

Parallel Sentence Mining Without Parallel Sentences: Scaling to Low-resource Languages

Anonymous ACL submission

Abstract

Parallel sentence mining aims to find translation pairs from comparable monolingual corpora and creates a valuable dataset for downstream tasks. Yet, it has been mostly considered for high-resource language pairs, and endeavours for low-resource languages still require a significant number of parallel sentences for success, although it remains challenging to gather them. Recent works on multilingual sentence representation focused on three techniques which could address this data constraint: enhancing the isotropy of the embeddings, using contrastive learning, and knowledge distillation. In this study, we hence assess the robustness of the three techniques in a low-resource context. We only use monolingual data in the low-resource language or parallel sentences in relevant but different languages to obtain a sentence representation. We extend and create a benchmark to cover sixteen language pairs with eight low-resource languages from three families. While all three methods improve the representation quality by tackling the underlying cross-lingual misalignment, monolingual pre-training and language proximity are essential factors that lead to better performance. We show a significant increase in mining quality, even in the most difficult language pairs.

1 Introduction

Parallel corpora remain a crucial resource for multilingual Natural Language Processing (NLP), as they still help cross-lingual performance of language models (Xu et al., 2024; Lin et al., 2025). Since they are costly to obtain, an automatic approach is parallel sentence mining: it aims to find translation pairs from comparable monolingual corpora (Zweigenbaum et al., 2017). Its performance notably relies on the quality of the sentence representations in the language. Low-resource languages suffer here twice: not only are parallel corpora scarcely available or small, their representations in widely-used models are poorer, leading to

worse mining quality. If high-resource languages (and pairs) have been extensively studied (e.g., the BUCC benchmark, Zweigenbaum et al., 2017) in general, it is less so for low-resource languages. Previous works either involved a large number of *parallel* sentences (>40k) to support an unseen language (Kvapilíková and Bojar, 2023; Tan et al., 2023) or did not achieve usable mining quality without parallel sentences (Okabe et al., 2025).

In this work, we therefore focus on parallel sentence mining for low-resource languages and consider three approaches to improve their language representation *without direct parallel sentences*. To do so, we extend the existing BELOPSEM benchmark of parallel sentence mining (Okabe et al., 2025) to cover eight low-resource languages (Occitan, Upper and Lower Sorbian, Chuvash, Corsican, Mingrelian, Gilaki, and Mazandarani), each paired with two higher-resourced target languages (i.e., sixteen language pairs). Our benchmark extends and improves the original by considering four writing systems and three language families in total. We then apply three types of techniques to obtain better multilingual sentence representation *without direct* parallel sentences, i.e., including the low-resource languages: isotropy-enhancing transformation, contrastive learning, and knowledge distillation. We evaluate and analyse their cross-lingual capabilities on our created mining benchmark. Moreover, we perform case studies for four selected languages in realistic conditions.

Our contributions are as follows: (i) we expand and release a parallel sentence mining benchmark to cover eight typologically diverse low-resource source languages,¹ (ii) we evaluate three successful approaches to improve language representation *without* requiring any parallel sentences featuring a low-resource language, and (iii) we further analyse the effect of each approach on the language repre-

¹Anonymous link for the benchmark and creation code.

083 sentation and perform two case studies. We also
084 release the best performing trained model².

085 2 Related works

086 Performance of parallel sentence mining is inter-
087 twined with the quality of multilingual sentence
088 representation, which explains why mined corpora
089 for low-resource languages are generally noisier.
090 Language models have better cross-lingual repre-
091 sentation when they have been trained on parallel
092 corpora. Yet, as the need for such resources is a
093 major constraint, especially for low-resource lan-
094 guages, several strategies have been devised.

095 Cross-lingual alignment can be improved with-
096 out requiring any specific pre-training by using
097 post-processing methods. Hämmerl et al. (2023)
098 consider ZCA whitening (Huang et al., 2021) and
099 cluster-based isotropy enhancement (Rajae and
100 Pilehvar, 2021), where both tackle the inherent
101 anisotropy in language representation.

102 Relatedly, contrastive learning (Chopra et al.,
103 2005; Hadsell et al., 2006) has been applied to
104 improve English sentence representation, such as
105 in SimCSE (Gao et al., 2021), without requiring
106 parallel sentences. This framework can be super-
107 vised using the Natural Language Inference (NLI)
108 sentence relationship. A positive and a negative
109 example can be compared for a given sentence
110 to bring similar sentences closer together. Wang
111 et al. (2022) extended it to multilingual sentence
112 embeddings with mSimCSE, showing that using
113 (English or multilingual) NLI datasets could im-
114 prove the overall cross-lingual generalisability of
115 the representation. This method proved to be more
116 efficient than similar methods without contrastive
117 learning, such as (Goswami et al., 2021). Such con-
118 trastive learning methods remain more accessible
119 for low-resource languages as they do not require
120 parallel sentences for training at all. This contrasts
121 with established sentence encoders such as LASER
122 (Artetxe and Schwenk, 2019b) or LaBSE (Feng
123 et al., 2022), trained with large parallel corpora.

124 Moreover, knowledge distillation can also help
125 extend a model to other languages (Reimers and
126 Gurevych, 2020). However, it has mostly been
127 used with *direct* parallel sentences in the studied
128 language pair (Tan et al., 2023), limiting its appli-
129 cation to better-resourced languages. Still, for low-
130 resource languages, Heffernan et al. (2022) notably
131 improve LASER using distillation with a multilin-

132 gual teacher and monolingual student (LASER3).
133 The training relies on both parallel corpora and
134 monolingual sentences of the source low-resource
135 language. Tan et al. (2023) extend this work by
136 using contrastive learning with large parallel cor-
137 pora (>40k sentences). For a given sentence pair,
138 the English translation is considered as a positive
139 example. We do not consider similar approaches
140 because of the small size of the available parallel
141 sentences for some of our language pairs (e.g., Min-
142 grelian–English), and instead consider using only
143 parallel sentences from *related* language pairs.

144 Finally, few works focused on parallel sen-
145 tence mining for low-resource languages. Okabe
146 et al. (2025) released a benchmark for this very
147 task, BELOPSEM, covering three language pairs.
148 They gained consistent improvement, notably with
149 anisotropy enhancement, but failed to outperform
150 LaBSE for two of the pairs. We build upon this
151 benchmark for more linguistic diversity (language
152 branches, families, and script) and show better min-
153 ing, still *without* parallel sentences.

154 3 Languages and corpora

155 3.1 Language pairing

156 For notation purposes, we will refer to the *low-*
157 *resource* language of a pair as ‘source’, and the
158 *higher-resourced* language as ‘target’, reflecting a
159 larger availability of sentences for the latter. We
160 study the following eight *source* languages: Occi-
161 tan, Upper Sorbian, Lower Sorbian, Chuvash, Cor-
162 sican, Mingrelian, Gilaki, and Mazandarani. They
163 are all classified as ‘scraping-by’ (1 on a scale from
164 0 to 5) at most in terms of available data in the
165 taxonomy of Joshi et al. (2020) and can hence be
166 considered as low-resource. Besides, Ethnologue
167 (Eberhard et al., 2025) considers all but three lan-
168 guages as endangered, while the UNESCO (2010)
169 lists five of the languages as definitely endangered.
170 Appendix A reports the exact status for all eight
171 languages. We pair each source language with a
172 well (or better)-resourced language (classified as 3
173 or above in the taxonomy from Joshi et al. (2020)).
174 We choose them from a geographical and cultural
175 standpoint and refer to them as *main target*
176 languages (§3.2). We also pair all eight languages
177 with English as a second target (§3.3).

178 3.2 Main language pairs

179 **Occitan–Spanish (OCI–ES)** Occitan (ISO code:
180 oci) is a Romance language spoken in southern

²Anonymous link for the best performing model.

181 France, Spain, and Italy. Both languages in the pair
182 are part of the same branch in the language family
183 (Ibero-Romance) and are written in the Latin script.
184 This is the closest language pair in our benchmark.

185 **Upper Sorbian and Lower Sorbian–German**
186 **(HSB–DE and DSB–DE)** Upper Sorbian (hsb)
187 and Lower Sorbian (dsb) are two Slavic languages
188 spoken in Germany. They constitute the Sorbian
189 branch and are related to Czech or Polish, respec-
190 tively. We hence have two pairs of Indo-European
191 languages, but from different branches, written in
192 the Latin script.

193 **Chuvash–Russian (CHV–RU)** Chuvash (chv) is
194 a Turkic language spoken in the Chuvash Republic
195 in Russia. It is quite distant from other related
196 languages, as it belongs to its own branch. Both
197 languages in the pair use the Cyrillic script but
198 belong to different language families.

199 **Corsican–French (COS–FR)** Corsican (cos) is
200 a language spoken on the islands of Corsica in
201 France and Sardinia in Italy. The language pair is
202 hence close, as both are from the Romance branch
203 of the broader Indo-European language family and
204 written in the Latin script.

205 **Mingrelian–Georgian (XMF–KA)** Mingrelian
206 (xmf) is a Kartvelian language (where Georgian
207 also belongs), spoken in Western Georgia. Both
208 languages are written with the Georgian script.

209 **Gilaki and Mazandarani–Persian (GLK–FA and**
210 **MZN–FA)** Gilaki (g1k) and Mazandarani (mzn)
211 are two Caspian languages from the Indo-Iranian
212 branch, spoken in northern Iran. They hence belong
213 to the larger Iranian branch as Persian, and all three
214 are written in the (Perso-)Arabic script.

215 3.3 Pairing with English

216 We have also paired all source languages with En-
217 glish for two reasons. First, English is the language
218 with the most resources, and language models have
219 been extensively trained on it. Second, it introduces
220 a different language distance for all language pairs,
221 as it is a Germanic language written in the Latin
222 script. Close pairs such as Mingrelian–Georgian
223 will then have a more challenging counterpart from
224 the script and language family perspective. Ap-
225 pendix B classifies the language pair difficulty.

226 3.4 Corpus creation

227 We create synthetic corpora for parallel sentence
228 mining, following the BUCC Shared Task method-

229 ology (Zweigenbaum et al., 2017) used in BE-
230 LOPSEM (Okabe et al., 2025). We mix gold parallel
231 sentence pairs into monolingual corpora in each
232 language, and the goal is to retrieve them back. We
233 insist here on the importance of the quality of the
234 parallel corpus in curating our benchmark and lan-
235 guage pairs, limiting the set of studied languages.

236 **English side** For the first five languages, we
237 have relied on machine translation using Google
238 Translate.³ As we translate from high-resource lan-
239 guages (Spanish, German, Russian, and French)
240 and into English, we deem the output to be of
241 decent quality.⁴ On the other hand, the parallel
242 sentences of the last three languages come from a
243 three-way parallel corpus with English. We hence
244 used actual parallel corpora between the main tar-
245 get language and English (e.g., Georgian–English)
246 and treated them separately as a ‘monolingual’ cor-
247 pus to create their benchmarks (XMF–KA and EN).

248 **Data source** For the three language pairs in BE-
249 LOPSEM (OCI–ES, HSB–DE, and CHV–RU), we use
250 the same *original* datasets. Besides, the HSB–DE,
251 DSB–DE, and CHV–RU pairs were considered in
252 the WMT Shared Tasks in Unsupervised MT and
253 Very Low Resource Supervised MT (Libovický and
254 Fraser, 2021; Weller-Di Marco and Fraser, 2022),
255 which give us both parallel and monolingual source
256 sentences. The target German and Russian sides
257 were taken from the Leipzig corpora (Goldhahn
258 et al., 2012). For COS–FR and OCI–ES, we use par-
259 allel corpora from OPUS (Tiedemann, 2012) and
260 monolingual sentences from the Leipzig corpora.

261 For the last three source languages, we use three-
262 way parallel corpora with English as the third lan-
263 guage. For Mingrelian, we use the Megrelian Lan-
264 guage Corpus (Gersamia and Lobzhanidze, 2022),
265 and for Gilaki and Mazandarani, we rely on the par-
266 allel sentences collected in (Ahmadi et al., 2025b).
267 For their target and English monolingual corpora,
268 we thus use an actual parallel corpus instead of ma-
269 chine translation. From the OPUS collection, we
270 select the OpenSubtitles corpus (Lison and Tiede-
271 mann, 2016) for Georgian and the TEP corpus (Pil-
272 evar et al., 2011) for Persian. While creating our
273 BUCC-style corpus, we ensure that the shuffling
274 and position of the parallel sentences are identical
275 across the two variants of the dataset (main and
276 English). Appendix D lists the exact corpora used.

³With the googletrans library.

⁴Appendix C verifies the MT quality of Google Translate.

lang. pair	train			test		
	source	target	par.	source	target	par.
OCI-ES/EN*	5,981	6,828	331	17,923	20,489	996
HSB-DE/EN*	10,948	13,029	546	32,850	39,091	1,642
DSB-DE/EN*	9,996	11,899	499	29,992	35,699	1,499
CHV-RU/EN*	6,716	6,957	372	20,152	20,877	1,120
COS-FR/EN*	2,749	3,276	136	8,260	9,835	413
XMF-KA/EN	1,023	1,218	50	3,075	3,659	153
GLK-FA/EN	2,161	2,568	124	6,486	7,708	374
MZN-FA/EN	4,248	5,056	211	12,750	15,172	637

Table 1: Size of the training and test datasets, where source and target include the injected parallel sentences (par.). * denotes the use of MT for English.

Data preparation We pre-process both monolingual and parallel sentences using GlotLID (Kargaran et al., 2023) for language identification. For parallel sentences, we also filter out sentence pairs with a large difference in length and with a high proportion of overlapping words. After pre-processing, we inject the gold sentence pairs into the monolingual corpora. Then, we split each created corpus into training and test sets following a 25:75 ratio. The smaller training set is used to find the best model (and parameter setting) to mine sentences from a larger corpus (i.e., test set). Table 1 presents the datasets of the eight language pairs. We note that some datasets (e.g., XMF-KA or COS-FR) remain small due to the original parallel dataset size.

4 Sentence representation

4.1 Baseline multilingual models

We use the standard approach of averaging the word embeddings⁵ to get the sentence-level representation from a language model. We mainly compare two models: XLM-RoBERTa (base) or XLM-R (Conneau et al., 2020) and Glot500-m (Imani et al., 2023), which we will denote G500. The latter extends the former towards more than 500 languages, with a particular focus on low-resource languages, through pre-training on monolingual data. We also compare our systems with LaBSE (Feng et al., 2022), a state-of-the-art sentence encoder which was trained using a large amount of *parallel* sentences. Appendix E summarises the language coverage of the three models.

4.2 Isotropy improving transformation

Multilingual sentence embeddings can be improved by tackling the anisotropy of the vectors in the mul-

tilingual space. Without relying on any additional data, operations such as whitening can transform embeddings from language models before downstream usage. To mine low-resourced language pairs, Okabe et al. (2025) explored a cluster-based isotropy enhancement or CBIE (Rajaei and Pilehvar, 2021). This technique first clusters the vectors and then uses a Principal Component Analysis to remove the top 12 principal components, as in (Hämmerl et al., 2023).⁶ As it was shown to be a simple yet effective method to improve mean-pooled representations on BELOPSEM, we also apply it to our sentence-level representation from Glot500-m. This setting will be called G500+CBIE.

4.3 Contrastive learning

For contrastive learning, we use the mSimCSE framework (Wang et al., 2022), which extends the English-focused SimCSE (Gao et al., 2021) in the multilingual space. We consider two types of datasets to train the models: NLI datasets and parallel sentences between English and the main target languages (i.e., no *source* language involved).

With NLI datasets When using a NLI dataset for contrastive learning, sentences with an ‘entailment’ relationship are considered as positive pairs, while the ‘contradiction’ relation is its hard negative. We consider two NLI datasets: an English-only dataset (Conneau et al., 2017) and the multilingual XNLI dataset (Conneau et al., 2018) to further foster cross-lingual transfer. In this work, we apply both batch contrastive learning methods to XLM-R for comparison: we denote the former mSC-XLM-en and the latter mSC-XLM-multi. The first approach led to significant improvement on both sentence matching and mining in (Wang et al., 2022), while it only used (monolingual) English sentences. The second approach led to additional gains on retrieval tasks comparatively. Yet, these results concerned high(er)-resourced languages.

With translation datasets Finally, we perform contrastive learning with parallel data as positive pairs and a random sentence in another language as a hard negative (mSC-XLM-tr). We use the Parallel Sentences dataset collection (Reimers and Gurevych, 2019) for the six *target* languages paired with English (e.g., Spanish-English). In this way, we still do *not* use *direct* parallel sentences with the low-resource source language.

⁵We use the 8th layer for all language models, following (Imani et al., 2023; Hämmerl et al., 2023).

⁶<https://github.com/KathyHaem/outliers>.

Changing the base model Additionally, we switch the base model from XLM-R to Glot500-m and carry out contrastive learning on the same three datasets. These newly trained models will be denoted as mSC-G500-`{enmultitr}` variants. These models enable us to identify the influence of monolingual training for the source languages.

4.4 Knowledge distillation

The third strategy we consider is a knowledge distillation approach. We choose LaBSE as the teacher model, based on its overall mining performance, and Glot500-m as the main student model, due to its better language coverage. We use the Sentence Transformer framework for distillation (Reimers and Gurevych, 2020). We train the model on *target* language pairs: English paired with the six target languages (de, es, fa, fr, ka, ru). We use the same parallel dataset as for contrastive learning, with a maximum of 100k sentence pairs per language. We call the model G500-LaBSE. To understand the impact of monolingual pre-training, we also use XLM-R as a student model (XLM-LaBSE).

5 Experimental setting

5.1 Mining pipeline

We use an established mining pipeline (Okabe et al., 2025), which updated the system of (Hangya and Fraser, 2019) with contextual embeddings. It consists of two steps: first, it converts each sentence into embeddings in the same multilingual space using the systems previously described. Then, sentences are compared according to a dedicated similarity metric, CSLS (Artetxe and Schwenk, 2019a):

$$\text{CSLS}(x, y) = 2 \cos(x, y) - \sum_{z \in \text{NN}_k(x)} \frac{\cos(x, z)}{k} - \sum_{z \in \text{NN}_k(y)} \frac{\cos(y, z)}{k}, \quad (1)$$

where $\text{NN}_k(x)$ indicates all k -nearest neighbours of vector x to avoid the hubness problem (Dinu et al., 2015). We use $k = 20$ in our experiments. We select our similarity threshold based on the training set performance. The experiments are mainly evaluated using the F-score. Appendix F lists experimental details for reproducibility.

5.2 Main mining results

The left side of the columns in Table 2 displays the sentence mining results for the eight main language pairs (i.e., presented in Section 3.2).

Baseline performance We first notice that pre-training on the language is crucial when using averaged word embeddings for sentence representation. Glot500-m consistently outperforms XLM-R on seen languages, while it similarly struggles with the unseen dsb and glk. The lower score for mzn could be due to the fewer pre-training sentences (74k) than for the other languages, which all comprise more than 100k sentences. Both approaches are also worse than the sentence encoder LaBSE, except for chv and xmf for G500. Explicit training on parallel sentences with a contrastive objective helps LaBSE represent unseen but related languages, whereas mean-pooled embeddings seem to suffer from structural cross-lingual misalignment.

CBIE Tackling the inherent anisotropy of the language representation with CBIE leads to consistent improvement for all languages for G500. This closes the gap with LaBSE and widens it for chv and xmf. We note here that the least performing cases are still for DSB-DE and GLK-FA, which constitute the two language pairs where the source language has not been seen by G500.

Contrastive learning Methods based on contrastive learning lead to noticeable improvement compared to their base model overall. If the English NLI dataset (-en) brings uneven effect, the cross-lingual training (-multi) leads to systematically better F-scores for all main language pairs and models. Our extension using Glot500-m brings significant improvement to the baseline G500 but also to mSC-XLM-en, in several language pairs. This shows that contrastive learning is also more effective, when the base model has been pre-trained on the language. Besides, we see that on average, using XNLI over English NLI is worse with XLM-R (17 vs. 21), while for Glot500-m, cross-lingual NLI benefits more (44 vs. 26), showing that the competitiveness found in (Wang et al., 2022) seems to mostly hold, when both languages have been seen during pre-training.

Finally, using translation datasets for contrastive learning (-tr) leads to rather comparable scores to the -multi models overall, with a few exceptions. A notable difference is that with XLM-R as a base model, this setting can potentially lead to a slightly worse score, whereas we see an improvement with G500, albeit sometimes small.

Knowledge distillation When we use knowledge distillation from LaBSE, we observe that it

source target	OCI		HSB		DSB		CHV		COS		XMF		GLK		MZN	
	ES	EN	DE	EN	DE	EN	RU	EN	FR	EN	KA	EN	FA	EN	FA	EN
XLM-R	47.44	35.72	0.40	1.10	0.42	0.20	0.64	0.35	16.06	12.75	4.55	0.00	1.64	0.00	3.84	0.00
G500	79.96	62.59	21.33	14.03	2.97	1.45	39.59	27.61	49.94	38.34	41.25	2.44	4.06	0.25	9.92	0.31
G500+CBIE	92.43	87.10	48.01	33.51	20.29	14.51	48.77	34.13	69.80	65.13	57.37	8.89	19.54	1.50	43.25	4.42
mSC-XLM-en	62.64	50.24	13.45	10.46	11.91	8.55	7.11	2.11	52.44	43.31	13.01	2.92	20.46	4.39	30.37	3.31
mSC-XLM-multi	61.93	39.45	18.40	5.82	11.70	3.08	3.53	1.24	47.95	27.99	7.84	1.96	21.77	1.86	24.05	0.41
mSC-XLM-tr	67.95	63.69	21.28	20.12	15.79	13.61	2.50	1.47	44.18	43.85	5.46	2.16	18.64	7.62	24.77	10.01
mSC-G500-en	73.95	34.54	22.38	4.65	11.32	2.52	36.57	5.18	51.66	26.01	56.12	5.24	32.89	4.91	47.07	5.72
mSC-G500-multi	83.24	74.83	53.34	41.10	27.02	18.42	57.58	46.68	68.93	61.44	57.54	29.91	31.77	5.23	48.50	5.32
mSC-G500-tr	86.13	81.40	63.15	55.28	27.80	24.57	64.90	55.18	69.17	60.79	62.75	30.63	33.23	18.04	49.63	21.83
LaBSE	98.20	97.04	74.87	77.43	69.60	74.69	30.85	31.35	90.32	90.39	28.47	27.62	45.49	28.68	56.64	32.63
XLM-LaBSE	96.05	95.69	43.64	44.20	39.06	40.89	24.57	23.34	69.71	70.53	18.05	21.05	44.32	26.41	58.93	31.69
G500-LaBSE	97.11	96.92	78.31	77.98	54.57	54.56	80.68	79.25	83.70	85.58	71.11	59.61	54.96	32.51	71.71	39.61

Table 2: F-scores (%) on the test set of the eight source languages paired with the main target language (left column for targets) and English (right). Results in **bold** indicate the best score.

leads to a positive improvement for both base models. However, using G500 as a student model is a better strategy, with some drastic differences, such as +43 points for XMF–KA. This stresses again the importance of monolingual pre-training in cross-lingual alignment in yet another strategy.

Across strategies Comparing G500+CBIE and mSC-G500-tr, we observe that the latter outperforms the former, except for the two closest (and easier) language pairs of OCI–ES and COS–FR. Overall, G500-LaBSE leads to the best score, with a similar resource cost to the -tr models. This suggests that specific training (contrastive learning or distillation) helps address deeper misalignment than isotropy enhancement. CBIE is, however, easily scalable (no training) and language independent. We also see that parallel sentences, even without the source language, remain more efficient in improving multilingual alignment than NLI datasets.

5.3 Results when paired with English

The right side of the columns in Table 2 shows the mining score when paired with the English translation of the target side. This allows for a direct comparison with the left side in a sometimes more challenging context: English can introduce a difference in script (e.g., CHV–EN for Cyrillic) and language family (XMF–KA to XMF–EN).

We observe that most of the trends seen earlier still hold. Applying CBIE, contrastive learning, or knowledge distillation leads to consistent improvement over the averaged word representation from the base models, and monolingual pre-training gives a noticeable advantage for cross-lingual alignment. Interestingly, NLI contrastive learning with

only English, which is the target language here, can lead to lower scores than for the main language pair (e.g., 22 vs. 5 for mSC-G500-en in hsb), showing some limits to this otherwise robust approach (Wang et al., 2022). This underlines the importance of multilingual contrastive learning, whether with NLI or translation data, depending on availability.

Challenging pairs When comparing the mining scores for the same source language, we observe that they are lower on average when paired with English, which is very likely due to the larger distance between the languages. The most notable drop is for xmf, where the F-score decreases by 22 points on average for contrastive learning and 48 for CBIE. This stark difference also affects both Caspian languages. This is mainly due to the divergence in script and language family, underlining the language or script ‘cluster’ phenomenon in multilingual language models (Wang et al., 2022; Liu et al., 2024). This negative effect, however, is reduced for LaBSE and distilled models, due to their more isotropic representation, with relatively little difference. Overall, the script and language family barrier has not been fully overcome yet, as three of the four most difficult language pairs with English (XMF–, GLK–, and MZN–EN) significantly perform worse, showing that the three methods improved still largely *within* a language and script cluster.

6 Analyses

6.1 Language pairs for knowledge distillation

Beyond the setting of G500-LaBSE, we use two other sets of languages. First, a more comprehensive set of languages than G500-LaBSE with six additional languages which are *related* to the *source*

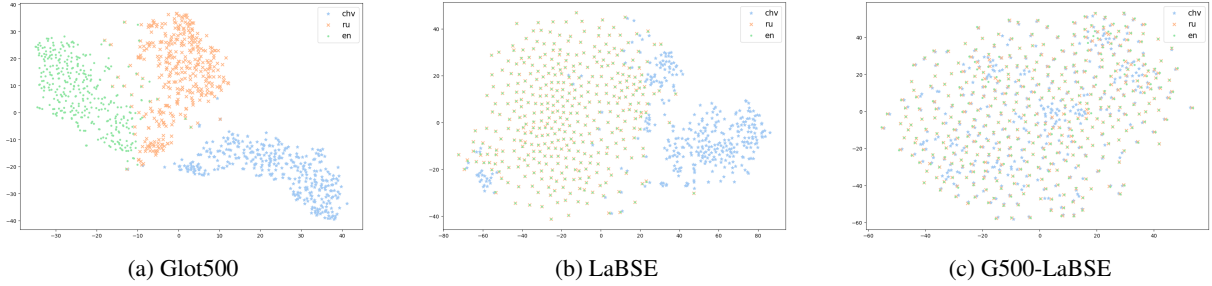


Figure 1: t-SNE visualisation for three-way parallel sentences in Chuvash, Russian, and English.

ones (ca, cs, pl, tr, it, ku; a total of twelve language pairs).⁷ Second, we simulate a more generic and multilingual model by using five language pairs representing the United Nations languages: ar, en, es, fr, ru, zh-cn. The former, denoted G500-LaBSE-all, approximates the upper bound of distillation with cross-lingual transfer from better-resourced languages. The latter (G500-LaBSE-un) explores a more general model with widely spoken and institutional languages. For our benchmark, we note that some languages are covered by this -un set from the target side, from the script (e.g., CHV-RU or GLK-FA), but not all (e.g., XMF-KA).

	main pair	English pair
LaBSE	61.81	57.48
G500-LaBSE	74.02	65.75
G500-LaBSE-all	78.98	69.72
G500-LaBSE-un	73.40	63.66

Table 3: Averaged F-score for the three knowledge distillation settings on the 16 language pairs.

Table 3 compares the average scores of the three models (and LaBSE for reference); full results are in Appendix G. We note that providing related languages (-all) improves the mining quality noticeably. Besides, using some unrelated languages (-un setting) still gives better results than the baseline LaBSE, underlining the benefits of Glot500-m as a student model, but lags behind the other two settings. We notice lower but similar results for most language pairs, with the exceptions of dsb, glk, and mzn, showing low transfer for the Arabic script.

6.2 Impact on isotropy

Anisotropy representation Figure 1 displays the t-SNE representation (van der Maaten and Hinton, 2008) of the *three-way parallel* sentences in Chuvash, Russian, and English for three models.

⁷Kurdish is for both Caspian languages, and Georgian is already in the first set of training languages. See Appendix B.

We observe that for Glot500, we easily distinguish three clusters that barely overlap, suggesting a high anisotropy. Monolingual pre-training, especially for languages from different families and scripts, struggles to properly align the sentences. We also notice a gap between the two high-resource languages. LaBSE, on the other hand, manages to better align Russian-English sentences, as we cannot distinguish between both sentence representations any more in the figure. Chuvash sentences seem to surround the isotropic bilingual space, suggesting poor relative cross-lingual alignment between Chuvash and the two target languages.

Finally, the knowledge distillation between Glot500 and LaBSE seems to alleviate the anisotropy problem, as the three language representations are now closer, with no visible language cluster. We stress here that no parallel sentences for CHV-RU were needed to obtain this representation. The pre-training focused on the low-resource language with monolingual sentences and widely-used Russian-English parallel sentences for distillation.

	Before	After CBIE
XLM-R	7.82	20.79
G500	24.75	40.54
LaBSE	59.64	47.84
mSC-G500-multi	44.43	34.92
G500-LaBSE	69.89	52.02

Table 4: Averaged F-score over the 16 language pairs before and after applying CBIE.

Applying CBIE Another evidence that the sentence representations are more isotropic is that applying CBIE worsens the mining quality. Table 4 shows the averaged F-score over the 16 language pairs before and after applying CBIE on the sentence representations. We note that the isotropy enhancement technique mainly improves the mean-pooled word representations, while it significantly degrades the performance of LaBSE or of the mod-

els obtained with contrastive learning or knowledge distillation. Hämmerl et al. (2023) observed that such a transformation benefits little to models trained on parallel data; we confirm that already isotropic representations do not benefit from it. We thus observe little to no synergy between CBIE and the other two approaches.

6.3 Case studies: for one language branch

One real-life scenario is to perform parallel sentence mining for a limited set of *related* language pairs rather than typologically diverse ones. We focus here on mining for two language branches, Sorbian and Caspian languages, as both feature a language not seen by Glot500 (dsb and glk).

We start from XLM-R and pre-train it in two ways: with the relevant languages seen by Glot500 (base setting)⁸ and additionally with the unseen language (+ setting). We use 100k monolingual sentences maximum for each language, drawn from the Glot500 corpus. For the two unseen languages, we use different sentences than those in our benchmark. Based on the results, we compare the simple mean-pooled representation (PT), a CBIE transformed version, as it scales easily (no training), and the best knowledge distillation with related languages (as in the -all setting in Section 6.1), denoted as PT-LaBSE. We note that models are developed separately per language branch.

language branch	Sorbian				Caspian			
	hsb		dsb		glk		mzn	
source	de	en	de	en	fa	en	fa	en
G500	21.33	14.03	21.33	1.45	4.06	0.25	9.92	0.31
LaBSE	74.87	77.43	74.87	74.69	45.49	28.68	56.64	32.63
G500-LaBSE	78.31	77.98	54.57	54.56	54.96	32.51	71.71	39.61
G500-LaBSE-all	89.67	88.26	67.29	66.41	59.30	32.06	70.39	40.16
PT	41.89	4.00	2.48	1.82	2.90	0.00	37.11	0.96
CBIE	44.29	11.25	17.33	9.45	17.31	1.05	56.73	4.02
PT-LaBSE	87.39	83.18	66.42	63.47	47.62	26.23	77.44	43.42
PT+	49.24	17.04	59.64	18.08	2.86	0.00	34.62	0.46
CBIE+	47.40	26.32	51.93	28.10	20.48	0.34	55.72	3.13
PT-LaBSE+	89.63	86.10	93.09	89.61	48.64	24.95	79.56	43.99

Table 5: Mining scores for models specifically pre-trained for Sorbian or Caspian languages.

Table 5 reports the results for all eight language pairs. We still notice relatively poor F-scores for the two languages not supported in the base setting PT, dsb and glk, despite the branch-specific pre-training. This shows that G500 benefits also from cross-lingual transfer from less related languages seen during pre-training. Interestingly, if using

⁸de, cs, pl, hsb for Sorbian languages and fa, ckb, mzn for Caspian languages. We exclude English for comparison.

monolingual data for Lower Sorbian (+ setting) improves its representation with a similar trend to hsb across settings (+26 points from distillation), Gilaki benefits far less from the additional pre-training. Nonetheless, we notice that the distillation approach remains reliable, as it systematically outperforms LaBSE. Besides, monolingual pre-training helped distillation, as PT-LaBSE+ leads to better mining than PT-LaBSE. Finally, G500-LaBSE-all appears as the best overall system.

We additionally perform another case study to integrate a realistically and commonly available bilingual source for low-resource language pairs. We consider for this a bilingual lexicon for the two pairs involving Chuvash in Appendix H.

7 Conclusion

We created a benchmark to evaluate parallel sentence mining quality for eight low-resource languages, paired with a main target language (a locally relevant high-resource language) and its English translation, allowing for direct comparison. We strove to represent various levels of language distance, both in terms of language family and script. We applied three types of approaches to improve multilingual sentence embeddings: isotropy enhancement (CBIE), which doesn’t need any external resources, contrastive learning, and knowledge distillation, using NLI datasets or parallel corpora of *related* languages. None of our strategies relied on *parallel* data involving our eight source languages; yet, we observed significant improvement for all three methods across the 16 language pairs, compared to a simply mean-pooled sentence representation. The best approach also competes or notably outperforms LaBSE. Our analyses stress the importance of monolingual pre-training to fully benefit from better cross-lingual alignment.

Most benchmarks focus on languages paired with English, which only obscures the extent of (mis-)alignment in the multilingual space. Models can still struggle when paired with English (and not their main target language), i.e., when language distance is higher. However, our results also suggest that isotropy improvement is still possible ‘locally’, for culturally relevant language pairs.

Future work will extend the study to more language pairs, notably on the source side with low-resource languages. Moreover, we will focus on improving the alignment between distant language pairs such as Gilaki–English.

666 Limitations

667 First, our sentence encoding baseline, LaBSE, is
668 already very robust, despite having never seen most
669 of the source languages we considered (none but
670 Corsican). Relatedly, we observe that for both Ro-
671 mance languages (Occitan and Corsican), we only
672 managed to get close to and not above the score
673 reached by the baseline LaBSE. It has, nevertheless,
674 been extensively pre-trained on related languages
675 (cf. Section E; e.g., Catalan, French, Spanish for
676 Occitan) or on the language itself (for Corsican), us-
677 ing parallel sentences, which can explain its ‘zero-
678 shot’ performance for mining. Besides, if we see
679 higher performance for Indo-European languages
680 written in the Latin script, other languages are com-
681 paratively lagging behind. Our approach, on the
682 contrary, scales to a variety of languages with a
683 similar degree of improvement, even for the most
684 challenging language pairs, such as CHV–EN or
685 XMF–EN. The harder the source language and the
686 further the language pair, the better our approach
687 performs compared to LaBSE for mining. It ef-
688 fectively improves the sentence representation of
689 a low-resource language using only monolingual
690 pre-training in the language (as in Glot500) and
691 parallel datasets from *related* parallel corpora.

692 Moreover, despite the high F-scores achieved by
693 G500-LaBSE-all (more than 70 for 10 language
694 pairs), we still see a relatively low mining qual-
695 ity for some pairs (e.g., Caspian languages paired
696 with English, between 30–40), suggesting a poor
697 representation. The poorest scores primarily oc-
698 cur when the two languages are distant in terms of
699 script or language family, and with smaller or no
700 monolingual pre-training.

701 Finally, the choice of language pairs (and cor-
702 responding datasets) is bound by the availability
703 of resources (both monolingual and parallel). The
704 main bottleneck is the number of *high-quality* par-
705 allel sentences, which restricts the overall dataset
706 size and explains its variation. It also reduces the
707 number of possible language pairs to study.

708 Ethical considerations

709 We create our benchmark from openly accessible
710 corpora, respecting their licences. For that reason,
711 each language pair has a different licence in our
712 released version. Moreover, the machine transla-
713 tion we use to generate the English side for four
714 *target* languages is deemed to be of decent quality,
715 due to the nature of the translation direction and in-

716 volved high-resource languages. We acknowledge,
717 however, the risk it entails and try to measure it in
718 Section C.

References

- 719 Sina Ahmadi, Razhan Hameed, and Rico Sennrich. 720
2025a. [Literary translations and synthetic data for 721
machine translation of low-resourced Middle Eastern 722
languages](#). In *Proceedings of the 22nd International 723
Conference on Spoken Language Translation (IWSLT 724
2025)*, pages 110–118, Vienna, Austria (in-person and 725
online). Association for Computational Linguistics. 726
727
- Sina Ahmadi, Rico Sennrich, Erfan Karami, Ako 728
Marani, Parviz Fekrazad, Gholamreza Akbarzadeh 729
Baghban, Hanah Hadi, Semko Heidari, Mahîr Dogan, 730
Pedram Asadi, Dashne Bashir, Mohammad Amin 731
Ghodrati, Kouros Amini, Zeynab Ashourinezhad, 732
Mana Baladi, Farshid Ezzati, Alireza Ghasemifar, 733
Daryoush Hosseinpour, Behrooz Abbaszadeh, and 734
14 others. 2025b. [PARME: Parallel corpora for low- 735
resourced Middle Eastern languages](#). In *Proceedings 736
of the 63rd Annual Meeting of the Association for 737
Computational Linguistics (Volume 1: Long Papers)*, 738
pages 30032–30053, Vienna, Austria. Association 739
for Computational Linguistics. 740
- Mikel Artetxe and Holger Schwenk. 2019a. [Margin- 741
based parallel corpus mining with multilingual sen- 742
tence embeddings](#). In *Proceedings of the 57th An- 743
nual Meeting of the Association for Computational 744
Linguistics*, pages 3197–3203, Florence, Italy. Asso- 745
ciation for Computational Linguistics. 746
- Mikel Artetxe and Holger Schwenk. 2019b. [Mas- 747
sively multilingual sentence embeddings for zero- 748
shot cross-lingual transfer and beyond](#). *Transactions 749
of the Association for Computational Linguistics*, 750
7:597–610. 751
- Sumit Chopra, Raia Hadsell, and Yann LeCun. 2005. 752
[Learning a similarity metric discriminatively, with 753
application to face verification](#). In *2005 IEEE Com- 754
puter Society Conference on Computer Vision and 755
Pattern Recognition (CVPR’05)*, volume 1, pages 756
539–546 vol. 1. 757
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, 758
Vishrav Chaudhary, Guillaume Wenzek, Francisco 759
Guzmán, Edouard Grave, Myle Ott, Luke Zettle- 760
moyer, and Veselin Stoyanov. 2020. [Unsupervised 761
cross-lingual representation learning at scale](#). In *Pro- 762
ceedings of the 58th Annual Meeting of the Asso- 763
ciation for Computational Linguistics*, pages 8440– 764
8451, Online. Association for Computational Lin- 765
guistics. 766
- Alexis Conneau, Douwe Kiela, Holger Schwenk, Loïc 767
Barrault, and Antoine Bordes. 2017. [Supervised 768
learning of universal sentence representations from 769
natural language inference data](#). In *Proceedings of 770*

885	Scaling multilingual corpora and language models to 500 languages.	NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Hefernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, and 20 others. 2024. Scaling neural machine translation to 200 languages.	940
886			941
887			942
888			943
889			944
890			945
891	Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019. Billion-scale similarity search with GPUs. <i>IEEE Transactions on Big Data</i> , 7(3):535–547.	Shu Okabe, Katharina Hämmerl, and Alexander Fraser. 2025. Improving parallel sentence mining for low-resource and endangered languages.	946
892			947
893			948
894	Pratik Joshi, Sebastin Santy, Amar Budhiraja, Kalika Bali, and Monojit Choudhury. 2020. The state and fate of linguistic diversity and inclusion in the NLP world.		949
895			950
896			951
897			952
898			953
899			954
900			955
901	Amir Hossein Kargaran, Ayyoob Imani, François Yvon, and Hinrich Schuetze. 2023. GlotLID: Language identification for low-resource languages.	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation.	956
902			957
903			958
904			959
905			960
906			961
907	Ivana Kvapilíková and Ondřej Bojar. 2023. Boosting unsupervised machine translation with pseudo-parallel data.	Mohammad Taher Pilevar, Hesham Faili, and Abdol Hamid Pilevar. 2011. Tep: Tehran english-persian parallel corpus.	962
908			963
909			964
910			965
911			966
912			967
913	Jindřich Libovický and Alexander Fraser. 2021. Findings of the WMT 2021 shared tasks in unsupervised MT and very low resource supervised MT.		968
914			969
915			970
916			971
917			972
918			973
919	Peiqin Lin, Andre Martins, and Hinrich Schuetze. 2025. A recipe of parallel corpora exploitation for multilingual large language models.	Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation.	974
920			975
921			976
922			977
923			978
924			979
925	Pierre Lison and Jörg Tiedemann. 2016. OpenSubtitles2016: Extracting large parallel corpora from movie and TV subtitles.	Matt Post. 2018. A call for clarity in reporting BLEU scores.	980
926			981
927			982
928			983
929			984
930			985
931			986
932	Yihong Liu, Chunlan Ma, Haotian Ye, and Hinrich Schuetze. 2024. TransliCo: A contrastive learning framework to address the script barrier in multilingual pretrained language models.	Sara Rajae and Mohammad Taher Pilehvar. 2021. A cluster-based approach for improving isotropy in contextual embedding space.	987
933			988
934			989
935			990
936			991
937			992
938			993
939			994
		Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence embeddings using Siamese BERT-networks.	995
			996
		Nils Reimers and Iryna Gurevych. 2020. Making monolingual sentence embeddings multilingual using knowledge distillation.	997

2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 4512–4525, Online. Association for Computational Linguistics.

Weiting Tan, Kevin Heffernan, Holger Schwenk, and Philipp Koehn. 2023. [Multilingual representation distillation with contrastive learning](#). In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 1477–1490, Dubrovnik, Croatia. Association for Computational Linguistics.

Jörg Tiedemann. 2012. [Parallel data, tools and interfaces in OPUS](#). In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC’12)*, pages 2214–2218, Istanbul, Turkey. European Language Resources Association (ELRA).

UNESCO. 2010. *Atlas of the world’s languages in danger*, 3rd edition. Paris, France.

Laurens van der Maaten and Geoffrey Hinton. 2008. [Visualizing data using t-SNE](#). *Journal of Machine Learning Research*, 9(86):2579–2605.

Yaushian Wang, Ashley Wu, and Graham Neubig. 2022. [English contrastive learning can learn universal cross-lingual sentence embeddings](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9122–9133, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Marion Weller-Di Marco and Alexander Fraser. 2022. [Findings of the WMT 2022 shared tasks in unsupervised MT and very low resource supervised MT](#). In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 801–805, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Hassan Awadalla. 2024. [A paradigm shift in machine translation: Boosting translation performance of large language models](#). In *The Twelfth International Conference on Learning Representations*.

Pierre Zweigenbaum, Serge Sharoff, and Reinhard Rapp. 2017. [Overview of the second BUCC shared task: Spotting parallel sentences in comparable corpora](#). In *Proceedings of the 10th Workshop on Building and Using Comparable Corpora*, pages 60–67, Vancouver, Canada. Association for Computational Linguistics.

A Language status

Table 6 reports the status of the eight source languages according to (Joshi et al., 2020), Ethnologue (Eberhard et al., 2025), and UNESCO’s *Atlas of the world’s languages in danger* (UNESCO, 2010). For Ethnologue, E stands for endangered, S for stable, and I for institutional, in order of language

ISO	Glottocode	NLP resource	Ethnologue	UNESCO
oci	occi1239	1	E	SE-DE
hsb	uppe1395	1	E	DE
dsb	lowe1385	1	E	DE
chv	chuv1255	1	E	VU
cos	cors1241	1	E	DE
xmf	ming1252	1	S	DE
glk	gila1241	1	I	NE
mzn	maza1291	0	S	NE

Table 6: Classification of languages for NLP resource availability and endangerment.

vitality. UNESCO classifies as follows, ordered from the most critical state to safe: SE: severely endangered, DE: definitely endangered, VU: vulnerable, and NE: safe.⁹ For Occitan, the classification depends on the considered dialect. We also report the Glottocodes (Hammarström and Forkel, 2022) of each language.

Besides, only two languages are present in the high-quality and multilingual FLORES+ benchmark for machine translation (NLLB Team et al., 2024): Chuvash and Occitan. The remaining six languages have thus not been considered by the popular multilingual sentence matching task.

B Language pair

source	main target	English	closest
Occitan	Spanish: 1	2	Catalan
Upper Sorbian	German: 2	2	Czech
Lower Sorbian	German: 2	2	Polish
Chuvash	Russian: 3	4	Turkish
Corsican	French: 1	2	Italian
Mingrelian	Georgian: 2	4	Georgian
Gilaki	Persian: 2	4	Kurdish
Mazandarani	Persian: 2	4	Kurdish

Table 7: Language pair difficulty when paired with the main target and English, from the easiest (1) to the hardest (4).

Table 7 summarises the language pair difficulty into four levels. Level 1 represents the easiest pairs with the same language branch, script, and relatedness to several high-resource languages (e.g., OCI-ES). Level 2 indicates a pair where the language family is the same but from a different branch (e.g., HSB-DE) or from the same branch but with a non-

⁹There are two other levels below SE.

Latin script (e.g., XMF–KA). Level 3 is only for CHV–RU, as they are from different language families but with the same script. Finally, the hardest is level 4, where the script and broader language branch or family differ (e.g., GLK–EN).

We also report the *closest high-resource* language according to the taxonomy of (Joshi et al., 2020). For Mingrelian, no language other than Georgian was deemed to be close enough. We use these languages paired with English during the additional distillation of LaBSE into Glot500 (G500-LaBSE-all) in Section 6.1.

C MT quality

Since part of our benchmark relies on Machine Translation into English, we evaluate the quality of Google Translate on the multilingual Flores+¹⁰ benchmark (NLLB Team et al., 2024). We evaluate the translations according to three widely-used metrics: BLEU (Papineni et al., 2002), chrF++ (Popović, 2015), and xCOMET(-XL) (Guerreiro et al., 2024). We use sacreBLEU (Post, 2018) for the implementation of the first two. The third was ranked best in the WMT 2024 metrics Shared Task (Freitag et al., 2024) and supports the five considered languages.

lang. pair	BLEU	chrF++	xCOMET
DE→EN	48.77	69.46	98.32
ES→EN	33.73	60.37	97.14
FR→EN	51.97	71.18	97.21
RU→EN	41.29	64.73	97.09

Table 8: BLEU, chrF++, and xCOMET scores for Google Translate on Flores+ for translations from German, Spanish, French, and Russian into English.

Table 8 presents the translation quality for the four language pairs where we used MT into English. While BLEU and chrF++ are fairly average, the xCOMET score shows very high quality semantically for all four language pairs. As it is the metric with the highest correlation with human evaluation in (Freitag et al., 2024), we deem the MT quality from Google Translate to be sufficient for our use.

D Corpora details

Below are the resources that we use to create our synthetic corpus for parallel sentence mining.

¹⁰https://huggingface.co/datasets/openlanguagedata/flores_plus.

Occitan-Spanish Following (Okabe et al., 2025), we use the Wikimedia corpus from OPUS (Tiedemann, 2012) for the parallel data. Similarly, the monolingual sentences for both languages come from the Leipzig corpus (Goldhahn et al., 2012), more specifically, the Wikipedia data from 2021 (30k sentences).

Upper Sorbian-German For Upper Sorbian sentences (monolingual and bilingual), we rely on the WMT 2020 Shared Tasks in Unsupervised MT and Very Low Resource Supervised MT (Fraser, 2020). The monolingual sentences in German come from the 2020 Wikipedia data of the Leipzig corpus (Goldhahn et al., 2012) with 300k sentences.

Lower Sorbian-German We use the data from the WMT 2022 Shared Tasks in Unsupervised MT and Very Low Resource Supervised MT¹¹ (Weller-Di Marco and Fraser, 2022) for both monolingual and parallel sentences in Lower Sorbian. More precisely, we use the 66408_DSB_monolingual.txt.gz file for the monolingual part and the 2022 training data files (40194_train_dsb_de) for the parallel corpus. For German, we use the news data from the Leipzig corpora (Goldhahn et al., 2012) in German (2022, 100k sentences).

Chuvash-Russian This language pair was studied in the WMT 2021 Shared Tasks in Unsupervised MT and Very Low Resource Supervised MT (Libovický and Fraser, 2021). We rely on the development dataset (in devtest.chv-ru.tgz) to have the parallel sentences. For both Chuvash and Russian monolingual data, we use the Leipzig corpora (Wikipedia 2021; 30k sentences)

Corsican-French For this language pair, we use parallel sentences from OPUS (Tiedemann, 2012). Given the dataset size, we combine the Wikimedia dataset and the eight sentences from Tatoeba and filter sentences with two or fewer words. We also manually corrected some sentences that were misaligned. On the monolingual side, we rely on the Leipzig corpora and use the Wikipedia corpus for both Corsican and French. We namely combine all three available corpora for Corsican (Wikipedia 2014, 2016, and 2021, each with 10k sentences). We use the 2021 30k Wikipedia corpus for French.

¹¹https://www.statmt.org/wmt22/unsup_and_very_low_res.html.

Mingrelian–Georgian/English The three-way parallel dataset from the Megrelian Language Corpus¹² is released under a CC BY-NC-SA 4.0 licence. We use the 2021 10k Wikipedia corpus for Mingrelian. For the Georgian or English monolingual side, we use a Georgian–English parallel corpus from OPUS: OpenSubtitles¹³ (Lison and Tiedemann, 2016).

Gilaki and Mazandarani–Persian/English For both Gilaki and Mazandarani, we use the three-way parallel dataset collected in the PARME dataset (Ahmadi et al., 2025b). It is licensed with a MIT licence. For the source monolingual sentences, we use the 2016 10k Wikipedia corpus for Gilaki and for Mazandarani, the 2021 30k Wikipedia corpus from the Leipzig corpus. For the target side, we use an actual parallel corpus for Persian–English: the TEP corpus (Pilevar et al., 2011) from OPUS, which can be used for research and non-commercial purposes. We note here that both source languages were namely studied for mining in (Ahmadi et al., 2025a), paired with English only.

E Language coverage of the models

Table 9 summarises the languages present in the pre-training dataset of the three language models that we compare, XLM-R, Glot500-m, and LaBSE. Glot500-m extends XLM-R (base) to more than 500 languages, including six of our source languages. Corsican has the most monolingual sentences among them (3,015,055), while Mazandarani has the fewest comparatively (73,719). We note that LaBSE is trained on more than 109 languages, including Corsican (both monolingual and parallel sentences). All three models have seen English and the six target languages (French, Georgian, German, Persian, Russian, and Spanish, in alphabetical order).

F Experimental details

F.1 Computational details

The parallel sentence mining pipeline relies on the creation of the sentence-level representation and the mining itself, both scaling with the dataset size (i.e., the longest experiments being for HSB–DE, the shortest for XMF–KA). The mining part carries out similarity search using Faiss (Johnson et al., 2019), which can also run with GPUs for faster

¹²<https://xmf.iliauni.edu.ge/>.

¹³<https://www.opensubtitles.org>.

	XLM-R	Glots500-m	LaBSE
oci	✗	✓ (1.4M)	✗
hsb	✗	✓ (104k)	✗
dsb	✗	✗	✗
chv	✗	✓ (860k)	✗
cos	✗	✓ (3.0M)	✓
xmf	✗	✓ (175k)	✗
glk	✗	✗	✗
mzn	✗	✓ (74k)	✗

Table 9: Languages seen during pre-training for the three back-end multilingual language models.

results. We used 1 GPU (NVIDIA A100 or H100) for all our experiments. The whole pipeline runs in less than one hour for one source language paired with the two target languages (a few minutes for the smallest pair).

We have also trained or retrained models using the mSimCSE framework (Wang et al., 2022). To ensure comparability, we use the same pre-training parameters as the default implementation. Using the same GPU resources as above, the training took a few hours; the longest computation time was for multilingual contrastive learning with XLNI or translation datasets, which reached around one million instances.

Finally, distillation with LaBSE took around eight hours for the largest dataset (-all setting) with 3 epochs. For our distillation experiments, convergence can happen earlier, depending on the language pairs used.

F.2 External datasets for model training

We list the datasets used to train or pre-train our models. For contrastive learning, we use the English NLI dataset (276k sentences) for the -en variants and the XNLI dataset (2.2M sentences) for the -multi variant. The -tr variant uses the same dataset as the base dataset used for knowledge distillation: each language pair (English paired with: French, Georgian, German, Persian, Russian, and Spanish) is represented through 100k sentences (maximum).

For the additional knowledge distillation settings (Section 6.1), the same 100k threshold is applied for each language pair, from the Parallel Sentences dataset. For the -all setting, this amounts to 1.0M sentences with the language pairs added (English paired with Catalan, Czech, Italian, Kurdish, Polish, and Turkish). For the -un setting, this leads to a dataset of 500k instances for the En-

	CHV-RU	CHV-EN
G500-LaBSE-all	85.17	85.07
dict-word	90.61	90.59
dict-sent	86.43	83.94

Table 10: Effect of a bilingual dictionary for Chuvash paired with Russian or English.

ble 10, the ‘code-switched’ sentences only brought a slight boost in CHV-RU and degraded the score when paired with the unseen English.

1279
1280
1281

1238 glish–Arabic, simplified Chinese, French, Russian,
1239 and Spanish combined language pairs.

1240 We use the same 100k threshold for our spe-
1241 cific pre-training experiments (Section 6.3). We
1242 consider the Glot500 corpus as source for the lan-
1243 guages covered by the model. For the Sorbian
1244 branch, there are 400k monolingual sentences in
1245 Czech, German, Polish, and Upper Sorbian. The
1246 monolingual sentences for Lower Sorbian come
1247 from the 2021 edition of the WMT Shared Task
1248 (the benchmark relies on the 2022 edition). For
1249 the Caspian branch, we have 271,710 sentences for
1250 Central Kurdish, Mazandarani (less than 100k), and
1251 Persian. For Gilaki, we use the other dataset avail-
1252 able from the Leipzig corpus: the 2017 community
1253 dataset. We make sure to pre-process the datasets
1254 before applying the 100k threshold to avoid inter-
1255 nal duplicates, but also overlaps with our created
1256 benchmark sentences.

1257 G Complete mining results

1258 Tables 11 and 12 display the F-scores for all models
1259 tested on the eight main language pairs and English
1260 versions, respectively.

1261 H Case study with a bilingual lexicon

1262 We also study strategies to integrate a realistically
1263 and commonly available bilingual source for low-
1264 resource language pairs: a dictionary between the
1265 source and target languages. We consider G500-
1266 LaBSE as a student model and LaBSE as a teacher
1267 model. We use the Chuvash-Russian dictionary
1268 (63k entries) provided during the WMT Shared
1269 Task (Libovický and Fraser, 2021) and evaluate its
1270 impact on two language pairs, CHV-RU and CHV-
1271 EN. We try two approaches: either with words
1272 (and phrases) given directly as a pair (dict-word) or
1273 with ‘code-switched’ Russian sentences, where we
1274 replace words covered by the dictionary in Chuvash
1275 (dict-sent). This allows us to find which method
1276 leads to a better representation from a lexicon.

1277 While using the word lexicon directly improves
1278 the mining score for both language pairs in Ta-

LM	OCI-ES	HSB-DE	DSB-DE	CHV-RU	COS-FR	XMF-KA	GLK-FA	MZN-FA
XLM-R	47.44	0.40	0.42	0.64	16.06	4.55	1.64	3.84
G500	79.96	21.33	2.97	39.59	49.94	41.25	4.06	9.92
XLM-R+CBIE	81.66	16.11	10.18	2.46	39.62	14.29	12.17	28.54
G500+CBIE	92.43	48.01	20.29	48.77	69.80	57.37	19.54	43.25
mSC-XLM-en	62.64	13.45	11.91	7.11	52.44	13.01	20.46	30.37
mSC-XLM-multi	61.93	18.40	11.70	3.53	47.95	7.84	21.77	24.05
mSC-XLM-tr	67.95	21.28	15.79	2.50	44.18	5.46	18.64	24.77
mSC-G500-en	73.95	22.38	11.32	36.57	51.66	56.12	32.89	47.07
mSC-G500-multi	83.24	53.34	27.02	57.58	68.93	57.54	31.77	48.50
mSC-G500-tr	86.13	63.15	27.80	64.90	69.17	62.75	33.23	49.63
LaBSE	98.20	74.87	69.60	30.85	90.32	28.47	45.49	56.64
XLM-LaBSE	96.05	43.64	39.06	24.57	69.71	18.05	44.32	58.93
G500-LaBSE	97.11	78.31	54.57	80.68	83.70	71.11	54.96	71.71
G500-LaBSE-all	97.07	89.67	67.29	85.17	85.99	76.92	59.30	70.39
G500-LaBSE-un	97.63	87.83	56.98	82.56	84.76	73.84	39.20	64.43

Table 11: All F-scores (%) on the test set of the eight main language pairs for all models. mSC stands for mSimCSE, while G500 designates Glot500-m.

LM	OCI-EN	HSB-EN	DSB-EN	CHV-EN	COS-EN	XMF-EN	GLK-EN	MZN-EN
XLM-R	35.72	1.10	0.20	0.35	12.75	0.00	0.00	0.00
G500	62.59	14.03	1.45	27.61	38.34	2.44	0.25	0.31
XLM-R+CBIE	71.77	9.82	6.97	0.47	34.30	1.62	1.31	1.29
G500+CBIE	87.10	33.51	14.51	34.13	65.13	8.89	1.50	4.42
mSC-XLM-en	50.24	10.46	8.55	2.11	43.31	2.92	4.39	3.31
mSC-XLM-multi	39.45	5.82	3.08	1.24	27.99	1.96	1.86	0.41
mSC-XLM-tr	63.69	20.12	13.61	1.47	43.85	2.16	7.62	10.01
mSC-G500-en	34.54	4.65	2.52	5.18	26.01	5.24	4.91	5.72
mSC-G500-multi	74.83	41.10	18.42	46.68	61.44	29.91	5.23	5.32
mSC-G500-tr	81.40	55.28	24.57	55.18	60.79	30.63	18.04	21.83
LaBSE	97.04	77.43	74.69	31.35	90.39	27.62	28.68	32.63
XLM-LaBSE	95.69	44.20	40.89	23.34	70.53	21.05	26.41	31.69
G500-LaBSE	96.92	77.98	54.56	79.25	85.58	59.61	32.51	39.61
G500-LaBSE-all	97.12	88.26	66.41	85.07	85.75	62.92	32.06	40.16
G500-LaBSE-un	97.24	83.72	52.97	79.80	83.81	60.70	20.41	30.60

Table 12: All F-scores (%) on the test set of the eight source languages paired with **English** for all models. mSC stands for mSimCSE, while G500 designates Glot500-m.