

MULTILINGUAL ALIGNMENT OF CONTEXTUAL WORD REPRESENTATIONS

Anonymous authors

Paper under double-blind review

ABSTRACT

We propose procedures for evaluating and strengthening contextual embedding alignment and show that they are useful in understanding and improving multilingual BERT. In particular, after our proposed alignment procedure, BERT exhibits significantly improved zero-shot performance on XNLI compared to the base model, remarkably matching fully-supervised models for Bulgarian and Greek. Further, using non-contextual and contextual versions of word retrieval, we show that BERT outperforms fastText while being able to distinguish between multiple uses of a word, suggesting that pre-training subsumes word vectors for learning cross-lingual signals. Finally, we use the contextual word retrieval task to gain a better understanding of the strengths and weaknesses of multilingual pre-training.

1 INTRODUCTION

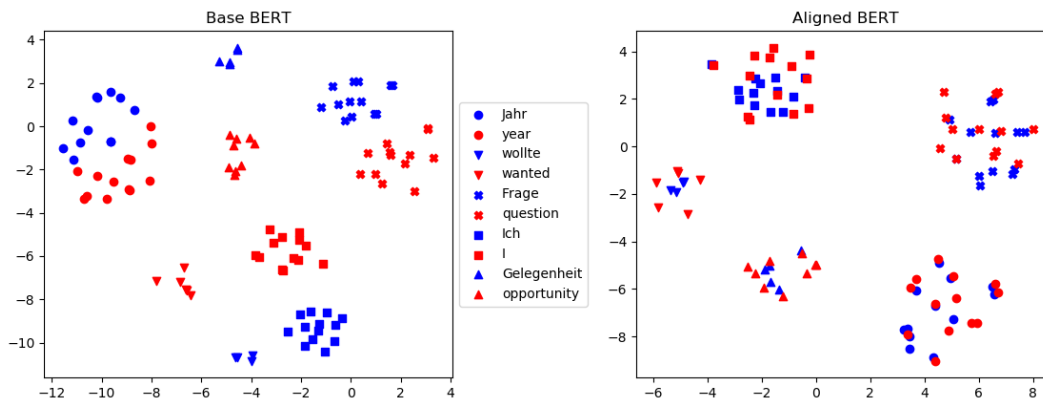


Figure 1: t-SNE (Maaten & Hinton, 2008) visualization of the embedding spaces of BERT pre- and post-alignment. Each point is a different instance of the word within a sentence in the Europarl corpus. This figure suggests that BERT begins already somewhat aligned out-of-the-box but becomes much more aligned after our proposed procedure.

Embedding alignment was originally studied for word vectors with the goal of enabling cross-lingual transfer, where the embeddings for two languages are in alignment if word translations, e.g. *cat* and *Katze*, have similar representations (Mikolov et al., 2013; Smith et al., 2017). Recently, large pre-trained models have largely subsumed word vectors based on their accuracy on downstream tasks, partly due to the fact that their word representations are context-dependent, allowing them to more richly capture the meaning of a word (Peters et al., 2018; Howard & Ruder, 2018; Radford et al., 2018; Devlin et al., 2018). Therefore, with the same goal of cross-lingual transfer but for these more complex models, we might consider contextual embedding alignment, where we observe whether word pairs within parallel sentences, e.g. *cat* in “*The cat sits*” and *Katze* in “*Die Katze sitzt*,” have similar representations.

One model relevant to these questions is multilingual BERT, a version of BERT pre-trained on 104 languages that achieves remarkable transfer on downstream tasks. For example, after the model is

fine-tuned on the English MultiNLI training set, it achieves 74.3% accuracy on the test set in Spanish, which is only 7.1% lower than the English accuracy (Devlin et al., 2018; Conneau et al., 2018b). Furthermore, while the model transfers better to languages similar to English, it still achieves reasonable accuracies even on languages with different scripts.

However, given the way that multilingual BERT was pre-trained, it is unclear why we should expect such high zero-shot performance. Compared to monolingual BERT which exhibits no zero-shot transfer, multilingual BERT differs only in that (1) during pre-training (i.e. masked word prediction), each batch contains sentences from all of the languages, and (2) it uses a single shared vocabulary, formed by WordPiece on the concatenated monolingual corpora (Devlin et al., 2019). Therefore, we might wonder: (1) How can we better understand BERT’s multilingualism? (2) Can we further improve BERT’s cross-lingual transfer?

In this paper, we show that contextual embedding alignment is a useful concept for addressing these questions. First, we propose a contextual version of word retrieval to evaluate the degree of alignment, where a model is presented with two parallel corpora, and given a word within a sentence in one corpus, it must find the correct word and sentence in the other. Using this metric of alignment, we show that multilingual BERT achieves zero-shot transfer because its embeddings are partially aligned, as depicted in Figure 1, with the degree of alignment predicting the degree of downstream transfer.

Next, using the Europarl corpus as parallel data (Koehn, 2005), we propose an alignment procedure and show that it significantly improves BERT as a multilingual model. Specifically, on zero-shot XNLI, where the model is trained on English MultiNLI and tested on other languages (Conneau et al., 2018b), the aligned model improves accuracies by 2.78% on average over the base model, and it remarkably matches fully-supervised models for Bulgarian and Greek.

To compare contextual alignment to its non-contextual counterpart, we also compare multilingual BERT to fastText augmented with sentence vectors (Devlin et al., 2018; Bojanowski et al., 2017; Rücklé et al., 2018). Interestingly, we find that aligned BERT matches fastText in non-contextual word retrieval. Furthermore, BERT is able to distinguish between multiple uses of the same word, a task that is impossible for non-contextual methods. These experiments provide evidence that pre-trained language models largely subsume word vectors for learning cross-lingual signals.

Finally, we present a finer-grained analysis of multilingual BERT using the contextual word retrieval task, with the goal of better understanding its strengths and shortcomings. Specifically, we find that base BERT has trouble matching open-class compared to closed-class parts-of-speech, suggesting insight into the pre-training procedure that we explore in Section 5. Together, these experiments support contextual alignment as an important task that provides useful insight into large multilingual pre-trained models.

2 RELATED WORK

Word vector alignment. A series of works has shown that after word vectors are trained for each language separately, to achieve alignment, a rotation matrix can be applied post-hoc and can be learned with or without supervision (Mikolov et al., 2013; Smith et al., 2017; Artetxe et al., 2017; 2018; Conneau et al., 2018a; Hoshen & Wolf, 2018; Xu et al., 2018; Chen & Cardie, 2018). Because alignment is evaluated using bilingual dictionary word retrieval, these papers also propose ways to mitigate the hubness problem in nearest neighbors, e.g. by using alternate similarity functions.

Incorporating context into alignment. One key challenge in making alignment context aware is that the embeddings are now different across multiple occurrences of the same word. Past papers have handled this issue by removing context and aligning the “average sense” of a word. In one such study, Schuster et al. (2019) learn a rotation to align contextual ELMo embeddings (Peters et al., 2018) with the goal of improving zero-shot multilingual dependency parsing, and they handle context by taking the average embedding for a word in all of its contexts. In another paper, Aldarmaki & Diab (2019) learn a rotation on sentence vectors, produced by taking the average word vector over the sentence, and they show that the resulting alignment also works well for word-level tasks. In this paper, we depart from these past works by (1) aligning not only the word but also the context, and (2) using more expressive alignment methods than rotation.

Incorporating parallel texts into pre-training. Instead of performing alignment post-hoc, another line of works proposes contextual pre-training procedures that are more cross-lingually-aware. Wieting et al. (2019) pre-train sentence embeddings using parallel texts by maximizing similarity between sentence pairs while minimizing similarity with negative examples. Lample & Conneau (2019) propose a cross-lingual pre-training objective that incorporates large amounts of parallel data in addition to monolingual corpora, leading to improved downstream cross-lingual transfer.

Analyzing multilingual BERT. Pires et al. (2019) present a series of probing experiments to better understand multilingual BERT. They find that transfer is possible even between dissimilar languages, but that it works better on languages that are typologically similar. They conclude that BERT is remarkably multilingual but falls short for certain language pairs.

3 METHODS

3.1 MULTILINGUAL PRE-TRAINING

We first briefly describe multilingual BERT (Devlin et al., 2018). Like monolingual BERT, multilingual BERT is pre-trained on sentences from Wikipedia to perform two tasks: masked word prediction, where it must predict words that are masked within a sentence, and next sentence prediction, where it must predict whether the second sentence follows the first one. The model is trained on 104 languages, with each batch containing training sentences from each language, and it uses a shared vocabulary formed by WordPiece on the 104 Wikipedias concatenated (Wu et al., 2016).

3.2 DEFINING AND EVALUATING CONTEXTUAL ALIGNMENT

In the following sections, we describe how to define, evaluate, and improve contextual alignment. Given two languages, a model is in contextual alignment if it has similar representations for word pairs within parallel sentences. More precisely, suppose we have N parallel sentences $C = \{(\mathbf{s}^1, \mathbf{t}^1), \dots, (\mathbf{s}^N, \mathbf{t}^N)\}$, where (\mathbf{s}, \mathbf{t}) is a source-target sentence pair. Also, let each sentence pair (\mathbf{s}, \mathbf{t}) have word pairs, denoted $a(\mathbf{s}, \mathbf{t}) = \{(i_1, j_1), \dots, (i_m, j_m)\}$, containing position tuples (i, j) such that the words s_i and t_j are translations of each other.¹ We will use f to represent a pre-trained model such that $f(i, \mathbf{s})$ is the contextual embedding for the i th word in \mathbf{s} .

As an example, we might have the following sentence pair:

$$\begin{aligned} \mathbf{s} &= \{I \overset{0}{a} \overset{1}{t} \overset{2}{e} \overset{3}{t} \overset{4}{h} \overset{5}{e} \overset{6}{.}\} & \mathbf{t} &= \{I \overset{0}{c} \overset{1}{h} \overset{2}{a} \overset{3}{b} \overset{4}{e} \overset{5}{d} \overset{6}{e} \overset{7}{n} \overset{8}{A} \overset{9}{p} \overset{10}{f} \overset{11}{e} \overset{12}{l} \overset{13}{g} \overset{14}{e} \overset{15}{s} \overset{16}{s} \overset{17}{e} \overset{18}{n} \overset{19}{.}\} \\ a(\mathbf{s}, \mathbf{t}) &= \{(0, 0), (1, 4), (2, 2), (3, 3), (4, 5)\} \end{aligned}$$

Then, we define the *contextual alignment* of the model f with respect to the parallel corpus C as its *accuracy in contextual word retrieval*. In this task, the model is presented with two parallel corpora, and given a word within a sentence in one corpus, it must find the correct word and sentence in the other. Specifically, we can define a nearest neighbor retrieval function

$$\text{neighbor}(i, \mathbf{s}; f, C) = \underset{\mathbf{t} \in C, 0 \leq j \leq \text{len}(\mathbf{t})}{\text{argmax}} \text{sim}(f(i, \mathbf{s}), f(j, \mathbf{t})),$$

where i and j denote the position within a sentence and sim is a similarity function. The accuracy is then given by the percentage of exact matches over the entire corpus, or

$$A(f; C) = \frac{1}{N} \sum_{(\mathbf{s}, \mathbf{t}) \in C} \sum_{(i, j) \in a(\mathbf{s}, \mathbf{t})} \mathbb{I}(\text{neighbor}(i, \mathbf{s}; f, C) = (j, \mathbf{t})),$$

where \mathbb{I} represents the indicator function. We can perform the same procedure in the other direction, where we pull target words given source words, so we report the average between the two directions. As our similarity function, we use CSLS, a modified version of cosine similarity that mitigates the hubness problem, with neighborhood size 10 (Conneau et al., 2018a).

¹These pairs are called word alignments in the machine translation community, but we use the term “word pairs” or “parallel words” to avoid confusion with embedding alignment.

Given parallel data, these word pairs can be procured in an unsupervised fashion using standard techniques developed by the machine translation community (Brown et al., 1993). While these methods can be noisy, by running the algorithm in both the source-target and target-source directions and only keeping word pairs in their intersection, we can trade-off coverage for accuracy, producing a reasonably high-precision dataset (Och & Ney, 2003).

One additional point is that this procedure can be made more or less contextual based on the corpus. In particular, a corpus with more occurrences for each word type requires better representations of context. Therefore, we also test non-contextual word retrieval by removing all but the first occurrence of each word type.

3.3 ALIGNING PRE-TRAINED CONTEXTUAL EMBEDDINGS

To improve the alignment of the model f with respect to the corpus C , we can encapsulate alignment in the loss function

$$L(f; C) = - \sum_{(\mathbf{s}, \mathbf{t}) \in C} \sum_{(i, j) \in a(\mathbf{s}, \mathbf{t})} \text{sim}(f(i, \mathbf{s}), f(j, \mathbf{t})),$$

where we sum the similarities between parallel words. Because the CSLS metric is not easily optimized, we instead use the squared error loss, or $\text{sim}(f(i, \mathbf{s}), f(j, \mathbf{t})) = -\|f(i, \mathbf{s}) - f(j, \mathbf{t})\|_2^2$.

However, note that this loss function does not account for the informativity of f ; for example, it is zero if f is constant. Therefore, at a high level, we would like to minimize $L(f; C)$ while maintaining some aspect of f that makes it useful, e.g. its high accuracy when fine-tuned on downstream tasks. Letting f_0 denote the initial pre-trained model before alignment, we achieve this goal by defining a regularization term

$$R(f; C) = \sum_{\mathbf{t} \in C} \sum_{i=1}^{\text{len}(\mathbf{t})} \|f(j, \mathbf{t}) - f_0(j, \mathbf{t})\|_2,$$

which imposes a penalty if the target language embeddings stray from their initialization. Then, we sample minibatches $B \subset C$ and take gradient steps of the function $L(f; B) + \lambda R(f; B)$ directly on the weights of f , which moves the source embeddings toward the target embeddings while preventing the latter from drifting too far. In our experiments, we set $\lambda = 1$.

In the multilingual case, suppose we have k parallel corpora C^1, \dots, C^k , where each corpus has a different source language with the target language as English. Then, we sample equal-sized batches $B^i \subset C^i$ from each corpus and take gradient steps on $\sum_{i=1}^k L(f; B^i) + \lambda R(f; B^i)$, which moves all of the non-English embeddings toward English.

Note that this alignment method departs from prior work, in which each non-English language is rotated to match the English embedding space through individual learned matrices. Rotation requires the strong assumption that the embedding spaces are roughly isometric, an assumption that may not hold for contextual pre-trained models given their increased complexity. Another advantage of our method is that the alignment for all languages is done simultaneously.

As our dataset, we use the Europarl corpora for English paired with Bulgarian, German, Greek, Spanish, and French, the languages represented in both Europarl and XNLI (Koehn, 2005). After tokenization (Koehn et al., 2007), we produce word pairs using fastAlign and keep the one-to-one pairs in the intersection (Dyer et al., 2013). We use the most recent 1024 sentences as the test set, the previous 1024 sentences as the development set, and the following 250K sentences as the training set. Furthermore, we modify the test set accuracy calculation to only include word pairs not seen in the training set. We also remove any exact matches, e.g. punctuation and numbers, because BERT is already aligned for these pairs due to its shared vocabulary.

3.4 SENTENCE-AUGMENTED NON-CONTEXTUAL BASELINE

Given that there has been a long line of work on word vector alignment (Artetxe et al., 2018; Conneau et al., 2018a; Smith et al., 2017, *inter alia*), we also compare BERT to a sentence-augmented fastText baseline (Bojanowski et al., 2017). Following Artetxe et al. (2018), we first normalize, then

	bg-en	de-en	el-en	es-en	fr-en	Average
Aligned fastText + sentence	44.0	46.4	42.0	48.6	44.5	45.1
Base BERT	19.5	26.1	13.9	32.5	28.3	24.1
Aligned BERT	50.7	51.3	49.8	51.0	48.6	50.3

Table 1: Contextual word retrieval accuracy for the aligned sentence-augmented fastText baseline and BERT pre- and post-alignment. Across languages, base BERT has variable accuracy while aligned BERT is consistently effective.

mean-center, then normalize the word vectors. Next, we learn a rotation W applied to the source vectors to minimize the distance between parallel word pairs, or

$$\min_W \sum_{(s,t) \in C} \sum_{(i,j) \in a(s,t)} \|\vec{\sigma}_{s_i} - W\vec{\tau}_{t_j}\|_2 \quad s.t. \quad W^\top W = I,$$

where $\vec{\sigma}_{s_i}$ and $\vec{\tau}_{t_j}$ are the source and target fastText vectors for the words s_i and t_j . This problem is known as the Procrustes problem and can be solved in closed form (Schonemann, 1966).

We also strengthen this baseline by including sentence information. Specifically, during word retrieval, we concatenate each word vector with a vector representing its sentence. Following Rücklé et al. (2018), we compute the sentence vector by concatenating the average, maximum, and minimum vector over all of the words in the sentence, a method that was shown to be state-of-the-art for a suite of cross-lingual tasks. We also experimented with other methods, such as first retrieving the sentence and then the word, but found this method resulted in the highest accuracy. Note that as a result, the fastText vectors are 1200-dimensional, while the BERT vectors are 768-dimensional.

3.5 TESTING ZERO-SHOT TRANSFER

The next step is to determine whether better alignment improves cross-lingual transfer. As our downstream task, we use the XNLI dataset, where the English MultiNLI development and test sets are human-translated into multiple languages (Conneau et al., 2018b; Williams et al., 2018). Given a pair of sentences, the task is to predict whether the first sentence implies the second, where there are three labels: entailment, neutral, or contradiction. Starting from either the base or aligned multilingual BERT, we train on English and evaluate on Bulgarian, German, Greek, Spanish, and French, the XNLI languages represented in Europarl.

As our architecture, following Devlin et al. (2018), we apply a linear layer followed by softmax on the [CLS] embedding of the sentence pair, producing scores for each of the three labels. The model is trained using cross-entropy loss and selected based on its development set accuracy averaged across all of the languages. As a fully-supervised ceiling, we also compare to models trained and tested on the same language. For the non-English training data, we use the machine translations of the English MultiNLI training data provided by Conneau et al. (2018b). While the quality of the training data is affected by the quality of the MT system, this comparison nevertheless serves as a good approximation for the fully-supervised setting.

4 RESULTS

4.1 CONTEXTUAL WORD RETRIEVAL

Table 1 displays the contextual word retrieval accuracies for the aligned sentence-augmented fastText baseline and BERT pre- and post-alignment. Unsurprisingly, aligned BERT outperforms the other two approaches. Furthermore, before alignment, BERT’s performance varies greatly between languages, while after alignment it is consistently effective. In particular, Bulgarian and Greek initially have very low accuracies. This phenomenon is also reflected in the XNLI numbers (Table 3), where Bulgarian and Greek receive the largest boosts from alignment.

These results are consistent with a hypothesis by Pires et al. (2019) to explain BERT’s multilingualism. They posit that due to the shared vocabulary, shared words between languages, e.g. numbers

	bg-en	de-en	el-en	es-en	fr-en	Average
Aligned fastText + sentence	61.3	65.4	61.6	71.1	64.8	64.8
Base BERT	29.1	37.0	22.3	46.5	41.8	35.3
Aligned BERT	62.8	64.3	67.5	68.4	66.3	65.9

Table 2: Non-contextual word retrieval accuracy, where there is only one occurrence per word type. Interestingly, even though the task is non-contextual, aligned BERT outperforms the fastText baseline by a small margin.

	English	Bulgarian	German	Greek	Spanish	French	Average
Fully-Supervised (Translate-Train)							
Base BERT	81.9	73.6	75.9	71.6	77.8	76.8	76.3
Zero-Shot							
Base BERT ²	80.4	68.7	70.4	67.0	74.5	73.4	72.4
Aligned BERT	80.1	73.4	73.1	71.4	75.5	74.5	74.7

Table 3: Accuracy on the XNLI test set. After alignment, Bulgarian and Greek match the fully-supervised ceiling, while German, Spanish, and French close between one-third and one-half of the gap.

and names, are forced to have the same representation. Then, due to the masked word prediction task, other words that co-occur with these shared words also receive similar representations. If this hypothesis is true, then languages with higher lexical overlap with English are likely to experience higher transfer. As an extreme form of this phenomenon, Bulgarian and Greek have completely different scripts and should experience worse transfer than the common-script languages, an intuition that is confirmed by the word retrieval and XNLI accuracies. The fact that all languages are equally aligned with English post-alignment suggests that the pre-training procedure is needlessly suboptimal for languages that share little lexical overlap but are otherwise similar to English. In particular, the model benefits greatly from additional alignment methods that do not rely on lexical overlap.

4.2 NON-CONTEXTUAL WORD RETRIEVAL

Table 2 shows the non-contextual word retrieval accuracies, where there is only one occurrence per word type. Even though there is no longer a need to disambiguate between contexts, aligned BERT still outperforms fastText on average by a small margin. Given that BERT is significantly better in the contextual case, this result suggests that incorporating context and larger scale pre-training improve a model’s ability to capture cross-lingual signals. In particular, cross-lingual alignment is still possible with these more complex models, making them unequivocally better than non-contextual methods for multilingual tasks.

4.3 ZERO-SHOT XNLI TRANSFER

To test whether alignment improves downstream transfer, we also apply the models to the XNLI task, displayed in Table 3. Looking at the unaligned zero-shot numbers, the word retrieval accuracies are highly correlated with downstream zero-shot performance (Figure 2). Furthermore, alignment greatly improves accuracies, with all languages seeing a gain of at least 1%. These results give credence to alignment as a predictor of cross-lingual transfer. In particular, the Bulgarian and Greek zero-shot numbers are boosted by almost 5% each and end up matching the fully-supervised numbers, suggesting that the alignment procedure is especially effective for languages that are ini-

²Note that the zero-shot Base BERT numbers are slightly different from those reported in Devlin et al. (2019) because we select a single model using the average accuracy across the six languages.

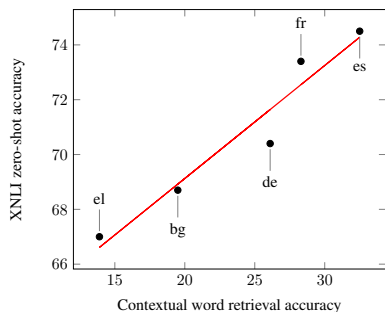


Figure 2: XNLI zero-shot versus contextual word retrieval for base BERT, where each point is a different language paired with English. This plot suggests that the degree of alignment correlates well with downstream transfer.

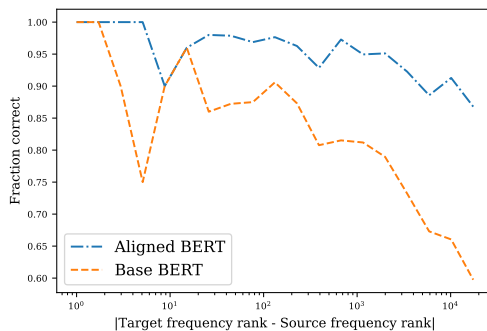


Figure 3: BERT word retrieval accuracy plotted against difference in frequency rank between the source and target. Base BERT accuracy decreases for larger differences, suggesting that its alignment depends on the word pairs having similar usage statistics.

tially difficult for BERT. The accuracies matching is especially notable given that there is one fully-supervised model for each language, while there is only a single model in the zero-shot setting.

4.4 CONTEXT-AWARE RETRIEVAL

In this section, we qualitatively show that aligned BERT is able to disambiguate between different occurrences of a word. The example shows two meanings of the word “like” occurring in the English-German Europarl test set, where BERT is able to find the correct match. Note also that in the second and third example, the two senses of “like” occur in the same sentence. We provide more examples in the appendix.

- This empire did not look for colonies far from home or overseas, **like** most Western European States, but close by.
Dieses Reich suchte seine Kolonien nicht weit von zu Hause und in bersee **wie** die meisten westeuropäischen Staaten, sondern in der unmittelbaren Umgebung.
- **Like** other speakers, I would like to support the call for the arms embargo to remain.
Wie andere Sprecher, so möchte auch ich den Aufruf zur Aufrechterhaltung des Waffenembargos untersttzen.
- Like other speakers, I would **like** to support the call for the arms embargo to remain.
Wie andere Sprecher, so **möchte** auch ich den Aufruf zur Aufrechterhaltung des Waffenembargos untersttzen.
- I would also **like**, although they are absent, to mention the Commission and the Council.
Ich **möchte** mir sogar erlauben, die Kommission und den Rat zu nennen, auch wenn sie nicht anwesend sind.

5 ANALYSIS

In the following sections, we use the word retrieval task to gain a better understanding multilingual BERT. Specifically, we examine which groups of words are aligned out-of-the-box to reveal shortcomings and gain insight into the pre-training procedure.

5.1 WORD RETRIEVAL PART-OF-SPEECH ANALYSIS

We first analyze the accuracy broken down by part-of-speech using the Universal Part-of-Speech Tagset (Petrov et al., 2012), annotated using polyglot (Al-Rfou et al., 2013) for Bulgarian and

Lexical Overlap	Numeral	Punctuation	Proper Noun			Average	
Base BERT	0.90	0.88	0.80			0.86	
Aligned BERT	0.97	0.96	0.95			0.96	
Closed-Class	Determiner	Preposition	Conjunction	Pronoun	Auxiliary	Average	
Base BERT	0.76	0.72	0.71	0.70	0.61	0.70	
Aligned BERT	0.91	0.86	0.89	0.89	0.84	0.88	
Open-Class	Noun	Adverb	Adjective	Verb			Average
Base BERT	0.61	0.57	0.50	0.49			0.54
Aligned BERT	0.90	0.88	0.90	0.89			0.89

Table 4: Accuracy by part-of-speech tag for non-contextual word retrieval. To achieve better word type coverage, we do not remove word pairs seen in the training set. The tags are grouped into lexically overlapping, closed-class, and open-class groups. The “Particle,” “Symbol,” “Interjection,” and “Other” tags are omitted.

spaCy (Honnibal & Montani, 2017) for the other languages, as displayed in Table 4. Unsurprisingly, multilingual BERT has high alignment out-of-the-box for groups with high lexical overlap, e.g. numerals, punctuation, and proper nouns, due to its shared vocabulary.

We further divide the remaining tags into closed-class and open-class, where closed-class parts-of-speech correspond to fixed sets of words serving grammatical functions (e.g. determiner, preposition, conjunction, pronoun, and auxiliary), while open-class parts-of-speech correspond to lexical words and are open in the sense that new words are constantly added (e.g. noun, adverb, adjective, verb). Interestingly, we see that base BERT has consistently lower accuracy for closed-class versus open-class categories (0.70 vs 0.54), but that this discrepancy disappears after alignment (0.88 vs 0.89).

5.2 USAGE HYPOTHESIS FOR ALIGNMENT

From this closed-class vs open-class difference, we hypothesize that BERT’s alignment of a particular word pair is influenced by the similarity of their usage statistics. Specifically, given that BERT is trained through masked word prediction, its embeddings are in large part determined by the co-occurrences between words. Therefore, two words that are more consistently used in similar contexts should be better aligned. This hypothesis provides an explanation of the closed-class vs open-class difference: closed-class words are typically grammatical, so they are used in similar ways across typologically similar languages. Furthermore, these words are not substitutable for one another due to their grammatical function. Therefore, their usage statistics are a strong signature that can be used for alignment. On the other hand, open-class words can be substituted for one another: for example, in most sentences, the noun tokens could be replaced by a wide range of semantically dissimilar nouns with the sentence remaining syntactically well-formed. By this effect, many nouns have similar co-occurrences, making them difficult to align through masked word prediction alone.

To further test this hypothesis, we plot the word retrieval accuracy as a function of the difference between the frequency rank of the target and source word, where this difference measures discrepancies in usage, as depicted in Figure 2. We see that base BERT drops off significantly as this difference increases. Furthermore, this shortcoming is remedied by alignment, revealing another systematic deficiency of multilingual pre-training.

6 CONCLUSION

Given that the degree of alignment is causally predictive of downstream cross-lingual transfer, contextual alignment proves to be a useful concept for understanding and improving multilingual pre-trained models. Contextual word retrieval also provides useful new insights into the pre-training procedure, opening up new avenues for analysis.

REFERENCES

- Rami Al-Rfou, Bryan Perozzi, and Steven Skiena. Polyglot: Distributed word representations for multilingual nlp. In *Proceedings of the Seventeenth Conference on Computational Natural Language Learning*, pp. 183–192, Sofia, Bulgaria, August 2013. Association for Computational Linguistics. URL <http://www.aclweb.org/anthology/W13-3520>.
- Hanan Aldarmaki and Mona Diab. Context-aware cross-lingual mapping. 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)*, pp. 3906–3911, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1391. URL <https://www.aclweb.org/anthology/N19-1391>.
- Mikel Artetxe, Gorka Labaka, and Eneko Agirre. Learning bilingual word embeddings with (almost) no bilingual data. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 451–462, Vancouver, Canada, July 2017. Association for Computational Linguistics. doi: 10.18653/v1/P17-1042. URL <https://www.aclweb.org/anthology/P17-1042>.
- Mikel Artetxe, Gorka Labaka, and Eneko Agirre. A robust self-learning method for fully unsupervised cross-lingual mappings of word embeddings. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 789–798, Melbourne, Australia, July 2018. Association for Computational Linguistics. URL <https://www.aclweb.org/anthology/P18-1073>.
- Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics*, 5:135–146, 2017. doi: 10.1162/tacl.a.00051. URL <https://www.aclweb.org/anthology/Q17-1010>.
- Peter F. Brown, Vincent J. Della Pietra, Stephen A. Della Pietra, and Robert L. Mercer. The mathematics of statistical machine translation: Parameter estimation. *Comput. Linguist.*, 19(2): 263–311, June 1993. ISSN 0891-2017. URL <http://dl.acm.org/citation.cfm?id=972470.972474>.
- Xilun Chen and Claire Cardie. Unsupervised multilingual word embeddings. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 261–270, Brussels, Belgium, October–November 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1024. URL <https://www.aclweb.org/anthology/D18-1024>.
- Alexis Conneau, Guillaume Lample, Marc’Aurelio Ranzato, Ludovic Denoyer, and Herve J’egou. Word translation without parallel data. In *Proceedings of the 6th International Conference on Learning Representations (ICLR 2018)*, 2018a. URL <https://arxiv.org/pdf/1710.04087.pdf>.
- Alexis Conneau, Ruty Rinott, Guillaume Lample, Adina Williams, Samuel Bowman, Holger Schwenk, and Veselin Stoyanov. XNLI: Evaluating cross-lingual sentence representations. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 2475–2485, Brussels, Belgium, October–November 2018b. Association for Computational Linguistics. doi: 10.18653/v1/D18-1269. URL <https://www.aclweb.org/anthology/D18-1269>.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv:1810.04805 [cs.CL]*, October 2018. URL <http://arxiv.org/abs/1810.04805>.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. <https://github.com/google-research/bert/blob/master/multilingual.md>, 2019.
- Chris Dyer, Victor Chahuneau, and Noah A. Smith. 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*, pp. 644–648, Atlanta, Georgia, June 2013. Association for Computational Linguistics. URL <https://www.aclweb.org/anthology/N13-1073>.
- Matthew Honnibal and Ines Montani. spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing. To appear, 2017.
- Yedid Hoshen and Lior Wolf. Non-adversarial unsupervised word translation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 469–478, Brussels, Belgium, October–November 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1043. URL <https://www.aclweb.org/anthology/D18-1043>.
- Jeremy Howard and Sebastian Ruder. Universal language model fine-tuning for text classification. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 328–339. Association for Computational Linguistics, 2018. URL <http://aclweb.org/anthology/P18-1031>.
- Philipp Koehn. Europarl: A parallel corpus for statistical machine translation. In *Conference Proceedings: The Tenth Machine Translation Summit*, pp. 79–86, Phuket, Thailand, 2005. AAMT.
- Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondřej Bojar, Alexandra Constantin, and Evan Herbst. Moses: Open source toolkit for statistical machine translation. In *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics Companion Volume Proceedings of the Demo and Poster Sessions*, pp. 177–180, Prague, Czech Republic, June 2007. Association for Computational Linguistics. URL <https://www.aclweb.org/anthology/P07-2045>.
- Guillaume Lample and Alexis Conneau. Cross-lingual language model pretraining. 2019. URL <https://arxiv.org/pdf/1901.07291.pdf>.
- Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of Machine Learning Research*, 9:2579–2605, 2008. URL <http://www.jmlr.org/papers/v9/vandermaaten08a.html>.
- Tomas Mikolov, Quoc V Le, and Ilya Sutskever. Exploiting similarities among languages for machine translation. 2013. URL <https://arxiv.org/pdf/1309.4168.pdf>.
- Franz Josef Och and Hermann Ney. A systematic comparison of various statistical alignment models. *Comput. Linguist.*, 29(1):19–51, March 2003. ISSN 0891-2017. doi: 10.1162/089120103321337421. URL <http://dx.doi.org/10.1162/089120103321337421>.
- Matthew Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. Deep contextualized word representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pp. 2227–2237, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-1202. URL <https://www.aclweb.org/anthology/N18-1202>.
- Slav Petrov, Dipanjan Das, and Ryan McDonald. A universal part-of-speech tagset. In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC-2012)*, pp. 2089–2096, Istanbul, Turkey, May 2012. European Languages Resources Association (ELRA). URL http://www.lrec-conf.org/proceedings/lrec2012/pdf/274_Paper.pdf.
- Telmo Pires, Eva Schlinger, and Dan Garrette. How multilingual is multilingual BERT? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 4996–5001, Florence, Italy, July 2019. Association for Computational Linguistics. URL <https://www.aclweb.org/anthology/P19-1493>.
- Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training. 2018. URL https://s3-us-west-2.amazonaws.com/openai-assets/research-covers/language-unsupervised/language_understanding_paper.pdf.

- Andreas Rücklé, Steffen Eger, Maxime Peyrard, and Iryna Gurevych. Concatenated p-mean word embeddings as universal cross-lingual sentence representations. *arXiv:1803.01400 [cs.CL]*, 2018. URL <http://arxiv.org/abs/1803.01400>.
- Peter H. Schonemann. A generalized solution of the orthogonal procrustes problem. *Psychometrika*, 31(1):1–10, 1966.
- Tal Schuster, Ori Ram, Regina Barzilay, and Amir Globerson. Cross-lingual alignment of contextual word embeddings, with applications to zero-shot dependency parsing. 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)*, pp. 1599–1613, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1162. URL <https://www.aclweb.org/anthology/N19-1162>.
- Samuel L. Smith, David H. P. Turban, Steven Hamblin, and Nils Y. Hammerla. Offline bilingual word vectors, orthogonal transformations and the inverted softmax. In *Proceedings of the 5th International Conference on Learning Representations (ICLR 2017)*, 2017. URL <https://openreview.net/pdf?id=r1Aab85gg>.
- John Wieting, Kevin Gimpel, Graham Neubig, and Taylor Berg-Kirkpatrick. Simple and effective paraphrastic similarity from parallel translations. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 4602–4608, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1453. URL <https://www.aclweb.org/anthology/P19-1453>.
- Adina Williams, Nikita Nangia, and Samuel Bowman. A broad-coverage challenge corpus for sentence understanding through inference. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pp. 1112–1122, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-1101. URL <https://www.aclweb.org/anthology/N18-1101>.
- Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, Jeff Klingner, Apurva Shah, Melvin Johnson, Xiaobing Liu, Lukasz Kaiser, Stephan Gouws, Yoshikiyo Kato, Taku Kudo, Hideto Kazawa, Keith Stevens, George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals, Greg Corrado, Macduff Hughes, and Jeffrey Dean. Google’s neural machine translation system: Bridging the gap between human and machine translation. *arXiv:1609.08144 [cs.CL]*, 2016.
- Ruochen Xu, Yiming Yang, Naoki Otani, and Yuexin Wu. Unsupervised cross-lingual transfer of word embedding spaces. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 2465–2474, Brussels, Belgium, October–November 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1268. URL <https://www.aclweb.org/anthology/D18-1268>.

A APPENDIX

A.1 OPTIMIZATION HYPERPARAMETERS

For both alignment and XNLI optimization, we use a learning rate of 5×10^{-5} with Adam hyperparameters $\beta = (0.9, 0.98)$, $\epsilon = 10^{-9}$ and linear learning rate warmup for the first 10% of the training data. For alignment, the model is trained for one epoch, with each batch containing 2 sentence pairs per language. For XNLI, each model is trained for 3 epochs with 32 examples per batch, and 10% dropout is applied to the BERT embeddings.

A.2 MORE EXAMPLES OF CONTEXT-AWARE RETRIEVAL

In this section, we provide a couple hand-picked examples where aligned BERT successfully performs context-aware retrieval.

Multiple meanings of “order”:

- Moreover, the national political elite had to make a detour in Ambon in **order** to reach the civil governor’s residence by warship.
In Ambon mußte die politische Spitze des Landes auch noch einen Umweg machen, **um** mit einem Kriegsschiff die Residenz des Provinzgouverneurs zu erreichen.
- Although the European Union has an interest in being surrounded by large, stable regions, the tools it has available in **order** to achieve this are still very limited.
Der Europäischen Union ist zwar an großen stabilen Regionen in ihrer Umgebung gelegen, aber sie verfügt nach wie vor nur über recht begrenzte Instrumente, **um** das zu erreichen.
- We could reasonably expect the new Indonesian government to take action in three fundamental areas: restoring public **order**, prosecuting and punishing those who have blood on their hands and entering into a political dialogue with the opposition.
Von der neuen indonesischen Regierung darf man mit Fug und Recht drei elementare Maßnahmen erwarten: die Wiederherstellung der öffentlichen **Ordnung**, die Verfolgung und Bestrafung derjenigen, an deren Händen Blut klebt, und die Aufnahme des politischen Dialogs mit den Gegnern.
- Firstly, I might mention the fact that the army needs to be reformed, secondly that a stable system of law and **order** needs to be introduced.
Ich nenne hier an erster Stelle die notwendige Reform der Armee, ferner die Einführung eines stabilen Systems rechtsstaatlicher **Ordnung**.

Multiple meanings of “support”:

- Financial **support** is needed to enable poor countries to take part in these court activities.
Arme Länder müssen finanziell **unterstützt** werden, damit auch sie sich an der Arbeit des Gerichtshofs beteiligen können.
- We must help them and ensure that a proper action plan is implemented to **support** their work.
Es gilt einen wirklichen Aktionsplan auf den Weg zu bringen, um die Arbeit dieser Organisationen zu **unterstützen**.
- So I hope that you will all **support** this resolution condemning the abominable conditions of prisoners and civilians in Djibouti.
Ich hoffe daher, daß Sie alle diese Entschließung **befürworten**, die die entsetzlichen Bedingungen von Inhaftierten und Zivilpersonen in Dschibuti verurteilt.
- It would be difficult to **support** a subsidy scheme that channelled most of the aid to the large farms in the best agricultural regions.
Es wäre auch problematisch, ein Beihilfesystem zu **befürworten**, das die meisten Beihilfen in die großen Betriebe in den besten landwirtschaftlichen Gebieten lenkt.

Multiple meanings of “close”:

- This empire did not look for colonies far from home or overseas, like most Western European States, but **close** by.
Dieses Reich suchte seine Kolonien nicht weit von zu Hause und in bersee wie die meisten westeuropäischen Staaten, sondern in der unmittelbaren **Umgebung**.
- In addition, if we are to shut down or refuse investment from every company which may have an association with the arms industry, then we would have to **close** virtually every American and Japanese software company on the island of Ireland with catastrophic consequences.
Wenn wir zudem jedes Unternehmen, das auf irgendeine Weise mit der Rüstungsindustrie verbunden ist, schließen oder Investitionen dieser Unternehmen unterbinden, dann müßten wir so ziemlich alle amerikanischen und japanischen Softwareunternehmen auf der irischen Insel **schließen**, was katastrophale Auswirkungen hätte.

- On the other hand, the deployment of resources left over in the Structural Funds from the programme planning period 1994 to 1999 is hardly worth considering as the available funds have already been allocated to specific measures, in this case in **close** collaboration with the relevant French authorities.

Die Verwendung verbliebener Mittel der Strukturfonds aus dem Programmplanungszeitraum 1994 bis 1999 ist dagegen kaum in Erwägung zu ziehen, da die verfügbaren Mittel bereits bestimmten Maßnahmen zugewiesen sind, und zwar im konkreten Fall im **engen** Zusammenwirken mit den zuständigen französischen Behörden.

- This is particularly justified given that, as already stated, many Member States have very **close** relations with Djibouti.

Zumal, wie erwähnt, viele Mitgliedstaaten sehr **enge** Beziehungen zu Dschibuti unterhalten.

- Mr President, it is regrettable that, at the **close** of the 20th century, a century symbolised so positively by the peaceful women's revolution, there are still countries, such as Kuwait and Afghanistan, where half the population, women that is, is still denied fundamental human rights.

Herr Präsident! Es ist wirklich bedauerlich, daß es am **Ende** des 20. Jahrhunderts, eines so positiv von der friedlichen Revolution der Frauen geprägten Jahrhunderts, noch immer Länder wie Kuwait und Afghanistan gibt, in denen der Hälfte der Bevölkerung, den Frauen, die elementaren Menschenrechte verweigert werden.