

How Language-Neutral is Multilingual BERT?

Jindřich Libovický¹ and Rudolf Rosa² and Alexander Fraser¹

¹Center for Information and Language Processing, LMU Munich, Germany

²Faculty of Mathematics and Physics, Charles University, Prague, Czech Republic

{libovicky, fraser}@cis.lmu.de rosa@ufal.mff.cuni.cz

Abstract

Multilingual BERT (mBERT) provides sentence representations for 104 languages, which are useful for many multi-lingual tasks. Previous work probed the cross-linguality of mBERT using zero-shot transfer learning on morphological and syntactic tasks. We instead focus on the semantic properties of mBERT. We show that mBERT representations can be split into a language-specific component and a language-neutral component, and that the language-neutral component is sufficiently general in terms of modeling semantics to allow high-accuracy word-alignment and sentence retrieval but is not yet good enough for the more difficult task of MT quality estimation. Our work presents interesting challenges which must be solved to build better language-neutral representations, particularly for tasks requiring linguistic transfer of semantics.

1 Introduction

Multilingual BERT (mBERT; [Devlin et al. 2019](#)) is gaining popularity as a contextual representation for various multilingual tasks, such as dependency parsing ([Kondratyuk and Straka, 2019](#); [Wang et al., 2019](#)), cross-lingual natural language inference (XNLI) or named-entity recognition (NER) ([Pires et al., 2019](#); [Wu and Dredze, 2019](#); [Kudugunta et al., 2019](#)).

[Pires et al. \(2019\)](#) present an exploratory paper showing that mBERT can be used cross-lingually for zero-shot transfer in morphological and syntactic tasks, at least for typologically similar languages. They also study an interesting semantic task, sentence-retrieval, with promising initial results. Their work leaves many open questions in terms of how good the cross-lingual mBERT representation is for semantics, motivating our work.

In this paper, we directly assess the semantic cross-lingual properties of mBERT. To avoid

methodological issues with zero-shot transfer (possible language overfitting, hyper-parameter tuning), we selected tasks that only involve a direct comparison of the representations: cross-lingual sentence retrieval, word alignment, and machine translation quality estimation (MT QE). Additionally, we explore how the language is represented in the embeddings by training language identification classifiers and assessing how the representation similarity corresponds to phylogenetic language families.

Our results show that the mBERT representations, even after language-agnostic fine-tuning, are not very language-neutral. However, the identity of the language can be approximated as a constant shift in the representation space. An even higher language-neutrality can still be achieved by a linear projection fitted on a small amount of parallel data.

Finally, we present attempts to strengthen the language-neutral component via fine-tuning: first, for multi-lingual syntactic and morphological analysis; second, towards language identity removal via a adversarial classifier.

2 Related Work

Since the publication of mBERT ([Devlin et al., 2019](#)), many positive experimental results were published.

[Wang et al. \(2019\)](#) reached impressive results in zero-shot dependency parsing. However, the representation used for the parser was a bilingual projection of the contextual embeddings based on word-alignment trained on parallel data.

[Pires et al. \(2019\)](#) recently examined the cross-lingual properties of mBERT on zero-shot NER and part-of-speech (POS) tagging but the success of zero-shot transfer strongly depends on how typologically similar the languages are. Similarly,

Wu and Dredze (2019) trained good multilingual models for POS tagging, NER, and XNLI, but struggled to achieve good results in the zero-shot setup.

Pires et al. (2019) assessed mBERT on cross-lingual sentence retrieval between three language pairs. They observed that if they subtract the average difference between the embeddings from the target language representation, the retrieval accuracy significantly increases. We systematically study this idea in the later sections.

Many experiments show (Wu and Dredze, 2019; Kudugunta et al., 2019; Kondratyuk and Straka, 2019) that downstream task models can extract relevant features from the multilingual representations. But these results do not directly show language-neutrality, i.e., to what extent are similar phenomena are represented similarly across languages. The models can obtain the task-specific information based on the knowledge of the language, which (as we show later) can be easily identified. Our choice of evaluation tasks eliminates this risk by directly comparing the representations. Limited success in zero-shot setups and the need for explicit bilingual projection in order to work well (Pires et al., 2019; Wu and Dredze, 2019; Rönqvist et al., 2019) also shows limited language neutrality of mBERT.

3 Centering mBERT Representations

Following Pires et al. (2019), we hypothesize that a sentence representation in mBERT is composed of a language-specific component, which identifies the language of the sentence, and a language-neutral component, which captures the meaning of the sentence in a language-independent way. We assume that the language-specific component is similar across all sentences in the language.

We thus try to remove the language-specific information from the representations by centering the representations of sentences in each language so that their average lies at the origin of the vector space. We do this by estimating the language centroid as the mean of the mBERT representations for a set of sentences in that language and subtracting the language centroid from the contextual embeddings.

We then analyze the semantic properties of both the original and the centered representations using a range of probing tasks. For all tasks, we test all layers of the model. For tasks utilizing a

single-vector sentence representation, we test both the vector corresponding to the `[cls]` token and mean-pooled states.

4 Probing Tasks

We employ five probing tasks to evaluate the language neutrality of the representations.

Language Identification. With a representation that captures all phenomena in a language-neutral way, it should be difficult to determine what language the sentence is written in. Unlike other tasks, language identification does require fitting a classifier. We train a linear classifier on top of a sentence representation to try to classify the language of the sentence.

Language Similarity. Experiments with POS tagging (Pires et al., 2019) suggest that similar languages tend to get similar representations on average. We quantify that observation by measuring how languages tend to cluster by the language families using V-measure over hierarchical clustering of the language centroid (Rosenberg and Hirschberg, 2007).

Parallel Sentence Retrieval. For each sentence in a multi-parallel corpus, we compute the cosine distance of its representation with representations of all sentences on the parallel side of the corpus and select the sentence with the smallest distance.

Besides the plain and centered `[cls]` and mean-pooled representations, we evaluate explicit projection into the “English space”. For each language, we fit a linear regression projecting the representations into English representation space using a small set of parallel sentences.

Word Alignment. While sentence retrieval could be done with keyword spotting, computing bilingual alignment requires resolving detailed correspondence on the word level.

We find the word alignment as a minimum weighted edge cover of a bipartite graph. The graph connects the tokens of the sentences in the two languages and edges between them are weighted with the cosine distance of the token representation. Tokens that get split into multiple subwords are represented using the average of the embeddings of the subwords. Note that this algorithm is invariant to representation centering which would only change the edge weights by a constant offset.

We evaluate the alignment using the F_1 score over both sure and possible alignment links in a manually aligned gold standard.

MT Quality Estimation. MT QE assesses the quality of an MT system output without having access to a reference translation.

The standard evaluation metric is the correlation with the Human-targeted Translation Error Rate which is the number of edit operations a human translator would need to do to correct the system output. This is a more challenging task than the two previous ones because it requires capturing more fine-grained differences in meaning.

We evaluate how cosine distance of the representation of the source sentence and of the MT output reflects the translation quality. In addition to plain and centered representations, we also test trained bilingual projection, and a fully supervised regression trained on training data.

5 Experimental Setup

We use a pre-trained mBERT model that was made public with the BERT release¹. The model dimension is 768, hidden layer dimension 3072, self-attention uses 12 heads, the model has 12 layers. It uses a vocabulary of 120k wordpieces that is shared for all languages.

To train the language identification classifier, for each of the BERT languages we randomly selected 110k sentences of at least 20 characters from Wikipedia, and keep 5k for validation and 5k for testing for each language. The training data are also used for estimating the language centroids.

For parallel sentence retrieval, we use a multi-parallel corpus of test data from the WMT14 evaluation campaign (Bojar et al., 2014) with 3,000 sentences in Czech, English, French, German, Hindi, and Russian. The linear projection experiment uses the WMT14 development data.

We use manually annotated word alignment datasets to evaluate word alignment between English on one side and Czech (2.5k sent.; Mareček, 2016), Swedish (192 sent.; Holmqvist and Ahrenberg, 2011), German (508 sent.), French (447 sent.; Och and Ney, 2000) and Romanian (248 sent.; Mihalcea and Pedersen, 2003) on the other side. We compare the results with FastAlign (Dyer et al., 2013) that was provided with 1M additional parallel sentences from ParaCrawl (Esplà

¹<https://github.com/google-research/bert>

| | mBERT | UDify | Ing-free |
|------------------|-------|-------|----------|
| [cls] | .935 | .938 | .796 |
| [cls], cent. | .867 | .851 | .337 |
| mean-pool | .919 | .896 | .230 |
| mean-pool, cent. | .285 | .243 | .247 |

Table 1: Accuracy of language identification, values from the best-scoring layers.



Figure 1: Language centroids of the mean-pooled representations from the 8th layer of cased mBERT on a tSNE plot with highlighted language families.

et al., 2019) in addition to the test data.

For MT QE, we use English-German data provided for the WMT19 QE Shared Task (Fonseca et al., 2019) consisting training and test data with source sentences, their automatic translations, and manually corrections.

6 Results

Language Identification. Table 1 shows that centering the sentence representations considerably decreases the accuracy of language identification, especially in the case of mean-pooled embeddings. This indicates that the proposed centering procedure does indeed remove the language-specific information to a great extent.

Language Similarity. Figure 1 is a tSNE plot (Maaten and Hinton, 2008) of the language centroids, showing that the similarity of the centroids tends to correspond to the similarity of the lan-

| cased | uncased | UDify | lng-free | random |
|-------|---------|-------|----------|--------|
| 82.42 | 82.09 | 80.03 | 80.59 | 62.14 |

Table 2: V-Measure for hierarchical clustering of language centroids and grouping languages into genealogical families for families with at least three languages covered by mBERT.

| | mBERT | UDify | lng-free |
|------------------|-------|-------|----------|
| [cls] | .639 | .462 | .549 |
| [cls], cent. | .684 | .660 | .686 |
| [cls], proj. | .915 | .933 | .697 |
| mean-pool | .776 | .314 | .755 |
| mean-pool, cent. | .838 | .564 | .828 |
| mean-pool, proj. | .983 | .906 | .983 |

Table 3: Average accuracy for sentence retrieval over all 30 language pairs.

guages. Table 2 confirms that the hierarchical clustering of the language centroids mostly corresponds to the language families.

Parallel Sentence Retrieval. Results in Table 3 reveal that the representation centering dramatically improves the retrieval accuracy, showing that it makes the representations more language-neutral. However, an explicitly learned projection of the representations leads to a much greater improvement, reaching a close-to-perfect accuracy, even though the projection was fitted on relatively small parallel data. The accuracy is higher for mean-pooled states than for the [cls] embedding and varies according to the layer of mBERT used (see Figure 2).

Word Alignment. Table 4 shows that word-alignment based on mBERT representations surpasses the outputs of the standard FastAlign tool even if it was provided large parallel corpus. This suggests that word-level semantics are well captured by mBERT contextual embeddings. For this task, learning an explicit projection had a negligible effect on the performance.²

MT Quality Estimation. Qualitative results of MT QE are tabulated in Table 5. Unlike sentence retrieval, QE is more sensitive to subtle differences

²We used an expectation-maximization approach that alternately aligned the words and learned a linear projection between the representations. This algorithm only brings a negligible improvement of .005 F₁ points.

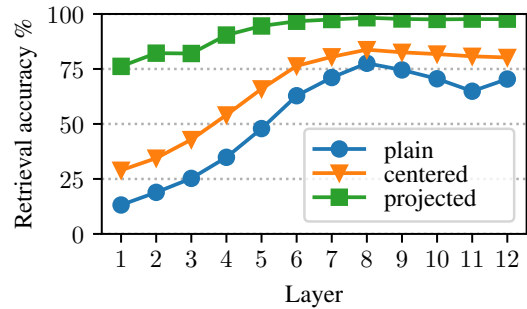


Figure 2: Accuracy of sentence retrieval for mean-pooled contextual embeddings from BERT layers.

| en- | FastAlign | mBERT | UDify | lng-free |
|-----|-----------|-------|-------|----------|
| cs | .692 | .738 | .708 | .744 |
| sv | .438 | .478 | .459 | .468 |
| de | .471 | .767 | .731 | .768 |
| fr | .583 | .612 | .581 | .607 |
| ro | .690 | .703 | .696 | .704 |

Table 4: Maximum F₁ score for word alignment across layers compared with FastAlign baseline.

between sentences. Measuring the distance of the non-centered sentence vectors does not correlate with translation quality at all. Centering or explicit projection only leads to a mild correlation, much lower than a supervisedly trained regression;³ and even better performance is possible (Fonseca et al., 2019). The results show that the linear projection between the representations only captures a rough semantic correspondence, which does not seem to be sufficient for QE, where the most indicative feature appears to be sentence complexity.

7 Fine-tuning mBERT

We also considered model fine-tuning towards stronger language neutrality. We evaluate two fine-tuned versions of mBERT: *UDify*, tuned for a multi-lingual dependency parser, and *lng-free*, tuned to jettison the language-specific information from the representations.

7.1 UDify

The UDify model (Kondratyuk and Straka, 2019) uses mBERT to train a single model for dependency parsing and morphological analysis of 75

³Supervised regression using either only the source or only MT output also shows a respectable correlation, which implies that structural features of the sentences are more useful than the comparison of the source sentence with MT output.

| BERT | cente- red | glob. proj. | supervised | | |
|----------|---------------|----------------|------------|------|------|
| | | | src | MT | both |
| cased | .005 | .163 | .362 | .352 | .419 |
| uncased | .027 | .204 | .367 | .390 | .425 |
| UDify | .039 | .167 | .368 | .375 | .413 |
| lng-free | .026 | .136 | .349 | .343 | .411 |

Table 5: Correlation of estimated MT quality with HTER for English-to-German translation on WMT19 data.

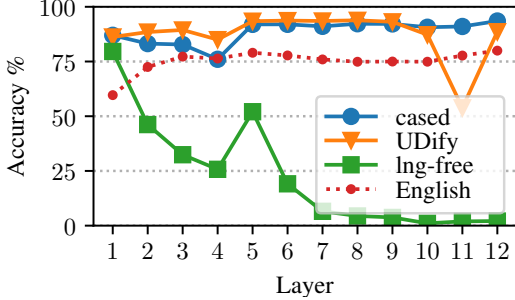


Figure 3: Language ID accuracy for different layers of mBERT.

languages. During the parser training, mBERT is fine-tuned, which improves the parser accuracy. Results on zero-shot parsing suggest that the fine-tuning leads to more cross-lingual representations with respect to morphology and syntax.

However, our analyses show that fine-tuning mBERT for multilingual dependency parsing does not remove the language identity information from the representations and actually makes the representations less semantically cross-lingual.

7.2 lng-free

In this experiment, we try to make the representations more language-neutral by removing the language identity from the model using an adversarial approach. We continue training mBERT in a multi-task learning setup with the masked LM objective with the same sampling procedure (Devlin et al., 2019) jointly with adversarial language ID classifiers (Elazar and Goldberg, 2018). For each layer, we train one classifier for the [cls] token and one for the mean-pooled hidden states with the gradient reversal layer (Ganin and Lempitsky, 2015) between mBERT and the classifier.

The results reveal that the adversarial removal of language information succeeds in dramatically decreasing the accuracy of the language identifica-

tion classifier; the effect is strongest in deeper layers for which the standard mBERT tend to perform better (see Figure 3). However, other tasks are not affected by the adversarial fine-tuning.

8 Conclusions

Using a set of semantically oriented tasks that require explicit semantic cross-lingual representations, we showed that mBERT contextual embeddings do not represent similar semantic phenomena similarly and therefore they are not directly usable for zero-shot cross-lingual tasks.

Contextual embeddings of mBERT capture similarities between languages and cluster the languages by their families. Neither cross-lingual fine-tuning nor adversarial language identity removal breaks this property. A part of language information is encoded by the position in the embedding space, thus a certain degree of cross-linguality can be achieved by centering the representations for each language. Exploiting this property allows a good cross-lingual sentence retrieval performance and bilingual word alignment (which is invariant to the shift). A good cross-lingual representation can be achieved by fitting a supervised projection on a small parallel corpus.

References

- Ondrej Bojar, Christian Buck, Christian Federmann, Barry Haddow, Philipp Koehn, Johannes Leveling, Christof Monz, Pavel Pecina, Matt Post, Herve Saint-Amand, Radu Soricut, Lucia Specia, and Aleš Tamchyna. 2014. [Findings of the 2014 workshop on statistical machine translation](#). In *Proceedings of the Ninth Workshop on Statistical Machine Translation*, pages 12–58, Baltimore, Maryland, USA. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Chris Dyer, Victor Chahuneau, and Noah A. Smith. 2013. A simple, fast, and effective reparameterization of IBM model 2. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 644–648, Atlanta, Georgia. Association for Computational Linguistics.

- Yanai Elazar and Yoav Goldberg. 2018. [Adversarial removal of demographic attributes from text data](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 11–21, Brussels, Belgium. Association for Computational Linguistics.
- Miquel Esplà, Mikel Forcada, Gema Ramírez-Sánchez, and Hieu Hoang. 2019. ParaCrawl: Web-scale parallel corpora for the languages of the EU. In *Proceedings of Machine Translation Summit XVII Volume 2: Translator, Project and User Tracks*, pages 118–119, Dublin, Ireland. European Association for Machine Translation.
- Erick Fonseca, Lisa Yankovskaya, André F. T. Martins, Mark Fishel, and Christian Federmann. 2019. [Findings of the WMT 2019 shared tasks on quality estimation](#). In *Proceedings of the Fourth Conference on Machine Translation (Volume 3: Shared Task Papers, Day 2)*, pages 1–10, Florence, Italy. Association for Computational Linguistics.
- Yaroslav Ganin and Victor Lempitsky. 2015. [Unsupervised domain adaptation by backpropagation](#). In *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 1180–1189, Lille, France. PMLR.
- Maria Holmqvist and Lars Ahrenberg. 2011. [A gold standard for English-Swedish word alignment](#). In *Proceedings of the 18th Nordic Conference of Computational Linguistics (NODALIDA 2011)*, pages 106–113, Riga, Latvia. Northern European Association for Language Technology (NEALT).
- Dan Kondratyuk and Milan Straka. 2019. [75 languages, 1 model: Parsing universal dependencies universally](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2779–2795, Hong Kong, China. Association for Computational Linguistics.
- Sneha Kudugunta, Ankur Bapna, Isaac Caswell, and Orhan Firat. 2019. [Investigating multilingual NMT representations at scale](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1565–1575, Hong Kong, China. Association for Computational Linguistics.
- Laurens van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-sne. *Journal of machine learning research*, 9(Nov):2579–2605.
- David Mareček. 2016. [Czech-english manual word alignment](#). LINDAT/CLARIN digital library at the Institute of Formal and Applied Linguistics (ÚFAL), Faculty of Mathematics and Physics, Charles University.
- Rada Mihalcea and Ted Pedersen. 2003. [An evaluation exercise for word alignment](#). In *Proceedings of the HLT-NAACL 2003 Workshop on Building and Using Parallel Texts: Data Driven Machine Translation and Beyond*, pages 1–10.
- Franz Josef Och and Hermann Ney. 2000. [Improved statistical alignment models](#). In *Proceedings of the 38th Annual Meeting of the Association for Computational Linguistics*, pages 440–447, Hong Kong. Association for Computational Linguistics.
- Telmo Pires, Eva Schlinger, and Dan Garrette. 2019. [How multilingual is multilingual BERT?](#) In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4996–5001, Florence, Italy. Association for Computational Linguistics.
- Samuel Rönqvist, Jenna Kanerva, Tapio Salakoski, and Filip Ginter. 2019. [Is multilingual BERT fluent in language generation?](#) In *Proceedings of the First NLP Workshop on Deep Learning for Natural Language Processing*, pages 29–36, Turku, Finland. Linköping University Electronic Press.
- Andrew Rosenberg and Julia Hirschberg. 2007. [V-measure: A conditional entropy-based external cluster evaluation measure](#). In *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*, pages 410–420, Prague, Czech Republic. Association for Computational Linguistics.
- Yuxuan Wang, Wanxiang Che, Jiang Guo, Yijia Liu, and Ting Liu. 2019. [Cross-lingual BERT transformation for zero-shot dependency parsing](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5725–5731, Hong Kong, China. Association for Computational Linguistics.
- Shijie Wu and Mark Dredze. 2019. [Beto, bentz, becas: The surprising cross-lingual effectiveness of BERT](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 833–844, Hong Kong, China. Association for Computational Linguistics.