

Joint Encoding of Linguistic Information in BERT-based Language Models

Apostolos Panagiotopoulos (#14021900), **Konstantinos Papakostas** (#13914332),
Matteo Rosati (#13858149), **Giulio Starace** (#13010840)

{ apostolos.panagiotopoulos, konstantinos.papakostas,
matteo.rosati, giulio.starace }@student.uva.nl

Abstract

BERT-based Language Models (BLMs) have gained considerable popularity due to their impressive results across diverse NLP tasks. However, their representations are still not interpretable. In this work, we study the information-sharing mechanisms in BLMs and highlight the presence of joint encoding across linguistic entities, such as part-of-speech tags and dependency labels. Furthermore, we find that the way this information is shared resembles the linguistic relations between entities and is consistent across languages. Our results are validated with RoBERTa and XLM-R in English, Italian, and Greek, and take a step toward a deeper understanding of information flow in BLMs. Our code is available at <https://github.com/thesofakillers/bert-infoshare>

1 Introduction

In recent years, contextualized word embeddings obtained from Large Language Models (LLMs) have consistently demonstrated state-of-the-art performance across multiple NLP tasks (Devlin et al., 2019; Liu et al., 2019; Wu and Dredze, 2020). However, these representations have proven hard to interpret with regard to the information they contain, leading to a series of works (Tenney et al., 2018; Hewitt and Manning, 2019; Ravishankar et al., 2019; Libovický et al., 2020; Choenni and Shutova, 2022) that probe these models to extract insights. Choenni and Shutova (2022) showed that multilingual LLMs jointly encode information across languages, which inspired us to investigate whether they also jointly encode information across different linguistic entities, such as *nouns* and *verbs*.

We approach this problem by studying how two state-of-the-art models, RoBERTa (Liu et al., 2019) and XLM-R (Conneau et al., 2020), encode information about part-of-speech (POS) tags and dependency relations, as well as how this information

may be shared between classes. We hypothesize that these BERT-based Language Models (BLMs) learn to jointly encode words from related classes in a way that mirrors linguistic relationships. We expect to see that learned representations for linguistic features are not completely independent, with information sharing between them increasing with higher co-occurrences or functional dependencies (e.g. nouns and determinants or nominal subject and clausal subject dependencies). These overlapping relationships can vary across languages, though we anticipate that lexical entities, i.e., certain parts-of-speech such as nouns, share some degree of information in a language-agnostic manner.

Our methodology leverages a cross-neutralizing method based on centroid estimation (Libovický et al., 2020), akin to the multilingual approach in Choenni and Shutova (2022). In particular, we use a pre-trained instance of RoBERTa (Liu et al., 2019) as the encoder architecture and probe various layers to perform the POS tagging and dependency labeling tasks. Then, we apply the cross-neutralizing method to uncover which entities share information. We extend our experiments to the multilingual setting using XLM-R (Conneau et al., 2020) for the English, Greek, and Italian languages.

This work bridges the gap between analyses of multilingual information sharing and attempts to localize granular semantic and syntactic information within LLMs. Tenney et al. (2019) and Clark et al. (2019) have evaluated *where* information is encoded by probing different layers within BERT. However, their work lacked an analysis of *how* this information is shared. On the other hand, papers that have explored a direction similar to ours focus on how information is shared *across languages* (Choenni and Shutova, 2022), rather than *across linguistic features*. By combining these two approaches we can arrive at a more comprehensive understanding of how information about linguistic

units is represented in BLMs.

We find that similarly to typological properties of languages (Choenni and Shutova, 2022), information is shared across different linguistic features in BLMs. We report that this phenomenon extends to the multilingual setting in the languages tested, with variations in which labels are jointly encoded that are language-specific. Additionally, cross-lingual self-neutralization (i.e., a NOUN centroid from Italian used to cross-neutralize the NOUN class in Greek) also yields a decrease in performance, suggesting that centroids can encapsulate language-agnostic features.

2 Related Work

Prior work has focused on two main approaches to inspect the data contained in contextualized embeddings; through descriptive “probing tasks” (Conneau et al., 2018) and through informative subspace representations that are used to decompose word vectors (Libovický et al., 2020).

At the sentence embedding level, a prominent curated suite of probing tasks was introduced by Conneau et al. (2018). These tasks examine whether sentence embeddings encode properties such as sentence length, tree depth, etc. This work is extended by Hewitt and Manning (2019) by proposing methods to discover linear transformations that encode word distance and tree depth in sentence parse trees using word representations. Tenney et al. (2018) goes further with “edge probing” tasks, which probe the token-level representations directly and expand to a broader range of syntactic and semantic tasks at the sub-sentence level.

In later work, Tenney et al. (2019) apply these probing methods to BERT’s hidden states to quantify where linguistic information is captured within the network. They find that the model represents the steps of the traditional NLP pipeline, with lower-level, syntactic features appearing earlier than more complex semantic roles and structure. This work has been extended in the multilingual setting to probe for the previously mentioned linguistic information (Ravishankar et al., 2019) or novel information, such as determining where typological properties of certain languages are encoded (Choenni and Shutova, 2020). Additionally, Clark et al. (2019) explore BERT’s attention mechanism to identify what linguistic characteristics it attends to, showing that they correspond to linguistic notions of syntax and co-reference.

Other than identifying which information is encoded in word embeddings and hidden states, substantial effort has been put into trying to identify how this encoded information is shared across linguistic categories and across languages (Şahin et al., 2020; Blevins et al., 2018), with most of the work being focused on the multilingual setting (Chi et al., 2020). In particular, Libovický et al. (2020) developed a method that successfully removes language-specific features without affecting encoded language-agnostic information (i.e., semantic meaning) by leveraging what they refer to as “language centroids”. Our work draws inspiration from that of Choenni and Shutova (2022), who probe for joint encoding of typological features of different languages by extending the previous method. They find that these feature values are encoded jointly across languages and are localizable in their respective language centroids.

Compared to this existing work, we apply the cross-neutralizing method to probe BLM’s hidden representations for joint encoding of syntactic classes instead of languages. In sum, we apply the multilingual methods to a multi-label paradigm. Using the centroid method introduced by Libovický et al. (2020) and extended in Choenni and Shutova (2022), we explore what class information is independent and what is jointly encoded, as well as where in the model this information emerges. We also explore how these results vary from monolingual to multi-lingual models, as well as across languages.

3 Methods

3.1 Extracting Word Representations

We use the base versions of RoBERTa (Liu et al., 2019) and XLM-R (Conneau et al., 2020) as our encoders, which produce a contextualized embedding for each sub-word token that is generated by their tokenizers. To extract a set of representations on the word level, we need to aggregate the sub-word token representations. We do so by either (i) taking the representation of the first sub-word token, (ii) taking the max-pooling over the sub-word tokens along the 768 dimensions, or (iii) taking the mean of the sub-word tokens.

3.2 NLP Tasks

To study information sharing in the aforementioned models we probe them for two well-studied NLP tasks: POS tagging and dependency labeling. For

POS tagging, we encode each sentence and probe its word-level representations from a given layer l using a shallow MLP. For dependency labeling, we follow a similar scheme and extract word-level representations for the child and the head of each dependency in the sentence. We concatenate them¹ to form a feature vector of size $2 \times H$, where H is the hidden state size of the encoder model, and feed it to a shallow MLP to predict their dependency relation. We include training details for the POS tagging and dependency labeling probes in Appendix A.

3.3 Information Sharing

We examine whether BLMs encode information about certain linguistic categories, such as POS tags (nouns, adjectives, etc.) or syntactic dependency relations (nominal subjects, objects, etc.) in common subspaces of their representation space.

To achieve this, we first localize the subspace of the model that corresponds to each of these categories by obtaining their mean vector representation, similarly to Libovický et al. (2020). For *POS tagging*, this corresponds to the centroid of each target POS tag t , which, given the word representations \mathbf{v} , is formally defined as:

$$\mathbf{u}_t^{(POS)} = \frac{1}{|V_t|} \sum_{\mathbf{v} \in V_t} \mathbf{v} \quad (1)$$

where V_t is the set of the representations of the words in our validation set that were predicted as tag t when probing the l -th layer of our encoder. For *dependency labeling*, we slightly modify the definition of the centroid, as each prediction depends on both the representation of the head \mathbf{h} and the child \mathbf{c} of the dependency relation. Here, the centroid of each target dependency label t is defined as:

$$\mathbf{u}_t^{(DEP)} = \frac{1}{|P_t|} \sum_{(\mathbf{h}, \mathbf{c}) \in P_t} [\mathbf{h} ; \mathbf{c}] \quad (2)$$

where P_t is the set of the (head, child) representation pairs that were predicted as dependency t when probing the l -th layer of our encoder, and $[\mathbf{h} ; \mathbf{c}]$ is the concatenation of the two vectors.

After obtaining the mean vector \mathbf{u}_t for each linguistic category t , we investigate information sharing with the cross-neutralizing method from

¹We tested different concatenation configurations, such as including the mean vector or the absolute difference of the pair, but found the simpler approach to perform equally well.

Choenni and Shutova (2022). To cross-neutralize with \mathbf{u}_t , we encode each word in the sentence using our pre-trained encoder and subtract \mathbf{u}_t from its representation. We pass the resulting representations through our probing network to obtain a new distribution over the linguistic categories and re-classify words, in case of POS tagging, or words with their heads, in case of dependency labeling.

To interpret our results, we group them in (neutralizer, target) pairs, meaning that we observe the change in classification accuracy for the class `target` when neutralized with the centroid of the class `neutralizer`. In the special case where the same linguistic class is both the target and the neutralizer, we refer to our method as *self-neutralizing*, and otherwise as *cross-neutralizing*. We argue that pairs of linguistic categories that result in substantial drops in performance after cross-neutralizing are jointly encoded, and hence information is shared among them.

3.4 Layer and Aggregation Selection

The choice of the layer l from which to extract the token embeddings is not arbitrary. The same non-triviality presents itself in the choice of aggregation function mentioned in Section 3.1. An ideal selection should capture the greatest amount of information for the words and the task at hand. For each aggregation function, we self-neutralize using embeddings from different layers and compare the relative drop in self-neutralizing accuracy.

We hypothesize that the configurations where the self-neutralizing percentage drop is the greatest are those closest to an ideal selection, as a higher decrease in performance signals in a higher amount of information being captured in the centroids. We also consider that this drop should be relative to a high baseline accuracy, to achieve more accurate centroid definitions. Therefore, we take the top quartile of our configurations in terms of baseline accuracy, rank them in descending order of relative drop in accuracy, and select the first to perform cross-neutralization.

4 Experiments

4.1 Datasets

We use the Universal Dependencies framework (Nivre et al., 2020) to train our probing models in POS tagging and dependency labeling, as it offers a consistent annotation style across a collection of treebanks over multiple languages. In

particular, we use the GUM corpus (Zeldes, 2017) for English, which consists of 9,130 sentences with a total of 164,488 words, the VIT corpus (Delmonte et al., 2017) for Italian with 10,087 sentences and 279,723 total words, and the GDT corpus (Prokopenko and Papageorgiou, 2017) for Greek with 2,521 sentences and 63,441 total words.

All of the datasets contain token-level annotations with a total of 17 possible POS tags and 36 dependency relations. Additionally, for dependency relations, a unique head corresponds to each word, with the exception of the `root` relation, which has no head, and hence it is omitted. A detailed description of the pre-processing pipeline for all three treebanks can be found in Appendix B.

4.2 Baseline Models

We use RoBERTa and XLM-R² as our encoders, and we keep their weights frozen while training our probe classifiers on the aforementioned datasets for both tasks. We try various combinations of layers to probe ($l \in \{1, 3, 6, 9, 12\}$) and sub-word token aggregations (`first`, `mean`, `max-pooling`), and use the one that achieves a high baseline accuracy while resulting in the largest accuracy drop when self-neutralizing, as described in Section 3.4. We display the results we obtained per configuration in Appendix C. We report the classification accuracy for the best self-neutralizing combination for all our baseline models in Table 1, and we observe that we are in line with the current literature (Straka et al., 2019). As we only employ a pre-trained encoder without fine-tuning on our corpus, we expect to see a slightly lower accuracy compared to state-of-the-art approaches in each task.

	RoBERTa	XLM-R		
	en_gum	en_gum	it_vit	el_gdt
POS	95.6%	95.5%	97.4%	97.9%
DEP	90.9%	91.4%	93.9%	94.8%

Table 1: Classification accuracy for part-of-speech tagging and dependency labeling for our baseline models.

4.3 Cross-Neutralizing in English

To perform cross-neutralization, we subtract the centroid of a given label (`neutralizer`) from the embeddings of our encoder and compare the accuracy drop over the labels (`targets`). Although it is reasonable to expect that subtracting

²The pre-trained models are available on Hugging Face: <https://huggingface.co/roberta-base> <https://huggingface.co/xlm-roberta-base>

a fixed vector from our embeddings might always cause a performance drop, we experimented with subtracting random vectors and noticed that the performance remained unchanged. Thus, we argue that considerable changes in accuracy for a (`neutralizer`, `target`) pair suggest that the model learns some information about these two classes jointly.

While the labels for POS tagging are trivially bound to a single word, dependency labels refer to pairs of child and head words, leading to the definition of centroids being less straightforward. We experiment with considering only the embedding of the child as the centroid, versus a concatenation of both the child and the head. The latter results in more consistent self-neutralization across the task labels, which we assume suggests a better choice for centroid definition. Therefore, we use this concatenation configuration in all further dependency labeling cross-neutralization analyses.

4.4 Monolingual vs Multilingual Models

We also extend our experiments to the multilingual setting using XLM-R as the encoder and evaluating cross-neutralization in three different languages. This lets us both verify the consistency of our method across models, as we can directly compare the results between RoBERTa and XLM-R in English, but also across languages, as we can compare the effects of cross-neutralization on English, Italian, and Greek. By doing so, we can also assess the effect of cross-neutralization at a language-family level, as these are three Indo-European languages.

4.5 Cross-Lingual Cross-Neutralizing

In the multilingual setting, we test whether our findings are consistent across languages. A different direction worth exploring is whether there is a joint encoding of information between two related linguistic categories from two different languages. We test this hypothesis by cross-neutralizing every linguistic entity in one language with another entity from another language, e.g. all Italian POS tags with English nouns.

4.6 Cross-Task Cross-Neutralizing

To study joint sharing of linguistic information *across tasks*, we additionally subtract the POS centroids from the child embeddings in child-parent concatenations for dependency labelling. In this way we effectively cross-neutralise across tasks,

with POS neutralizers and dependency label targets. Similarly, we subtract the child portion of the dependency label centroids from the POS embeddings for POS tagging, to test whether any sort of joint learning happens in both directions.

5 Results & Analysis

5.1 Monolingual Cross-Neutralization

The decrease in accuracy when cross-neutralizing English POS tags and dependency labels is presented in Figure 1. Overall, we observe a pattern that verifies our hypothesis that related classes are jointly encoded in RoBERTa, for both the POS and DEP tasks. Subtracting the centroids associated with specific classes leads to a decrease in accuracy for other classes, indicating that information is shared across them. Notably, this cross-neutralizing phenomenon is not necessarily symmetric; in several instances, the accuracy change for a label y when using neutralizer x is not similar in magnitude to the accuracy change for x when y is the neutralizer. This may be due to the structure of language, which can be represented as a directed graph in which linguistic entities may modify others, but may not themselves be modified by their object. For example, adjectives modify nouns, and yet nouns do not modify adjectives. Further work could further clarify why certain entities cross-neutralize symmetrically and why others do not. In the rest of this subsection, we kept a subset of the linguistic categories both for the POS and DEP tasks, to highlight some task-specific patterns that emerge from our cross-neutralizing experiments. The figures containing the cross-neutralizing results between all classes can be found in Appendix D.

POS Tagging The top part of Figure 1 shows the impact of cross-neutralization on a selection of POS tags. The shown POS tags have been selected to ensure a fair distribution between open and closed-class words. We observe that while there is evidence of linguistically related units being jointly encoded, this phenomenon is not ubiquitous. For instance, we see that the auxiliaries (AUX) and verbs (VERB) are jointly encoded, as indicated by the relative decrease of 53% in VERB classification accuracy. This aligns with our hypothesis since auxiliaries can be considered function words acting on verbs. However, the same is not observed when neutralizing nouns (NOUN)

	NOUN	ADJ	VERB	PRON	DET	NUM	ADV	AUX
POS Neutralizer								
NOUN	-90	-21	-10	-13	-8	-32	-32	-15
ADJ	-8	-90	-11	-11	-3	-17	-29	-12
VERB	-6	-14	-90	-5	-3	-15	-19	-11
PRON	-3	-4	-1	-86	-3	-2	-9	-4
DET	-2	-16	-0	-12	-89	-21	-24	-5
NUM	-5	-10	-3	-6	-3	-94	-13	-1
ADV	-2	-21	-0	-3	-2	-14	-89	-2
AUX	-3	-12	-53	-2	-2	-5	-11	-89
DEP Neutralizer								
PUNCT	-99	4	-4	-11	-61	-4	-1	-4
NSUBJ	-0	-97	-12	-21	-45	-7	-1	-6
OBJ	-0	-25	-97	-61	-68	-13	-2	-9
OBL	-1	-10	-16	-97	-91	-15	-3	-7
ADVCL	-0	-1	-3	-17	-93	-4	-3	-5
CASE	0	-5	-10	-17	-26	-96	-2	-2
DET	0	-3	-6	-21	-29	-3	-98	-8
AMOD	0	-7	-7	-16	-22	-0	-2	-96
	PUNCT	NSUBJ	OBJ	OBL	ADVCL	CASE	DET	AMOD
Target								

Figure 1: Relative change in accuracy when cross-neutralizing 8 sampled POS (top) and DEP (bottom) tags using embeddings from RoBERTa in English.

with determiner (DET) centroids, where the relative percentage decrease is only 2%, suggesting that these tags are learned independently, despite determiners acting as a lexical modifier for nouns. Finally, we note how certain tags are particularly susceptible to (such as numerals, NUM) or adept at (such as NOUN) cross-neutralization, regardless of their relationship to the neutralizer or target.

Dependency Labeling The lower half of Figure 1 also shows the impact of cross-neutralization on a selection of dependency labels. These labels were selected to highlight the major trends in joint encoding within this task. We observe that information sharing is most pronounced when dependency classes are (1) closely related by definition, or (2)

		en_gum								it_vit								el_gdt							
POS Neutralizer		NOUN	ADJ	VERB	PRON	DET	NUM	ADV	AUX	NOUN	ADJ	VERB	PRON	DET	NUM	ADV	AUX	NOUN	ADJ	VERB	PRON	DET	NUM	ADV	AUX
	NOUN	-92	-29	-15	-13	-12	-39	-25	-20	-98	-25	-11	-16	-3	-39	-22	-20	-89	-34	-5	-16	-1	-28	-20	-13
	ADJ	-13	-91	-16	-8	-4	-27	-31	-18	-21	-96	-22	-17	-0	-26	-44	-18	-18	-93	-6	-9	0	-25	-40	-8
	VERB	-10	-25	-94	-7	-7	-33	-17	-13	-33	-63	-96	-9	-2	-43	-20	-10	-18	-18	-90	-9	-1	-29	-12	-6
	PRON	-7	-3	-3	-91	-3	-13	-6	-7	-13	-5	-9	-93	-3	-17	-7	-14	-4	-7	-0	-88	0	-13	-15	-2
	DET	-1	-8	-2	-18	-93	-19	-15	-6	-6	-7	-4	-32	-94	-35	-8	-23	-1	-28	-1	-27	-93	-48	-18	-16
	NUM	-13	-16	-8	-8	-4	-94	-23	-4	-59	-47	-12	-18	-1	-95	-27	-25	-17	-23	-4	-1	-0	-91	-10	-2
	ADV	-4	-22	-2	-2	-2	-27	-97	-7	-8	-27	-7	-8	-3	-27	-94	-8	-8	-11	-1	-5	0	-13	-94	-3
	AUX	-5	-12	-59	-3	-3	-22	-9	-94	-6	-16	-60	-7	-2	-19	-5	-96	-3	-2	-6	-5	0	-28	-16	-97
DEP Neutralizer		en_gum								it_vit								el_gdt							
	PUNCT	-100	-0	-13	-38	-40	-8	-5	-15	-99	-1	-15	-14	-52	-3	-0	-2	-100	-0	-3	-40	-19	-3	-0	-5
	NSUBJ	-0	-96	-13	-17	-32	-10	-3	-8	0	-98	-43	-64	-61	-4	-2	-11	-1	-95	-14	-15	-18	-3	-1	-5
	OBJ	-0	-25	-98	-66	-63	-19	-5	-14	0	-40	-98	-67	-42	-5	-3	-9	-1	-9	-94	-33	-12	-4	-1	-8
	OBL	0	-6	-18	-99	-90	-18	-7	-12	0	-13	-27	-98	-68	-6	-1	-17	-1	1	-7	-95	-40	-4	-0	-8
	ADVCL	-0	0	-5	-34	-91	-7	-3	-5	0	-2	-14	-9	-100	-1	0	-70	1	0	-3	-22	-99	-1	-0	-13
	CASE	0	-7	-16	-23	-18	-98	-2	-10	0	-4	-17	-12	-36	-98	-3	-15	2	-9	-16	-19	-43	-95	-5	-13
	DET	0	-3	-9	-18	-29	-4	-98	-14	0	-6	-15	-22	-43	-4	-96	-23	0	-6	-14	-14	-47	-2	-96	-26
	AMOD	-1	-8	-15	-21	-28	-5	-1	-97	0	-5	-14	-18	-63	-5	-1	-99	-1	-2	-6	-9	-32	0	-1	-94
		Target								Target								Target							
		PUNCT	NSUBJ	OBJ	OBL	ADVCL	CASE	DET	AMOD	PUNCT	NSUBJ	OBJ	OBL	ADVCL	CASE	DET	AMOD	PUNCT	NSUBJ	OBJ	OBL	ADVCL	CASE	DET	AMOD

Figure 2: Relative change in accuracy when cross-neutralizing 8 sampled POS (top) and DEP (bottom) tags using embeddings from XLM-R in English (left), Italian (center) and Greek (right).

hierarchically connected. Both of these aspects support our hypothesis of joint encoding of linguistically related dependencies and elevate important details of linguistic structure. As an example of (1), one can consider the classes NSUBJ, OBJ, and OBL. Each is jointly encoded with the others, as they help define the basic syntactic structure of the sentence in relation to the predicate.

On the other hand, the results of PUNCT and ADVCL lend support to (2). ADVCL is a clausal dependency, and is therefore significantly affected when removing information that is more related to word components of clauses such as OBJ, AMOD, CASE, and PUNCT. This latter dependency is quite noteworthy; it is a strong neutralizer across several classes, and yet is not affected by the removal of any other centroid. Punctuation plays a strong role in defining the syntactic tree of a sentence and therefore removing that information will distort its structure and thus the predicted dependency labels. For example, we can consider the striking differ-

ence between the sentences “let’s eat, grandma” and “let’s eat grandma”. On the other hand, one can read a sentence in a foreign language and still be able to point out the punctuation, indicating that, even with limited information, punctuation information is orthogonal to other linguistic features.

5.2 Multilingual Cross-Neutralization

In this section, we repeat the evaluation of Section 5.1 using embeddings from XLM-R trained on English, Italian and Greek datasets and report results in Figure 2. To facilitate comparison both across languages and with the results acquired from RoBERTa, we make the same selection of POS tags and dependency labels and present the full results in Appendix D.

Consistency Across Models When comparing RoBERTa and XLM-R on the English dataset, we observe similar patterns across linguistic categories for both POS tagging and dependency labeling. For example, numerals seem to be jointly encoded with

nouns and auxiliaries. This supports the idea that even in a multilingual training setting, the linguistic entities of each language are encoded similarly to a monolingual model.

Consistency Across Languages Similar patterns can also be observed when comparing across languages in XLM-R. For instance, we see that nouns consistently neutralize numerators by a significant factor. Likewise, in the dependency labeling setting, we note that the lack of neutralization of determiners (DET) is consistent across all three languages. This suggests that the tags involved in the cross-neutralization pairs where this observation holds are language-agnostic. On the same note, language-specific behavior is also evident. For instance, VERB POS tags cross-neutralize adjectives (ADJ) more prominently in Italian than in the other two languages. Similarly, while the neutralization of ADVCL by NSUBJ in English dependency labeling is significant, the effect is less pronounced in the other two languages. The results suggest that XLM-R discovers and leverages language-agnostic information when possible, while also learning language-specific information when necessary.

5.3 Cross-Lingual Cross-Neutralization

We extend the experiments of the previous section by cross-neutralizing using linguistic categories from different languages. Figure 3 shows the results of cross-neutralizing for Italian POS tagging and dependency labeling when using the corresponding English centroids from each linguistic category. We notice a substantial drop in performance across the diagonals of these graphs, which correspond to the self-neutralizing scenario. This indicates the presence of information sharing between languages but is limited to the same linguistic feature. We also observe a correlation between different features of different languages, such as Italian adverbs and English adjectives, or Italian nominal subjects and English objects. This behavior highlights the tremendous learning capacity of these models, as information can be shared between the representations for different languages.

Finally, we repeat the experiment by using English neutralizers on the Greek corpus, and we observed similar but milder trends. We hypothesize that this is due to English being linguistically closer to Italian than Greek, which leads to a higher degree of information sharing between these two languages. For more detailed plots on cross-

		NOUN	ADJ	VERB	PRON	DET	NUM	ADV	AUX
English POS Neutralizer	NOUN	-79	-23	-13	-13	-2	-48	-23	-29
	ADJ	-10	-70	-14	-9	-1	-29	-47	-19
	VERB	-6	-22	-67	-5	-2	-41	-13	-13
	PRON	-2	-6	-3	-73	-8	-24	-7	-1
	DET	-1	-4	-1	-17	-35	-35	-5	-1
	NUM	-19	-26	-10	-11	-1	-78	-25	-23
	ADV	-2	-20	-7	-5	-1	-34	-90	-10
	AUX	-1	-10	-32	-5	-1	-18	-8	-52
English DEP Neutralizer	CONJ	-96	-3	-8	-9	-23	-1	-0	-8
	NSUBJ	-13	-5	-5	-14	-53	-4	0	-7
	OBJ	-20	-82	-16	-18	-15	-4	-1	-6
	OBL	-11	-4	-13	-25	-23	-2	-7	-18
	ADVCL	-30	-1	-6	-17	-86	-5	0	-10
	CASE	-16	-3	-9	-7	-20	-74	-2	-5
	DET	-10	-2	-4	-7	-7	-2	-53	-1
	AMOD	-17	-3	-9	-11	-36	-2	-2	-86
		CONJ	NSUBJ	OBJ	OBL	ADVCL	CASE	DET	AMOD
		Italian Target							

Figure 3: Relative change in accuracy when cross-neutralizing Italian with 8 sampled POS (top) and DEP (bottom) English tags using embeddings from XLM-R.

lingual cross-neutralization, we refer the reader to Appendix D.

6 Conclusion

In this work, we studied the joint encoding of linguistic entities in BLMs. We approached the problem by cross-neutralizing pairs of entities in two well-studied tasks: POS tagging and dependency labeling. We first evaluated our methodology in a simple monolingual scenario using RoBERTa and found joint encoding of similar linguistic categories for both tasks. We extended our evaluation to include XLM-R and found our hypothesis to hold both across languages and to be persistent across two different BERT-based models. Finally, we noticed that related linguistic categories are jointly

encoded across different languages in XLM-R.

Our findings are consistent with recent work on information sharing in BLMs and shed light on the flow of information within them. They also give credit to the efficacy of multilingual models over the use of multiple monolingual ones. We hope that our work will aid the exploration and interpretation of BLMs and encourage further research on their information-sharing mechanisms. As a future direction, we would like to extend our experimental setting to include more downstream NLP tasks and evaluate our findings in more languages.

References

- Terra Blevins, Omer Levy, and Luke Zettlemoyer. 2018. [Deep RNNs Encode Soft Hierarchical Syntax](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 14–19, Melbourne, Australia. Association for Computational Linguistics.
- Ethan A. Chi, John Hewitt, and Christopher D. Manning. 2020. [Finding Universal Grammatical Relations in Multilingual BERT](#). *ACL*.
- Rochelle Choenni and Ekaterina Shutova. 2020. [What does it mean to be language-agnostic? Probing multilingual sentence encoders for typological properties](#). *arXiv:2009.12862 [cs]*. ArXiv: 2009.12862.
- Rochelle Choenni and Ekaterina Shutova. 2022. [Investigating Language Relationships in Multilingual Sentence Encoders Through the Lens of Linguistic Typology](#). *Computational Linguistics*, pages 1–38.
- Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. 2019. [What Does BERT Look at? An Analysis of BERT’s Attention](#). In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 276–286, Florence, Italy. Association for Computational Linguistics.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. [Unsupervised Cross-lingual Representation Learning at Scale](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.
- Alexis Conneau, German Kruszewski, Guillaume Lample, Loïc Barrault, and Marco Baroni. 2018. [What you can cram into a single \\$&!#* vector: Probing sentence embeddings for linguistic properties](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2126–2136, Melbourne, Australia. Association for Computational Linguistics.
- Rodolfo Delmonte, Antonella Bristot, and Sara Tonelli. 2017. *VIT – Venice Italian Treebank: Syntactic and Quantitative Features*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- John Hewitt and Christopher D. Manning. 2019. [A Structural Probe for Finding Syntax in Word Representations](#). 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 4129–4138, Minneapolis, Minnesota. Association for Computational Linguistics.
- Taku Kudo and John Richardson. 2018. [SentencePiece: A simple and language independent subword tokenizer and detokenizer for Neural Text Processing](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71, Brussels, Belgium. Association for Computational Linguistics.
- Jindřich Libovický, Rudolf Rosa, and Alexander Fraser. 2020. [On the Language Neutrality of Pre-trained Multilingual Representations](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1663–1674, Online. Association for Computational Linguistics.
- 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:1907.11692 [cs]*. ArXiv: 1907.11692.
- Ilya Loshchilov and Frank Hutter. 2019. [Decoupled Weight Decay Regularization](#). In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.
- Joakim Nivre, Marie-Catherine de Marneffe, Filip Ginter, Jan Hajič, Christopher D. Manning, Sampo Pyysalo, Sebastian Schuster, Francis Tyers, and Daniel Zeman. 2020. [Universal Dependencies v2: An Evergrowing Multilingual Treebank Collection](#). In *Proceedings of the 12th Language Resources and Evaluation Conference*, pages 4034–4043, Marseille, France. European Language Resources Association.
- Prokopis Prokopidis and Haris Papageorgiou. 2017. [Universal Dependencies for Greek](#). In *Proceedings of the NoDaLiDa 2017 Workshop on Universal Dependencies (UDW 2017)*, pages 102–106, Gothenburg, Sweden. Association for Computational Linguistics.

- Vinit Ravishankar, Memduh Gökırmak, Lilja Øvrelid, and Erik Velldal. 2019. [Multilingual Probing of Deep Pre-Trained Contextual Encoders](#). In *Proceedings of the First NLPL Workshop on Deep Learning for Natural Language Processing*, pages 37–47, Turku, Finland. Linköping University Electronic Press.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. [Neural Machine Translation of Rare Words with Subword Units](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.
- Milan Straka, Jana Straková, and Jan Hajič. 2019. [Evaluating Contextualized Embeddings on 54 Languages in POS Tagging, Lemmatization and Dependency Parsing](#). Technical Report arXiv:1908.07448, arXiv. ArXiv:1908.07448 [cs] version: 1 type: article.
- Ian Tenney, Dipanjan Das, and Ellie Pavlick. 2019. [BERT Rediscovered the Classical NLP Pipeline](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4593–4601, Florence, Italy. Association for Computational Linguistics.
- Ian Tenney, Patrick Xia, Berlin Chen, Alex Wang, Adam Poliak, R. Thomas McCoy, Najoung Kim, Benjamin Van Durme, Samuel R. Bowman, Dipanjan Das, and Ellie Pavlick. 2018. What do you learn from context? Probing for sentence structure in contextualized word representations. In *International Conference on Learning Representations*.
- Shijie Wu and Mark Dredze. 2020. [Are All Languages Created Equal in Multilingual BERT?](#) In *Proceedings of the 5th Workshop on Representation Learning for NLP*, pages 120–130, Online. Association for Computational Linguistics.
- Amir Zeldes. 2017. [The GUM corpus: creating multilayer resources in the classroom](#). *Language Resources and Evaluation*, 51(3):581–612.
- Gözde Gül Şahin, Clara Vania, Ilia Kuznetsov, and Iryna Gurevych. 2020. [LINSPECTOR: Multilingual Probing Tasks for Word Representations](#). *Computational Linguistics*, 46(2):335–385.

A Training details for the probes

For both tasks, we opt for a two-layer MLP probe with a tanh activation. During training, we keep the encoder’s parameters frozen and train the linear probes using the AdamW optimizer (Loshchilov and Hutter, 2019) with a learning rate of 10^{-3} and weight decay of 10^{-2} , and employ early stopping according to the validation set accuracy.

B Dataset pre-processing

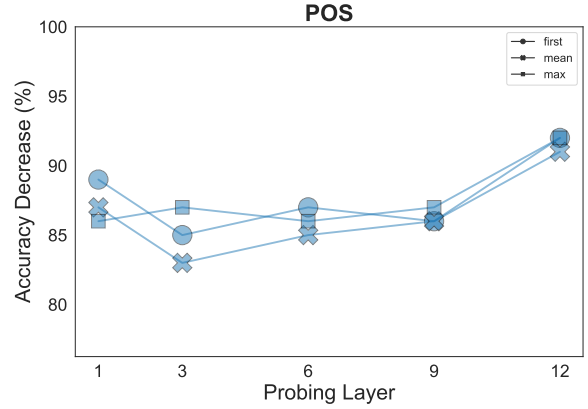
The sentences in the corpora from the Universal Dependencies framework are already tokenized to the word level and stored as lists of words in the `tokens` field. However, since we use sub-word tokenizers, namely the Byte Pair Encoding (Senrich et al., 2016) and SentencePiece (Kudo and Richardson, 2018) tokenizer, we further split the words into their sub-word tokens. Depending on the task, we also include either the `upos` field, which is a list of integers corresponding to one of the 17 universal POS tags available, or the `head` and `deprel` fields which contain the head and one of the 36 dependence relations for dependency labeling³. It should be noted that we only keep the language-independent relations, as some of them appear only with a language-specific modifier, and including them would make comparisons across languages less straightforward.

Furthermore, upon inspecting the datasets we observed that the annotators had split contractions into their parts and included them next to the original contraction for the Italian and Greek corpora. However, ground-truth labels were only provided for the sub-words, with the compound words annotated as a special class “_”. Hence, we filtered out the compound words from these datasets and retained their sub-parts. In addition to that, for dependency labeling, we ignored words with the root dependency label, as they have no head and their prediction is trivial.

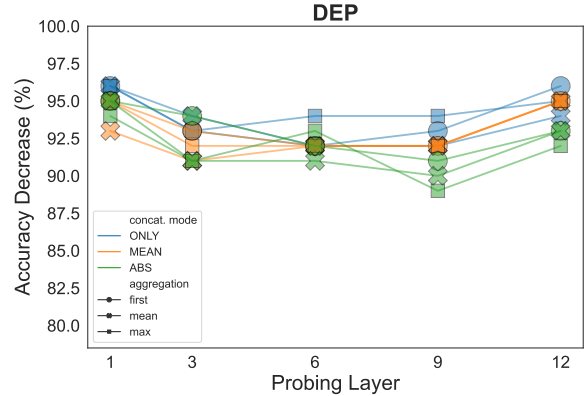
C Choosing a configuration for probing

Figures 4 and 5 showcase the decrease in accuracy when self-neutralizing in the POS tagging/dependency labeling task for RoBERTa and XLM-R accordingly. We report the best configurations for each model, language, and task combination in Table 2, and correspond to the ones we used for our cross-neutralization experiments.

³A full list of [POS tags](#) and [dependency relations](#) can be found on the Universal Dependencies website.



(a) Accuracy decrease for **POS tagging** when using different WordPiece aggregations.



(b) Accuracy decrease for **dependency labeling** when using different WordPiece aggregations and child-head concatenation configurations.

Figure 4: Decrease in performance for RoBERTa when **self-neutralizing** in the POS tagging (top) and dependency labeling (bottom) tasks using embeddings extracted from different layers and setup configurations.

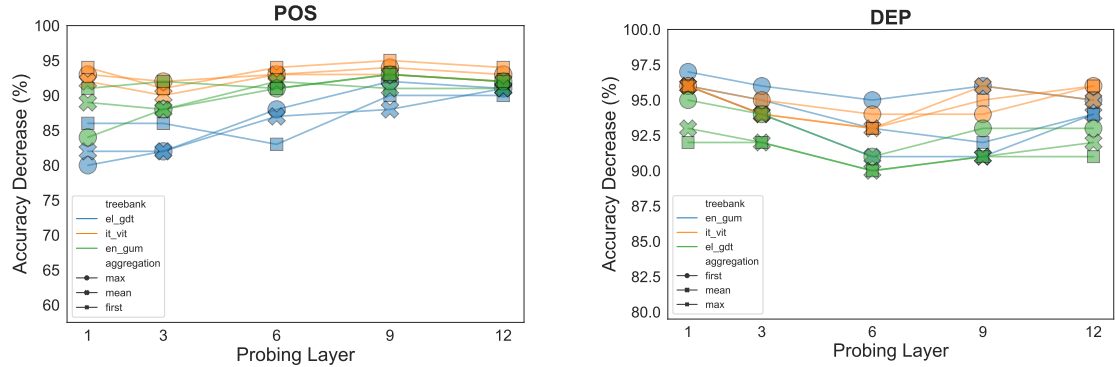
Figure 6 shows the results of cross-neutralizing in the dependency labeling task when the centroids are based exclusively on the child word’s representation. As mentioned in Section 5.1, this results in less prominent self-neutralizing. We hypothesize that the inclusion of information about the head within the centroid is likely leveraged by the classifier to retain good accuracy.

D Detailed cross-neutralizing results

We display the complete results for our cross-neutralization experiments in Figures 7 – 26. More specifically, Figures 7 and 8 correspond to the monolingual setting with RoBERTa on English, Figures 9 – 14 to the multilingual with XLM-R on English, Italian and Greek. Figures 15 – 26 show the cross-lingual setting for every possible pairing of our three languages with XLM-R.

POS			DEP		
Encoder / treebank	Layer	WP Aggregation	Layer	WP Aggregation	Concatenation
RoBERTa / en_gum	3	max	3	mean	[child ; head]
XLM-R / en_gum	9	max	9	first	[child ; head]
XLM-R / it_vit	9	first	9	mean	[child ; head]
XLM-R / el_gdt	12	mean	9	mean	[child ; head]

Table 2: The best probing configuration for each encoder model (RoBERTa & XLM-R), task (POS & DEP) and language (English, Italian & Greek) combination, chosen as outlined in Section 3.4. Ultimately, these are the configurations that we use for all of our cross-neutralization experiments.



(a) POS tagging Accuracy Decrease using different WordPiece (b) Dependency labeling Accuracy Decrease using different WordPiece aggregations.

Figure 5: Decrease in performance for XLM-R when **self-neutralizing** in the POS tagging (left) and dependency labeling (right) tasks using embeddings extracted from different layers and setup configurations, for each of for English (en_gum), Italian (it_vit) and Greek (el_gdt) treebanks. For dependency labeling, we use the best child-head concatenation mode (ONLY) based on the results we acquired with RoBERTa, as shown in Figure 4b.

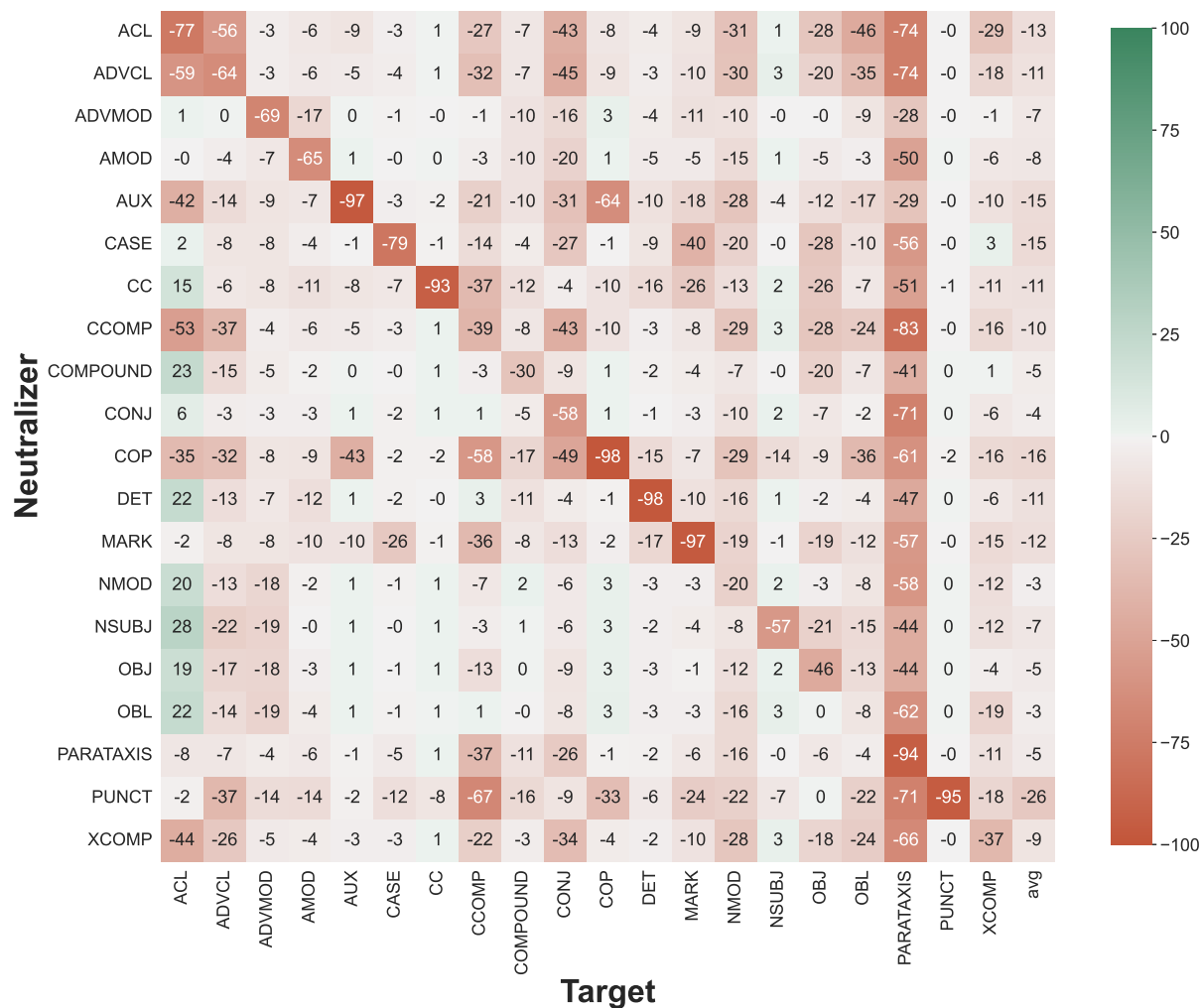


Figure 6: Relative change in accuracy when cross-neutralizing DEP tags using embeddings from RoBERTa on en_gum, in a setup configuration where we **only use the child embeddings** for the centroid calculation.

Neutralizer	ADJ	-91	-2	-31	-18	-6	-4	-13	-27	-10	-8	-1	-27	-16	-12
	ADP	-11	-94	-20	-11	-3	-3	-1	-22	-24	-9	-0	-37	-3	-10
	ADV	-22	-1	-97	-7	-3	-2	-4	-27	-4	-2	-0	-5	-2	-6
	AUX	-12	-3	-9	-94	-0	-3	-5	-22	-7	-3	-0	-14	-59	-11
	CCONJ	-1	-2	-21	-3	-97	-2	-4	-14	-9	-2	-0	-11	-0	-5
	DET	-8	-4	-15	-6	-2	-93	-1	-19	-17	-18	-0	-29	-2	-9
	NOUN	-29	-3	-25	-20	-6	-12	-92	-39	-31	-13	-1	-39	-15	-24
	NUM	-16	-12	-23	-4	-3	-4	-13	-94	-4	-8	-0	-47	-8	-10
	PART	-3	-2	-14	-3	-2	-2	-7	-5	-98	-1	-0	-8	-3	-4
	PRON	-3	-1	-6	-7	-0	-3	-7	-13	-7	-91	-1	-22	-3	-7
	PUNCT	-4	-4	-21	-12	-9	-10	-7	-18	-19	-5	-95	-35	-6	-17
	VERB	-5	-11	-40	-5	-3	-2	-5	-19	-11	-1	-0	-91	-2	-5
	avg	-25	-2	-17	-13	-4	-7	-10	-33	-12	-7	-0	-27	-94	-13
		ADJ	ADP	ADV	AUX	CCONJ	DET	NOUN	NUM	PART	PRON	PUNCT	SCONJ	VERB	avg
		Target													

Figure 7: Relative change in accuracy when cross-neutralizing POS tags using RoBERTa embeddings on en_gum.

Neutralizer	ACL	-72	-87	-12	-10	-6	-2	-11	-34	-4	-36	-13		-1	-9	-27	-10	-12	-18	-7	0	-29	-8
	ADVCL	2	-93	-6	-5	-2	-4	-4	-1	-5	-35	-24		-3	-5	-7	-1	-3	-17	-16	-0	-10	-5
	ADVMOD	-16	-40	-97	-9	-5	-1	-2	-14	-1	-19	-24	-47	-1	-7	-7	0	-3	-33	-42	0	-10	-9
	AMOD	15	-22	-25	-96	-9	-0	-5	-60	-43	-16	-35	-24	-2	-6	-34	-7	-7	-16	-53	0	-28	-14
	AUX	-5	-52	-10	-7	-95	-4	-2	-13	-19	-43	-24		-0	-10	-17	-2	-1	-32	-53	0	-24	-10
	CASE	8	-26	-28	-2	-15	-96	-6	-49	-22	-10	-15	-72	-2	-57	-34	-5	-10	-17	-45	0	-38	-17
	CC	-3	-8	-18	-3	-3	-1	-99	-5	1	-5	-1	-68	-2	-5	-12	-3	-2	-4	-65	0	-0	-7
	CCOMP	-2	-74	-10	-8	-5	-4	-9	-96	-9	-51	-31		-2	-9	-13	2	-8	-13	-67	-0	-23	-7
	COMPOUND	22	-16	-11	-13	-5	0	-5	-52	-94	-11	-26	-54	-2	-4	-22	-8	-3	-8	-43	0	-19	-9
	CONJ	32	-32	-10	-3	-6	-3	-10	-30	-1	-93	-22	-64	-1	-8	-17	-1	-1	-11	-60	-0	-19	-7
	COP	7	-34	-4	-11	-11	-1	-0	-31	-21	-39	-95		-0	-4	-6	-7	-6	-10	-18	0	-32	-6
	CSUBJ	0	-76	-10	-5	-4	-1	-4	-12	-15	-14	-6		-1	-5	-5	-25	-11	-10	-8	0	-13	-5
	DET	16	-29	-10	-8	-3	-3	-3	-68	-3	-13	2	-56	-98	-11	-52	-3	-6	-21	-42	0	-48	-15
	MARK	-29	-46	-33	-4	-15	-13	-4	-15	-7	-10	-10		-1	-95	-12	3	-5	-46	-55	0	-3	-10
	NMOD	18	-38	-12	-2	-5	-1	-8	-55	-9	-15	-17	-62	-1	-17	-96	-6	-7	-40	-65	0	-52	-11
	NSUBJ	-32	-45	-12	-6	-7	-7	-5	-15	-22	-19	-2		-1	-17	-14	-97	-12	-21	-49	-0	-19	-13
	OBJ	-27	-68	-17	-9	-10	-13	-5	-36	-36	-35	-33		-2	-17	-27	-25	-97	-61	-57	-0	-37	-17
	OBL	-23	-91	-19	-7	-12	-15	-3	-16	-38	-35	-41	-30	-3	-16	-69	-10	-16	-97	-38	-1	-11	-17
	PARATAXIS	1	-44	-12	-8	-11	-11	-3	-14	-8	-62	-30		-2	-11	-14	1	-2	-15	-100	-0	-9	-7
	PUNCT	4	-61	-22	-4	-10	-4	-5	-27	-6	-35	-9	-80	-1	-7	-20	4	-4	-11	-49	-99	-14	-24
	XCOMP	4	-62	-10	-6	-4	-5	-4	-6	-12	-41	-42		-2	-4	-10	1	-4	-20	-43	-0	-91	-6
		ACL	ADVCL	ADVMOD	AMOD	AUX	CASE	CC	CCOMP	COMPOUND	CONJ	COP	CSUBJ	DET	MARK	NMOD	NSUBJ	OBJ	OBL	PARATAXIS	PUNCT	XCOMP	avg
		Target																					

Figure 8: Relative change in accuracy when cross-neutralizing DEP tags using RoBERTa embeddings on en_gum. Note the relative increase in performance for the PARATAXIS relation; this is likely due to noise, as PARATAXIS labels make up less than 1% of the dataset.

Neutralizer	ADJ	ADP	ADV	AUX	CCONJ	DET	NOUN	NUM	PART	PRON	PUNCT	SCONJ	VERB	avg
	-91	-2	-31	-18	-6	-4	-13	-27	-10	-8	-1	-27	-16	-12
	-11	-94	-20	-11	-3	-3	-1	-22	-24	-9	-0	-37	-3	-10
	-22	-1	-97	-7	-3	-2	-4	-27	-4	-2	-0	-5	-2	-6
	-12	-3	-9	-94	-0	-3	-5	-22	-7	-3	-0	-14	-59	-11
	-1	-2	-21	-3	-97	-2	-4	-14	-9	-2	-0	-11	-0	-5
	-8	-4	-15	-6	-2	-93	-1	-19	-17	-18	-0	-29	-2	-9
	-29	-3	-25	-20	-6	-12	-92	-39	-31	-13	-1	-39	-15	-24
	-16	-12	-23	-4	-3	-4	-13	-94	-4	-8	-0	-47	-8	-10
	-3	-2	-14	-3	-2	-2	-7	-5	-98	-1	-0	-8	-3	-4
Target	-3	-1	-6	-7	-0	-3	-7	-13	-7	-91	-1	-22	-3	-7
	-4	-4	-21	-12	-9	-10	-7	-18	-19	-5	-95	-35	-6	-17
	-5	-11	-40	-5	-3	-2	-5	-19	-11	-1	-0	-91	-2	-5
	-25	-2	-17	-13	-4	-7	-10	-33	-12	-7	-0	-27	-94	-13

Figure 9: Relative change in accuracy when cross-neutralizing POS tags using XLM-R embeddings on en_gum.

Neutralizer	ADJ	-96	-1	-44	-18	-13	-0	-21	-26	-17	-27	0	-45	-22	-14
	ADP	-11	-97	-16	-17	-9	-2	-2	-37	-10	-31	0	-23	-6	-15
	ADV	-27	-4	-94	-8	-29	-3	-8	-27	-8	-35	0	-15	-7	-11
	AUX	-16	-1	-5	-96	-5	-2	-6	-19	-7	-30	0	-7	-60	-12
	CCONJ	-4	-1	-28	-8	-97	-1	-2	-26	-4	-6	0	-14	-6	-5
	DET	-7	-1	-8	-23	-5	-94	-6	-35	-32	-30	0	-48	-4	-16
	NOUN	-25	-0	-22	-20	-10	-3	-98	-39	-16	-26	0	-33	-11	-21
	NUM	-47	-1	-27	-25	-18	-1	-59	-95	-18	-7	0	-52	-12	-17
	PRON	-5	-1	-7	-14	-4	-3	-13	-17	-93	-23	0	-24	-9	-9
	PROPN	-17	-1	-18	-24	-6	-1	-43	-18	-17	-96	0	-42	-12	-15
VERB	PUNCT	-25	-4	-66	-21	-12	-3	-43	-13	-25	-30	-97	-24	-11	-29
	SCONJ	-8	-1	-24	-8	-19	-3	-6	-38	-23	-18	0	-98	-20	-7
	VERB	-63	-1	-20	-10	-7	-2	-33	-43	-9	-35	0	-19	-96	-19
	avg														
		Target													

Figure 10: Relative change in accuracy when cross-neutralizing POS tags using XLM-R embeddings on it_vit.

Neutralizer	ADJ	-93	-3	-40	-8	-7	0	-18	-25	-16	-9	-22	-4	-27	-6	-36	-16
	ADP	-4	-95	-25	-7	-4	0	-2	-16	-13	-6	-22	0	-19	-1	-55	-9
	ADV	-11	-1	-94	-3	-7	0	-8	-13	-5	-5	-27	-0	-7	-1	-29	-6
	AUX	-2	-4	-16	-97	-6	0	-3	-28	-14	-5	-14	-0	-18	-6	-23	-7
	CCONJ	-4	-1	-21	-4	-97	-0	-11	-32	-2	-2	-7	-2	-6	-1	-26	-8
	DET	-28	-6	-18	-16	-8	-93	-1	-48	-26	-27	-31	-0	-24	-1	-49	-20
	NOUN	-34	-6	-20	-13	-10	-1	-89	-28	-19	-16	-34	-4	-16	-5	-45	-24
	NUM	-23	-0	-10	-2	-1	-0	-17	-91	-7	-1	-27	-1	-4	-4	-48	-9
	PART	-2	-1	-31	-0	-2	0	-4	-12	-94	-4	-15	0	0	-1	-21	-3
	PRON	-7	-2	-15	-2	-3	0	-4	-13	-5	-88	-26	0	-19	-0	-27	-6
	PROPN	-21	-7	-18	-8	-6	-1	-24	-18	-4	-11	-88	-4	-8	-4	-23	-13
	PUNCT	-10	-5	-6	-21	-0	-0	-28	-31	0	-16	-18	-81	-12	-4	-48	-17
	SCONJ	-6	-1	-73	-1	-3	-0	-14	-18	-2	-10	-25	-1	-82	-1	-30	-8
	VERB	-18	-4	-12	-6	-7	-1	-18	-29	-11	-9	-30	-3	-8	-90	-38	-14
	X	-8	-2	-9	-9	-3	-1	-12	-18	-1	-9	-29	-3	-3	-3	-96	-8
		ADJ	ADP	ADV	AUX	CCONJ	DET	NOUN	NUM	PART	PRON	PROPN	PUNCT	SCONJ	VERB	X	avg
		Target															

Figure 11: Relative change in accuracy when cross-neutralizing POS tags using XLM-R embeddings on el_gdt.

Neutralizer	ACL	-72	-87	-12	-10	-6	-2	-11	-34	-4	-36	-13		-1	-9	-27	-10	-12	-18	-7	0	-29	-8
	ADVCL	2	-93	-6	-5	-2	-4	-4	-1	-5	-35	-24		-3	-5	-7	-1	-3	-17	-16	-0	-10	-5
	ADVMOD	-16	-40	-97	-9	-5	-1	-2	-14	-1	-19	-24	-47	-1	-7	-7	0	-3	-33	-42	0	-10	-9
	AMOD	15	-22	-25	-96	-9	-0	-5	-60	-43	-16	-35	-24	-2	-6	-34	-7	-7	-16	-53	0	-28	-14
	AUX	-5	-52	-10	-7	-95	-4	-2	-13	-19	-43	-24		-0	-10	-17	-2	-1	-32	-53	0	-24	-10
	CASE	8	-26	-28	-2	-15	-96	-6	-49	-22	-10	-15	-72	-2	-57	-34	-5	-10	-17	-45	0	-38	-17
	CC	-3	-8	-18	-3	-3	-1	-99	-5	1	-5	-1	-68	-2	-5	-12	-3	-2	-4	-65	0	-0	-7
	CCOMP	-2	-74	-10	-8	-5	-4	-9	-96	-9	-51	-31		-2	-9	-13	2	-8	-13	-67	-0	-23	-7
	COMPOUND	22	-16	-11	-13	-5	0	-5	-52	-94	-11	-26	-54	-2	-4	-22	-8	-3	-8	-43	0	-19	-9
	CONJ	32	-32	-10	-3	-6	-3	-10	-30	-1	-93	-22	-64	-1	-8	-17	-1	-1	-11	-60	-0	-19	-7
	COP	7	-34	-4	-11	-11	-1	-0	-31	-21	-39	-95		-0	-4	-6	-7	-6	-10	-18	0	-32	-6
	CSUBJ	0	-76	-10	-5	-4	-1	-4	-12	-15	-14	-6		-1	-5	-5	-25	-11	-10	-8	0	-13	-5
	DET	16	-29	-10	-8	-3	-3	-3	-68	-3	-13	2	-56	-98	-11	-52	-3	-6	-21	-42	0	-48	-15
	MARK	-29	-46	-33	-4	-15	-13	-4	-15	-7	-10	-10		-1	-95	-12	3	-5	-46	-55	0	-3	-10
	NMOD	18	-38	-12	-2	-5	-1	-8	-55	-9	-15	-17	-62	-1	-17	-96	-6	-7	-40	-65	0	-52	-11
	NSUBJ	-32	-45	-12	-6	-7	-7	-5	-15	-22	-19	-2		-1	-17	-14	-97	-12	-21	-49	-0	-19	-13
	OBJ	-27	-68	-17	-9	-10	-13	-5	-36	-36	-35	-33		-2	-17	-27	-25	-97	-61	-57	-0	-37	-17
	OBL	-23	-91	-19	-7	-12	-15	-3	-16	-38	-35	-41	-30	-3	-16	-69	-10	-16	-97	-38	-1	-11	-17
	PARATAXIS	1	-44	-12	-8	-11	-11	-3	-14	-8	-62	-30		-2	-11	-14	1	-2	-15	-100	-0	-9	-7
	PUNCT	4	-61	-22	-4	-10	-4	-5	-27	-6	-35	-9	-80	-1	-7	-20	4	-4	-11	-49	-99	-14	-24
	XCOMP	4	-62	-10	-6	-4	-5	-4	-6	-12	-41	-42		-2	-4	-10	1	-4	-20	-43	-0	-91	-6
		ACL	ADVCL	ADVMOD	AMOD	AUX	CASE	CC	CCOMP	COMPOUND	CONJ	COP	CSUBJ	DET	MARK	NMOD	NSUBJ	OBJ	OBL	PARATAXIS	PUNCT	XCOMP	avg
		Target																					

Figure 12: Relative change in accuracy when cross-neutralizing DEP tags using XLM-R embeddings on en_gum.

Neutralizer	ACL	-89	-95	-8	-48	-6	-5	-7	-53	-62	-4	-0	-4	-44	-50	-30	-57	-3	-26	-24	-12	-29	0	-42	-13
	ADVCL	-25	-100	-8	-70	-12	-1	-7	-10	-32	-29	0	-2	-51	-56	-18	-23	-2	-29	-14	-9	-11	0	-27	-9
	ADVMOD	-32	-47	-99	-33	-3	-1	-6	13	-17	-15	-1	-8	-13	-55	-4	-12	-0	-13	-6	-24	-24	0	-15	-11
	AMOD	-30	-63	-33	-99	-17	-5	-10	-15	-26	-18	-1	-13	-35	-37	-26	-54	-5	-18	-14	-18	-30	0	-37	-14
	AUX	-24	-59	-6	-6	-98	-4	-6	-46	-39	-17	-0	0	-18	-62	-11	-16	-2	-8	-1	-34	-41	-0	-21	-11
	CASE	-4	-36	-28	-15	-30	-98	-9	-16	-27	-15	-3	-16	-43	-40	-78	-42	-4	-33	-17	-12	-41	0	-58	-23
	CC	-7	-45	-28	-6	-4	-0	-99	3	-19	1	-1	-4	-7	-25	-10	-16	-3	-10	-10	-11	-45	-0	-41	-9
	CCOMP	-13	-87	-2	-22	-2	-2	-4	-73	-32	-15	0	0	-18	-44	-4	-7	1	-20	-5	-2	-21	0	-35	-6
	CONJ	-4	-25	-9	-16	-23	-1	-14	-7	-98	-31	-0	-15	-23	-25	-16	-24	-5	-30	-14	-25	-11	0	-36	-9
	COP	-10	-71	-6	-14	-14	-1	-5	-37	-36	-96	-0	-1	-29	-41	-5	-6	-8	-8	-12	-1	-27	0	-48	-7
	DET	-10	-43	-29	-23	-6	-4	-5	-16	-17	-3	-96	-11	-51	-47	-18	-77	-6	-29	-15	-22	-52	0	-61	-26
	EXPL	-40	-61	-26	-12	-1	-4	-2	-54	-32	-12	-8	-97	-19	-29	-20	-17	-36	-52	-82	-48	-29	-0	-56	-16
	FIXED	-24	-70	-14	-6	-30	-2	-2	11	-17	-3	-0	-1	-96	-6	-35	-75	-3	-7	-16	-62	-43	0	-22	-14
	FLAT	-23	-32	-8	-16	-26	-1	-5	-34	-26	-5	-0	-9	-24	-96	-11	-86	-3	-15	-10	-30	-30	0	-53	-12
	MARK	-44	-88	-20	-8	-6	-6	-9	-57	-18	-11	-0	-4	-8	-29	-98	-16	-8	-25	-10	-46	-24	0	-13	-11
	NMOD	-15	-28	-12	-5	-25	-1	-13	-19	-26	-15	-0	-22	-40	-32	-21	-99	-6	-16	-21	-59	-27	0	-60	-14
	NSUBJ	-46	-61	-8	-11	0	-4	-5	-57	-34	-16	-2	-5	-45	-41	-25	-35	-98	-28	-43	-64	-35	0	-49	-17
	NUMMOD	-15	-49	-11	-37	-20	-2	-10	-31	-23	-44	-1	-6	-36	-12	-21	-89	-8	-99	-21	-50	-40	0	-67	-15
	OBJ	-41	-42	-13	-9	-3	-5	-17	-44	-47	-41	-3	-8	-55	-36	-22	-46	-40	-31	-98	-67	-34	0	-68	-17
	OBL	-41	-68	-18	-17	-4	-6	-7	-31	-38	-46	-1	-15	-35	-48	-18	-89	-13	-18	-27	-98	-24	0	-36	-19
	PARATAXIS	-8	-75	-4	-22	-2	-2	-1	-8	-51	-2	-2	-2	-32	-35	-2	-9	-1	-33	-7	-3	-91	0	-21	-6
	PUNCT	-34	-52	-16	-2	-9	-3	-3	-24	-8	-40	-0	-0	-26	-42	-10	-16	-1	-2	-15	-14	-23	-99	-58	-23
	XCOMP	-11	-65	-8	-39	-13	-2	-17	-35	-21	-33	-0	-1	-19	-32	-7	-8	0	-26	-24	-9	-9	0	-98	-7
		ACL	ADVCL	ADVMOD	AMOD	AUX	CASE	CC	CCOMP	CONJ	COP	DET	EXPL	FIXED	FLAT	MARK	NMOD	NSUBJ	NUMMOD	OBJ	OBL	PARATAXIS	PUNCT	XCOMP	avg
		Target																							

Figure 13: Relative change in accuracy when cross-neutralizing DEP tags using embeddings from XLM-R on it_vit.

Neutralizer	ACL	-27	-85	-11	-7	-5	-2	-2	-29	-11	-8	-1	-0	-4	-23	-7	6	-19	-18	-1	-33	-7
	ADVCL	57	-99	-7	-13	0	-1	-1	-6	-31	-3	-0	-6	-1	-8	0	-2	-3	-22	1	-24	-4
	ADVMOD	-6	-46	-96	-2	0	0	-4	-30	5	2	-0	-3	-4	-5	-0	4	-3	-22	1	-23	-6
	AMOD	47	-32	-26	-94	-5	0	-4	-62	-0	-23	-1	-36	-13	-36	-2	-14	-6	-9	-1	-28	-13
	AUX	-8	-30	-30	-4	-98	-0	-4	0	-3	1	-1	-5	-9	-10	-11	-4	-15	-12	2	-12	-8
	CASE	47	-43	-66	-13	-5	-95	-2	-70	-2	-4	-5	-17	-34	-43	-9	-22	-16	-19	2	-31	-20
	CC	5	-11	-21	-2	-0	-1	-100	-28	-2	-0	-0	-15	-5	-12	-3	-9	-4	-9	2	-11	-6
	CCOMP	35	-54	-3	-8	-0	-7	-1	-95	-32	-4	-1	-10	0	-9	0	-2	-9	-11	2	-41	-4
	CONJ	48	-31	-5	-2	-8	-1	-0	-46	-91	-27	-0	-14	-1	-14	-2	-11	-3	-5	-1	-19	-4
	COP	40	-73	-9	-15	-1	0	0	-34	-13	-98	-3	-9	-1	-11	-35	-2	-17	-14	2	-31	-8
	DET	25	-47	-14	-26	-6	-2	-1	-54	-8	-1	-96	-20	-1	-49	-6	-46	-14	-14	0	-45	-29
	FLAT	47	-42	-5	-2	-5	-0	-0	-38	-10	-21	-0	-77	-0	-66	-2	-7	-8	-14	-0	-10	-8
	MARK	-5	-29	-56	-12	-2	-8	-6	-15	-6	-2	-2	-6	-98	-8	-6	-11	-2	-17	2	-31	-8
	NMOD	51	-39	-9	-1	-7	-1	-1	-59	-4	-19	0	-27	-3	-93	-2	-13	-8	-21	-1	-40	-10
	NSUBJ	-6	-18	-8	-5	0	-3	-0	-20	-6	2	-1	-0	-7	-5	-95	-17	-14	-15	-1	-15	-8
	NUMMOD	6	-24	-11	-9	-1	-0	-0	-34	-4	1	-0	-24	-6	-19	0	-86	-1	-9	-0	-17	-7
	OBJ	-3	-12	-10	-8	1	-4	-2	-16	-7	-0	-1	-1	-2	-10	-9	-12	-94	-33	-1	-51	-8
	OBL	-7	-40	-12	-8	0	-4	-1	-26	-6	1	-0	-22	-3	-35	1	-19	-7	-95	-1	-25	-10
	PUNCT	-7	-19	-17	-5	-0	-3	-2	-18	-14	-6	-0	-13	-5	-14	-0	-6	-3	-40	-100	-5	-18
	XCOMP	44	-29	-6	-10	0	-3	-1	-11	-18	-7	-1	-4	-0	-8	1	-1	-2	-17	2	-97	-3
		ACL	ADVCL	ADVMOD	AMOD	AUX	CASE	CC	CCOMP	CONJ	COP	DET	FLAT	MARK	NMOD	NSUBJ	NUMMOD	OBJ	OBL	PUNCT	XCOMP	avg
		Target																				

Figure 14: Relative change in accuracy when cross-neutralizing DEP tags using XLM-R embeddings on el_gdt.

English Neutralizer	ADJ	-59	-1	-19	-1	-3	-0	-7	-15	-8	-12	-28	-0	-10	-0	-20	-7
	ADP	-3	-23	-9	-2	-3	0	-1	-2	-2	-2	-10	0	-13	-0	-24	-2
	ADV	-5	-1	-75	-1	-3	0	-2	-6	-4	-5	-6	0	-10	-1	-22	-3
	AUX	-3	-1	-12	-56	-2	-0	-1	-5	0	-3	-0	-0	-11	-18	-8	-5
	CCONJ	-0	-1	-10	-3	-61	0	-3	-4	0	-1	-8	-0	-4	-1	-15	-2
	DET	-14	-1	-2	-7	-1	-9	0	-22	0	-8	-8	-0	-9	-2	-4	-4
	NOUN	-13	-0	-16	-6	-3	-0	-69	-20	0	-14	-33	-2	-2	-3	-25	-16
	NUM	-5	-1	-7	0	-1	-0	-19	-69	-4	-3	-12	-0	-3	-4	0	-6
	PART	-1	-1	-7	-7	-1	0	0	-3	0	-1	3	-0	0	-2	-17	-1
	PRON	-1	0	-8	-1	0	0	-0	-0	0	-37	-2	0	-7	1	3	-1
	PROPN	-5	-0	-9	0	-2	-0	-11	-8	1	-6	-22	0	-1	-2	-37	-4
	PUNCT	-1	-0	-11	-3	-4	0	-9	-17	-14	-6	-10	-46	-10	-1	-57	-9
	SCONJ	-4	-3	-63	-1	-2	0	-4	-9	1	-4	-4	-0	-9	-0	-36	-3
	VERB	-8	-0	-8	-6	-1	0	-6	-8	1	-7	-10	0	-2	-34	-8	-5
	X	-10	0	-6	0	0	0	-51	-23	1	-3	-23	0	1	-0	-24	-12
		ADJ	ADP	ADV	AUX	CCONJ	DET	NOUN	NUM	PART	PRON	PROPN	PUNCT	SCONJ	VERB	X	avg
		Greek Target															

Figure 15: Relative change in accuracy when cross-neutralizing Greek POS tags using English XLM-R embeddings.

English Neutralizer	ACL	13	-62	-7	-1	-19	-1	-0	-33	-4	-17	-0	-1	-24	-5	2	-17	-13	-1	-31	-5
	ADVCL	33	-70	-6	-9	-0	-2	-1	-14	-2	-4	-1	5	-4	-0	-3	-3	-18	-1	-15	-2
	ADVMOD	-8	-32	-90	-1	0	-0	-2	-21	6	2	-0	-3	-4	0	6	-2	-15	2	-21	-5
	AMOD	27	-23	-25	-77	-4	-0	-3	-51	3	-11	-1	-29	-26	-1	-6	-10	-11	-1	-32	-11
	AUX	21	-53	-10	-10	-10	0	-0	-10	-11	-5	-1	-2	-7	-13	-3	-5	-5	2	-21	-4
	CASE	33	-25	-22	-3	-5	-26	-0	-54	-4	-3	-1	-3	-25	-5	-3	-10	-10	2	-30	-7
	CC	-7	-10	-12	-3	-0	-1	-89	-8	8	-2	-0	-6	-5	-1	-2	0	-3	2	-5	-4
	CCOMP	49	-20	-3	-10	-1	-3	-1	-59	-5	-4	-1	-8	-6	0	-4	-4	-5	2	-10	-2
	CONJ	24	-26	-3	-4	-5	-1	-0	-20	-88	-10	-0	-3	-6	-2	-8	-2	-4	-1	-14	-3
	COP	34	-52	-5	-12	-1	0	0	-33	-6	-60	-1	0	-8	-36	-2	-19	-6	2	-44	-6
	DET	29	-37	-7	-9	-1	-0	0	-43	-5	-1	-39	-5	-19	-2	-10	-6	-8	0	-33	-13
	FLAT	9	-38	-11	-6	-1	-2	-0	-23	-1	-1	-0	-5	-13	-3	-10	-18	-64	1	-22	-6
	NMOD	-11	-15	-8	-8	1	-2	-0	-10	-7	1	-0	-10	-53	-1	-16	-44	-52	2	-27	-11
	NSUBJ	47	-23	-9	-20	-2	-7	-0	-7	-2	-7	0	-15	-31	-6	-16	-4	-9	2	-10	-6
	NUMMOD	-4	-16	-24	-5	-2	-4	-2	-6	1	2	-0	-4	-3	-1	-4	-1	-3	2	-6	-3
	OBJ	34	-24	-6	1	-5	-2	-0	-48	1	-5	0	-27	-69	-1	-7	-9	-21	0	-46	-8
	OBL	-5	-4	-3	-4	0	-1	-0	-9	4	2	-1	1	-5	-78	-4	-11	-8	0	-5	-6
	PUNCT	-5	-29	-11	-8	-0	-4	-0	-20	-3	2	-1	-18	-49	1	-19	-8	-86	-0	-25	-11
	XCOMP	-7	-16	-19	-5	-0	-1	-4	-17	-11	-8	-0	-12	-24	-4	-8	-4	-46	-100	-14	-20
		ACL	ADVCL	ADVMOD	AMOD	AUX	CASE	CC	CCOMP	CONJ	COP	DET	FLAT	NMOD	NSUBJ	NUMMOD	OBJ	OBL	PUNCT	XCOMP	avg
		Greek Target																			

Figure 16: Relative change in accuracy when cross-neutralizing Greek DEP tags using English XLM-R embeddings.

Italian Neutralizer	ADJ	-48	-0	-18	-1	-3	-0	-12	-14	-7	-36	-0	-13	0	-25	-7
	ADP	-6	-51	-12	-1	-1	0	-0	-8	-4	-12	0	-12	-0	-36	-5
	ADV	-11	-1	-85	0	-2	0	-2	-5	-7	-21	0	-7	0	-43	-5
	AUX	-5	-3	-6	-57	-1	-1	-1	-8	-4	-0	0	-9	-34	-17	-6
	CCONJ	-1	-0	-20	-3	-56	-0	-5	-12	-1	-11	-0	-9	-1	-11	-3
	DET	-20	-2	-5	-12	-3	-49	-3	-22	-14	-41	0	-9	-4	-17	-11
	NOUN	-12	0	-10	-1	-0	-0	-76	-17	-16	-43	-2	-2	-3	-32	-17
	NUM	-10	-1	-8	-3	-0	-0	-53	-84	-2	-29	-1	-6	-8	-11	-15
	PRON	-5	0	-5	-1	-1	0	-2	-3	-34	-17	0	-12	1	-19	-2
	PROPN	-4	0	-5	-1	-0	-0	-10	-6	-2	-16	0	-3	-3	-42	-3
	PUNCT	-6	-34	-74	-33	-13	-0	-14	-26	-7	-14	-100	-6	-10	-73	-23
	SCONJ	-3	-1	-68	-1	-3	-0	-4	-9	-5	-10	-0	-17	-0	-48	-4
	VERB	-17	-0	-8	-12	-2	0	-5	-12	-6	-26	0	-3	-42	-25	-6
	X	-8	-0	-5	0	-1	0	-16	-9	-7	-27	0	-2	0	-93	-5
		ADJ	ADP	ADV	AUX	CCONJ	DET	NOUN	NUM	PRON	PROPN	PUNCT	SCONJ	VERB	X	avg
		Greek Target														

Figure 17: Relative change in accuracy when cross-neutralizing Greek POS tags using Italian XLM-R embeddings.

Italian Neutralizer	ACL	-32	-80	-7	-4	-8	-2	-1	-26	-7	-11	-1	-7	-5	-39	-6	4	-23	-16	-1	-8
	ADVCL	43	-66	-7	-3	-10	-1	-1	-53	1	-26	-0	-13	-3	-14	-4	-0	-11	-8	-1	-4
	ADVMOD	-5	-39	-96	-0	-0	0	-3	-39	5	2	0	-2	-5	-3	-1	6	-4	-20	2	-6
	AMOD	55	-28	-22	-55	-8	0	-2	-65	-1	-24	-1	-25	-11	-55	-5	-14	-17	-17	-1	-12
	AUX	39	-60	-9	-13	-17	-0	-1	-10	-39	-9	-0	-9	-1	-9	-9	-3	-2	-6	2	-5
	CASE	40	-40	-37	-12	-0	-46	-1	-65	-3	1	-5	-13	-28	-51	-7	-12	-12	-17	2	-15
	CC	4	-19	-37	-1	-0	-0	-81	-34	1	-0	-0	-18	-7	-18	-2	-4	-4	-7	2	-6
	CCOMP	43	-37	-2	-7	-0	-2	-1	-8	-6	-8	-1	-1	-0	-3	-1	1	-3	-4	-1	-1
	CONJ	38	-27	-5	-3	-9	-0	-0	-47	-94	-20	-0	-17	-1	-21	-2	-14	-4	-7	-1	-5
	COP	46	-69	-7	-18	-2	-0	0	-29	-22	-64	-2	-11	0	-9	-22	-8	-17	-5	1	-7
	DET	31	-46	-7	-24	-9	-1	-1	-62	-4	-1	-68	-19	-1	-59	-4	-27	-10	-12	2	-24
	FLAT	8	-11	-7	-5	-0	-2	0	-35	5	2	-1	-6	-2	-39	-6	-9	-14	-40	3	-7
	MARK	-4	-31	-41	-5	1	0	-3	-8	-14	-0	-0	-4	-5	-11	-3	-5	-17	-12	2	-6
	NMOD	-5	-19	-41	-6	-7	-8	-2	-17	3	2	-0	1	-24	-5	-2	-4	0	-4	2	-4
	NSUBJ	42	-43	-6	-1	-12	0	0	-61	-7	-11	0	-30	-2	-90	-2	-10	-13	-32	-1	-11
	NUMMOD	-6	-9	-7	-5	0	-1	-0	-21	-5	2	-1	-1	-6	-4	-89	-6	-12	-13	-0	-8
	OBJ	3	-23	-5	-1	-4	-0	-1	-27	-5	-5	-0	-7	-2	-77	-2	-21	-6	-21	1	-10
	OBL	-5	-17	-9	-7	0	-3	-1	-12	-12	2	-2	-6	-1	-11	-8	-9	-86	-23	-1	-7
	PUNCT	55	-76	-5	-10	-1	-2	-0	-42	-35	-1	-0	-4	0	-5	-4	-2	-4	-6	-1	-3
		ACL	ADVCL	ADVMOD	AMOD	AUX	CASE	CC	CCOMP	CONJ	COP	DET	FLAT	MARK	NMOD	NSUBJ	NUMMOD	OBJ	OBL	PUNCT	avg
		Greek Target																			

Figure 18: Relative change in accuracy when cross-neutralizing Greek DEP tags using Italian XLM-R embeddings.

English Neutralizer	ADJ	-70	-1	-47	-19	-9	-1	-10	-29	-9	-12	0	-34	-14	-10
	ADP	-8	-63	-12	-13	-5	-2	-0	-44	-12	-22	0	-27	-5	-10
	ADV	-20	-1	-90	-10	-12	-1	-2	-34	-5	-19	0	-10	-7	-7
	AUX	-10	-0	-8	-52	-5	-1	-1	-18	-5	-8	0	-10	-32	-6
	CCONJ	-2	-1	-23	-6	-82	-1	-2	-27	-6	-2	0	-11	-7	-4
	DET	-4	-1	-5	-1	-3	-35	-1	-35	-17	-12	0	-37	-1	-6
	NOUN	-23	-0	-23	-29	-5	-2	-79	-48	-13	-19	0	-33	-13	-17
	NUM	-26	-0	-25	-23	-7	-1	-19	-78	-11	-2	0	-34	-10	-9
	PRON	-6	-1	-7	-1	-3	-8	-2	-24	-73	-4	0	-42	-3	-5
	PROPN	-14	-0	-16	-18	-5	-1	-25	-25	-14	-73	0	-30	-10	-10
	PUNCT	-5	-2	-22	-19	-4	-1	-9	-3	-15	-10	-14	-14	-11	-8
	SCONJ	-6	-8	-35	-5	-11	-3	-4	-45	-13	-10	0	-42	-11	-6
	VERB	-22	-0	-13	-13	-2	-2	-6	-41	-5	-15	0	-16	-67	-9
Italian Target		ADJ	ADP	ADV	AUX	CCONJ	DET	NOUN	NUM	PRON	PROPN	PUNCT	SCONJ	VERB	avg

Figure 19: Relative change in accuracy when cross-neutralizing Italian POS tags using English XLM-R embeddings.

English Neutralizer	ACL	-46	-83	-3	-17	-13	-3	-2	-19	-30	-4	-0	-4	-26	-26	-18	-41	-3	-29	-17	-9	-12	-8
	ADVCL	-7	-86	-6	-10	-3	-5	-2	-22	-30	-29	0	-2	-15	-28	-7	-8	-1	-25	-6	-17	-9	-5
	ADVMOD	-36	-33	-92	-16	-1	-1	-2	15	-16	-18	0	-14	-11	-31	-3	-8	1	-11	-2	-14	-19	-8
	AMOD	-9	-36	-31	-86	-11	-2	-9	-8	-17	-14	-2	-14	-35	-44	-16	-34	-3	-19	-9	-11	-31	-11
	AUX	-35	-38	-7	-6	-69	-2	-1	-25	-35	-14	-0	-1	-9	-14	-6	-11	-2	-6	-5	-9	-35	-7
	CASE	2	-20	-17	-5	-14	-74	-6	-1	-16	-1	-2	-19	-37	-12	-49	-23	-3	-36	-9	-7	-42	-15
	CC	-21	-39	-16	-2	1	-1	-93	12	-14	2	-0	-5	-9	-12	-13	-5	-2	-13	-3	-2	-39	-6
	CCOMP	-8	-70	-5	-11	-3	-2	-8	-69	-25	-34	0	-1	-19	-16	-8	-6	0	-30	-13	-1	-24	-5
	CONJ	2	-23	-8	-8	-16	-1	-6	3	-96	-34	-0	-7	-15	-13	-8	-8	-3	-30	-8	-9	-7	-6
	COP	-10	-53	-6	-14	-2	-1	-4	-26	-29	-85	-0	-0	-14	-12	-3	-3	-10	-6	-10	-3	-24	-6
	DET	1	-7	-7	-1	1	-2	-1	-6	-10	5	-53	-18	-23	-22	-4	-35	-2	-32	-4	-7	-44	-13
	EXPL	-36	-3	-23	-7	-0	-2	-1	-41	-6	3	-1	-7	-11	-14	-6	-9	-43	-3	-13	-11	-15	-7
	FIXED	-6	-62	-11	-3	-7	-1	-5	11	-9	-0	-0	-12	-69	-5	-10	-38	-7	-17	-11	-38	-46	-8
	FLAT	-20	-23	-6	-6	-15	-3	-2	-21	-22	-8	-0	-5	-31	-90	-18	-82	-8	-17	-12	-16	-21	-11
	MARK	-40	-19	-9	-7	-2	-6	-7	-25	-19	-38	-1	-2	-22	-9	-15	-40	-2	-19	-51	-35	-19	-10
	NMOD	-16	-34	-16	-17	-28	-2	-7	-34	-38	-13	-0	-9	-56	-33	-49	-99	-32	-50	-18	-32	-95	-17
	NSUBJ	-35	-53	-13	-7	-1	-4	-3	-24	-13	-2	0	-5	-6	-9	-84	-8	-5	-9	-5	-14	-29	-7
	NUMMOD	-5	-9	-4	3	-11	-1	-3	-14	-17	-3	-1	-5	-28	-16	-17	-84	-4	-13	-14	-31	-24	-10
	OBJ	-39	-15	-3	-6	1	-4	-2	3	-20	-5	-1	-2	-28	-18	-22	-15	-82	-18	-16	-18	-23	-9
	OBL	-8	-23	-8	-18	-1	-2	-5	-16	-11	-31	-7	-11	-29	-20	-22	-54	-4	-88	-13	-25	-22	-11
	PARATAXIS	-38	-50	-21	-15	-2	-8	-5	-20	-38	-45	-2	-14	-40	-20	-18	-75	-6	-17	-25	-89	-22	-17
		ACL	ADVCL	ADVMOD	AMOD	AUX	CASE	CC	CCOMP	CONJ	COP	DET	EXPL	FIXED	FLAT	MARK	NMOD	NSUBJ	NUMMOD	OBJ	OBL	PARATAXIS	avg
		Italian Target																					

Figure 20: Relative change in accuracy when cross-neutralizing Italian DEP tags using English XLM-R embeddings.

Greek Neutralizer	ADJ	ADP	ADV	AUX	CCONJ	DET	NOUN	NUM	PRON	PROPN	PUNCT	SCONJ	VERB	avg
	-73	-1	-24	-20	-6	-1	-21	-35	-18	-11	0	-48	-15	-11
	-9	-79	-16	-17	-17	-1	-1	-42	-3	-20	0	-30	-7	-12
	-15	-4	-85	-14	-14	-1	-5	-35	-4	-21	0	-9	-7	-8
	-5	-3	-5	-11	-1	-1	-2	-18	-3	-9	0	-6	-14	-3
	-4	-0	-18	-7	-40	-1	-4	-34	-2	-3	0	-19	-8	-4
	-6	-1	-7	-2	-7	-58	-1	-49	-23	-9	0	-47	-1	-9
	-26	-0	-23	-26	-9	-2	-89	-45	-16	-13	0	-45	-12	-19
	-38	-0	-23	-21	-7	-2	-19	-54	-20	-5	0	-48	-11	-9
	-6	-0	-7	-8	-3	-7	-3	-28	-74	-15	0	-43	-4	-6
Italian Target	-26	-0	-24	-33	-4	-3	-68	-35	-27	-55	0	-47	-15	-19
	-7	-2	-19	-17	-2	-1	-12	-3	-10	-11	-7	-11	-11	-7
	-8	-4	-36	-7	-30	-3	-11	-54	-18	-29	0	-81	-15	-9
	-31	-1	-15	-6	-3	-4	-16	-48	-20	-29	0	-17	-81	-14

Figure 21: Relative change in accuracy when cross-neutralizing Italian POS tags using Greek XLM-R embeddings.

Greek Neutralizer	ACL	-31	-66	-2	-11	0	-2	-4	-20	-33	1	0	-3	-30	-30	-13	-20	-3	-18	-10	-7	-19	-6
	ADVCL	-11	-96	-4	-16	0	-5	-1	-58	-48	-18	0	-1	-36	-54	-10	-14	1	-19	-5	-28	-22	-8
	ADVMOD	-34	-36	-93	-24	-7	0	-4	15	-11	-20	0	-10	-12	-48	-2	-13	1	-14	-4	-25	-18	-10
	AMOD	-11	-31	-15	-85	-12	-3	-5	-19	-16	-7	-2	-9	-33	-33	-21	-40	-4	-14	-5	-10	-18	-10
	AUX	-38	-38	-19	-11	-25	-3	-1	-13	-15	-18	0	1	-25	-29	-48	-12	-6	-4	-10	-24	-14	-8
	CASE	-0	-10	-24	-4	-18	-84	-7	-5	-20	-4	-4	-5	-55	-34	-60	-29	-5	-38	-11	-9	-39	-18
	CC	-0	-4	-18	0	-4	0	-70	13	-13	4	-1	-1	-5	-21	-7	-10	0	-18	-3	-2	-40	-5
	CCOMP	-17	-59	-5	-26	-1	-5	-5	-94	-39	-23	-1	-1	-32	-44	-13	-12	1	-31	-20	-12	-38	-8
	CONJ	-1	-9	-6	-8	-15	-1	-6	-2	-94	-27	0	-13	-25	-20	-16	-19	-5	-29	-9	-9	-7	-7
	COP	-19	-65	-18	-14	-19	-2	-7	-43	-32	-99	0	-2	-46	-50	-8	-9	-5	-8	-14	-4	-37	-9
	DET	-5	-14	-9	-7	1	-3	-4	-21	-9	5	-76	-12	-37	-14	-7	-42	-5	-47	-6	-9	-37	-17
	EXPL	-1	-92	-14	-1	-14	-1	-6	-32	-13	3	0	-4	-89	-7	-44	-57	-4	-18	-17	-59	-50	-13
	FIXED	-15	-21	-6	-4	-19	-1	-2	-26	-20	-19	0	-4	-25	-33	-21	-90	-5	-15	-7	-19	-19	-10
	FLAT	-33	4	-18	-13	-23	-1	-2	14	-49	-33	-38	-19	-12	-23	-7	-17	-40	-62	-20	-22	-18	-16
	MARK	-10	-17	-7	-2	-13	-1	-4	-24	-15	-11	0	-10	-35	-20	-22	-92	-6	-13	-10	-30	-22	-10
	NMOD	-44	-35	-8	-8	0	-5	-2	-25	-30	-15	-2	-6	-49	-28	-26	-32	-96	-19	-34	-46	-31	-14
	NSUBJ	-3	-18	-14	-32	-2	-2	-7	-17	-9	-24	-7	-13	-52	-34	-20	-42	-4	-93	-10	-19	-14	-10
	NUMMOD	-42	-36	-10	-7	-5	-5	-8	-44	-38	-38	-2	-11	-50	-34	-23	-44	-33	-23	-96	-56	-32	-15
	OBJ	-41	-41	-18	-20	-5	-7	-4	-20	-29	-42	-1	-12	-37	-31	-23	-84	-8	-12	-25	-93	-17	-17
	OBL	-6	-15	-15	-2	-23	-9	-6	4	-17	-2	-1	-3	-37	-7	-29	-100	-12	-22	-34	-99	-10	-17
	PARATAXIS	-29	-26	-8	1	-16	-2	-2	-11	-5	-31	0	-2	-23	-34	-10	-13	0	-7	-16	-5	-20	-21
		ACL	ADVCL	ADVMOD	AMOD	AUX	CASE	CC	CCOMP	CONJ	COP	DET	EXPL	FIXED	FLAT	MARK	NMOD	NSUBJ	NUMMOD	OBJ	OBL	PARATAXIS	avg
		Italian Target																					

Figure 22: Relative change in accuracy when cross-neutralizing Italian DEP tags using Greek XLM-R embeddings.

Italian Neutralizer	ADJ	-63	-3	-25	-13	-6	-3	-11	-19	-4	-1	-33	-13	-10
	ADP	-4	-61	-14	-3	-2	-2	-0	-14	-3	-0	-47	-2	-6
	ADV	-15	-1	-88	-5	-2	-1	-6	-18	-1	-0	-7	-1	-6
	AUX	-9	-6	-7	-50	-1	-0	-5	-9	-1	-0	-8	-42	-7
	CCONJ	-2	-2	-25	-2	-78	-5	-4	-5	-2	-0	-16	-0	-4
	DET	-2	-2	-4	-2	0	-9	-1	-6	-7	-0	-18	-2	-2
	NOUN	-21	-4	-25	-14	-9	-5	-79	-30	-13	-1	-46	-14	-20
	NUM	-12	-7	-23	-5	-3	-5	-30	-61	-9	-0	-52	-10	-13
	PRON	-2	-1	-6	-2	-2	-1	-7	-2	-22	-0	-17	-5	-3
	PUNCT	-13	-3	-35	-6	-12	-6	-23	-14	-3	-100	-61	-8	-22
VERB	SCONJ	-2	-2	-32	-3	-3	-0	-5	-7	-1	-0	-61	-1	-3
	VERB	-20	-3	-15	-19	-7	-3	-6	-11	-3	-1	-28	-81	-11
		ADJ	ADP	ADV	AUX	CCONJ	DET	NOUN	NUM	PRON	PUNCT	SCONJ	VERB	avg
		English Target												

Figure 23: Relative change in accuracy when cross-neutralizing English POS tags using Italian XLM-R embeddings.

Italian Neutralizer	ACL	-83	-67	-17	-12	-2	-8	-7	-40	-5	-38	-5	-0	-21	-39	-9	-17	-11
	ADVCL	-47	-59	-9	-21	-2	-2	-4	-37	-3	-11	-13	-2	-7	-13	-6	-12	-6
	ADVMOD	-1	-39	-98	-16	-8	-1	-4	-33	-1	-6	-19	0	-3	-11	-0	-39	-9
	AMOD	4	-28	-32	-78	-12	-5	-4	-59	-20	-4	-13	-1	-19	-51	-5	-20	-14
	AUX	5	-52	-14	-11	-85	-4	-2	-32	-25	-47	-23	-1	-6	-11	-18	-13	-10
	CASE	1	-28	-35	-8	-1	-85	-9	-51	-53	-3	-15	-2	-56	-45	-7	-23	-22
	CC	2	-26	-40	-0	0	-2	-92	-31	-10	-1	2	-2	-6	-23	-3	-16	-9
	CCOMP	-0	-44	-5	-5	-4	-3	-4	-14	-27	-35	-19	-1	-11	-7	-4	-3	-6
	COMPOUND	10	-32	-23	-7	-12	-7	-1	-42	-82	-7	-17	-1	-18	-49	-6	-18	-13
	CONJ	16	-36	-11	-1	-7	-4	-7	-31	-19	-96	-13	-2	-19	-27	-2	-19	-9
	COP	-6	-46	-12	-14	-17	-3	-1	-27	-27	-44	-72	-2	-4	-7	-16	-11	-8
	DET	0	-38	-15	-20	0	-11	-3	-65	-32	-5	0	-36	-18	-69	-5	-32	-16
	MARK	-28	-48	-14	-7	-14	-13	-4	-40	-36	-10	-22	-3	-12	-17	-87	-36	-14
	NMOD	-9	-44	-33	-15	-34	-17	-3	-21	-43	-7	-17	-2	-42	-65	-6	-36	-16
	NSUBJ	-14	-49	-20	-12	-12	-17	-7	-35	-63	-24	-28	-4	-13	-28	-19	-47	-17
	OBL	13	-71	-12	-5	-1	-3	-6	-18	-25	-51	-2	-1	-9	-11	-12	-3	-7
		ACL	ADVCL	ADVMOD	AMOD	AUX	CASE	CC	CCOMP	COMPOUND	CONJ	COP	DET	MARK	NMOD	NSUBJ	OBL	avg
		English Target																

Figure 24: Relative change in accuracy when cross-neutralizing English DEP tags using Italian XLM-R embeddings.

Greek Neutralizer	ADJ	ADP	ADV	AUX	CCONJ	DET	NOUN	NUM	PART	PRON	PUNCT	SCONJ	VERB	avg
	-77	-2	-24	-27	-6	-6	-9	-40	-5	-6	-0	-18	-29	-12
	-3	-60	-16	-6	-3	-2	-1	-8	-12	-3	-0	-40	-1	-6
	-18	-2	-78	-13	-5	-3	-5	-26	-5	-2	-0	-8	-5	-6
	-5	-0	-16	-2	-1	0	-3	-3	0	-1	-0	-4	-16	-3
	-0	-3	-16	-3	-50	-2	-4	-2	-3	-1	-0	-19	-2	-3
	-3	-4	-7	-3	-0	-13	-1	-3	-0	-5	-0	-13	-3	-2
	-29	-2	-28	-25	-10	-15	-73	-35	-32	-21	-0	-33	-21	-23
	-30	-4	-23	-11	-1	-3	-10	-39	-2	-6	-1	-42	-8	-8
	-3	-1	-35	-1	-0	0	-9	1	-36	-5	-0	-9	-1	-5
English Target	-0	-0	-5	-6	-1	-1	-6	-2	-1	-26	-0	-15	-3	-4
	-4	-2	-22	-7	-5	-12	-7	-29	-20	-6	-91	-27	-6	-17
	-2	-2	-32	-8	-0	-1	-6	-4	-2	-0	-0	-63	-4	-4
	-29	-2	-15	-17	-6	-14	-10	-35	-7	-7	-0	-21	-94	-14

Figure 25: Relative change in accuracy when cross-neutralizing English POS tags using Greek XLM-R embeddings.

Greek Neutralizer	ACL	-30	-49	-13	-4	-1	-6	-6	-25	-19	-6	-0	-16	-11	-12	-9	-10	21	-0	-6
	ADVCL	-3	-87	-6	-6	-12	-4	-3	-14	-26	-22	-2	-6	-6	1	-2	-20	32	0	-5
	ADVMOD	-12	-46	-95	-10	-9	0	-4	-27	-3	-31	0	-3	-15	-0	-3	-39	45	0	-9
	AMOD	10	-28	-27	-87	-12	-7	-4	-67	-4	-11	-1	-25	-49	-5	-9	-17	-27	-1	-15
	AUX	-15	-31	-21	-6	-45	-3	-2	-35	-5	-24	-0	-35	-9	-4	-4	-11	71	0	-7
	CASE	0	-20	-38	-4	-7	-95	-7	-57	-9	-14	-1	-63	-39	-7	-13	-21	-39	0	-21
	CC	-4	-13	-28	-1	-0	-2	-79	-18	-5	-2	-2	-3	-17	-1	-1	-8	-11	0	-8
	CCOMP	12	-41	-7	-1	-6	-6	-7	-84	-33	-23	-1	-19	-6	-3	-16	-5	23	0	-5
	CONJ	23	-31	-10	-1	-9	-3	-4	-42	-85	-15	-2	-20	-24	-1	-3	-9	-77	-0	-7
	COP	-36	-43	-7	-8	-17	-3	-2	-35	-27	-89	-1	-7	-9	-11	-13	-9	-43	0	-8
	DET	8	-25	-16	-7	-0	-9	-3	-63	-13	-2	-37	-11	-57	-3	-8	-19	-25	0	-12
	MARK	0	-28	-30	-40	-30	-4	-5	-57	-6	-18	-2	-24	-51	-4	-6	-23	-1	-0	-14
	NMOD	-23	-48	-16	-11	-12	-11	-6	-40	-18	-32	-4	-12	-31	-18	-94	-44	-23	-0	-15
	NSUBJ	-25	-72	-18	-10	-12	-11	-5	-31	-12	-39	-5	-13	-77	-5	-13	-91	-12	-0	-16
	OBJ	12	-19	-7	-1	-2	-2	-3	-13	-5	-0	-2	-10	-3	-7	-5	-2	-61	0	-3
	OBL	6	-17	-15	-7	1	-3	-2	-17	-6	-7	-3	-1	-19	-0	-3	-16	-10	-49	-15
	PARATAXIS	-36	-63	-16	-12	-9	-9	-3	-20	-19	-22	-12	-25	-90	-10	-14	-72	45	0	-17
PUNCT	7	-46	-8	-5	-9	-5	-7	-25	-18	-32	-3	-10	-7	-1	-5	-8	48	0	-5	
	ACL	ADVCL	ADVMOD	AMOD	AUX	CASE	CC	CCOMP	CONJ	COP	DET	MARK	NMOD	NSUBJ	OBJ	OBL	PARATAXIS	PUNCT	avg	
	English Target																			

Figure 26: Relative change in accuracy when cross-neutralizing English DEP tags using Greek XLM-R embeddings.