–8451, Online. Association for Computational Linguistics.	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Léo Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. <i>Mistral 7b</i> . Preprint, arXiv:2310.06825.	734
682			735
683			736
684			737
685			738
686			739
687			740
688			741
689			742
690	DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, and 181 others. 2025. <i>Deepseek-v3 technical report</i> . Preprint, arXiv:2412.19437.	Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, Léo Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, and 7 others. 2024. <i>Mixtral of experts</i> . Preprint, arXiv:2401.04088.	743
691			744
692			745
693			746
694			747
695			748
696			749
697	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. <i>Bert: Pre-training of deep bidirectional transformers for language understanding</i> . Preprint, arXiv:1810.04805.	Julie Kallini, Dan Jurafsky, Christopher Potts, and Martijn Bartelds. 2025. <i>False Friends are not foes: Investigating vocabulary overlap in multilingual language models</i> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2025</i> , pages 21138–21154, Suzhou, China. Association for Computational Linguistics.	750
698			751
699			752
700			753
701	Pedro J Chamizo Dominguez and Brigitte Nerlich. 2002. False friends: their origin and semantics in some selected languages. <i>Journal of pragmatics</i> , 34(12):1833–1849.	Tomasz Limisiewicz, Jiří Balhar, and David Mareček. 2023. <i>Tokenization impacts multilingual language modeling: Assessing vocabulary allocation and overlap across languages</i> . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 5661–5681, Toronto, Canada. Association for Computational Linguistics.	754
702			755
703			756
704			757
705	María del Pilar Durán Escribano. 2004. Exploring cognition processes in second language acquisition: the case of cognates and false-friends in est. <i>Ibérica (Madrid)</i> , 7(1):87–106.	Nikola Ljubešić and Darja Fišer. 2013. Identifying false friends between closely related languages. In <i>Proceedings of the 4th Biennial International Workshop on Balto-Slavic Natural Language Processing</i> , pages 69–77.	758
706			759
707			760
708			761
709	Tiago Vidal Figueroa. 1995. Presuntos falsos amigos entre portugués e galego. I. <i>Viceversa. Revista galega de tradución</i> , pages 145–152.	Federico Martelli, Najla Kalach, Gabriele Tola, and Roberto Navigli. 2021. <i>SemEval-2021 task 2: Multilingual and cross-lingual word-in-context disambiguation (MCL-WiC)</i> . In <i>Proceedings of the 15th International Workshop on Semantic Evaluation (SemEval-2021)</i> , pages 24–36, Online. Association for Computational Linguistics.	762
710			763
711			764
712	Marcos Garcia, Carlos Gómez-Rodríguez, and Miguel A Alonso. 2018. New treebank or repurposed? on the feasibility of cross-lingual parsing of romance languages with universal dependencies. <i>Natural Language Engineering</i> , 24(1):91–122.		765
713			766
714			767
715			768
716			769
717	Aitor Gonzalez-Agirre, Egoitz Laparra, and German Rigau. 2012. Multilingual central repository version 3.0. In <i>LREC</i> , pages 2525–2529.		770
718			771
719			772

779	George A. Miller, Richard Beckwith, Christiane D. Fellbaum, Derek Gross, and Katherine J. Miller. 1990. Introduction to wordnet: An on-line lexical database . <i>International Journal of Lexicography</i> , 3:235–244.	
780		
781		
782		
783	Kanishka Misra. 2022. minicons: Enabling flexible behavioral and representational analyses of transformer language models. <i>arXiv preprint arXiv:2203.13112</i> .	
784		
785		
786	Ruslan Mitkov, Viktor Pekar, Dimitar Blagoev, and Andrea Mulloni. 2007. Methods for extracting and classifying pairs of cognates and false friends . <i>Machine Translation</i> , 21(1):29–53.	
787		
788		
789		
790	Roberto Navigli. 2009. Word sense disambiguation: A survey. <i>ACM computing surveys (CSUR)</i> , 41(2):1–69.	
791		
792		
793	Mitko Nikov, Žan Tomaž Šprajc, and Žan Bedrač. 2024. Cross-Lingual False Friend Classification via LLM-based Vector Embedding Analysis , page 33–36. Univerzitetna založba Univerze v Mariboru.	
794		
795		
796		
797	OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, and 262 others. 2024. Gpt-4 technical report . <i>Preprint</i> , arXiv:2303.08774.	
798		
799		
800		
801		
802		
803		
804		
805	Lluís Padró and Evgeny Stanilovsky. 2012. FreeLing 3.0: Towards wider multilinguality . In <i>Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC’12)</i> , pages 2473–2479, Istanbul, Turkey. European Language Resources Association (ELRA).	
806		
807		
808		
809		
810		
811	Carlos Perrián-Pascual and Nicolás José Fernández Martínez. 2025. Detección de cognados verdaderos y falsos amigos con word embeddings . <i>Revista Signos. Estudios de Lingüística</i> , 58(118).	
812		
813		
814		
815	Mohammad Taher Pilehvar and Jose Camacho-Collados. 2019. Wic: the word-in-context dataset for evaluating context-sensitive meaning representations. In <i>Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: Human language technologies, volume 1 (Long and short papers)</i> , pages 1267–1273.	
816		
817		
818		
819		
820		
821		
822	Alessandro Raganato, Tommaso Pasini, Jose Camacho-Collados, and Mohammad Taher Pilehvar. 2020. XL-WiC: A multilingual benchmark for evaluating semantic contextualization . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 7193–7206, Online. Association for Computational Linguistics.	
823		
824		
825		
826		
827		
828		
829	Roberto Samartim. 2012. Língua somos: A construção da ideia de língua e da identidade coletiva na galiza (pré-) constitucional. In <i>Novas achegas ao estudo da cultura galega II: enfoques socio-históricos e lingüístico-literarios</i> , pages 27–36. Universidade de Santiago de Compostela, Santiago de Compostela.	
830		
831		
832		
833		
834		
	Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, Louis Rouillard, Thomas Mesnard, Geoffrey Cideron, Jean bastien Grill, Sabela Ramos, Edouard Yvenc, Michelle Casbon, Etienne Pot, Ivo Penchev, and 197 others. 2025. Gemma 3 technical report . <i>Preprint</i> , arXiv:2503.19786.	835
		836
		837
		838
		839
		840
		841
		842
	Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, Johan Ferret, Peter Liu, Pouya Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos, Ravin Kumar, Charline Le Lan, Sammy Jerome, and 179 others. 2024. Gemma 2: Improving open language models at a practical size . <i>Preprint</i> , arXiv:2408.00118.	843
		844
		845
		846
		847
		848
		849
		850
		851
	Lianet Sepúlveda Torres and Sandra Aluísio. 2011. Using machine learning methods to avoid the pitfall of cognates and false friends in spanish-portuguese word pairs. In <i>Proceedings of the 8th Brazilian Symposium in Information and Human Language Technology</i> .	852
		853
		854
		855
		856
		857
	Ana Sabina Uban and Liviu P. Dinu. 2020. Automatically building a multilingual lexicon of false Friends with no supervision . In <i>Proceedings of the Twelfth Language Resources and Evaluation Conference</i> , pages 3001–3007, Marseille, France. European Language Resources Association.	858
		859
		860
		861
		862
		863
	Ana Sabina Uban, Liviu P Dinu, Ioan-Bogdan Iordache, Simona Georgescu, and Claudia Vlad. 2025. Friend or Foe? A Computational Investigation of Semantic False Friends across Romance Languages . In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 15309–15323, Suzhou, China. Association for Computational Linguistics.	864
		865
		866
		867
		868
		869
		870
		871
	Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019. Superglue: A stickier benchmark for general-purpose language understanding systems . <i>CoRR</i> , abs/1905.00537.	872
		873
		874
		875
		876
	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, and 3 others. 2020. Transformers: State-of-the-art natural language processing . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 38–45, Online. Association for Computational Linguistics.	877
		878
		879
		880
		881
		882
		883
		884
		885
		886
		887
	Yuemei Xu, Ling Hu, Jiayi Zhao, Zihan Qiu, Kexin Xu, Yuqi Ye, and Hanwen Gu. 2025. A survey on multilingual large language models: corpora, alignment, and bias . <i>Frontiers of Computer Science</i> , 19(11).	888
		889
		890
		891

892 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,
893 Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao,
894 Chengen Huang, Chenxu Lv, Chujie Zheng, Dayi-
895 heng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge,
896 Haoran Wei, Huan Lin, Jialong Tang, and 41 oth-
897 ers. 2025. Qwen3 technical report. *arXiv preprint*
898 *arXiv:2505.09388*.

899 Xiang Zhang, Senyu Li, Bradley Hauer, Ning Shi, and
900 Grzegorz Kondrak. 2023. [Don't trust ChatGPT when](#)
901 [your question is not in English: A study of multilin-](#)
902 [gual abilities and types of LLMs](#). In *Proceedings of*
903 *the 2023 Conference on Empirical Methods in Natu-*
904 *ral Language Processing*, pages 7915–7927, Singa-
905 pore. Association for Computational Linguistics.

906 José Ángel González, Ian Borrego Obrador, Álvaro
907 Romo Herrero, Areg Mikael Sarvazyan, Mara
908 China-Ríos, Angelo Basile, and Marc Franco-
909 Salvador. 2026. [IberBench: LLM evaluation on](#)
910 [Iberian languages](#). *Computer Speech & Language*,
911 96:101899.

A Complete Dataset Statistics

Dataset	Pair	COG	TFFs	PFFs	Tot.
Auto	ES-PT	152	66	87	304
Auto	GL-ES	50	30	20	100
Auto	GL-PT	65	22	43	130
Human	ES-PT	30	11	19	60
Human	GL-ES	25	16	9	50
Human	GL-PT	30	10	20	60

Table 8: Distribution of cognates (COG), total false friends (TFF), and partial false friends (PFF) for each language pair and data source (automatic or human created). The table also reports the total number of word pairs in each subset (Tot.).

B False Friends Resources

Online Resources:

- https://ec.europa.eu/translation/portuguese/magazine/documents/folha47_lista_pt.pdf
- <https://github.com/pln-fing-udelar/false-friends/tree/master>
- https://www.lingua.gal/c/document_library/get_file?file_path=/portal-lingua/curso/medio-administrativo/_13_falsosamigos.pdf

Academic Works and Books:

- Ramos, E. *Portugalizar. Portugués para galego-falantes.*
- Vidal Figueiroa, T. *Presuntos falsos amigos entre portugués e galego. (VOL. I-II-III)*

C Baseline experiments

Model	Pair	L.	Thr.	Acc.	F1
bert-base-multilingual-uncased					
	ES-PT	12	0.58	0.73	0.75
	GL-ES	12	0.62	0.79	0.76
	GL-PT	8	0.70	0.58	0.68
	H-ES-PT	12	0.56	0.73	0.77
	H-GL-ES	9	0.76	0.74	0.78
	H-GL-PT	0	0.76	0.57	0.58
bert-base-multilingual-cased					
	ES-PT	9	0.74	0.69	0.72
	GL-ES	8	0.62	0.71	0.76
	GL-PT	10	0.74	0.53	0.66
	H-ES-PT	8	0.68	0.72	0.76
	H-GL-ES	12	0.52	0.76	0.76
	H-GL-PT	7	0.70	0.60	0.71
pierluigi/xl-lexeme					
	ES-PT	13	0.82	0.81	0.83
	GL-ES	14	0.82	0.83	0.83
	GL-PT	16	0.80	0.68	0.72
	H-ES-PT	20	0.92	0.82	0.83
	H-GL-ES	15	0.82	0.80	0.82
	H-GL-PT	14	0.84	0.73	0.76
xlm-roberta-base					
	ES-PT	6	0.84	0.69	0.74
	GL-ES	3	0.88	0.70	0.69
	GL-PT	12	0.98	0.54	0.68
	H-ES-PT	6	0.86	0.68	0.72
	H-GL-ES	0	0.22	0.66	0.74
	H-GL-PT	5	0.80	0.60	0.71
xlm-roberta-large					
	ES-PT	14	0.86	0.73	0.74
	GL-ES	13	0.86	0.72	0.75
	GL-PT	13	0.84	0.63	0.71
	H-ES-PT	13	0.88	0.75	0.78
	H-GL-ES	11	0.88	0.68	0.75
	H-GL-PT	13	0.88	0.68	0.75

Table 9: Results for baseline experiments for all the models. We report the optimal transformer layer (*L*) and the decision threshold (*Thr.*), alongside classification Accuracy (*Acc.*) and F1-score for each language pairs. Bold values indicate the highest F1-score.

D Explained results

Encoder models										
Method	Family	Model	ES-PT	GL-ES	GL-PT	H-ES-PT	H-GL-ES	H-GL-PT	AVG	
Baseline	XLM (large)	XL-Lexeme	81.25	83.00	68.46	81.67	80.00	73.33	77.80	
Log. Regression	XLM (large)	XL-Lexeme	73.36	67.00	70.00	78.30	64.00	73.30	70.16	
Open-weight LLMs										
Size	Family	Model	ES-PT	GL-ES	GL-PT	H-ES-PT	H-GL-ES	H-GL-PT	AVG	
Small [<10B]	Google	gemma-2-9b-it	71.40	68.00	67.70	73.30	70.00	70.00	70.00	
		gemma-3-4b	70.40	84.00	67.70	76.70	72.00	75.00	74.30	
	Meta	llama-3.1-8b	65.10	62.00	64.60	68.30	68.00	66.70	65.78	
		llama-3.2-3b	55.90	52.00	52.30	48.30	48.00	56.70	52.20	
	Microsoft	phi-3-mini-128k	64.50	67.00	48.50	65.00	68.00	70.00		
	Mistral AI	ministral-3b	73.70	75.00	60.00	78.30	70.00	63.30	70.05	
mixtral-8x7b		69.40	76.00	59.20	68.30	76.00	65.00	68.98		
Medium [10B-50B]	Google	gemma-3-12b-it	81.60	81.00	76.20	90.00	84.00	81.70	82.41	
		gemma-3-27b-it	85.90	87.00	84.60	90.00	88.00	83.30	86.46	
	Microsoft	phi-3-medium-128k	71.40	70.00	70.80	71.70	62.00	75.00	70.15	
		phi-4-multimodal	80.30	84.00	80.00	93.30	88.00	85.00	85.10	
		phi-4-reasoning-plus	81.90	83.00	72.30	86.70	84.00	86.70	82.43	
	Qwen	qwen3-30b-a3b	85.50	83.00	76.20	93.30	80.00	81.70	83.28	
		mistral-small-24b	86.50	87.00	76.90	88.30	86.00	78.30	83.83	
	Mistral AI	mistral-small-3.2-24b	87.80	85.00	82.30	95.00	86.00	80.00	86.01	
		DeepSeek	deepseek-chat-v3.1	88.80	85.00	80.00	93.30	82.00	81.70	85.13
			Qwen	qwen3-235b	87.50	88.00	83.10	93.30	92.00	95.00
Meta	qwen3-next-80b	85.20	87.00	81.50	90.00	86.00	91.70	86.90		
	llama-3.3-70b	88.50	87.00	80.00	95.00	82.00	81.70	85.70		
Closed-weight LLMs										
Google		gemini-2.0-flash	84.90	91.00	78.50	96.70	84.00	85.00	86.68	
		gemini-2.5-flash	91.80	92.00	86.90	93.30	92.00	81.70	89.61	
		gemini-2.5-flash-lite	80.60	87.00	75.40	85.00	96.00	81.70	84.28	
OpenAI		gpt-4o-mini	85.20	79.00	73.10	86.70	80.00	81.70	80.95	
		gpt-4.1	86.80	90.00	86.20	90.00	96.00	88.30	89.55	
		gpt-4.1-mini	90.80	86.00	78.50	91.70	92.00	85.00	87.33	
xAI		grok-3-mini	93.10	93.00	89.20	93.30	98.00	88.30	92.48	
		grok-4-fast	95.40	93.00	89.20	96.70	90.00	90.00	92.38	
		grok-4	94.10	96.00	93.10	96.70	96.00	95.00	95.15	

Table 10: Accuracy(%) results across models and language pairs for semi-automatic and manual constructed datasets (H). LLM performance is based on few-shot prompting. The final column (AVG) denotes the macro-average accuracy across all datasets.

E Zero-shot prompt

Task: Decide whether the target words have the same sense in both sentences.
 TARGET WORD 1: {w1} SENTENCE 1: {s1}
 TARGET WORD 2: {w2} SENTENCE 2: {s2}

Rules:

- Analyze only the meaning of the target words in their respective contexts.
- Ignore grammatical or syntactic differences unless they affect the meaning.
- If the meaning of the target words is the same, respond with: "Label: YES".
- If the meaning of the target words is different, respond with: "Label: NO".
- Provide a concise explanation for your decision in the format: "Reason: <short sentence>".
- Ensure the output is strictly in plain text and follows the exact format specified.

Example 1:

TARGET WORD 1: fabricar SENTENCE 1: A sua empresa fabrica cadeiras de madeira
 TARGET WORD 2: fabricar SENTENCE 2: Eles fabricam brinquedos para as crianças.
 Label: YES Reason: Both refers to the action of making. Example 2:
 TARGET WORD 1: brinco SENTENCE 1: Maria deu un brinco de alegria cando lle dixeron que ia ser nai. TARGET WORD 2: brinco SENTENCE 2: Ela colocou um brinco dourado para combinar com o vestido.
 Label: NO Reason: The first refers to a jump of joy, while the second refers to a piece of jewelry.

Rules:

- Analyze only the meaning of the target words in their respective contexts.
- Ignore grammatical or syntactic differences unless they affect the meaning.
- If the meaning of the target words is the same, respond with: "Label: YES"
- If the meaning of the target words is different, respond with: "Label: NO"
- Provide a concise explanation for your decision in the format: "Reason: <short sentence>"
- Ensure the output is strictly in

F Few-shot prompt

Task: Decide whether the target words have the same sense in both sentences.
 TARGET WORD 1: w1 SENTENCE 1: s1
 TARGET WORD 2: w2 SENTENCE 2: s2
 Examples:

988 plain text and follows the exact format
989 specified.
990

991 G Mean Zero vs. Few-shot

Dataset	ZERO mean	FEW mean
ES-PT	81.00	82.50
GL-ES	78.70	82.10
GL-PT	72.30	76.30
H-ES-PT	84.60	85.80
H-GL-ES	82.60	82.60
H-GL-PT	76.60	80.70
Overall	79.30	81.60

Table 11: Comparison between zero-shot and few-shot mean performance (%) across the six datasets.

992 H Deviation

Model	Pair	Dev.	Acc.
llama-3.2-3b-instruct	GL-ES	0.75	50.00
mixtral-8x7b-instruct	GL-PT	0.67	62.10
gemma-2-9b-it	GL-ES	0.64	69.00
llama-3.2-3b-instruct	GL-PT	0.63	54.50
mixtral-8x7b-instruct	ES-PT	0.57	68.85
ministral-3b	GL-PT	0.57	61.65
llama-3.2-3b-instruct	ES-PT	0.54	52.10
gemma-2-9b-it	GL-PT	0.48	68.85
gemma-2-9b-it	ES-PT	0.48	71.35
llama-3.1-8b-instruct	ES-PT	0.47	66.70
mixtral-8x7b-instruct	GL-ES	0.40	72.50
ministral-3b	GL-ES	0.37	76.00

Table 12: Models with the highest prediction deviation (Dev.) across language pairs (and corresponding overall accuracy). Higher values indicate a systematic bias toward one class.

I Part-of-Speech distribution

	A	A/N	ART/N	C	PRON	N	N/NUM	N/P	N/R	N/V	R	V	TOTAL
ES-PT	90	18	0	5	1	199	1	1	1	1	4	43	364
GL-ES	14	4	0	0	0	89	0	2	0	0	0	41	150
GL-PT	41	17	2	0	0	101	0	0	0	3	0	26	190

Table 13: Distribution of Part-of-Speech (POS) tags across all the datasets. Abbreviations include Adjectives (A), Articles (ART), Conjunctions (C), Pronouns (PRON), Nouns (N), Numerals (NUM), Preposition (P), Adverbs (R), and Verbs (V). Combined tags indicate lexical categories that vary between the two languages.

J Part-of-Speech error rates

Pair	Model	A (%)	N (%)	V (%)
ES-PT	llama-3.1-8b	23.33	43.22	13.95
	gemma-2-9b-it	50.00	16.58	37.21
	gemma-3-27b-it	12.22	12.56	25.58
	qwen3-235b	10.00	11.56	16.28
	grok-4	4.44	6.53	6.98
GL-ES	llama-3.1-8b	28.57	33.71	43.90
	gemma-2-9b-it	50.00	34.83	21.95
	gemma-3-27b-it	21.43	8.99	17.07
	qwen3-235b	7.14	10.11	14.63
	grok-4	7.14	3.37	4.88
GL-PT	llama-3.1-8b	12.20	38.61	42.31
	gemma-2-9b-it	48.78	30.69	34.62
	gemma-3-27b-it	19.51	17.82	15.38
	qwen3-235b	12.20	17.82	7.69
	grok-4	12.20	6.93	00.00

Table 14: Error rates (%) per Part-of-Speech category (A = adjectives, N = nouns, V = verbs). Each value represents the percentage of misclassified instances within each category for each model and language pair.

K Translation models

- Qwen3-235B-A22B-2507⁴

- GPT-4o-mini⁵

- Grok-4.1-Fast⁶

⁴<https://huggingface.co/Qwen/Qwen3-235B-A22B-Instruct-2507>

⁵<https://platform.openai.com/docs/models/gpt-4o-mini>

⁶<https://x.ai/news/grok-4-1-fast>