

# When is BERT Multilingual? Isolating Crucial Ingredients for Cross-lingual Transfer

**Ameet Deshpande**  
Department of Computer Science  
Princeton University, USA  
asd@cs.princeton.edu

**Partha Talukdar**  
Google Research  
India  
partha@google.com

**Karthik Narasimhan**  
Department of Computer Science  
Princeton University, USA  
karthikn@cs.princeton.edu

## Abstract

While recent work on multilingual language models has demonstrated their capacity for cross-lingual zero-shot transfer, there is a lack of consensus in the community as to what shared properties between languages enable transfer on downstream tasks. Analyses involving pairs of natural languages are often inconclusive and contradictory since languages simultaneously differ in many linguistic aspects. In this paper, we perform a large-scale empirical study to isolate the effects of various linguistic properties by measuring zero-shot transfer between four diverse natural languages and their counterparts constructed by modifying aspects such as the script, word order, and syntax. Among other things, our experiments show that the absence of sub-word overlap significantly affects zero-shot transfer when languages differ in their word order, and there is a strong correlation between transfer performance and word embedding alignment between languages (e.g.,  $\rho_s = 0.94$  on the task of NLI). Our results call for focus in multilingual models on explicitly improving word embedding alignment between languages rather than relying on its implicit emergence.<sup>1</sup>

## 1 Introduction

Multilingual language models like XLM (Conneau et al., 2020a) and Multilingual-BERT (Devlin, 2019) are trained with masked-language modeling (MLM) objective on a combination of raw text from multiple languages. Surprisingly, these models exhibit decent cross-lingual zero-shot transfer, where fine-tuning on a task in a source language translates to good performance for a different language (target).

**Requirements for zero-shot transfer** Recent studies have provided inconsistent explanations for properties required for zero-shot transfer (hereon,

<sup>1</sup>Code is available at <https://github.com/princeton-nlp/MultilingualAnalysis>

transfer). For example, while Wu and Dredze (2019) conclude that sub-word overlap is vital for transfer, K et al. (2020) demonstrate that it is not crucial, although they consider only English as the source language. While Pires et al. (2019) suggest that typological similarity (e.g., similar SVO order) is essential for transfer, other works (Kakwani et al., 2020; Conneau et al., 2020a) successfully build multilingual models for dissimilar languages.

**Need for systematic analysis** A major cause of these discrepancies is a large number of varying properties (e.g., syntax, script, and vocabulary size) between languages, which make isolating crucial ingredients for transfer difficult. Some studies alleviate this issue by creating synthetic languages which differ from natural ones only in specific linguistic properties like script (K et al., 2020; Dufter and Schütze, 2020). However, their focus is only on English as a source language, and the scale of their experiments is small (in number of tasks or pre-training corpora size), thus limiting the scope of their findings to their settings alone.

**Our approach** We perform a systematic study of cross-lingual transfer on bilingual language models trained on a natural language and a systematically *derived* counterpart. We choose four diverse natural languages (English, French, Arabic, and Hindi) and create *derived* variants using four different transformations on structural properties such as inverting or permuting word order, altering scripts, or varying syntax (Section 3.2). We train models on each of the resulting sixteen language pairs, and evaluate zero-shot transfer on four downstream tasks – natural language inference (NLI), named-entity recognition (NER), part-of-speech tagging (POS), and question-answering (QA).

**Our experiments reveal the following:**

1. Contrary to previous belief, the absence of sub-word overlap degrades transfer when languages

differ in their word order (e.g., by more than 40 F1 points on POS tagging, (§ 4.1)).

2. There is a strong correlation between token embedding alignment and zero-shot transfer across different tasks (e.g.,  $\rho_s = 0.94, p < .005$  for XNLI, Fig 4).
3. Using pre-training corpora from similar sources for different languages (e.g., Wikipedia) boosts transfer when compared to corpora from different sources (e.g., 17 F1 points on NER, Fig 3).

To our knowledge, we are the first study to quantitatively show that zero-shot transfer between languages is strongly correlated with token embedding alignment ( $\rho_s = 0.94$  for NLI). We also show that the current multilingual pre-training methods (Dodapaneni et al., 2021) fall short of aligning embeddings even between simple natural and derived language pairs, leading to failure in zero-shot transfer. Our results call for training objectives that explicitly improve alignment using either supervised (e.g., parallel corpora or bilingual dictionaries) or unsupervised data.

## 2 Related work

### Multilingual pre-training for Transformers

The success of monolingual Transformer language models (Devlin et al., 2019; Radford et al., 2018) has driven studies that learn a multilingual language-model (LM) on several languages. Multilingual-BERT (M-BERT) (Devlin, 2019) is a single neural network pre-trained using the masked language-modeling (MLM) objective on a corpus of text from 104 languages. XLM (Conneau and Lample, 2019) introduced translation language-modeling, which performs MLM on pairs of parallel sentences, thus encouraging alignment between their representations. These models exhibit surprising zero-shot cross-lingual transfer performance (Conneau and Lample, 2019; K et al., 2020), a setup where the model is fine-tuned on a source language and evaluated on a different target language.

**Analysis of cross-lingual transfer** While Pires et al. (2019), Conneau et al. (2020b), and K et al. (2020) showed that transfer works even without a shared vocabulary between languages, Wu and Dredze (2019) discovered a correlation between sub-word overlap and zero-shot performance. Conneau et al. (2020b) and Artetxe et al. (2020a) showed that shared parameters for languages with

different scripts were crucial for transfer. Pires et al. (2019) and (Wu and Dredze, 2019) observed that transfer for NER and POS tagging works better between typologically similar languages. However, a study conducted by Lin et al. (2019) showed that there is no simple rule of thumb to gauge when transfer works between languages. Hsu et al. (2019) observed that changing the syntax (SOV) order of the source to match that of the target does not improve performance.

### Transfer between real and synthetic Languages

K et al. (2020) create a synthetic language by changing English’s script and find that transfer between it and Spanish works even without common sub-words. However, they use only English as their source language, test only on two tasks, and use a single natural-synthetic language pair. Dufter and Schütze (2020) study transfer between English and *synthetic* English obtained by changing the script, word order, or model delimiters. However, they use a small corpus (228K words) compared to current standards (we use 3 orders more) and measure only embedding similarity and not zero-shot transfer. A contemporary work (Wu et al., 2022) uses synthetic transformations to modify the GLUE dataset (Wang et al., 2018) and analyze properties required for good zero-shot transfer, but they perform their experiments only on English and do not perform token embedding alignment analysis. We show that the latter is crucial for good transfer.

## 3 Approach

We first provide some background on bilingual language models (Section 3.1), followed by descriptions of our transformations (Section 3.2), and our training and evaluation setup (Section 3.3).

### 3.1 Background

**Bilingual pre-training** The standard setup (Conneau and Lample, 2019) trains a bilingual language model (*Bi-LM*) on raw text corpora from two languages simultaneously. *Bi-LM* uses the masked language-modeling loss ( $\mathcal{L}_{\text{MLM}}$ ) on the corpora from the two languages ( $\mathcal{C}_1, \mathcal{C}_2$ ) separately with no explicit cross-lingual signal:

$$\mathcal{L}_{\text{Bi-LM}}^\theta(\mathcal{C}_1 + \mathcal{C}_2) = \mathcal{L}_{\text{MLM}}^\theta(\mathcal{C}_1) + \mathcal{L}_{\text{MLM}}^\theta(\mathcal{C}_2)$$

A shared byte pair encoding tokenizer (Sennrich et al., 2015) is trained on  $\mathcal{C}_1 + \mathcal{C}_2$ . A single batch contains instances from both languages, but each instance belongs to a single language.

Transformation	Instance ( $s$ )	Transformed instance ( $\mathcal{T}(s)$ )
<b>Inversion</b> ( $\mathcal{T}_{\text{inv}}$ )	Welcome to NAACL at Seattle	Seattle at NAACL to Welcome
<b>Permutation</b> ( $\mathcal{T}_{\text{perm}}$ )	This is a conference	a This conference is
<b>Transliteration</b> ( $\mathcal{T}_{\text{trans}}$ )	I am Sam . I am	♣ <sub>(I)</sub> ♥ <sub>(am)</sub> ◇ <sub>(Sam)</sub> ♠ <sub>(.)</sub> ♣ <sub>(I)</sub> ♥ <sub>(am)</sub>
<b>Syntax</b> ( $\mathcal{T}_{\text{syn}}$ )	Sara (S) ate (V) apples (O)	Sara (S) apples (O) ate (V)
	Une table (N) ronde (A)	Une ronde (A) table (N)

Table 1: Examples of our transformations applied to different sentences (without sub-word tokenization). *Inversion* inverts the tokens, *Permutation* samples a random reordering, and *Transliteration* changes the script. We use symbols (♣) to denote words in the new script and mention the corresponding original word in brackets. *Syntax* stochastically modifies the syntactic structure. In the first example for *Syntax*, the sentence in Subject-Verb-Object (SVO) order gets transformed to SOV order, and in the second, the sentence in Noun-Adjective (NA) order gets transformed to the AN order. The examples are high probability re-orderings and other ones might be sampled too.

**Zero-shot transfer evaluation** Consider a bilingual model (*Bi-LM*) pre-trained on two languages, *source* and *target*. Zero-shot transfer involves fine-tuning *Bi-LM* on downstream task data from *source* and evaluating on test data from *target*. This is considered zero-shot because *Bi-LM* is not fine-tuned on any data belonging to *target*.

### 3.2 Generating language variants with systematic transformations

Natural languages typically differ in several ways, like the script, word order, and syntax. To isolate the affect of these properties on zero-shot transfer, we obtain *derived* language corpora (hereon, *derived* corpora) from *original* (natural) language corpora by performing sentence level transformations ( $\mathcal{T}$ ) which change particular properties. For example, an “*inversion*” transformation could be used to invert each sentence in the corpus ( $Welcome_1$  to<sub>2</sub> NAACL<sub>3</sub>  $\Rightarrow$  NAACL<sub>3</sub> to<sub>2</sub> Welcome<sub>1</sub>). Since the transformation ( $\mathcal{T}$ ) is applied on each sentence of the *original* corpus, the size of the *original* and the *derived* corpus is the same. In the following sections, we will use the following notation:

$$\begin{aligned} \mathcal{C}_{\text{orig}} &\equiv \text{Original corpus} \\ &= \{s_i \mid i = 1 : N, s_i = \text{sentence}\} \\ \mathcal{T} &\equiv \text{Sentence-level transformation} \\ \mathcal{C}_{\text{deriv}} &\equiv \text{Derived corpus} \\ &= \{\mathcal{T}(\text{sent}) \mid \forall \text{sent} \in \mathcal{C}_{\text{orig}}\} \end{aligned}$$

**Types of transformations** We consider four transformations which modify different aspects of sentences (examples in Table 1):

1. **Inversion** ( $\mathcal{T}_{\text{inv}}$ ): Invert the order of *tokens* in the sentence, like in [Dufter and Schütze](#)

(2020). The first token becomes the last, and vice versa.

2. **Permutation** ( $\mathcal{T}_{\text{perm}}$ ): Permute the order of tokens in a sentence uniformly at random. For a sentence of  $n$  tokens, we sample a random ordering with probability  $\frac{1}{n!}$ .
3. **Transliteration** ( $\mathcal{T}_{\text{trans}}$ ): Change the script of all tokens other than the special tokens (like [CLS]). This creates a *derived* vocabulary ( $\mathcal{V}_{\text{deriv}}$ ) with a one-to-one correspondence with the original vocabulary ( $\mathcal{V}_{\text{orig}}$ ).
4. **Syntax** ( $\mathcal{T}_{\text{syn}}$ ): Modify a sentence to match the syntactic properties of a different natural language by re-ordering the dependents of nouns and verbs in the dependency parse. These transformations are stochastic because of the errors in parsing and sampling over possible re-orderings ([Wang and Eisner, 2016](#)).

These transformations allow us to systematically evaluate the effect of corresponding properties on zero-shot transfer. We also consider composed transformations (§4.2) which consecutively apply two transformations. We note that while real languages typically differ in more than one or two properties considered in our transformations, our methodology remains useful in isolating crucial properties that enable good transfer and can be extended to more transformations.

**Transformations for downstream tasks** We obtain the downstream corpus in the *derived* language ( $\mathcal{D}_{\text{deriv}}$ ) by applying the same transformation ( $\mathcal{T}$ ) used during pre-training on the *original* downstream corpus ( $\mathcal{D}_{\text{orig}}$ ). Unlike pre-training corpora which contain raw sentences, instances in downstream tasks contain one or more sentences with annotated labels. For text classification tasks like

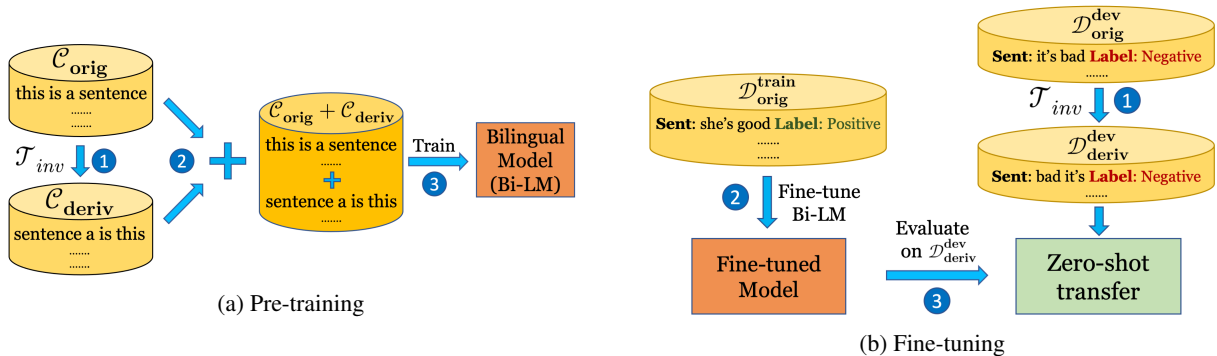


Figure 1: (a) During pre-training, we ① obtain the *derived* language corpus ( $\mathcal{C}_{\text{deriv}}$ ) by *transforming* the *original* language corpus ( $\mathcal{C}_{\text{orig}}$ ). ② The two corpora are combined and, ③ a bilingual model (*Bi-LM*) is learned using the MLM objective. (b) During fine-tuning, we ① obtain the *derived dev* dataset ( $\mathcal{D}_{\text{deriv}}^{\text{dev}}$ ) by *transforming* the original dev dataset ( $\mathcal{D}_{\text{orig}}^{\text{dev}}$ ). ② *Bi-LM* is fine-tuned on the original train dataset ( $\mathcal{D}_{\text{orig}}^{\text{train}}$ ), and ③ evaluated on  $\mathcal{D}_{\text{deriv}}^{\text{dev}}$ , which is the standard zero-shot cross lingual setup.

Evaluation	Corpus source		
	Pre-train	Fine-tune (train)	Fine-tune (dev)
<b>BZ</b>	$\mathcal{C}_{\text{orig}} + \mathcal{C}_{\text{deriv}}$	$\mathcal{D}_{\text{orig}}$	$\mathcal{D}_{\text{deriv}}$
<b>BS</b>	$\mathcal{C}_{\text{orig}} + \mathcal{C}_{\text{deriv}}$	$\mathcal{D}_{\text{deriv}}$	$\mathcal{D}_{\text{deriv}}$
<b>MZ</b>	$\mathcal{C}_{\text{orig}}$	$\mathcal{D}_{\text{orig}}$	$\mathcal{D}_{\text{deriv}}$
	$\Delta_{(\text{BZ}-\text{BS})} = (\text{BZ} - \text{BS})$		
	$\Delta_{(\text{MZ}-\text{BS})} = (\text{MZ} - \text{BS})$		

Table 2: Summary of evaluation metrics defined in § 3.3.  $\mathcal{C}$  and  $\mathcal{D}$  denote the pre-training and downstream corpus respectively, and their subscript indicates their source (*original* or *derived*). **BZ** and **MZ** represent bilingual and monolingual zero-shot transfer scores, and **BS** is the supervised learning baseline on *derived*. The differences in the setting of **BZ** and other scores are typeset in blue. We use  $\Delta_{(\text{BZ}-\text{BS})}$  and  $\Delta_{(\text{MZ}-\text{BS})}$  (defined in the last two rows) throughout our paper.

NLI, we apply the transformation on each sentence in every dataset instance. For token classification tasks (e.g., NER, POS), any transformation which changes the order of the tokens also changes the order of the labels. We present the mathematical specification in Appendix A.

### 3.3 Model Training and Evaluation

We now describe our pre-training and zero-shot transfer evaluation setup. Figure 1 provides an overview of pre-training and fine-tuning, and Table 2 summarizes the evaluation metrics we use.

**Pre-training** Let  $\mathcal{C}_{\text{orig}}$  and  $\mathcal{C}_{\text{deriv}}$  be the *original* and *derived* language pre-training corpora. We train two models for each *original-derived* pair:

1. **Bilingual Model (Bi-LM)**: A bilingual model pre-trained on the combined corpus ( $\mathcal{C}_{\text{orig}} +$

$\mathcal{C}_{\text{deriv}}$ ) (Figure 1a).

2. **Monolingual Model (Mono-LM)**: A monolingual model trained only on  $\mathcal{C}_{\text{orig}}$  for the same number of steps as *Bi-LM*'s. *Mono-LM* is used as a baseline to measure zero-shot transfer of a model not pre-trained on *derived*.

**Evaluation** Let  $\mathcal{D}_{\text{orig}}^{\text{train}}$  and  $\mathcal{D}_{\text{orig}}^{\text{dev}}$  be the *original* language training and development sets for a downstream task, and  $\mathcal{D}_{\text{deriv}}^{\text{train}}$  and  $\mathcal{D}_{\text{deriv}}^{\text{dev}}$  be the corresponding *derived* language datasets. For evaluation, we first fine-tune the pre-trained models on a downstream training set and evaluate the resulting model on a development set (Figure 1b). Since our goal is to investigate the extent of zero-shot transfer, we require appropriate lower and upper bounds to make informed conclusions. To this end, we compute three metrics, all on the same development set (summarized in Table 2):

- **Bilingual zero-shot transfer (BZ)**: This is the standard zero-shot transfer score (Conneau and Lample, 2019) which measures how well a bilingual model fine-tuned on  $\mathcal{D}_{\text{orig}}^{\text{train}}$  zero-shot transfers to the other language ( $\mathcal{D}_{\text{deriv}}^{\text{dev}}$ ).
- **Bilingual supervised synthetic (BS)**: This is the supervised learning performance on the *derived* language obtained by fine-tuning *Bi-LM* on  $\mathcal{D}_{\text{deriv}}^{\text{train}}$  and evaluating it on  $\mathcal{D}_{\text{deriv}}^{\text{dev}}$ .
- **Monolingual zero-shot transfer (MZ)**: This measures the zero-shot performance of the baseline *Mono-LM*, which is not pre-trained on the *derived* language, by fine-tuning *Mono-LM* on  $\mathcal{D}_{\text{orig}}^{\text{train}}$  and evaluating it on  $\mathcal{D}_{\text{deriv}}^{\text{dev}}$ .

*BS* uses fine-tuning train data from the *derived* language and serves as an upper-bound on *BZ* and *MZ*

which don't use it. *MZ* doesn't pre-train on the *derived* language and serves as a lower-bound on *BZ* which does pre-train on it. For easier comparison of *BZ* and *MZ* with *BS* (upper-bound), we report the following score differences (Table 2), which are both negative in our experiments.

$$\Delta_{(BZ-BS)} = (BZ - BS) \quad (1)$$

$$\Delta_{(MZ-BS)} = (MZ - BS) \quad (2)$$

*BZ* alone cannot capture the quality of the zero-shot transfer. A large and negative  $\Delta_{(BZ-BS)}$  implies that bilingual zero-shot transfer is much worse than supervised fine-tuning on *derived*. Concurrently,  $\Delta_{(BZ-BS)} \approx \Delta_{(MZ-BS)}$  implies that *Bi-LM* transfers as poorly as *Mono-LM*. **Thus, good zero-shot transfer is characterized by  $\Delta_{(BZ-BS)} \approx 0$  and  $\Delta_{(BZ-BS)} \gg \Delta_{(MZ-BS)}$ .**

### 3.4 Experimental Setup

**Languages** We choose four diverse natural languages: English (Indo-European, Germanic), French (Indo-European, Romance), Hindi (Indo-European, Indo-Iranian), and Arabic (Afro-Asiatic, Semitic), which are represented in the multilingual XTREME benchmark (Hu et al., 2020). For each language, we consider four transformations (Section 3.2) to create *derived* counterparts, giving us 16 different original-derived pairs in total. For the *Syntax* transformation, we use Qi et al. (2020) for parsing. We modify the syntax of FR, HI, and AR to that of EN, and the syntax of EN to that of FR.

**Datasets** For the pre-training corpus ( $\mathcal{C}_{\text{orig}}$ ), we use a 500MB (uncompressed) subset of Wikipedia ( $\approx 100\text{M}$  tokens) for each language. This matches the size of WikiText-103 (Merity et al., 2016), a standard language-modeling dataset. For downstream evaluation, we choose four tasks from the XTREME benchmark (Hu et al., 2020). Table 4 lists all the datasets and their evaluation metrics.

**Implementation Details** We use a variant of RoBERTa (Liu et al., 2019) which has 8 layers, 8 heads, and a hidden dimensionality of 512. We train each model on 500K steps, a batch size of 128, and a learning rate of  $1e-4$  with a linear warmup of 10K steps. We use an *original* language vocabulary size of 40000 for all the models and train on 8 Cloud TPU v3 cores for 32-48 hours. For fine-tuning, we use standard hyperparameters (Appendix F) from the XTREME benchmark and report our scores on the development sets.

## 4 Results

Our experiments reveal several interesting findings for bilingual models including the situational importance of sub-word overlap for zero-shot transfer (§ 4.1, 4.2), the effect of domain mismatch between languages (§ 4.3), and correlation of zero-shot performance with embedding alignment (§ 4.4). We connect our findings to zero-shot transfer results between natural languages in Section 4.5.

### 4.1 Sub-word overlap is not strictly necessary for strong zero-shot transfer

Sub-word overlap is the number of common tokens between two different language corpora. If  $\mathcal{E}_1$  and  $\mathcal{E}_2$  are sets of tokens which appear in the two corpora, then: Sub-word overlap =  $|\mathcal{E}_1 \cap \mathcal{E}_2| / |\mathcal{E}_1 \cup \mathcal{E}_2|$  (Pires et al., 2019). The *Transliteration* transformation ( $\mathcal{T}_{\text{trans}}$ ) creates *original-derived* language pairs that have 0% sub-word overlap (equivalently, different scripts), but follow the same word order.

Table 3 displays  $\Delta_{(BZ-BS)}$  scores for  $\mathcal{T}_{\text{trans}}$ , averaged over four languages (Appendix B contains a breakdown). We observe that  $\Delta_{(BZ-BS)} \approx 0$  for all tasks while  $\Delta_{(MZ-BS)}$  is highly negative, implying that zero-shot transfer is strong and on par with supervised learning. This result indicates that zero-shot transfer is possible even when languages with different scripts have similar word orders (in line with K et al. (2020)). However, it is unrealistic for natural languages to differ only in their script and not other properties (e.g., word order).

### 4.2 Absence of sub-word overlap significantly hurts zero-shot performance when languages differ in their word-orders

To simulate a more realistic scenario, we create *original* and *derived* language pairs which differ both in their scripts (0% sub-word overlap) and in word order. We achieve this by composing two transformations on the *original* language corpus, one of which is *Transliteration* ( $\mathcal{T}_{\text{trans}}$ ). We experiment with three different compositions, (a)  $\mathcal{T}_{\text{trans}} \circ \mathcal{T}_{\text{inv}}$ , (b)  $\mathcal{T}_{\text{trans}} \circ \mathcal{T}_{\text{perm}}$ , and (c)  $\mathcal{T}_{\text{trans}} \circ \mathcal{T}_{\text{syn}}$ . Here,  $\alpha \circ \beta$  means that transformation  $\beta$  is applied before  $\alpha$ . A composed transformation ( $\mathcal{T}_{\text{trans}} \circ$

<sup>4</sup>XQuAD is a question-answering task where the correct answer is a *contiguous* span. We do not report scores on XQuAD for  $\mathcal{T}_{\text{perm}}$  and  $\mathcal{T}_{\text{syn}}$  because they can potentially reorder individual words in the contiguous answer, thus distributing them throughout the transformed sentence and making the question unanswerable. On the other hand,  $\mathcal{T}_{\text{inv}}$  and  $\mathcal{T}_{\text{trans}}$  do not have this issue because they maintain the spans.

Task	Inversion ( $\mathcal{T}_{\text{inv}}$ )			Permutation ( $\mathcal{T}_{\text{perm}}$ )			Syntax ( $\mathcal{T}_{\text{syn}}$ )			Transliteration ( $\mathcal{T}_{\text{trans}}$ )		
	$\Delta_{(\text{BZ}-\text{BS})}$	$\Delta_{(\text{MZ}-\text{BS})}$	BZ	$\Delta_{(\text{BZ}-\text{BS})}$	$\Delta_{(\text{MZ}-\text{BS})}$	BZ	$\Delta_{(\text{BZ}-\text{BS})}$	$\Delta_{(\text{MZ}-\text{BS})}$	BZ	$\Delta_{(\text{BZ}-\text{BS})}$	$\Delta_{(\text{MZ}-\text{BS})}$	BZ
XNLI	-10.2	-13.0	58.4	-3.6	-8.6	62.6	-0.9 *	-1.1	67.8	-1.0 *	-36.7	69.3
NER	-49.1	-46.7	37.9	-26.3	-35.4	47.3	-14.6	-16.6	62.9	-1.9 *	-82.6	83.7
POS	-30.2	-36.2	64.2	-11.2	-25.2	73.6	-4.4	-7.6	89.4	-0.4 *	-95.0	95.4
XQuAD <sup>4</sup>	-32.8	-31.0	22.8	— <sup>4</sup>	—	—	— <sup>4</sup>	—	—	0.0 *	-55.9	61.2

Table 3: **(1) Evaluation:** We report  $\Delta_{(\text{BZ}-\text{BS})}$  and  $\Delta_{(\text{MZ}-\text{BS})}$  (§ 3.3 and Table 2) for transformations on different tasks, averaged over four languages (EN, FR, HI, AR). We report the breakdown for different languages in Appendix B. *BZ*, the bilingual zero-shot performance, is reported for reference. **(2) Interpreting scores:** Smaller (more negative)  $\Delta_{(\text{BZ}-\text{BS})}$  implies worse bilingual zero-shot transfer, whereas  $\Delta_{(\text{BZ}-\text{BS})} \approx 0$  implies strong transfer.  $\Delta_{(\text{BZ}-\text{BS})} \gg \Delta_{(\text{MZ}-\text{BS})}$  implies that bilingual pre-training is extremely useful. Scores are highlighted based on their value (lower scores have a higher intensity of red). Cases with strong zero-shot transfer ( $\Delta_{(\text{BZ}-\text{BS})} \approx 0$ ) are marked with an asterisk. **(3) Trends:**  $\mathcal{T}_{\text{trans}}$  exhibits strong transfer on all tasks and languages (high  $\Delta_{(\text{BZ}-\text{BS})}$  scores), and bilingual pre-training is extremely useful ( $\Delta_{(\text{BZ}-\text{BS})} \gg \Delta_{(\text{MZ}-\text{BS})}$ ), implying that zero-shot transfer is possible between languages with different scripts but the same word order.  $\mathcal{T}_{\text{inv}}$  and  $\mathcal{T}_{\text{perm}}$  suffer on all tasks (small  $\Delta_{(\text{BZ}-\text{BS})}$  scores) whereas  $\mathcal{T}_{\text{syn}}$  suffers significantly lesser, which provides evidence that local changes to the word order made by *Syntax* ( $\mathcal{T}_{\text{syn}}$ ) hurts zero-shot transfer significantly lesser than global changes made by *Inversion* ( $\mathcal{T}_{\text{inv}}$ ) and *Permutation* ( $\mathcal{T}_{\text{perm}}$ ).

Dataset	Task	Metric
XNLI (Conneau et al., 2018)	NLI	Accuracy
Wikiann (Pan et al., 2017)	NER	F1
UD v2.5 (Nivre et al., 2018)	POS	F1
XQuAD (Artetxe et al., 2020b)	QA	F1

Table 4: XTREME benchmark datasets used for zero-shot transfer evaluation. NLI=Natural Language Inference, NER=Named-entity recognition, POS=Part-of-speech tagging, QA=Question-Answering.

$\beta$ ) differs from its second constituent ( $\beta$ ) in that the former produces a *derived* language which has 0% sub-word overlap with the *original* language whereas the latter has a 100% sub-word overlap.

**Results** Our results (Figure 2, breakdown in Appendix C) show that zero-shot performance is significantly hurt for composed transformations when compared to its constituents.  $|\Delta_{(\text{BZ}-\text{BS})}|$  is much larger for  $\mathcal{T}_{\text{trans}} \circ \mathcal{T}_{\text{inv}}$  when compared to  $\mathcal{T}_{\text{trans}}$  or  $\mathcal{T}_{\text{inv}}$  individually. For example, for XNLI,  $|\Delta_{(\text{BZ}-\text{BS})}| = 19$  for the composed transformation and just 2 and 3 for  $\mathcal{T}_{\text{trans}}$  and  $\mathcal{T}_{\text{inv}}$  individually.  $\mathcal{T}_{\text{trans}} \circ \mathcal{T}_{\text{perm}}$  is worse by  $\approx 20$  points on XNLI and NER, and over 40 points on POS when compared to  $\mathcal{T}_{\text{perm}}$ .  $\mathcal{T}_{\text{trans}} \circ \mathcal{T}_{\text{syn}}$  suffers lesser than the other two composed transformations, but it is still worse than  $\mathcal{T}_{\text{syn}}$  by 3, 6, and 1 point on XNLI, NER, and POS. In conclusion, the absence of sub-word overlap significantly degrades zero-shot per-

formance in the realistic case of languages with different word orders.

### 4.3 Data from the same domain boosts bilingual performance

Previously, we considered transformations ( $\mathcal{T}$ ) that modified the *original* pre-training corpus to get a parallel corpus,  $\mathcal{C}_{\text{deriv}} = \mathcal{T}(\mathcal{C}_{\text{orig}})$ , such that there is a one-to-one correspondence between sentences in  $\mathcal{C}_{\text{orig}}$  and  $\mathcal{C}_{\text{deriv}}$  (we call this setting *parallel*). Since procuring large parallel corpora is expensive in practice, we consider two other settings which use different corpora for *original* and *derived*.

**Setup** Consider two text corpora of the same size,  $\mathcal{C}_{\text{orig}}^1$  and  $\mathcal{C}_{\text{orig}}^2$ . We compare two settings: (1) The *parallel* setting pre-trains a bilingual model on  $\mathcal{C}_{\text{orig}}^1 + \mathcal{T}(\mathcal{C}_{\text{orig}}^1)$ , whereas the (2) *non-parallel* corpus setting uses  $\mathcal{C}_{\text{orig}}^1 + \mathcal{T}(\mathcal{C}_{\text{orig}}^2)$ . We consider two variants of *non-parallel*, (1) *non-parallel (same)* which uses different splits of Wikipedia data (hence, *same* domain), and (2) *non-parallel (diff)* which uses Wikipedia data for the *original* and common crawl data (web text) for the *derived* language (hence, *different* domain). We use the *Transliteration* transformation ( $\mathcal{T}_{\text{trans}}$ ) to generate the *derived* language corpus and report  $|\Delta_{(\text{BZ}-\text{BS})}|$  averaged over all languages in Figure 3.

**Results** We observe consistently on all tasks that the *parallel* setting (blue bar) performs better than both the non-parallel settings. *Non-parallel (same)*

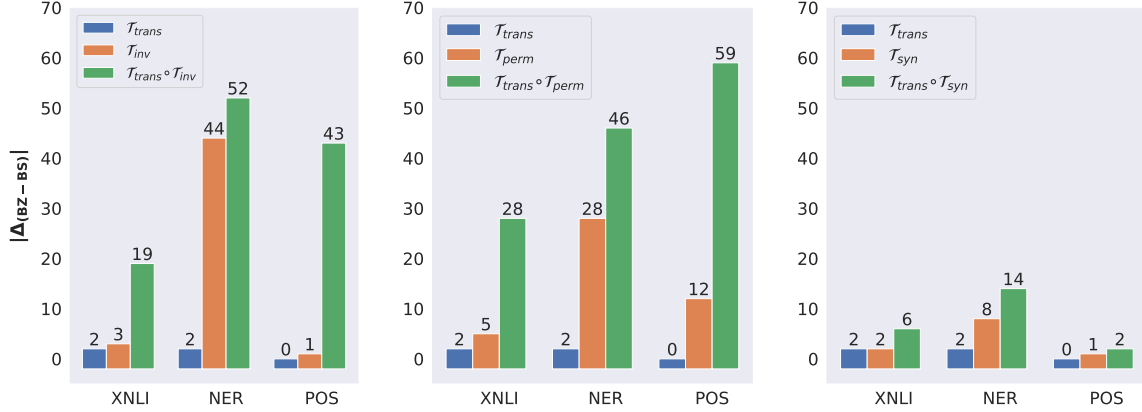


Figure 2:  $|\Delta_{(\text{BZ-BS})}|$  for composed transformations (§ 4.2) applied on EN as the *original* language. Larger scores imply worse zero-shot transfer.  $\mathcal{T}_{\text{trans}}$  = *Transliteration*,  $\mathcal{T}_{\text{inv}}$  = *Inversion*,  $\mathcal{T}_{\text{perm}}$  = *Permutation*, and  $\mathcal{T}_{\text{syn}}$  = *Syntax*. Sub-word overlap between the *original* and *derived* language is 0% when composed transformations are used (e.g.  $\mathcal{T}_{\text{trans}} \circ \mathcal{T}_{\text{inv}}$ ) and 100% when the second constituent is used (here,  $\mathcal{T}_{\text{inv}}$ ). We observe that the composed transformations (green bars) do significantly worse than their constituents (blue and orange), showing that  $\mathcal{T}_{\text{trans}} \circ \mathcal{T}_{\text{inv}}$  is worse than  $\mathcal{T}_{\text{inv}}$  by over 16 points on XNLI and 42 points on POS, with similar trends for  $\mathcal{T}_{\text{trans}} \circ \mathcal{T}_{\text{perm}}$ .  $\mathcal{T}_{\text{trans}} \circ \mathcal{T}_{\text{syn}}$  doesn’t suffer as much, but its performance degradation when compared to *Syntax* is still large (ranges between 1 point on POS to 6 points on NER). **absence of sub-word overlap significantly hurts performance when languages differ in their word orders.**

performs better than *non-parallel (diff)*, with gains ranging between 2 points on XQuAD to 17 points on NER. This result shows that even for *original* and *derived* language pairs which differ only in their script, having parallel pre-training corpora leads to the best zero-shot transfer. Since large-scale parallel unsupervised data is hard to procure, the best alternative is to use corpora from similar domains (Wikipedia) rather than different ones (Wikipedia v.s. web text).

#### 4.4 Zero-shot performance is strongly correlated with embedding alignment

Our previous results (§ 4.2, 4.3) showed cases where zero-shot transfer between languages is poor when there is no sub-word overlap. To investigate this further, we analyze the static word embeddings learned by bilingual models and find that zero-shot transfer between languages is strongly correlated with the alignment between word embeddings for the *original* and *derived* languages.

**Setup** The *original* and the *derived* languages have a one-to-one correspondence between their sub-word vocabularies when we use *transliteration* ( $\mathcal{T}_{\text{trans}}$ ). For a token embedding in the *original*-language embedding matrix, its alignment score is 100% if it retrieves the corresponding token embedding in the *derived* language when a nearest-neighbor search is performed, and 0% otherwise.

We average the alignment score over all the tokens and call it *alignment*.

**Results** We measure the *alignment* of bilingual models pre-trained on different *original- derived* language pairs created using *transliteration*, namely the composed transformations (§ 4.2), *parallel*, and *non-parallel* (§ 4.3). We plot the *alignment* along with the corresponding  $\Delta_{(\text{BZ-BS})}$  scores for XNLI in Figure 4. Results for other tasks are in Appendix E.

We observe that higher *alignment* is associated with lower  $\Delta_{(\text{BZ-BS})}$ , implying better zero-shot transfer. *Alignment* is lower for composed transformations like  $\mathcal{T}_{\text{trans}} \circ \mathcal{T}_{\text{inv}}$  and  $\mathcal{T}_{\text{trans}} \circ \mathcal{T}_{\text{perm}}$  which have large and negative  $\Delta_{(\text{BZ-BS})}$ . *Alignment* also explains the results in Section 4.3, with *non-parallel* variants having lower alignment scores than *parallel*, which is in line with their lower  $\Delta_{(\text{BZ-BS})}$ . Overall, we find a strong and significant Spearman’s rank correlation between *alignment* and  $\Delta_{(\text{BZ-BS})}$ , with  $\rho = 0.94, p < .005$  for XNLI,  $\rho = 0.93, p < .005$  for NER, and  $\rho = 0.89, p < .01$  for POS, indicating that increasing the embedding alignment between languages helps improve zero-shot transfer.

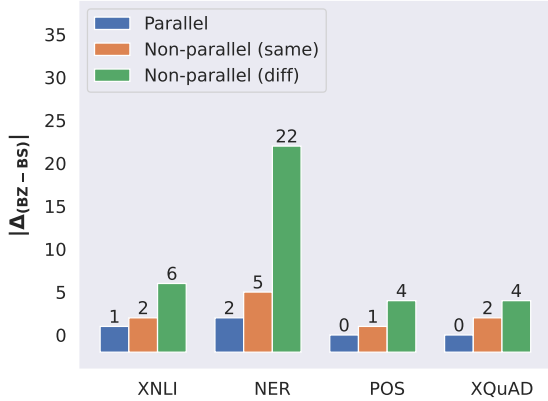


Figure 3:  $|\Delta_{(BZ-BS)}|$  for  $\mathcal{T}_{\text{trans}}$  under different conditions on the source of *original* and *derived* language pre-training corpora (hereon, corpora) (§ 4.3), averaged over four languages. Larger values imply worse zero-shot transfer. The breakdown of scores for different languages is in Appendix D. (1) *Non-parallel (diff)* (green bar), which uses corpora from different domains is worse than (2) *Non-parallel (same)* (orange bar), which uses different sets of sentences sampled from the same domain, which is in turn worse than (3) *Parallel*, which uses the same sentences. Having pre-training corpora from the same domain like Wikipedia (*Non-parallel (same)*) gives performance boosts between 2 points for QA to 17 points for NER when compared to *Non-parallel (diff)*.

#### 4.5 Connections to results on natural language pairs

**Effect of sub-word overlap** In § 4.2, we showed that when languages have different scripts (0% sub-word overlap), zero-shot transfer significantly degrades when they additionally have different word orders. However, the zero-shot transfer is good when languages differ only in the script and have similar or the same word order. This is in line with anecdotal evidence in Pires et al. (2019), where zero-shot transfer works well between *English* and *Bulgarian* (different script but same subject-verb-object order – SVO), but is poor between *English* and *Japanese* (different script *and* word order – SVO v.s. SOV). Our result also corroborates findings in Conneau et al. (2020b) that artificially increasing sub-word overlap between natural languages (which have different word orders) improves performance (e.g., 3 points on XNLI).

**Effect of token embedding alignment** In § 4.4, we showed that zero-shot transfer is strongly correlated with word embedding alignment between languages. This explains the usefulness of recent stud-

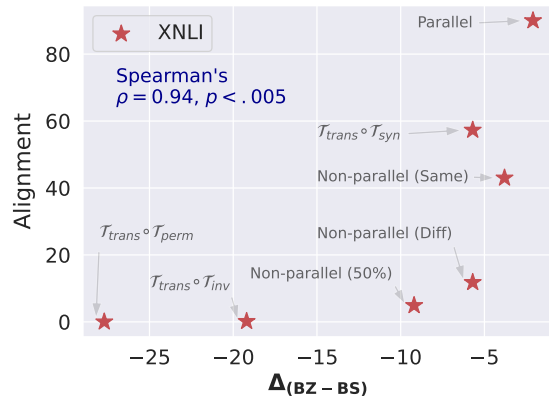


Figure 4:  $\Delta_{(BZ-BS)}$  for *Transliteration* ( $\mathcal{T}_{\text{trans}}$ ) variants on XNLI. Larger values (less negative) imply better zero-shot transfer. We see that alignment (§ 4.4) between token embeddings of different languages is correlated with  $\Delta_{(BZ-BS)}$ , and hence with better zero-shot transfer. For example,  $\mathcal{T}_{\text{trans}} \circ \mathcal{T}_{\text{inv}}$  (bottom left) which has poor zero-shot transfer also has lower alignment, whereas *Parallel* (top right) which has strong transfer is accompanied with higher alignment. We find a strong and statistically significant Spearman’s correlation of  $\rho_s = 0.94, p < .005$  on XNLI,  $\rho_s = 0.93, p < .005$  on NER, and  $\rho_s = 0.89, p < .01$  on POS. Plots for other tasks are in Appendix E.

ies which try to improve multilingual pre-training with the help of auxiliary objectives, which improve word or sentence embedding alignment.

DICT-MLM (Chaudhary et al., 2020) and ReLateLM (Khemchandani et al., 2021) require the model to predict cross-lingual synonyms as an auxiliary objective, thus indirectly improving word-embedding alignment and the zero-shot performance on multiple tasks. Hu et al. (2021) add an auxiliary objective that implicitly improves word embedding alignment and show that they can achieve performance similar to larger models. Cao et al. (2019) explicitly improve contextual word embedding alignment with the help of word-level alignment information in machine-translated cross-lingual sentence pairs. Since they apply this post hoc and not during pre-training, the improvement, albeit significant, is small (2 points on XNLI). While these studies do not fully utilize word and sentence embedding alignment information, our results lead us to posit that they are a step in the right direction and that baking alignment information more explicitly into pre-training will be beneficial.



## 5 Conclusion

Through a systematic study of zero-shot transfer between four diverse natural languages and their counterparts created by modifying specific properties like the script, word order, and syntax, we showed that (1) absence of sub-word overlap hurts zero-shot performance when languages differ in their word order, and (2) zero-shot performance is strongly correlated with word embedding alignment between languages. Some recent studies have implicitly or unknowingly attempted to improve alignment and have shown slight improvements in zero-shot transfer performance. However, our results lead us to posit that explicitly improving word embedding alignment during pre-training by using either supervised (e.g., parallel sentences and translation dictionaries) or unsupervised data will significantly improve zero-shot transfer. Although real languages typically differ in more ways than the set of properties considered in our transformations, our methodology is still useful to help isolate crucial properties for transfer. Future work can experiment with more sophisticated transformations and investigate closer connections with human language pairs.

## Acknowledgments

This work was funded through a grant from the Chadha Center for Global India at Princeton University. We thank Shunyu Yao, Vishvak Murahari, Tianyu Gao, Sadhika Malladi, and Jens Tuyls for reading our draft and providing valuable comments. We also thank the Google Cloud Research program for computational support in running our experiments.

## References

- Mikel Artetxe, Sebastian Ruder, and Dani Yogatama. 2020a. [On the cross-lingual transferability of monolingual representations](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 4623–4637. Association for Computational Linguistics.
- Mikel Artetxe, Sebastian Ruder, and Dani Yogatama. 2020b. [On the cross-lingual transferability of monolingual representations](#). *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*.
- Samuel Bowman, Christopher Potts, and Christopher D Manning. 2015. Recursive neural networks can learn logical semantics. In *Proceedings of the 3rd workshop on continuous vector space models and their compositionality*, pages 12–21.
- Steven Cao, Nikita Kitaev, and Dan Klein. 2019. Multilingual alignment of contextual word representations. In *International Conference on Learning Representations*.
- Aditi Chaudhary, Karthik Raman, Krishna Srinivasan, and Jiecao Chen. 2020. [DICT-MLM: improved multilingual pre-training using bilingual dictionaries](#). *CoRR*, abs/2010.12566.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020a. [Unsupervised cross-lingual representation learning at scale](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 8440–8451. Association for Computational Linguistics.
- Alexis Conneau and Guillaume Lample. 2019. [Cross-lingual language model pretraining](#). In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 7057–7067.
- Alexis Conneau, Ruty Rinott, Guillaume Lample, Adina Williams, Samuel R. Bowman, Holger Schwenk, and Veselin Stoyanov. 2018. [Xnli: Evaluating cross-lingual sentence representations](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics.
- Alexis Conneau, Shijie Wu, Haoran Li, Luke Zettlemoyer, and Veselin Stoyanov. 2020b. [Emerging cross-lingual structure in pretrained language models](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6022–6034.
- Jacob Devlin. 2019. Multilingual-BERT. <https://github.com/google-research/bert/blob/master/multilingual.md>. [Online; accessed 20-December-2020].
- 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.
- Sumanth Doddapaneni, Gowtham Ramesh, Anoop Kunchukuttan, Pratyush Kumar, and Mitesh M Khapra. 2021. A primer on pretrained multilingual language models. *arXiv preprint arXiv:2107.00676*.

- Philipp Dufter and Hinrich Schütze. 2020. Identifying elements essential for bert’s multilinguality. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4423–4437.
- Tsung-Yuan Hsu, Chi-Liang Liu, and Hung-Yi Lee. 2019. Zero-shot reading comprehension by cross-lingual transfer learning with multi-lingual language representation model. 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 5933–5940.
- Junjie Hu, Melvin Johnson, Orhan Firat, Aditya Siddhant, and Graham Neubig. 2021. Explicit alignment objectives for multilingual bidirectional encoders. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3633–3643.
- Junjie Hu, Sebastian Ruder, Aditya Siddhant, Graham Neubig, Orhan Firat, and Melvin Johnson. 2020. Xtreme: A massively multilingual multi-task benchmark for evaluating cross-lingual generalisation. In *International Conference on Machine Learning*, pages 4411–4421. PMLR.
- Karthikeyan K, Zihan Wang, Stephen Mayhew, and Dan Roth. 2020. [Cross-lingual ability of multilingual BERT: an empirical study](#). In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- Divyanshu Kakwani, Anoop Kunchukuttan, Satish Golla, NC Gokul, Avik Bhattacharyya, Mitesh M Khapra, and Pratyush Kumar. 2020. [inlpsuite: Monolingual corpora, evaluation benchmarks and pre-trained multilingual language models for indian languages](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: Findings*, pages 4948–4961.
- Yash Khemchandani, Sarvesh Mehtani, Vaidehi Patil, Abhijeet Awasthi, Partha Talukdar, and Sunita Sarawagi. 2021. [Exploiting language relatedness for low web-resource language model adaptation: An Indic languages study](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1312–1323, Online. Association for Computational Linguistics.
- Yu-Hsiang Lin, Chian-Yu Chen, Jean Lee, Zirui Li, Yuyan Zhang, Mengzhou Xia, Shruti Rijhwani, Junxian He, Zhisong Zhang, Xuezhe Ma, et al. 2019. Choosing transfer languages for cross-lingual learning. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3125–3135.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. 2016. Pointer sentinel mixture models. *arXiv preprint arXiv:1609.07843*.
- Joakim Nivre, Mitchell Abrams, Željko Agić, Lars Ahrenberg, Lene Antonsen, Maria Jesus Aranzabe, Gashaw Arutie, Masayuki Asahara, Luma Ateyah, and Mohammed Attia et al. 2018. [Universal dependencies 2.2](#). LINDAT/CLARIAH-CZ digital library at the Institute of Formal and Applied Linguistics (ÚFAL), Faculty of Mathematics and Physics, Charles University.
- Xiaoman Pan, Boliang Zhang, Jonathan May, Joel Nothman, Kevin Knight, and Heng Ji. 2017. Cross-lingual name tagging and linking for 282 languages. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1946–1958.
- 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.
- Peng Qi, Yuhao Zhang, Yuhui Zhang, Jason Bolton, and Christopher D Manning. 2020. Stanza: A python natural language processing toolkit for many human languages. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 101–108.
- Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving language understanding by generative pre-training.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2015. Neural machine translation of rare words with subword units. *arXiv preprint arXiv:1508.07909*.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. 2018. Glue: A multi-task benchmark and analysis platform for natural language understanding. In *International Conference on Learning Representations*.
- Dingquan Wang and Jason Eisner. 2016. The galactic dependencies treebanks: Getting more data by synthesizing new languages. *Transactions of the Association for Computational Linguistics*, 4:491–505.
- 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 2019, Hong Kong, China, November 3-7, 2019*, pages 833–844. Association for Computational Linguistics.

Zhengxuan Wu, Isabel Papadimitriou, and Alex Tamkin. 2022. Oolong: Investigating what makes crosslingual transfer hard with controlled studies. *arXiv preprint arXiv:2202.12312*.

## Appendices

### A Mathematical Specification for Transformation of Downstream Datasets

**Text classification** Text classification tasks like news classification or sentiment analysis typically have instances which contain a single sentence and a label. Instances in other classification tasks like natural language inference (NLI) (Bowman et al., 2015) contain two sentences and one label. For such tasks, we apply the transformation ( $\mathcal{T}$ ) on each sentence within every instance, and leave the annotated label as is. Therefore, for a dataset of size  $n$  which contains  $m$  sentences per instance, we have:

$$\mathcal{D}_{\text{orig}} = \{(s_{i1}, \dots, s_{im}, y_i) \mid i = 1 : N\}$$
$$\mathcal{D}_{\text{deriv}} = \{(\mathcal{T}(s_{i1}), \dots, \mathcal{T}(s_{im}), y_i) \mid i = 1 : N\}$$

**Token-classification tasks** Tasks like named-entity recognition (NER) and part-of-speech tagging (POS tagging) have labels associated with *each* token in the sentence. For these datasets, we ensure that any transformation ( $\mathcal{T}$ ) that changes the order of the tokens also changes the order of the corresponding labels.

We define a few quantities to express the transformation mathematically. Let  $s_i = (w_{i1}, \dots, w_{ik})$  be a sentence comprised of  $k$  tokens and  $y_i = (y_{i1}, \dots, y_{ik})$  be labels corresponding to the tokens in the sentence. We define a new transformation ( $\mathcal{T}_{\text{aug}}$ ) which operates on the label augmented sentence,  $s_i^{\text{aug}} = ((w_{i1}, y_{i1}), \dots, (w_{ik}, y_{ik}))$ . Let  $s_i^{\text{aug}}[j]$  correspond to the  $j^{\text{th}}$  element in the sequence, and  $s_i^{\text{aug}}[j][\text{word}]$  and  $s_i^{\text{aug}}[j][\text{label}]$  correspond to the word and label of the  $j^{\text{th}}$  element. Let  $\mathcal{T}_{\text{aug}}(s_i^{\text{aug}})[j][\text{orig}]$  denote the index of the  $j^{\text{th}}$  element in the transformed sequence with respect to the original sequence  $s_i^{\text{aug}}$ . Then, the new transformation  $\mathcal{T}_{\text{aug}}$  is such that,

$$\mathcal{T}_{\text{aug}}(s_i^{\text{aug}})[j][\text{orig}] = \mathcal{T}(s_i)[j][\text{orig}]$$
$$\text{Let } \text{orig\_j} = \mathcal{T}_{\text{aug}}(s_i^{\text{aug}})[j][\text{orig}]$$
$$\mathcal{T}_{\text{aug}}(s_i^{\text{aug}})[j][\text{label}] = s_i^{\text{aug}}[\text{orig\_j}][\text{label}]$$

We transform the dataset using  $\mathcal{T}_{\text{aug}}$ :

$$\mathcal{D}_{\text{orig}} = \{s_i^{\text{aug}} \mid i = 1 : N\}$$
$$\mathcal{D}_{\text{deriv}} = \{\mathcal{T}_{\text{aug}}(s_i^{\text{aug}}) \mid i = 1 : N\}$$

### B Zero-shot transfer results for different transformations

Table 5 in the appendix is the extended version of Table 3 in the main paper with a breakdown for all languages. It reports  $\Delta_{(\text{BZ}-\text{BS})}$ ,  $\Delta_{(\text{MZ}-\text{BS})}$ , and  $\text{BZ}$  for different languages and transformations considered.

### C Composed Transformations

Table 6 in the appendix presents the breakdown of results in Figure 2 of the main paper. It reports  $\Delta_{(\text{BZ}-\text{BS})}$  scores for composed transformations and their constituents.

### D Comparing different sources for original and derived language corpora

Table 8 in the appendix contains the breakdown of results in Figure 3 of the main paper. It reports  $\Delta_{(\text{BZ}-\text{BS})}$  for different languages on different tasks for the settings mentioned in Section 4.3.

### E Alignment Correlation

We present alignment results (Section 4.4) for all XNLI, NER, and POS in Figure 5. We observe strong correlations between alignment and zero-shot transfer, with  $\rho_s = 0.94, p < .005$  on XNLI,  $\rho_s = 0.93, p < .005$  on NER, and  $\rho_s = 0.89, p < .01$  on POS. We present the raw scores in Table 7.

### F Hyperparameters for XTREME

- XNLI: Learning rate –  $2e-5$ , maximum sequence length – 128, epochs – 5, batch size – 32.
- NER: Learning rate –  $2e-5$ , maximum sequence length – 128, epochs – 10, batch size – 32.
- POS: Learning rate –  $2e-5$ , maximum sequence length – 128, epochs – 10, batch size – 32.
- Tatoeba: Maximum sequence length – 128, pooling strategy – representations from the middle layer ( $\frac{n}{2}$ ) of the model.
- XQuAD: Learning rate –  $3e-5$ , maximum sequence length – 384, epochs – 2, document stride – 128, warmup steps – 500, batch size – 16, weight decay – 0.0001.

Task	Language	Inversion			Permutation			Syntax			Transliteration		
		BZ	$\Delta_{(BZ-BS)}$	$\Delta_{(MZ-BS)}$	BZ	$\Delta_{(BZ-BS)}$	$\Delta_{(MZ-BS)}$	BZ	$\Delta_{(BZ-BS)}$	$\Delta_{(MZ-BS)}$	BZ	$\Delta_{(BZ-BS)}$	$\Delta_{(MZ-BS)}$
XNLI	English	73.2	-3.4	-14.9	68.6	-5	-7.7	74.1	-1.8	-1.5	74.1	-1.7	-42.5
	French	62.5	-9.5	-8.8	68.4	-1	-7.6	69.6	-2.2	-1.4	71.6	-1.6	-39.9
	Hindi	43.9	-15.7	-15.8	51.2	-6.2	-13.1	61.6	-0.3	-1.6	63.4	-0.1	-29.4
	Arabic	54	-12.3	-12.5	62.1	-2.3	-6	65.9	0.7	0.3	68	-0.4	-35.1
	Avg.	58.4	<b>-10.2</b>	<b>-13</b>	62.6	<b>-3.6</b>	<b>-8.6</b>	67.8	<b>-0.9</b>	<b>-1.1</b>	69.3	<b>-1.0</b>	<b>-36.7</b>
NER	English	39.8	-44.5	-35.9	40.2	-28.5	-33.2	61.1	-7.8	-10.3	78	-2.1	-70.2
	French	54.5	-34.4	-51.3	44.4	-36.0	-39.8	59.6	-21.9	-25.9	84.3	-3.1	-87.4
	Hindi	19.4	-63.9	-63.2	38.5	-21.9	-37.4	64.8	-8.4	-7.3	84.4	-0.5	-82.9
	Arabic	37.8	-53.6	-36.3	66.2	-18.8	-31.1	66.1	-20.1	-23	88	-1.9	-89.9
	Avg.	37.9	<b>-49.1</b>	<b>-46.7</b>	47.3	<b>-26.3</b>	<b>-35.4</b>	62.9	<b>-14.6</b>	<b>-16.6</b>	83.7	<b>-1.9</b>	<b>-82.6</b>
POS	English	94.4	-0.7	-24.3	78.3	-11.9	-17.6	92.9	-0.9	-2.2	94.6	-0.5	-95.1
	French	74.3	-22.7	-22.9	82	-12.2	-20.9	93.5	-3.2	-5.2	97.2	-0.2	-97.4
	Hindi	19	-74.5	-74.5	51	-14	-41.8	91.6	-3.3	-11.3	96.5	-0.1	-96.6
	Arabic	69.2	-23	-23	83.1	-6.5	-20.6	79.4	-10	-11.5	93.2	-0.8	-90.9
	Avg.	64.2	<b>-30.2</b>	<b>-36.2</b>	73.6	<b>-11.2</b>	<b>-25.2</b>	89.4	<b>-4.4</b>	<b>-7.6</b>	95.4	<b>-0.4</b>	<b>-95.0</b>
XQuAD	English	30.4	-43.2	-35.5	-	-	-	-	-	-	72.4	-4	-73
	French	25.2	-29.5	-29.6	-	-	-	-	-	-	60.9	-1	-55.5
	Hindi	14.5	-27.3	-27.3	-	-	-	-	-	-	57.3	10.6	-43.5
	Arabic	21	-31.2	-31.4	-	-	-	-	-	-	54	-0.5	-51.7
	Avg.	22.8	<b>-32.8</b>	<b>-31.0</b>	-	-	-	-	-	-	61.2	<b>1.3</b>	<b>-55.9</b>

Table 5: This table is an extended version of Table 3 in the main paper. Smaller (more negative)  $\Delta_{(BZ-BS)}$  implies worse bilingual zero-shot transfer, whereas  $\Delta_{(BZ-BS)} \approx 0$  implies strong transfer.  $\Delta_{(BZ-BS)} \gg \Delta_{(MZ-BS)}$  implies that bilingual pre-training is extremely useful. Scores are highlighted based on their value (lower scores have a higher intensity of red). **(1) Discussing  $\Delta_{(BZ-BS)}$ :**  $\mathcal{T}_{trans}$  exhibits strong transfer on all tasks and languages (high  $\Delta_{(BZ-BS)}$  scores), and bilingual pre-training is extremely useful ( $\Delta_{(BZ-BS)} \gg \Delta_{(MZ-BS)}$ ), implying that zero-shot transfer is possible between languages with different scripts but the same word order.  $\mathcal{T}_{inv}$  and  $\mathcal{T}_{perm}$  suffer on all tasks (small  $\Delta_{(BZ-BS)}$  scores) whereas  $\mathcal{T}_{syn}$  suffers significantly lesser, which provides evidence that local changes to the word order made by *Syntax* ( $\mathcal{T}_{syn}$ ) hurts zero-shot transfer significantly lesser than global changes made by *Inversion* ( $\mathcal{T}_{inv}$ ) and *Permutation* ( $\mathcal{T}_{perm}$ ). **(1) Discussing  $\Delta_{(MZ-BS)}$ :**  $\Delta_{(BZ-BS)}$  is much larger than  $\Delta_{(MZ-BS)}$  for  $\mathcal{T}_{trans}$ , implying that bilingual pre-training (hereon, pre-training) is extremely useful.  $\Delta_{(BZ-BS)}$  and  $\Delta_{(MZ-BS)}$  are similar for  $\mathcal{T}_{inv}$  and  $\mathcal{T}_{syn}$ , implying that pre-training is not beneficial for these transformations.  $\Delta_{(BZ-BS)}$  is slightly larger than  $\Delta_{(MZ-BS)}$  for  $\mathcal{T}_{perm}$ , which means that pre-training is moderately useful.

$\mathcal{T}$	XNLI		NER		POS	
	BZ	$\Delta_{(BZ-BS)}$	BZ	$\Delta_{(BZ-BS)}$	BZ	$\Delta_{(BZ-BS)}$
$\mathcal{T}_{trans}$	74.1	-2.1	78	-2.3	94.6	-0.5
$\mathcal{T}_{inv}$	73.2	-3.4	39.8	-44.5	94.4	-0.7
$\mathcal{T}_{trans} \circ \mathcal{T}_{inv}$	55.7	-19.2	32.5	-51.5	52.2	-42.7
$\mathcal{T}_{perm}$	68.6	-5	40.2	-28.5	78.3	-11.9
$\mathcal{T}_{trans} \circ \mathcal{T}_{perm}$	44	-27.7	17.1	-46.3	29.5	-59
$\mathcal{T}_{syn}$	74.1	-1.8	61.1	-7.8	92.9	-0.9
$\mathcal{T}_{trans} \circ \mathcal{T}_{syn}$	69.8	-5.7	53.5	-14.2	91.5	-2

Table 6: Breakdown of results in Figure 2 of the main paper. *BZ* is the zero-shot performance.  $\Delta_{(BZ-BS)}$ ,  $\Delta_{(MZ-BS)}$ , and *BZ* are described in Section 3.3 and Table 2. Composing transformations always hurts  $\Delta_{(BZ-BS)}$  when compared to individual transformations.

Transliteration Variant	$\Delta_{(BZ-BS)}$ ( $\uparrow$ )			Alignment ( $\uparrow$ )
	XNLI	NER	POS	
Parallel	-2.1	-2.3	-0.5	90.0
Trans $\circ$ Syntax	-5.7	-14.2	-2	57.3
Non-parallel (Same)	-3.8	-4.1	-0.7	43.0
Non-parallel (Diff)	-5.7	-14.3	-1.5	11.8
Trans $\circ$ Inv	-19.2	-51.5	-42.7	0.16
Trans $\circ$ Perm	-27.7	-46.3	-59	0.01

Table 7:  $\Delta_{(BZ-BS)}$  and *alignment* scores for different *Transliteration* variants. The table contains raw scores for results in Section 4.4 of the main paper. Rows are sorted in descending order based on *alignment*. We observe strong correlations between alignment and zero-shot transfer, with  $\rho_s = 0.94, p < .005$  on XNLI,  $\rho_s = 0.93, p < .005$  on NER, and  $\rho_s = 0.89, p < .01$  on POS.

Task	Language	XNLI	NER	POS	XQuAD
		$\Delta_{(BZ-BS)}$	$\Delta_{(BZ-BS)}$	$\Delta_{(BZ-BS)}$	$\Delta_{(BZ-BS)}$
Parallel	English	-1.7	-2.1	-0.5	-4
	French	-1.6	-3.1	-0.2	-1
	Hindi	-0.1	-0.5	-0.1	10.6
	Arabic	-0.4	-1.9	-0.8	-0.5
	<b>Avg.</b>	<b>-1.0</b>	<b>-1.9</b>	<b>-0.4</b>	<b>1.3</b>
Non-parallel (Same)	English	-3.8	-4.1	-0.7	-6.9
	French	-1	-6.3	-0.5	-0.9
	Hindi	-0.4	-3.1	-0.2	4.5
	Arabic	-2	-6.1	-1.5	0.7
	<b>Avg.</b>	<b>-1.8</b>	<b>-4.9</b>	<b>-0.7</b>	<b>-0.6</b>
Non-parallel (Diff)	English	-5.7	-14.3	-1.5	-9.3
	French	-10.9	-30.3	-10.5	-5.2
	Hindi	-0.5	-8.6	-1	5
	Arabic	-6.3	-34.7	-3.7	-1.9
	<b>Avg.</b>	<b>-5.9</b>	<b>-22.0</b>	<b>-4.2</b>	<b>-2.9</b>

Table 8:  $|\Delta_{(BZ-BS)}|$  for  $\mathcal{T}_{\text{trans}}$  under different conditions on the source of *original* and *derived* language pre-training corpora (§ 4.3). Larger values imply worse zero-shot transfer. For all languages: (1) *Non-parallel (diff)*, which uses corpora from different domains is worse than (2) *Non-parallel (same)*, which uses different sets of sentences sampled from the same domain, which is in turn worse than (3) *Parallel*, which uses the same sentences.

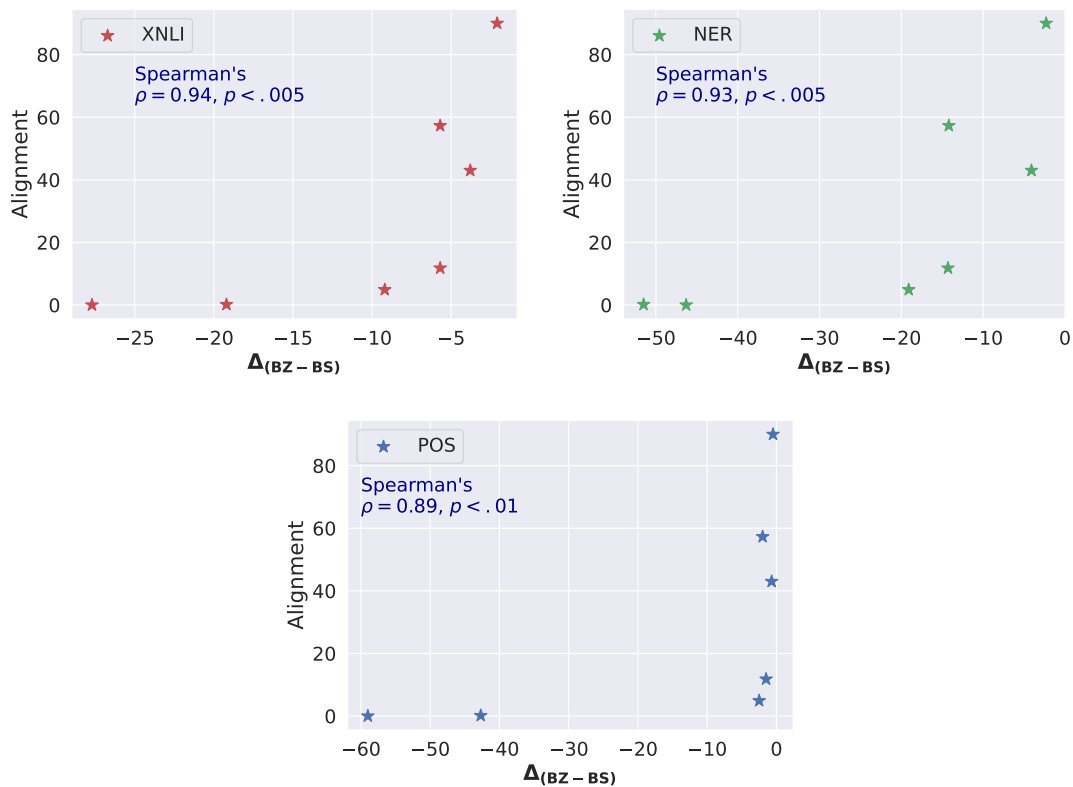


Figure 5: Alignment v.s.  $\Delta_{(BZ-BS)}$  plots for XNLI, NER, and POS. We observe strong correlations between alignment and zero-shot transfer, with  $\rho_s = 0.94, p < .005$  on XNLI,  $\rho_s = 0.93, p < .005$  on NER, and  $\rho_s = 0.89, p < .01$  on POS.