

WikiMatrix: Mining 135M Parallel Sentences in 1620 Language Pairs from Wikipedia

Holger Schwenk Facebook AI **Vishrav Chaudhary** Facebook AI **Shuo Sun** Johns Hopkins University **Hongyu Gong** Univ. of Illinois Urbana-Champaign **Francisco Guzmán** Facebook AI

schwenk@fb.com

vishrav@fb.com

ssun32@jhu.edu

hgong6@illinois.edu

fguzman@fb.com

Abstract

We present an approach based on multilingual sentence embeddings to automatically extract parallel sentences from the content of Wikipedia articles in 96 languages, including several dialects or low-resource languages. We systematically consider all possible language pairs. In total, we are able to extract 135M parallel sentences for 1620 different language pairs, out of which only 34M are aligned with English. This corpus is freely available.¹

To get an indication on the quality of the extracted bitexts, we train neural MT baseline systems on the mined data only for 1886 languages pairs, and evaluate them on the TED corpus, achieving strong BLEU scores for many language pairs. The WikiMatrix bitexts seem to be particularly interesting to train MT systems between distant languages without the need to pivot through English.

1 Introduction

Most of the current approaches in Natural Language Processing are data-driven. The size of the resources used for training is often the primary concern, but the quality and a large variety of topics may be equally important. Monolingual texts are usually available in huge amounts for many topics and languages. However, multilingual resources, i.e. sentences which are mutual translations, are more limited, in particular when the two languages do not involve English. An important source of parallel texts is from international organizations like the European Parliament (Koehn, 2005) or the United Nations (Ziems et al., 2016). Several projects rely on volunteers to provide translations for public texts, e.g. news commentary (Tiedemann, 2012), OpensubTitles (Lison and Tiedemann, 2016) or the TED corpus (Qi et al., 2018)

Wikipedia is probably the largest free multilingual resource on the Internet. The content of Wikipedia is very diverse and covers many topics. Articles exist in more than 300 languages. Some content on Wikipedia was human translated from an existing article into another language, not necessarily from or into English. Eventually, the translated articles have been later independently edited and are not parallel anymore. Wikipedia strongly discourages the use of unedited machine translation,² but the existence of such articles cannot be totally excluded. Many articles have been written independently, but may nevertheless contain sentences which are mutual translations. This makes Wikipedia a very appropriate resource to mine for parallel texts for a large number of language pairs. To the best of our knowledge, this is the first work to process the entire Wikipedia and systematically mine for parallel sentences in all language pairs.

In this work, we build on a recent approach to mine parallel texts based on a distance measure in a joint multilingual sentence embedding space (Schwenk, 2018; Artetxe and Schwenk, 2018a), and a freely available encoder for 93 languages. We approach the computational challenge to mine in almost six hundred million sentences by using fast indexing and similarity search algorithms.

The paper is organized as follows. In the next section, we first discuss related work. We then summarize the underlying mining approach. Section 4 describes in detail how we applied this approach to extract parallel sentences from Wikipedia in 1620 language pairs. In section 5, we assess the quality of the extracted bitexts by training NMT systems for a subset of language pairs and evaluate them on the TED corpus (Qi et al., 2018) for 45 languages. The paper concludes with a discussion of future research directions.

¹<https://github.com/facebookresearch/LASER/tree/master/tasks/WikiMatrix>

²<https://en.wikipedia.org/wiki/Wikipedia:Translation>

2 Related work

There is a large body of research on mining parallel sentences in monolingual texts collections, usually named “*comparable corpora*”. Initial approaches to bitext mining have relied on heavily engineered systems often based on metadata information, e.g. (Resnik, 1999; Resnik and Smith, 2003). More recent methods explore the textual content of the comparable documents. For instance, it was proposed to rely on cross-lingual document retrieval, e.g. (Utiyama and Isahara, 2003; Munteanu and Marcu, 2005) or machine translation, e.g. (Abdul-Rauf and Schwenk, 2009; Bouamor and Sajjad, 2018), typically to obtain an initial alignment that is then further filtered. In the shared task for bilingual document alignment (Buck and Koehn, 2016), many participants used techniques based on n -gram or neural language models, neural translation models and bag-of-words lexical translation probabilities for scoring candidate document pairs. The STACC method uses seed lexical translations induced from IBM alignments, which are combined with set expansion operations to score translation candidates through the Jaccard similarity coefficient (Etchegoyhen and Azpeitia, 2016; Azpeitia et al., 2017, 2018). Using multilingual noisy web-crawls such as ParaCrawl³ for filtering good quality sentence pairs has been explored in the shared tasks for high resource (Koehn et al., 2018) and low resource (Koehn et al., 2019) languages.

In this work, we rely on massively multilingual sentence embeddings and margin-based mining in the joint embedding space, as described in (Schwenk, 2018; Artetxe and Schwenk, 2018a,b). This approach has also proven to perform best in a low resource scenario (Chaudhary et al., 2019; Koehn et al., 2019). Closest to this approach is the research described in España-Bonet et al. (2017); Hassan et al. (2018); Guo et al. (2018); Yang et al. (2019). However, in all these works, only bilingual sentence representations have been trained. Such an approach does not scale to many languages, in particular when considering all possible language pairs in Wikipedia. Finally, related ideas have been also proposed in Bouamor and Sajjad (2018) or Grégoire and Langlais (2017). However, in those works, mining is not solely based on multilingual sentence embeddings, but they are part of a larger system. To the best of our knowledge, this work is the first one that applies the same mining approach

to all combinations of many different languages, written in more than twenty different scripts. In follow up work, the same underlying mining approach was applied to a huge collection of Common Crawl texts (Schwenk et al., 2019). Hierarchical mining in Common Crawl texts was performed by El-Kishky et al. (2020).⁴

Wikipedia is arguably the largest comparable corpus. One of the first attempts to exploit this resource was performed by Adafre and de Rijke (2006). An MT system was used to translate Dutch sentences into English and to compare them with the English texts, yielding several hundreds of Dutch/English bitexts. Later, a similar technique was applied to Persian/English (Mohammadi and GhasemAghaee, 2010). Structural information in Wikipedia such as the topic categories of documents was used in the alignment of multilingual corpora (Otero and López, 2010). In another work, the mining approach of Munteanu and Marcu (2005) was applied to extract large corpora from Wikipedia in sixteen languages (Smith et al., 2010). Otero et al. (2011) measured the comparability of Wikipedia corpora by the translation equivalents on three languages Portuguese, Spanish, and English. Patry and Langlais (2011) came up with a set of features such as Wikipedia entities to recognize parallel documents, and their approach was limited to a bilingual setting. Tufis et al. (2013) proposed an approach to mine bitexts from Wikipedia textual content, but they only considered high-resource languages, namely German, Spanish and Romanian paired with English. Tsai and Roth (2016) grounded multilingual mentions to English Wikipedia by training cross-lingual embeddings on twelve languages. Gottschalk and Demidova (2017) searched for parallel text passages in Wikipedia by comparing their named entities and time expressions. Finally, Aghaebrahimian (2018) propose an approach based on bilingual BiLSTM sentence encoders to mine German, French and Persian parallel texts with English. Parallel data consisting of aligned Wikipedia titles have been extracted for twenty-three languages.⁵ We are not aware of other attempts to systematically mine for parallel sentences in the textual content of Wikipedia for a large number of languages.

⁴<http://www.statmt.org/cc-aligned/>

⁵<https://linguatools.org/tools/corpora/wikipedia-parallel-titles-corpora/>

³<http://www.paracrawl.eu/>

3 Distance-based mining approach

The underlying idea of the mining approach used in this work is to first learn a multilingual sentence embedding. The distance in that space can be used as an indicator of whether two sentences are mutual translations or not. Using a simple absolute threshold on the cosine distance was shown to achieve competitive results (Schwenk, 2018). However, it has been observed that an absolute threshold on the cosine distance is globally not consistent, e.g. (Guo et al., 2018). This is particularly true when mining bitexts for many different language pairs.

3.1 Margin criterion

The alignment quality can be substantially improved by using a margin criterion (Artetxe and Schwenk, 2018a). The margin between two candidate sentences x and y is defined as the ratio between the cosine distance between the two sentence embeddings, and the average cosine similarity of its nearest neighbors in both directions:

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

where $\text{NN}_k(x)$ denotes the k unique nearest neighbors of x in the other language, and analogously for $\text{NN}_k(y)$. We used $k = 4$ in all experiments.

We follow the “*max*” strategy of Artetxe and Schwenk (2018a): the margin is first calculated in both directions for all sentences in language L_1 and L_2 . We then create the union of these forward and backward candidates. Candidates are sorted and pairs with source or target sentences that were already used are omitted. We then apply a threshold on the margin score to decide whether two sentences are mutual translations or not.

The complexity of a distance-based mining approach is $O(N \times M)$, where N and M are the number of sentences in each monolingual corpus. This makes a brute-force approach with exhaustive distance calculations intractable for large corpora. The languages with the largest Wikipedia are English and German with 134M and 51M sentences, respectively. This would require 6.8×10^{15} distance calculations.⁶ We show in Section 3.3 how to tackle this computational challenge.

⁶Strictly speaking, Cebuano and Swedish are larger than German, yet mostly consist of template/machine translated text

3.2 Multilingual sentence embeddings

Distance-based bitext mining requires a joint sentence embedding for all the considered languages. One may be tempted to train a bi-lingual embedding for each language pair, e.g. (Española-Bonet et al., 2017; Hassan et al., 2018; Guo et al., 2018; Yang et al., 2019), but this is difficult to scale to thousands of language pairs present in Wikipedia. Instead, we chose to use one single massively multilingual sentence embedding for all languages, namely the one proposed by the open-source LASER toolkit (Artetxe and Schwenk, 2018b). Training one joint multilingual embedding on many languages at once also has the advantage that low-resource languages can benefit from the similarity to other languages in the same language family. For example, we were able to mine parallel data for several Romance (minority) languages like Aragonese, Lombard, Mirandese or Sicilian although data in those languages was not used to train the multilingual LASER embeddings. The reader is referred to Artetxe and Schwenk (2018b) for a detailed description how LASER was trained.

3.3 Fast similarity search

In this work, we use the open-source FAISS library⁷ which implements highly efficient algorithms to perform similarity search on billions of vectors (Johnson et al., 2017). Our sentence representations being 1024-dimensional, all English sentences require $134 \cdot 10^6 \times 1024 \times 4 = 536$ GB of memory. Therefore, dimensionality reduction and data compression are needed for efficient search. We chose a rather aggressive compression based on a 64-bit product-quantizer (Jégou et al., 2011), and partitioning the search space in 32k cells.⁸ We build one FAISS index for each language.

The compressed FAISS index for English requires only 9.2GB, i.e. more than fifty times smaller than the original sentences embeddings. This makes it possible to load the whole index on a standard GPU and to run the search in a very efficient way on multiple GPUs in parallel, without the need to shard the index. The overall mining process for German/English requires less than 3.5 hours on 8 GPUs, including the nearest neighbor search in both direction and scoring all candidates.

⁷<https://github.com/facebookresearch/faiss>

⁸FAISS index type OPQ64, IVF32768, PQ64, see <https://github.com/facebookresearch/faiss/wiki/Faiss-indexes>

4 Bitext mining in Wikipedia

For each Wikipedia article, it is possible to get the link to the corresponding article in other languages. This could be used to mine sentences limited to the respective articles. On one hand, this **local mining** has several advantages: 1) mining is very fast since each article usually has a few hundreds of sentences only; 2) it seems reasonable to assume that a translation of a sentence is more likely to be found in the same article than anywhere in the whole Wikipedia. On the other hand, we hypothesize that the margin criterion will be less efficient since one article usually has few sentences which are similar. This may lead to many sentences in the overall mined corpus of the type “*NAME was born on DATE in CITY*”, “*BUILDING is a monument in CITY built on DATE*”, etc. Although those alignments may be correct, we hypothesize that they are of limited use to train an NMT system.

The other option is to consider the whole Wikipedia for each language: for each sentence in the source language, we mine in all target sentences. This **global mining** has the advantages that we can try to align two languages even though there are only a few articles in common. A drawback is a potentially increased risk of misalignment. In this work, we chose the global mining option.

4.1 Corpus preparation

Extracting the textual content of Wikipedia articles in all languages is a rather challenging task, i.e. removing all tables, citations, footnotes or formatting markup. There are several ways to download Wikipedia content. In this study, we use the so-called `CirrusSearch` dumps since they directly provide the textual content without any meta information.⁹ We downloaded this dump in March 2019. A total of about 300 languages are available, but the size obviously varies a lot between languages. We applied the following processing: 1) extract the textual content; 2) split the paragraphs into sentences; 3) remove duplicate sentences; and 4) perform language identification and remove sentences which are not in the expected language.

It should be pointed out that sentence segmentation is not a trivial task. Some languages do not use specific symbols to mark the end of a sentence, namely Thai. We are not aware of a freely available sentence segmenter for Thai and we had to exclude

⁹<https://dumps.wikimedia.org/other/cirrussearch/>

L_1 (French)	<i>Ceci est une très grande maison</i>
L_2 (German)	<i>Das ist ein sehr großes Haus</i> <i>This is a very big house</i> <i>Ez egy nagyon nagy ház</i> <i>Ini rumah yang sangat besar</i>

Table 1: Illustration how sentences in the wrong language can hurt the alignment process with a margin criterion. See text for a detailed discussion.

it. We used a freely available Python tool¹⁰ to detect sentence boundaries. Regular expressions were used for most of the Asian languages, falling back to English for the remaining languages. This gives us 879 million sentences in 300 languages. The margin criterion to mine for parallel data requires that the texts do not contain duplicates. This removes about 25% of the sentences.¹¹

LASER’s sentence embeddings are totally language agnostic. This has the side effect that the sentences in other languages (e.g. citations or quotes) may be considered closer in the embedding space than a potential translation in the target language. Table 1 illustrates this problem. The algorithm would not select the German sentence although it is a perfect translation. The sentences in the other languages are also valid translations which would yield a very small margin. To avoid this problem, we perform language identification (LID) on all sentences and remove those which are not in the expected language. LID is performed with `fasttext`¹² (Joulin et al., 2016). `Fasttext` does not support all the 300 languages present in Wikipedia and we disregarded the missing ones (which typically have only a few sentences anyway). After deduplication and LID, we dispose of 595M sentences in 182 languages. English accounts for 134M sentences, and German with 51M sentences is the second largest language. The sizes for all languages are given in Tables 3 and 5 (in the appendix).

4.2 Threshold optimization

Artetxe and Schwenk (2018a) optimized their mining approach for each language pair on a provided corpus of gold alignments. This is not possible when mining Wikipedia, in particular when con-

¹⁰<https://pypi.org/project/sentence-splitter/>

¹¹The Cebuano and Waray Wikipedia were largely created by a bot and contain more than 65% of duplicates.

¹²<https://fasttext.cc/docs/en/language-identification.html>

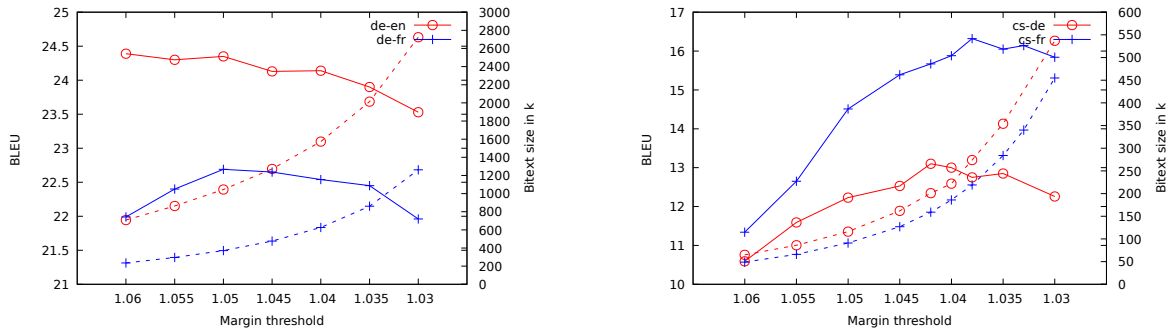


Figure 1: BLEU scores (continuous lines) for several NMT systems trained on bitexts extracted from Wikipedia for different margin thresholds. The size of the mined bitexts are depicted as dashed lines.

sidering many language pairs. In this work, we use an evaluation protocol inspired by the WMT shared task on parallel corpus filtering for low-resource conditions (Koehn et al., 2019): an NMT system is trained on the extracted bitexts – for different thresholds – and the resulting BLEU scores are compared. We choose *newstest2014* of the WMT evaluations since it provides an N -way parallel test sets for English, French, German and Czech. We favoured the translation between two morphologically rich languages from different families and considered the following language pairs: German/English, German/French, Czech/German and Czech/French. The size of the mined bitexts is in the range of 100k to more than 2M (see Table 2 and Figure 1). We did not try to optimize the architecture of the NMT system to the size of the bitexts and used the same architecture for all systems: the encoder and decoder are 5-layer transformer models as implemented in *fairseq* (Ott et al., 2019). The goal of this study is not to develop the best performing NMT system for the considered languages pairs, but to compare different mining parameters.

The evolution of the BLEU score in function of the margin threshold is given in Figure 1. Decreasing the threshold naturally leads to more mined data – we observe an exponential increase of the data size. The performance of the NMT systems trained on the mined data seems to change as expected: the BLEU score first improves with increasing amounts of available training data, reaches a maximum and then decreases since the additional data gets more and more noisy, i.e. contains wrong translations. It is also not surprising that a careful choice of the margin threshold is more important in a low-resource setting. Every additional parallel sentence is important. According to Figure 1, the optimal value of the margin threshold seems to be

Bitexts	de-en	de-fr	cs-de	cs-fr
Europarl	1.9M	1.9M	568k	627k
	21.5	23.6	14.9	21.5
Mined Wikipedia	1.0M	370k	200k	220k
	21.2	21.1	12.6	19.2
Europarl + Wikipedia	3.0M	2.3M	768k	846k
	25.5	25.6	17.7	24.0

Table 2: Comparison of NMT systems trained on the Europarl corpus and on bitexts automatically mined in Wikipedia by our approach at a threshold of 1.04. We give the number of sentences (first line) and the BLEU score (second line of each bloc) on *newstest2014*.

1.05 when many sentences can be extracted, in our case German/English and German/French. When less parallel data is available, i.e. Czech/German and Czech/French, a value in the range of 1.03–1.04 seems to be a better choice. Aiming at one threshold for all language pairs, we chose a value of 1.04. It seems to be a good compromise for most language pairs. However, for the open release of this corpus, we provide all mined sentence with a margin of 1.02 or better. This enables end users to choose an optimal threshold for their particular applications. However, it should be emphasized that we do not expect that many sentence pairs with a margin as low as 1.02 are good translations.

For comparison, we also trained NMT systems on the Europarl corpus V7 (Koehn, 2005), i.e. professional human translations, first on all available data, and then on the same number of sentences as the mined ones (see Table 2). With the exception of Czech/French, we were able to achieve better

BLEU scores with the mined bitexts in Wikipedia than with Europarl of the same size. Adding the mined bitexts to the full Europarl corpus, leads to further improvements of 1.1 to 3.1 BLEU.

5 Result analysis

We run the alignment process for all possible combinations of languages in Wikipedia. This yielded 1620 language pairs for which we were able to mine at least ten thousand sentences. Remember that mining $L_1 \rightarrow L_2$ is identical to $L_2 \rightarrow L_1$, and is counted only once. We propose to analyze and evaluate the extracted bitexts in two ways. First, we discuss the amount of extracted sentences (Section 5.1). We then turn to a qualitative assessment by training NMT systems for all language pairs with more than twenty-five thousand mined sentences (Section 5.2).

5.1 Quantitative analysis

Due to space limits, Table 3 summarizes the number of extracted parallel sentences only for languages which have a total of at least five hundred thousand parallel sentences (with all other languages at a margin threshold of 1.04). Additional results are given in Table 5 in the Appendix.

There are many reasons which can influence the number of mined sentences. Obviously, the larger the monolingual texts, the more likely it is to mine many parallel sentences. Not surprisingly, we observe that more sentences could be mined when English is one of the two languages. Let us point out some languages for which it is usually not obvious to find parallel data with English, namely Indonesian (1M), Hebrew (545k), Farsi (303k) or Marathi (124k sentences). The largest mined texts not involving English are Russian/Ukrainian (2.5M), Catalan/Spanish (1.6M), or between the Romance languages French, Spanish, Italian and Portuguese (480k–923k), and German/French (626k).

It is striking to see that we were able to mine more sentences when Galician and Catalan are paired with Spanish than with English. On one hand, this could be explained by the fact that LASER’s multilingual sentence embeddings may be better since the involved languages are linguistically very similar. On the other, it could be that the Wikipedia articles in both languages share a lot of content, or are obtained by mutual translation.

Services from the European Commission provide human translations of (legal) texts in all the

24 official languages of the European Union. This N-way parallel corpus enables training of MT system to directly translate between these languages, without the need to pivot through English. This is usually not the case when translating between other major languages, for example in Asia. Some interesting language pairs for which we were able to mine more than 100k sentences include: Korean/Japanese (222k), Russian/Japanese (196k), Indonesian/Vietnamese (146k), or Hebrew/Romance languages (120–150k sentences).

Overall, we were able to extract at least ten thousand parallel sentences for 96 different languages. For several low-resource languages, we were able to extract more parallel sentences with other languages than English. These include, among others, Aragonese with Spanish, Lombard with Italian, Breton with several Romance languages, Western Frisian with Dutch, Luxembourgish with German or Egyptian Arabic and Wu Chinese with the respective major language.

Finally, Cebuano (ceb) falls clearly apart: it has a rather huge Wikipedia (17.9M filtered sentences), but most of it was generated by a bot, as for the Waray language.¹³ This certainly explains that only a very small number of parallel sentences could be extracted. Although the same bot was also used to generate articles in the Swedish Wikipedia, our alignments seem to be better for that language.

5.2 Qualitative evaluation

Aiming to perform a large-scale assessment of the quality of the extracted parallel sentences, we trained NMT systems on the bitexts. We identified a publicly available dataset which provides test sets for many language pairs: translations of TED talks as proposed in the context of a study on pretrained word embeddings for NMT¹⁴ (Qi et al., 2018). We would like to emphasize that we did not use the training data provided by TED – we only trained on the mined sentences from Wikipedia. The goal of this study is not to build state-of-the-art NMT system for for the TED task, but to get an estimate of the quality of our extracted data, for many language pairs. In particular, there may be a mismatch in the topic and language style between Wikipedia texts and the transcribed and translated TED talks.

NMT systems are trained with a transformer model from fairseq (Ott et al., 2019) with the

¹³<https://en.wikipedia.org/wiki/Lsjobot>

¹⁴<https://github.com/neulab/word-embeddings-for-nmt>

References

- Sadaf Abdul-Rauf and Holger Schwenk. 2009. **On the Use of Comparable Corpora to Improve SMT performance**. In *EACL*, pages 16–23.
- Sisay Fissaha Adafre and Maarten de Rijke. 2006. Finding similar sentences across multiple languages in Wikipedia. In *Proceedings of the Workshop on NEW TEXT Wikis and blogs and other dynamic text sources*.
- Ahmad Aghaebrahimian. 2018. Deep neural networks at the service of multilingual parallel sentence extraction. In *Coling*.
- Mikel Artetxe and Holger Schwenk. 2018a. Margin-based Parallel Corpus Mining with Multilingual Sentence Embeddings. <https://arxiv.org/abs/1811.01136>.
- Mikel Artetxe and Holger Schwenk. 2018b. Massively multilingual sentence embeddings for zero-shot cross-lingual transfer and beyond. In <https://arxiv.org/abs/1812.10464>.
- Andoni Azpeitia, Thierry Etchegoyhen, and Eva Martínez García. 2017. **Weighted Set-Theoretic Alignment of Comparable Sentences**. In *BUCC*, pages 41–45.
- Andoni Azpeitia, Thierry Etchegoyhen, and Eva Martínez García. 2018. Extracting Parallel Sentences from Comparable Corpora with STACC Variants. In *BUCC*.
- Houda Bouamor and Hassan Sajjad. 2018. H2@BUCC18: Parallel Sentence Extraction from Comparable Corpora Using Multilingual Sentence Embeddings. In *BUCC*.
- Christian Buck and Philipp Koehn. 2016. **Findings of the wmt 2016 bilingual document alignment shared task**. In *Proceedings of the First Conference on Machine Translation*, pages 554–563, Berlin, Germany. Association for Computational Linguistics.
- Vishrav Chaudhary, Yuqing Tang, Francisco Guzmán, Holger Schwenk, and Philipp Koehn. 2019. Low-resource corpus filtering using multilingual sentence embeddings. In *Proceedings of the Fourth Conference on Machine Translation (WMT)*.
- Ahmed El-Kishky, Vishrav Chaudhary, Francisco Guzman, and Philipp Koehn. 2020. A massive collection of cross-lingual web-document pairs. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5960–5969.
- Cristina España-Bonet, Ádam Csaba Varga, Alberto Barrón-Cedeño, and Josef van Genabith. 2017. An Empirical Analysis of NMT-Derived Interlingual Embeddings and their Use in Parallel Sentence Identification. *IEEE Journal of Selected Topics in Signal Processing*, pages 1340–1348.
- Thierry Etchegoyhen and Andoni Azpeitia. 2016. **Set-Theoretic Alignment for Comparable Corpora**. In *ACL*, pages 2009–2018.
- Simon Gottschalk and Elena Demidova. 2017. Multitwiki: Interlingual text passage alignment in Wikipedia. *ACM Transactions on the Web (TWEB)*, 11(1):6.
- Francis Grégoire and Philippe Langlais. 2017. **BUCC 2017 Shared Task: a First Attempt Toward a Deep Learning Framework for Identifying Parallel Sentences in Comparable Corpora**. In *BUCC*, pages 46–50.
- Mandy Guo, Qinlan Shen, Yinfei Yang, Heming Ge, Daniel Cer, Gustavo Hernandez Abrego, Keith Stevens, Noah Constant, Yun-Hsuan Sung, Brian Strope, and Ray Kurzweil. 2018. Effective Parallel Corpus Mining using Bilingual Sentence Embeddings. *arXiv:1807.11906*.
- Hany Hassan, Anthony Aue, Chang Chen, Vishal Chowdhary, Jonathan Clark, Christian Federmann, Xuedong Huang, Marcin Junczys-Dowmunt, William Lewis, Mu Li, Shujie Liu, Tie-Yan Liu, Renqian Luo, Arul Menezes, Tao Qin, Frank Seide, Xu Tan, Fei Tian, Lijun Wu, Shuangzhi Wu, Yingce Xia, Dongdong Zhang, Zhirui Zhang, and Ming Zhou. 2018. Achieving Human Parity on Automatic Chinese to English News Translation. *arXiv:1803.05567*.
- Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2017. Billion-scale similarity search with GPUs. *arXiv preprint arXiv:1702.08734*.
- Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. 2016. Bag of tricks for efficient text classification. <https://arxiv.org/abs/1607.01759>.
- H. Jégou, M. Douze, and C. Schmid. 2011. Product quantization for nearest neighbor search. *IEEE Trans. PAMI*, 33(1):117–128.
- Philipp Koehn. 2005. Europarl: A parallel corpus for statistical machine translation. In *MT summit*.
- Philipp Koehn, Francisco Guzmán, Vishrav Chaudhary, and Juan M. Pino. 2019. Findings of the wmt 2019 shared task on parallel corpus filtering for low-resource conditions. In *Proceedings of the Fourth Conference on Machine Translation, Volume 2: Shared Task Papers*, Florence, Italy. Association for Computational Linguistics.
- Philipp Koehn, Huda Khayrallah, Kenneth Heafield, and Mikel L. Forcada. 2018. **Findings of the wmt 2018 shared task on parallel corpus filtering**. In *Proceedings of the Third Conference on Machine Translation: Shared Task Papers*, pages 726–739, Belgium, Brussels. Association for Computational Linguistics.

- Taku Kudo and John Richardson. 2018. [Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71, Brussels, Belgium. Association for Computational Linguistics.
- P. Lison and J. Tiedemann. 2016. Opensubtitles2016: Extracting large parallel corpora from movie and tv subtitles. In *LREC*.
- Mehdi Zadeh Mohammadi and Nasser GhasemAghaee. 2010. Building bilingual parallel corpora based on Wikipedia. In *2010 Second International Conference on Computer Engineering and Applications*, pages 264–268.
- Dragos Stefan Munteanu and Daniel Marcu. 2005. [Improving Machine Translation Performance by Exploiting Non-Parallel Corpora](#). *Computational Linguistics*, 31(4):477–504.
- P Otero, I López, S Cilenis, and Santiago de Compostela. 2011. Measuring comparability of multilingual corpora extracted from Wikipedia. *Iberian Cross-Language Natural Language Processings Tasks (ICL)*, page 8.
- Pablo Gamallo Otero and Isaac González López. 2010. Wikipedia as multilingual source of comparable corpora. In *Proceedings of the 3rd Workshop on Building and Using Comparable Corpora, LREC*, pages 21–25.
- Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. [fairseq: A fast, extensible toolkit for sequence modeling](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pages 48–53, Minneapolis, Minnesota. Association for Computational Linguistics.
- Alexandre Patry and Philippe Langlais. 2011. Identifying parallel documents from a large bilingual collection of texts: Application to parallel article extraction in Wikipedia. In *Proceedings of the 4th Workshop on Building and Using Comparable Corpora: Comparable Corpora and the Web*, pages 87–95. Association for Computational Linguistics.
- Matt Post. 2018. [A call for clarity in reporting bleu scores](#). In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Belgium, Brussels. Association for Computational Linguistics.
- Ye Qi, Devendra Sachan, Matthieu Felix, Sarguna Padmanabhan, and Graham Neubig. 2018. [When and why are pre-trained word embeddings useful for neural machine translation?](#) In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 529–535, New Orleans, Louisiana. Association for Computational Linguistics.
- Philip Resnik. 1999. [Mining the Web for Bilingual Text](#). In *ACL*.
- Philip Resnik and Noah A. Smith. 2003. [The Web as a Parallel Corpus](#). *Computational Linguistics*, 29(3):349–380.
- Holger Schwenk. 2018. Filtering and mining parallel data in a joint multilingual space. In *ACL*, pages 228–234.
- Holger Schwenk, Guillaume Wenzek, Sergey Edunov, Edouard Grave, Armand Joulin, and Angela Fan. 2019. CCMatrix: Mining billions of high-quality parallel sentences on the web. In <https://arxiv.org/abs/1911.04944>.
- Jason R. Smith, Chris Quirk, and Kristina Toutanova. 2010. Extracting parallel sentences from comparable corpora using document level alignment. In *NAACL*, pages 403–411.
- J. Tiedemann. 2012. Parallel data, tools and interfaces in OPUS. In *LREC*.
- Chen-Tse Tsai and Dan Roth. 2016. Cross-lingual wikification using multilingual embeddings. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 589–598.
- Dan Tufis, Radu Ion, Ștefan Daniel, Dumitrescu, and Dan Ștefănescu. 2013. Wikipedia as an smt training corpus. In *RANLP*, pages 702–709.
- Masao Utiyama and Hitoshi Isahara. 2003. [Reliable Measures for Aligning Japanese-English News Articles and Sentences](#). In *ACL*.
- Yinfei Yang, Gustavo Hernández Ábrego, Steve Yuan, Mandy Guo, Qinlan Shen, Daniel Cer, Yun-Hsuan Sung, Brian Strope, and Ray Kurzweil. 2019. Improving multilingual sentence embedding using bi-directional dual encoder with additive margin softmax. In <https://arxiv.org/abs/1902.08564>.
- Michał Ziemski, Marcin Junczys-Dowmunt, and Bruno Pouliquen. 2016. The United Nations Parallel Corpus v1.0. In *LREC*.

A Appendix

Table 5 provides the amounts of mined parallel sentences for languages which have a rather small Wikipedia. Aligning those languages obviously yields to a very small amount of parallel sentences. Therefore, we only provide these results for alignment with high resource languages. It is also likely that several of these alignments are of low quality since the LASER embeddings were not directly trained on most these languages, but we still hope to achieve reasonable results since other languages of the same family may be covered.

ISO	Name	Language Family	size	ca	da	de	en	es	fr	it	nl	pl	pt	sv	ru	zh	total
an	Aragonese	Romance	222	24	7	12	23	33	16	13	9	10	14	9	11	6	324
arz	Egyptian Arabic	Arabic	120	7	6	11	18	12	12	10	8	9	10	8	12	7	278
as	Assamese	Indo-Aryan	124	8	6	11	7	11	12	10	9	9	8	8	9	3	216
azb	South Azer-Turkic bajani	Turkic	398	6	4	9	8	9	10	9	7	6	8	6	7	3	172
bar	Bavarian	Germanic	214	7	6	41	16	12	12	10	8	9	10	8	10	5	261
bpy	Bishnupriya	Indo-Aryan	128	2	1	4	4	3	4	2	2	3	2	2	3	1	71
br	Breton	Celtic	413			20	16	22	23	22			19	16	6	200	
ce	Chechen	Northeast Caucasian	315	2	1	2	2	2	2	2	2	2	2	2	2	1	56
ceb	Cebuano	Malayo-Polynesian	17919	14	9	22	29	27	24	24	15	17	20	55	21	9	594
ckb	Central Kurdish	Kur-Iranian	127	2	2	6	8	5	5	4	4	4	4	3	6	4	113
cv	Chuvash	Turkic	198	4	3	5	4	6	6	7	5	4	6	5	8	2	129
dv	Maldivian	Indo-Aryan	52	2	2	5	6	4	4	3	3	3	3	3	5	3	96
fo	Faroese	Germanic	114	13	12	14	32	21	18	15	11	11	17	12	13	6	335
fy	Western Frisian	Germanic	493	13	8	16	32	21	18	17	38	12	18	13	14	5	453
gd	Gaelic	Celtic	66	1	1	1	1	1	1	1	1	1	1	1	1	1	41
ga	Irish	Irish	216			2	3	4	3	3	2	2	3	2	3	1	70
gom	Goan Konkani	Indo-Aryan	69	9	7	10	8	13	13	13	9	9	11	9	10	4	240
ht	Haitian Creole	Creole	60	2	1	3	4	3	4	3	2	3	2	2	3	1	72
ilo	Iloko	Philippine	63	3	2	4	5	4	4	4	3	3	4	3	4	2	96
io	Ido	constructed	153	5	3	6	11	7	7	5	5	5	6	5	5	3	143
jv	Javanese	Malayo-Polynesian	220	8	5	8	13	12	10	11	8	7	11	8	8	3	219
ka	Georgian	Kartvelian	480	11	7	15	12	16	17	16	12	11	14	12	13	5	288
ku	Kurdish	Iranian	165	5	4	8	5	8	7	8	7	6	7	6	6	3	222
la	Latin	Romance	558	12	9	17	32	20	18	17	12	13	18	13	14	6	478
lb	Luxembourgish	Germanic	372	12	7	26	22	19	18	15	11	11	16	12	11	4	305
lmo	Lombard	Romance	147	6	3	7	10	7	7	11	6	5	7	5	5	3	144
mg	Malagasy	Malayo-Polynesian	263	6	5	9	13	9	12	8	7	7	7	8	7	4	199
mhr	Eastern Mari	Uralic	61	3	2	4	3	4	4	5	3	3	4	3	4	2	96
min	Minangkabau	Malayo-Polynesian	255	4	2	6	7	5	5	5	4	4	4	5	5	2	121
mn	Mongolian	Mongolic	255	4	3	7	5	6	6	7	6	5	5	5	5	3	197
mw1	Mirandese	Romance	64	6	3	4	10	8	6	5	3	4	34	3	4	2	154
nds nl	Low German/Saxon	Germanic	65	5	4	6	10	7	7	6	15	5	6	5	5	3	151
ps	Pashto	Iranian	89			2	3	2	3	3	3	3	3	3	3	1	73
rm	Romansh	Italic	57	2	2	10	5	4	4	3	2	3	3	3	3	1	86
sah	Yakut	Turkic/Sib	134	4	3	7	5	6	6	6	5	5	5	5	6	3	134
scn	Sicilian	Romance	81	5	3	6	9	7	7	11	5	5	6	5	5	2	143
sd	Sindhi	Iranian	115			3	9	8	8	7	7	6	7	5	8	5	152
su	Sundanese	Malayo-Polynesian	120	4	3	5	7	6	5	6	4	4	5	4	4	2	117
tk	Turkmen	Turkic	56	2	2	3	3	4	3	4	2	2	4	2	3	1	76
tg	Tajik	Iranian	248	5	4	11	15	9	9	8	8	7	8	6	10	6	192
ug	Uighur	Turkic	83	4	3	9	10	7	8	6	6	5	6	5	9	6	168
ur	Urdu	Indo-Aryan	150	2	2	3	5	3	3	3	3	3	3	3	3	2	123
wa	Walloon	Romance	56	3	2	4	5	5	4	4	3	3	4	3	3	2	93
wuu	Wu Chinese	Chinese	75	8	6	11	17	12	11	10	8	9	11	9	10	43	283
yi	Yiddish	Germanic	131	3	2	4	3	4	4	5	3	3	4	3	4	1	92

Table 5: WikiMatrix (part 2): number of extracted sentences (in thousands) for languages with a rather small Wikipedia. Alignments with other languages yield less than 5k sentences and are omitted for clarity.

Table 2 gives the detailed configuration which was used to train NMT models on the mined data in Section 5. An 5000 subword vocabulary was learnt using SentencePiece (Kudo and Richardson, 2018). Decoding was done with beam size 5 and length normalization 1.2.

```

--arch transformer
--share-all-embeddings
--encoder-layers 5
--decoder-layers 5
--encoder-embed-dim 512
--decoder-embed-dim 512
--encoder-ffn-embed-dim 2048
--decoder-ffn-embed-dim 2048
--encoder-attention-heads 2
--decoder-attention-heads 2
--encoder-normalize-before
--decoder-normalize-before
--dropout 0.4
--attention-dropout 0.2
--relu-dropout 0.2
--weight-decay 0.0001
--label-smoothing 0.2
--criterion label_smoothed_cross_entropy
--optimizer adam
--adam-betas '(0.9, 0.98)'
--clip-norm 0
--lr-scheduler inverse_sqrt
--warmup-update 4000
--warmup-init-lr 1e-7
--lr 1e-3 --min-lr 1e-9
--max-tokens 4000
--update-freq 4
--max-epoch 100
--save-interval 10

```

Figure 2: Model settings for NMT training with fairseq

Finally, Table 6 gives the BLEU scores on the TED corpus when translating into and from English for some additional languages.

Lang	xx → en	en → xx
et	15.9	14.3
eu	10.1	7.6
fa	16.7	8.8
fi	10.9	10.9
lt	13.7	10.0
hi	17.8	21.9
mr	2.6	3.5

Table 6: BLEU scores on the TED test set as proposed in (Qi et al., 2018). NMT systems were trained on bi-texts mined in Wikipedia only. No other resources were used.