

Positional Artefacts Propagate Through Masked Language Model Embeddings

Ziyang Luo^{1*}, Artur Kulmizev¹, Xiaoxi Mao²

¹ Department of Linguistics and Philology, Uppsala University, Sweden

² Fuxi AI Lab, NetEase Inc., Hangzhou, China

Ziyang.Luo.9588@student.uu.se, artur.kulmizev@lingfil.uu.se

maoxiaoxi@corp.netease.com

Abstract

In this work, we demonstrate that the contextualized word vectors derived from pretrained masked language model-based encoders share a common, perhaps undesirable pattern across layers. Namely, we find cases of persistent outlier neurons within BERT and RoBERTa’s hidden state vectors that consistently bear the smallest or largest values in said vectors. In an attempt to investigate the source of this information, we introduce a neuron-level analysis method, which reveals that the outliers are closely related to information captured by positional embeddings. We also pre-train the RoBERTa-base models from scratch and find that the outliers disappear without using positional embeddings. These outliers, we find, are the major cause of anisotropy of encoders’ raw vector spaces, and *clipping* them leads to increased similarity across vectors. We demonstrate this in practice by showing that clipped vectors can more accurately distinguish word senses, as well as lead to better sentence embeddings when mean pooling. In three supervised tasks, we find that *clipping* does not affect the performance.

1 Introduction

A major area of NLP research in the deep learning era has concerned the representation of words in low-dimensional, continuous vector spaces. Traditional methods for achieving this have included word embedding models such as Word2Vec (Mikolov et al., 2013), GloVe (Pennington et al., 2014), and FastText (Bojanowski et al., 2017). However, though influential, such approaches all share a uniform pitfall in assigning a single, static vector to a word type. Given that the vast majority of words are polysemous (Klein and Murphy, 2001), static word embeddings cannot possibly represent a word’s changing meaning in context.

In recent years, deep language models, like ELMo (Peters et al., 2018), BERT (Devlin et al., 2019) and RoBERTa (Liu et al., 2019b), have achieved great success across many NLP tasks. Such models introduce a new type of word vectors, deemed the *contextualized* variety, where the representation is computed with respect to the context of the target word. Since these vectors are sensitive to context, they can better address the polysemy problem that hinders traditional word embeddings. Indeed, studies have shown that replacing static embeddings (e.g. word2vec) with contextualized ones (e.g. BERT) can benefit many NLP tasks, including constituency parsing (Kitaev and Klein, 2018), coreference resolution (Joshi et al., 2019) and machine translation (Liu et al., 2020).

However, despite the major success in deploying these representations across linguistic tasks, there remains little understanding about information embedded in contextualized vectors and the mechanisms that generate them. Indeed, an entire research area central to this core issue — the interpretability of neural NLP models — has recently emerged (Linzen et al., 2018, 2019; Alishahi et al., 2020). A key theme in this line of work has been the use of linear probes in investigating the linguistic properties of contextualized vectors (Tenney et al., 2019; Hewitt and Manning, 2019). Such studies, among many others, show that contextualization is an important factor that sets these embeddings apart from static ones, the latter of which are unreliable in extracting features central to context or linguistic hierarchy. Nonetheless, much of this work likewise fails to engage with the raw vector spaces of language models, preferring instead to focus its analysis on the transformed vectors. Indeed, the fraction of work that has done the former has shed some curious insights: that untransformed BERT sentence representations still lag behind word embeddings across a variety

* Work partly done during internship at NetEase Inc..

of semantic benchmarks (Reimers and Gurevych, 2019) and that the vector spaces of language models are explicitly anisotropic (Ethayarajh, 2019; Li et al., 2020a). Certainly, an awareness of the patterns inherent to models’ untransformed vector spaces — even if shallow — can only benefit the transformation-based analyses outlined above.

In this work, we shed light on a persistent pattern that can be observed for contextualized vectors produced by BERT and RoBERTa. Namely, we show that, across all layers, select neurons in BERT and RoBERTa consistently bear extremely large values. We observe this pattern across vectors for all words in several datasets, demonstrating that these singleton dimensions serve as major outliers to the distributions of neuron values in both encoders’ representational spaces. With this insight in mind, the contributions of our work are as follows:

1. We introduce a neuron-level method for analyzing the origin of a model’s outliers. Using this, we show that they are closely related to positional information.
2. In investigating the effects of *clipping* the outliers (zeroing-out), we show that the degree of anisotropy in the vector space diminishes significantly.
3. We show that after *clipping* the outliers, the BERT representations can better distinguish between a word’s potential senses in the word-in-context (WiC) dataset (Pilehvar and Camacho-Collados, 2019), as well as lead to better sentence embeddings when mean pooling.

2 Finding outliers

In this section, we demonstrate the existence of large-valued vector dimensions across nearly all tokens encoded by BERT and RoBERTa. To illustrate these patterns, we employ two well-known datasets — SST-2 (Socher et al., 2013) and QQP¹. SST-2 (60.7k sentences) is a widely-employed sentiment analysis dataset of movie reviews, while QQP (727.7k sentences) is a semantic textual similarity dataset of Quora questions, which collects questions across many topics. We choose these datasets in order to account for a reasonably wide distributions of domains and topics, but note that

¹<https://www.quora.com/q/quoradata/First-Quora-Dataset-Release-Question-Pairs>

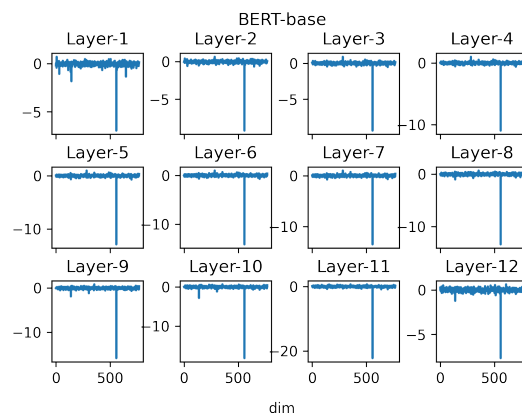


Figure 1: Average vectors for each layer of BERT-base.

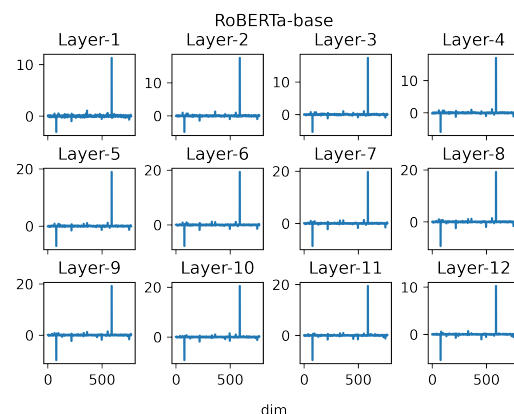


Figure 2: Average vectors for each layer of RoBERTa-base.

any dataset would illustrate our findings well. We randomly sample 10k sentences from the training sets of both SST-2 and QQP, tokenize them, and encode them via BERT-base and RoBERTa-base. All models are downloaded from the Huggingface Transformers Library (Wolf et al., 2020), though we replicated our results for BERT by loading the provided model weights via our own loaders.

When discounting the input embedding layers of each model, we are left with 3.68M and 3.59M contextualized token embeddings for BERT-base and RoBERTa-base, respectively. In order to illustrate the outlier patterns, we average all subword vectors for each layer of each model.

In examining BERT-base, we find that the minimum value of 96.60% of vectors lies in the 557th dimension. Figure 1 displays the averaged subword vectors for each layer of BERT-base, corroborating that these patterns exist across all layers. For RoBERTa-base, we likewise find that the maximum value of all vectors is the 588th element. Interestingly, the *minimum* element of 88.19% of vectors in

RoBERTa-base is the 77th element, implying that RoBERTa has two such outliers. Figure 2 displays the average vectors for each layer of RoBERTa-base.

Our observations here reveal a curious pattern that is present in the base versions of BERT and RoBERTa. We also corroborate the same findings for the large and distilled (Sanh et al., 2020) variants of these architectures, which can be found in the Appendix A. Indeed, it would be difficult to reach any sort of conclusion about the representational geometry of such models without understanding the outliers’ origin(s).

3 Where do outliers come from?

In this section, we attempt to trace the source of the outlier dimensions in BERT-base and RoBERTa-base (henceforth BERT and RoBERTa). Similarly to the previous section, we can corroborate the results of the experiments described here (as well as in the remainder of the paper) for the large and distilled varieties of each respective architecture. Thus, for reasons of brevity, we focus our forthcoming analyses on the base versions of BERT and RoBERTa and include results for the remaining models in the Appendix B.2 for the interested reader.

In our per-layer analysis in §2, we report that outlier dimensions exist across every layer in each model. Upon a closer look at the input layer (which features a vector sum of positional, segment, and token embeddings), we find that the same outliers also exist in positional embeddings. Figure 3 shows that the 1st positional embedding of BERT has two such dimensions, where the 557th element is likewise the minimum. Interestingly, this pattern does not exist in other positional embeddings, nor in segment or token embeddings. Furthermore, Figure 4 shows that the 4th positional embedding of RoBERTa has four outliers, which include the aforementioned 77th and 588th dimensions. We also find that, from the 4th position to the final position, the maximum element of 99.8% positional embeddings is the 588th element.

Digging deeper, we observe similar patterns in the Layer Normalization (LN, Ba et al. (2016)) parameters of both models. Recall that LN has two learnable parameters — *gain* (γ) and *bias* (β) — both of which are 768-dimension vectors (in the case of the base models). These are designed as an affine transformation over dimension-wise nor-

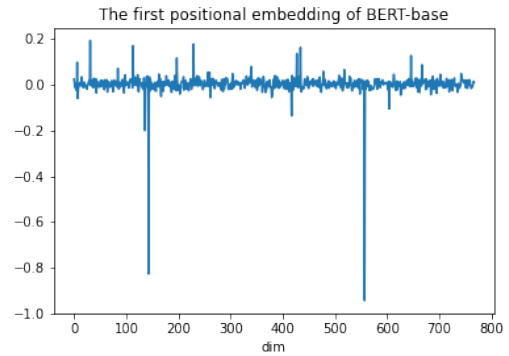


Figure 3: The first positional embedding of BERT-base.

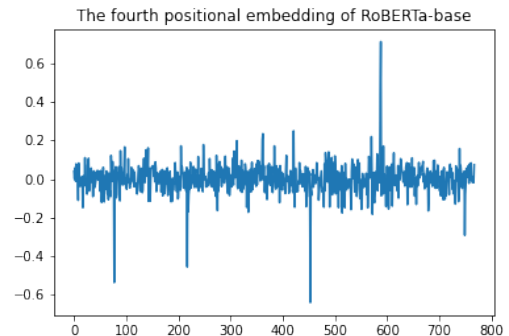


Figure 4: The fourth positional embedding of RoBERTa-base.

malized vectors in order to, like most normalization strategies, improve their expressive ability and to aid in optimization. Every layer of BERT and RoBERTa applies separate LNs post-attention and pre-output. For BERT, the 557th element of the γ vector is always among the top-6 largest values for the first ten layers’ first LN. Specifically, it is the largest value in the first three layers. For RoBERTa, the 588th element of the first LN’s β vector is always among the top-2 largest values for all layers — it is largest in the first five layers. Furthermore, the 77th element of the second LN’s γ are among the top-7 largest values from the second to the tenth layer.

It is reasonable to conclude that, after the vector normalization performed by LN, the outliers observed in the raw embeddings are lost. We hypothesize that these particular neurons are somehow important to the network, such that they retained after scaling the normalized vectors by the affine transformation involving γ and β . Indeed, we observe that, in BERT, only the 1st position’s embedding has such an outlier. However, it is subsequently observed in every layer and token

after the first LN is applied. Since LayerNorm is trained globally and is not token specific, it happens to rescale every vector such that the positional information is retained. We corroborate this by observing that all vectors share the same γ . This effectively guarantees the presence of outliers in the 1st layer, which are then propagated upward by means of the Transformer’s residual connection (He et al., 2015). Also, it is important to note that, in the case of BERT, the first position’s embedding is directly tied to the requisite [CLS] token, which is prepended to all sequences as part of the MLM training objective. This has been recently noted to affect e.g. attention patterns, where much of the probability mass is distributed to this particular token alone, despite it bearing the smallest norm among all other vectors in a given layer and head (Kobayashi et al., 2020).

Neuron-level analysis In order to test the extent to which BERT and RoBERTa’s outliers are related to positional information, we employ a probing technique inspired by Durrani et al. (2020). First, we train a linear probe $W \in R^{M \times N}$ without bias to predict the position of a contextualized vector in a sentence. In Durrani et al. (2020), the weights of the classifier are employed as a proxy for selecting the most relevant neurons to the prediction. In doing so, they assume that, the larger the absolute value of the weight, the more important the corresponding neuron. However, this method disregards the magnitudes of the values of neurons, as a large weights do not necessarily imply that the neuron has high contribution to the final classification result. For example, if the value of a neuron is close to zero, a large weight also leads to a small contribution. In order to address this issue, we define the contribution of the i^{th} neuron as $c(i) = abs(w_i * v_i)$ for $i = 1, 2, 3, \dots, n$, where w_i is the i^{th} weight and v_i is the i^{th} neuron in the contextualized word vector. We name $C = [c(1), c(2), \dots, c(n)]$ as a contribution vector. If a neuron has a high contribution, this means that this neuron is highly relevant to the final classification result.

We train, validate, and test our probe on the splits provided in the SST-2 dataset (as mentioned in §2, we surmise that any dataset would be adequate for demonstrating this). The linear probe is a 768×300 matrix, which we train separately for each layer. Since all SST-2 sentences are shorter than 300 tokens in length, we set $M = 300$. We

use a batch size of 128 and train for 10 epochs with a categorical cross-entropy loss, optimized by Adam (Kingma and Ba, 2017).

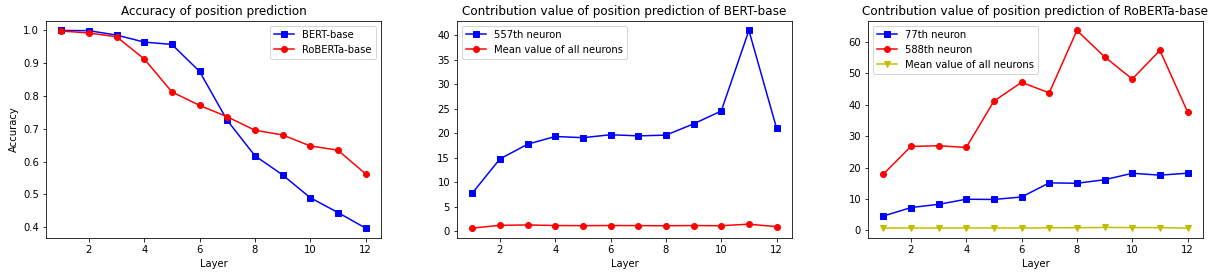
Figure 5a shows that, while it is possible to decode positional information from the lowest three layers with almost perfect accuracy, much of this information is gradually lost higher up in the model. Furthermore, it appears that the higher layers of RoBERTa contain more positional information than BERT. Looking at Figure 5b, we see that BERT’s outlier neuron has a higher contribution in position prediction than the average contribution of all neurons. We also find that the contribution values of the same neuron are the highest in all layers. Combined with the aforementioned pattern of the first positional embedding, we can conclude that the 557th neuron is related to positional information. Likewise, for RoBERTa, Figure 5c shows that the 77th and 588th neurons have the highest contribution for position prediction. We also find that the contribution values of the 588th neurons are always largest for all layers, which implies that these neurons are likewise related to positional information.²

Removing positional embeddings In order to isolate the relation between outlier neurons and positional information, we pre-train two RoBERTa-base models (with and without positional embeddings) from scratch using Fairseq (Ott et al., 2019). Our pre-training data is the English Wikipedia Corpus³, where we train for 200k steps with a batch size of 256, optimized by Adam. All models share the same hyper-parameters, which are listed in the Appendix C.1. We use four NVIDIA A100 GPUs to pre-train each model, costing about 35 hours per model.

We find that, without the help of positional embeddings, the validation perplexity of RoBERTa-base is very high at 354.0, which is in line with Lee et al. (2019)’s observation that the self-attention mechanism of Transformer Encoder is order-invariant. In other words, the removal of PEs from RoBERTa-base makes it a bag-of-word model, whose outputs do not contain any positional information. In contrast, the perplexity of RoBERTa equipped with standard positional embeddings is much lower at 4.3, which is likewise expected.

²We also use heatmaps to show the contribution values in Appendix B.1.

³We randomly select 158.4M sentences for training and 50k sentences for validation.



(a) Accuracy of position prediction.

(b) The contribution value of BERT-base’s outlier neuron on position prediction.

(c) The contribution value of RoBERTa-base’s outlier neurons on position prediction.

In examining outlier neurons, we employ the same datasets detailed in §2. For the RoBERTa-base model with PEs, we find that the maximum element of 82.56% of all vectors is the 81st dimension⁴, similarly to our findings above. However, we do not observe the presence of such outlier neurons in the RoBERTa-base model without PEs, which indicates that the outlier neurons are tied directly to positional information. Similar to §2, we display the averaged subword vectors for each layer of our models in Appendix C.2, which also corroborate our results.

4 Clipping the outliers

In §3, we demonstrated that outlier neurons are related to positional information. In this section, we investigate the effects of zeroing out these dimensions in contextualized vectors, a process which we refer to as clipping.

4.1 Vector space geometry

Anisotropy Ethayarajh (2019) observe that contextualized word vectors are anisotropic in all non-input layers, which means that the average cosine similarity between uniformly randomly sampled words is close to 1. To corroborate this finding, we randomly sample 2000 sentences from the SST-2 training set and create 1000 sentence-pairs. Then, we randomly select a token in each sentence, discarding all other tokens. This effectively sets the correspondence between the two sentences to two tokens instead. Following this, we compute the cosine similarity between these two tokens to measure the anisotropy of contextualized vectors.

In the left plot of Figure 6, we can see that contextualized representations of BERT and RoBERTa are more anisotropic in higher layers. This is espe-

cially true for RoBERTa, where the average cosine similarity between random words is larger than 0.5 after the first non-input layer. This implies that the internal representations in BERT and RoBERTa occupy a narrow cone in the vector space.

Since outlier neurons tend to be valued higher or lower than all other contextualized vector dimensions, we hypothesize that they are the main culprit behind the degree of observed anisotropy. To verify our hypothesis, we clip BERT and RoBERTa’s outliers by setting each neuron’s value to zero. The left plot in Figure 6 shows that, after clipping the outliers, their vector spaces become close to isotropic.

Self-similarity In addition to remarking upon the anisotropic characteristics of contextualized vector spaces, Ethayarajh (2019) introduce several measures to gauge the extent of “contextualization” inherent models. One such metric is *self-similarity*, which the authors employ to compare the similarity of a word’s internal representations in different contexts. Given a word w and n different sentences s_1, s_2, \dots, s_n which contain such word, $f_l^i(w)$ is the internal representation of w in sentence s_i in the l^{th} layer. The average self-similarity of w in the l^{th} layer is then defined as:

$$\text{SelfSim}_l(w) = \frac{\sum_{i=1}^n \sum_{j=i+1}^n \cos(f_l^i(w), f_l^j(w))}{n(n-1)} \quad (1)$$

Intuitively, a self-similarity score of 1 indicates that no contextualization is being performed by the model (e.g. static word embeddings), while a score of 0 implies that representations for a given word are maximally different given various contexts.

To investigate the effect of outlier neurons on a model’s self-similarity, we sample 1000 different words from SST-2 training set, all of which appear at least in 10 different sentences. We then com-

⁴Different initializations make our models have different outlier dimensions.

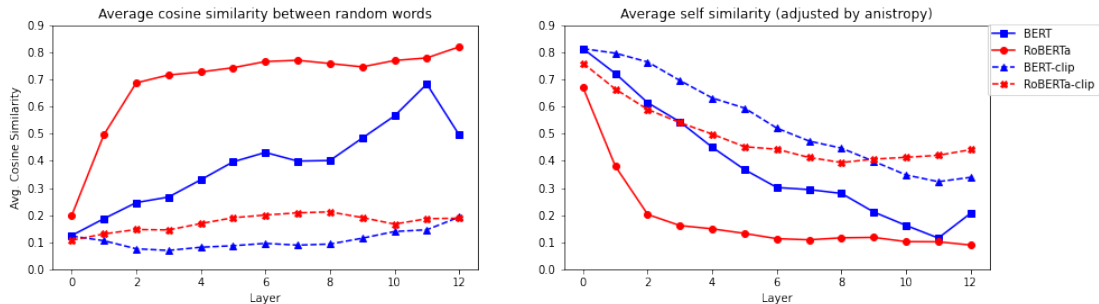


Figure 6: Left: anisotropy measurement of contextualized word vectors in BERT and RoBERTa before and after clipping the outlier dimensions. Right: self-similarity measurement of BERT and RoBERTa before and after clipping.

pute the average self-similarity of these words as contextualized by BERT and RoBERTa — before and after clipping the outliers. To adjust for the effect of anisotropy, we subtract the self-similarity from each layer’s anisotropy measurement, as in [Ethayarajh \(2019\)](#).

The right plot in Figure 6 shows that, similarly to the findings in [\(Ethayarajh, 2019\)](#), a word’s self-similarity is highest in the lower layers, but decreases in higher layers. Crucially, we also observe that, after clipping the outlier dimensions, the self-similarity increases, indicating that vectors become closer to each other in the contextualized space. This bears some impact on studies attempting to characterize the vector spaces of models like BERT and RoBERTa, as it is clearly possible to overstate the degree of “contextualization” without addressing the effect of positional artefacts.

4.2 Word sense

Bearing in mind the findings of the previous section, we now turn to the question of word sense, as captured by contextualized embeddings. Suppose that we have a target word w , which appears in two sentences. w has the same sense in these two sentences, but its contextualized representations are not identical due to the word appearing in (perhaps slightly) different contexts. In the previous few sections, we showed that outlier neurons are related to positional information and that clipping them can make a word’s contextualized vectors more similar. Here, we hypothesize that clipping such dimensions can likewise aid in intrinsic semantic tasks, like differentiating senses of a word.

To test our hypothesis, we analyze contextualized vectors using the word-in-context (WiC) dataset [\(Pilehvar and Camacho-Collados, 2019\)](#), which is designed to identify the meaning of words

Model	Layer	Threshold	Accuracy
Baseline	-	-	50.0%
Before clipping			
BERT	7	0.7	67.5%
RoBERTa	10	0.9	69.0%
After clipping			
BERT-clip	10	0.5	68.4%
RoBERTa-clip	11	0.6	69.9%

Table 1: The best accuracy scores on WiC dataset. **Bold** indicates that the best result increases after clipping.

in different contexts. WiC is a binary classification task, where, given a target word and two sentences which contain it, models must determine whether the word has the same meaning across the two sentences.

In order to test how well we can identify differences in word senses using contextualized vectors, we compute the cosine similarity between contextualized vectors of target words across pairs of sentences, as they appear in the WiC dataset. If the similarity value is larger than a specified threshold, we assign the true label to the sentence pair; otherwise, we assign the false label. We use this method to compare the accuracy of BERT and RoBERTa on WiC before and after clipping the outliers. Since this method does not require any training, we test our models on the WiC training dataset.⁵ We compare 9 different thresholds from 0.1 to 0.9, as well as a simple baseline model that assigns the true labels to all samples.

Table 1 shows that after clipping outliers, the best accuracy scores of BERT and RoBERTa increase about 1%.⁶ This indicates that these neurons

⁵The WiC test set does not provide labels and the size of validation set is too small (638 sentences pairs). We thus choose to use the training dataset (5428 sentences pairs).

⁶The thresholds are different due to the fact that the cosine

Dataset	STS-B	SICK-R	STS-12	STS-13	STS-14	STS-15	STS-16
Baseline							
Avg. GloVe	58.02	53.76	55.14	70.66	59.73	68.25	63.66
Before clipping							
BERT	58.61(3)	60.78(2)	48.00(1)	61.19(12)	50.10(12)	61.15(1)	62.38(12)
RoBERTa	56.60(11)	64.68(11)	40.00(1)	58.33(11)	49.79(8)	64.39(9)	64.82(11)
After clipping							
BERT-clip	63.06(2)	61.74(2)	50.40(1)	61.44(1)	54.52(2)	67.00(2)	64.18(2)
RoBERTa-clip	60.61(11)	64.82(11)	43.44(1)	59.72(11)	51.92(3)	66.15(3)	67.14(11)

Table 2: Experimental results on semantic textual similarity, where the baselines results are published in Reimers and Gurevych (2019). We show the best Spearman rank correlation between sentence embeddings’ cosine similarity and the golden labels. The results are reported as $r \times 100$. The number in the parenthesis denotes that this result belongs to the specific layer. **Bold** indicates that the best result increases after clipping.

are less related to word sense information and can be safely clipped for this particular task (if performed in an unsupervised fashion).

4.3 Sentence embedding

Venturing beyond the word-level, we also hypothesize that outlier clipping can lead to better sentence embeddings when relying on the cosine similarity metric. To test this, we follow Reimers and Gurevych (2019) in evaluating our models on 7 semantic textual similarity (STS) datasets, including the STS-B benchmark (STS-B) (Cer et al., 2017), the SICK-Relatedness (SICK-R) dataset (Bentivogli et al., 2016) and the STS tasks 2012-2016 (Agirre et al., 2012, 2013, 2014, 2015, 2016). Each sentence pair in these datasets is annotated with a relatedness score on a 5-point rating scale, as obtained from human judgments. We load each dataset using the SentEval toolkit (Conneau and Kiela, 2018).

Indeed, the most common approach for computing sentence embeddings from contextualized models is simply averaging all subword vectors that comprise a given sentence (Reimers and Gurevych, 2019). We follow this method in obtaining embeddings for each pair of sentences in the aforementioned tasks, between which we compute the cosine similarity. Given a set of similarity and gold relatedness scores, we then calculate the Spearman rank correlation. As a comparison, we also consider averaged GloVe embeddings as our baseline.

Table 2 shows that, after clipping the outliers, the best Spearman rank correlation scores for BERT and RoBERTa increase across all datasets, some by a large margin. This indicates that clipping the outlier neurons can lead to better sentence embeddings when mean pooling. However, like Li et al.

similarity is inflated in the presence of outlier neurons.

Model	SST-2	IMDB	SST-5
Before clipping			
BERT	85.9%(12)	86.8%(10)	46.2%(10)
RoBERTa	88.4%(8)	91.5%(9)	46.9%(7)
After clipping			
BERT-clip	85.4%(12)	86.4%(10)	46.1%(12)
RoBERTa-clip	88.7%(8)	91.6%(9)	47.0%(7)

Table 3: The best accuracy scores on different supervised tasks. The number in the parenthesis denotes that this result belongs to the specific layer.

(2020b), we also notice that averaged GloVe embeddings still manage outperform both BERT and RoBERTa on all STS 2012-16 tasks. This implies that the post-clipping reduction in anisotropy is only a partial explanation for why contextualized, mean-pooled sentence embeddings still lag behind static word embeddings in capturing the semantics of a given sentence.

4.4 Supervised tasks

In the previous sections, we analyzed the effects of clipping outlier neurons on various intrinsic semantic tasks. Here, we explore the effects of clipping in a supervised scenario, where we hypothesize that a model will learn to discard outlier information if it is not needed for a given task. We consider two binary classification tasks, SST-2 and IMDB (Maas et al., 2011), and a multi-class classification task, SST-5, which is a 5-class version of SST-2. First, we freeze all the parameters of the pre-trained models and use the same method in §4.3 to get the sentence embedding of each sentence. Then, we train a simple linear classifier $W \in R^{768 \times N}$ for each layer, where N is the number of classes. We use different batch sizes for different tasks, 768 for SST-2, 128 for IMDB and 1536 for SST-5. Then we train for 10 epochs with a categorical cross-entropy loss, optimized by Adam.

Table 3 shows that there is little difference in employing raw vs. clipped vectors in terms of task performance. This indicates that using vectors with clipped outliers does not drastically affect classifier accuracy when it comes to these common tasks.

5 Discussion

The experiments detailed in the previous sections point to the dangers of relying on metrics like cosine similarity when making observations about models’ representational spaces. This is particularly salient when the vectors being compared are taken off-the-shelf and their composition is not widely understood. Given the presence of model idiosyncrasies like the outliers highlighted here, mean-sensitive, L2 normalized metrics (e.g. cosine similarity or Pearson correlation) will inevitably weigh the comparison of vectors along the highest-valued dimensions. In the case of positional artefacts propagating through the BERT and RoBERTa networks, the basis of comparison is inevitably steered towards whatever information is captured in those dimensions. Furthermore, since such outlier values show little variance across vectors, proxy metrics of anisotropy like measuring the average cosine similarity across random words (detailed in §4.1) will inevitably return an exceedingly high similarity, no matter what the context. When cosine similarity is viewed primarily as means of semantic comparison between word or sentence vectors, the prospect of calculating cosine similarity for a benchmark like WiC or STS-B becomes erroneous. Though an examination of distance metrics is outside the scope of this study, we acknowledge similar points as having been addressed in regards to static word embeddings (Mimno and Thompson, 2017) as well as contextualized ones (Li et al., 2020b). Likewise, we would like to stress that our manual clipping operation was performed for illustrative purposes and that interested researchers should employ more systematic post-hoc normalization strategies, e.g. whitening (Su et al., 2021), when working with hidden states directly.

Relatedly, the anisotropic nature of the vector space that persists even after clipping the outliers suggests that positional artefacts are simply part of the explanation. Per this point, Gao et al. (2019) prove that, in training any sort of model with likelihood loss, the representations learned for tokens being predicted will be naturally be pushed away from most other tokens in order to achieve a higher like-

lihood. They relate this observation to the Zipfian nature of word distributions, where the vast majority of words are infrequent. Li et al. (2020a) extend this insight specifically to BERT and show that, while high frequency words concentrate densely, low frequency words are much more sparsely distributed. Though we do not attempt to dispute these claims with our findings, we do hope our experiments will highlight the important role that positional embeddings play in the representational geometry of Transformer-based models. Indeed, recent work has demonstrated that employing relative positional embeddings and untying them from the simultaneously learned word embeddings has lead to impressive gains for BERT-based architectures across common benchmarks (He et al., 2020; Ke et al., 2020). It remains to be seen how such procedures affect the representations of such models, however.

Beyond this, it is clear that LayerNorm is the reason positional artefacts propagate through model representations in the first place. Indeed, our experiments show that the outlier dimension observed for BERT is tied directly to the [CLS] token, which always occurs at the requisite 1st position — despite having no linguistic bearing on the sequence of observed tokens being modeled. However, the fact that RoBERTa (which employs a similar delimiter) retains outliers originating from different positions’ embeddings implies that the issue of artefact propagation is not simply a relic of task design. It is possible that whatever positional idiosyncrasies contribute to a task’s loss are likewise retained in their respective embeddings. In the case of BERT, the outlier dimension may be granted a large negative weight in order to differentiate the (privileged) 1st position between all others. This information being reconstructed by the LayerNorm parameters, which are shared for all positions in the sequence length, and then propagated up through the Transformer network is a phenomenon worthy of further attention.

6 Related work

In recent years, an explosion of work focused on understanding the inner workings of pretrained neural language models has emerged. One line of such work investigates the self-attention mechanism of Transformer-based models, aiming to e.g. characterize its patterns or decode syntactic structure (Raganato and Tiedemann, 2018; Vig, 2019;

Mareček and Rosa, 2018; Voita et al., 2019; Clark et al., 2019; Kobayashi et al., 2020). Another line of work analyzes models’ internal representations using probes. These are often linear classifiers that take representations as input and are trained with supervised tasks in mind, e.g. POS-tagging, dependency parsing (Tenney et al., 2019; Liu et al., 2019a; Lin et al., 2019; Hewitt and Manning, 2019; Zhao et al., 2020). In such work, high probing accuracies are often likened to a particular model having “learned” the task in question.

Most similar to our work, Ethayarajh (2019) investigate the extent of “contextualization” in models like BERT, ELMo, and GPT-2 (Radford et al., 2019). Mainly, they demonstrate that the contextualized vectors of all words are non-isotropic across all models and layers. However, they do not indicate why these models have such properties. Also relevant are the studies of Dalvi et al. (2018), who introduce a neuron-level analysis method, and Durrani et al. (2020), who use this method to analyze individual neurons in contextualized word vectors. Similarly to our experiment, Durrani et al. (2020) train a linear probe to predict linguistic information stored in a vector. They then employ the weights of the classifier as a proxy to select the most relevant neurons to a particular task. In a similar vein, Coenen et al. (2019) demonstrate the existence of syntactic and semantic subspaces in BERT representations.

7 Conclusion

In this paper, we called attention to sets of outlier neurons that appear in BERT and RoBERTa’s internal representations, which bear consistently large values when compared to the distribution of values of all other neurons. In investigating the origin of these outliers, we employed a neuron-level analysis method which revealed that they are artefacts derived from positional embeddings and Layer Normalization. Furthermore, we found that outliers are a major cause for the anisotropy of a model’s vector space (Ethayarajh, 2019). Clipping them, consequently, can make the vector space more directionally uniform and increase the similarity between words’ contextual representations. In addition, we showed that outliers can distort results when investigating word sense within contextualized representations as well as obtaining sentence embeddings via mean pooling, where removing them leads to uniformly better results. Lastly, we

find that “clipping” does not affect models’ performance on three supervised tasks.

It is important to note that the exact dimensions at which the outliers occur will vary pending different initializations and training procedures (as evidenced by our own RoBERTa model). As such, future work will aim at investigating strategies for mitigating the propagation of these artefacts when pretraining. Furthermore, given that both BERT and RoBERTa are masked language models, it will be interesting to investigate whether or not similar artefacts occur in e.g. autoregressive models like GPT-2 (Radford et al., 2019) or XLNet (Yang et al., 2019). Per the insights of Gao et al. (2019), it is very likely that the representational spaces of such models are anisotropic, but it is important to gauge the extent to which this can be traced to positional artefacts.

Authors’ Note We would like to mention Koval-eva et al. (2021)’s contemporaneous work, which likewise draws attention to BERT’s outlier neurons. While our discussion situates outliers in the context of positional embeddings and vector spaces, Kovaleva et al. (2021) offer an exhaustive analysis of LayerNorm parameterization and its impact on masked language modeling and finetuning. We refer the interested reader to that work for a thorough discussion of LayerNorm’s role in the outlier neuron phenomenon.

Acknowledgments We would like to thank Joakim Nivre and Daniel Dakota for fruitful discussions and the anonymous reviewers for their excellent feedback.

References

- Eneko Agirre, Carmen Banea, Claire Cardie, Daniel Cer, Mona Diab, Aitor Gonzalez-Agirre, Weiwei Guo, Iñigo Lopez-Gazpio, Montse Maritxalar, Rada Mihalcea, German Rigau, Larraitz Uria, and Janyce Wiebe. 2015. *SemEval-2015 task 2: Semantic textual similarity, English, Spanish and pilot on interpretability*. In *Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015)*, pages 252–263, Denver, Colorado. Association for Computational Linguistics.
- Eneko Agirre, Carmen Banea, Claire Cardie, Daniel Cer, Mona Diab, Aitor Gonzalez-Agirre, Weiwei Guo, Rada Mihalcea, German Rigau, and Janyce Wiebe. 2014. *SemEval-2014 task 10: Multilingual semantic textual similarity*. In *Proceedings of the 8th International Workshop on Semantic Evaluation*

- (*SemEval 2014*), pages 81–91, Dublin, Ireland. Association for Computational Linguistics.
- Eneko Agirre, Carmen Banea, Daniel Cer, Mona Diab, Aitor Gonzalez-Agirre, Rada Mihalcea, German Rigau, and Janyce Wiebe. 2016. [SemEval-2016 task 1: Semantic textual similarity, monolingual and cross-lingual evaluation](#). In *Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016)*, pages 497–511, San Diego, California. Association for Computational Linguistics.
- Eneko Agirre, Daniel Cer, Mona Diab, and Aitor Gonzalez-Agirre. 2012. [SemEval-2012 task 6: A pilot on semantic textual similarity](#). In **SEM 2012: The First Joint Conference on Lexical and Computational Semantics – Volume 1: Proceedings of the main conference and the shared task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation (SemEval 2012)*, pages 385–393, Montréal, Canada. Association for Computational Linguistics.
- Eneko Agirre, Daniel Cer, Mona Diab, Aitor Gonzalez-Agirre, and Weiwei Guo. 2013. [*SEM 2013 shared task: Semantic textual similarity](#). In *Second Joint Conference on Lexical and Computational Semantics (*SEM), Volume 1: Proceedings of the Main Conference and the Shared Task: Semantic Textual Similarity*, pages 32–43, Atlanta, Georgia, USA. Association for Computational Linguistics.
- Afra Alishahi, Yonatan Belinkov, Grzegorz Chrupała, Dieuwke Hupkes, Yuval Pinter, and Hassan Sajjad, editors. 2020. [Proceedings of the Third BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP](#). Association for Computational Linguistics, Online.
- Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. 2016. [Layer normalization](#).
- Luisa Bentivogli, Raffaella Bernardi, Marco Marelli, Stefano Menini, Marco Baroni, and Roberto Zamparelli. 2016. [Sick through the semeval glasses: Lesson learned from the evaluation of compositional distributional semantic models on full sentences through semantic relatedness and textual entailment](#). *Language Resources and Evaluation*, 50.
- Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. [Enriching word vectors with subword information](#). *Transactions of the Association for Computational Linguistics*, 5:135–146.
- Daniel Cer, Mona Diab, Eneko Agirre, Iñigo Lopez-Gazpio, and Lucia Specia. 2017. [SemEval-2017 task 1: Semantic textual similarity multilingual and crosslingual focused evaluation](#). In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages 1–14, Vancouver, Canada. Association for Computational Linguistics.
- 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.
- Andy Coenen, Emily Reif, Ann Yuan, Been Kim, Adam Pearce, Fernanda Viégas, and Martin Wattenberg. 2019. [Visualizing and measuring the geometry of bert](#). *arXiv preprint arXiv:1906.02715*.
- Alexis Conneau and Douwe Kiela. 2018. [SentEval: An evaluation toolkit for universal sentence representations](#). In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).
- Fahim Dalvi, Nadir Durrani, Hassan Sajjad, Yonatan Belinkov, Anthony Bau, and James Glass. 2018. [What is one grain of sand in the desert? analyzing individual neurons in deep nlp models](#).
- 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.
- Nadir Durrani, Hassan Sajjad, Fahim Dalvi, and Yonatan Belinkov. 2020. [Analyzing individual neurons in pre-trained language models](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4865–4880, Online. Association for Computational Linguistics.
- Kawin Ethayarajh. 2019. [How contextual are contextualized word representations? comparing the geometry of BERT, ELMo, and GPT-2 embeddings](#). 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 55–65, Hong Kong, China. Association for Computational Linguistics.
- Jun Gao, Di He, Xu Tan, Tao Qin, Liwei Wang, and Tiejun Liu. 2019. [Representation degeneration problem in training natural language generation models](#). *arXiv preprint arXiv:1907.12009*.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2015. [Deep residual learning for image recognition](#).
- Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2020. [Deberta: Decoding-enhanced bert with disentangled attention](#). *arXiv preprint arXiv:2006.03654*.

- 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.
- Mandar Joshi, Omer Levy, Luke Zettlemoyer, and Daniel Weld. 2019. [BERT for coreference resolution: Baselines and analysis](#). 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 5803–5808, Hong Kong, China. Association for Computational Linguistics.
- Guolin Ke, Di He, and Tie-Yan Liu. 2020. [Rethinking the positional encoding in language pre-training](#). *arXiv preprint arXiv:2006.15595*.
- Diederik P. Kingma and Jimmy Ba. 2017. [Adam: A method for stochastic optimization](#).
- Nikita Kitaev and Dan Klein. 2018. [Constituency parsing with a self-attentive encoder](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2676–2686, Melbourne, Australia. Association for Computational Linguistics.
- Devorah E. Klein and Gregory L. Murphy. 2001. [The representation of polysemous words](#). *Journal of Memory and Language*, 45(2):259 – 282.
- Goro Kobayashi, Tatsuki Kuribayashi, Sho Yokoi, and Kentaro Inui. 2020. [Attention is not only a weight: Analyzing transformers with vector norms](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7057–7075, Online. Association for Computational Linguistics.
- Olga Kovaleva, Saurabh Kulshreshtha, Anna Rogers, and Anna Rumshisky. 2021. [Bert busters: Outlier layernorm dimensions that disrupt bert](#).
- Juho Lee, Yoonho Lee, Jungtaek Kim, Adam Kosiosek, Seungjin Choi, and Yee Whye Teh. 2019. [Set transformer: A framework for attention-based permutation-invariant neural networks](#). In *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 3744–3753. PMLR.
- Bohan Li, Hao Zhou, Junxian He, Mingxuan Wang, Yiming Yang, and Lei Li. 2020a. [On the sentence embeddings from pre-trained language models](#). *arXiv preprint arXiv:2011.05864*.
- Bohan Li, Hao Zhou, Junxian He, Mingxuan Wang, Yiming Yang, and Lei Li. 2020b. [On the sentence embeddings from pre-trained language models](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9119–9130, Online. Association for Computational Linguistics.
- Yongjie Lin, Yi Chern Tan, and Robert Frank. 2019. [Open sesame: Getting inside bert’s linguistic knowledge](#).
- Tal Linzen, Grzegorz Chrupała, and Afra Alishahi, editors. 2018. *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*. Association for Computational Linguistics, Brussels, Belgium.
- Tal Linzen, Grzegorz Chrupała, Yonatan Belinkov, and Dieuwke Hupkes, editors. 2019. *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*. Association for Computational Linguistics, Florence, Italy.
- Nelson F. Liu, Matt Gardner, Yonatan Belinkov, Matthew E. Peters, and Noah A. Smith. 2019a. [Linguistic knowledge and transferability of contextual 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 1073–1094, Minneapolis, Minnesota. Association for Computational Linguistics.
- Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. 2020. [Multilingual denoising pre-training for neural machine translation](#).
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019b. [Roberta: A robustly optimized bert pretraining approach](#).
- Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. 2011. [Learning word vectors for sentiment analysis](#). In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 142–150, Portland, Oregon, USA. Association for Computational Linguistics.
- David Mareček and Rudolf Rosa. 2018. [Extracting syntactic trees from transformer encoder self-attentions](#). In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 347–349, Brussels, Belgium. Association for Computational Linguistics.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. [Efficient estimation of word representations in vector space](#).
- David Mimno and Laure Thompson. 2017. [The strange geometry of skip-gram with negative sampling](#). In *Empirical Methods in Natural Language Processing*.

- Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. [fairseq: A fast, extensible toolkit for sequence modeling](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pages 48–53, Minneapolis, Minnesota. Association for Computational Linguistics.
- Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. [GloVe: Global vectors for word representation](#). In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, Doha, Qatar. Association for Computational Linguistics.
- Matthew Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. [Deep contextualized word representations](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 2227–2237, New Orleans, Louisiana. Association for Computational Linguistics.
- Mohammad Taher Pilehvar and Jose Camacho-Collados. 2019. [WiC: the word-in-context dataset for evaluating context-sensitive meaning 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 1267–1273, Minneapolis, Minnesota. Association for Computational Linguistics.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. [Language models are unsupervised multitask learners](#).
- Alessandro Raganato and Jörg Tiedemann. 2018. [An analysis of encoder representations in transformer-based machine translation](#). In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 287–297, Brussels, Belgium. Association for Computational Linguistics.
- Nils Reimers and Iryna Gurevych. 2019. [Sentence-BERT: Sentence embeddings using Siamese BERT-networks](#). 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 3982–3992, Hong Kong, China. Association for Computational Linguistics.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2020. [Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter](#).
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013. [Recursive deep models for semantic compositionality over a sentiment treebank](#). In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1631–1642, Seattle, Washington, USA. Association for Computational Linguistics.
- Jianlin Su, Jiarun Cao, Weijie Liu, and Yangyiwen Ou. 2021. [Whitening sentence representations for better semantics and faster retrieval](#). *arXiv preprint arXiv:2103.15316*.
- Ian Tenney, Patrick Xia, Berlin Chen, Alex Wang, Adam Poliak, R Thomas McCoy, Najoung Kim, Benjamin Van Durme, Sam Bowman, Dipanjan Das, and Ellie Pavlick. 2019. [What do you learn from context? probing for sentence structure in contextualized word representations](#). In *International Conference on Learning Representations*.
- Jesse Vig. 2019. [Visualizing attention in transformer-based language representation models](#).
- Elena Voita, David Talbot, Fedor Moiseev, Rico Senrich, and Ivan Titov. 2019. [Analyzing multi-head self-attention: Specialized heads do the heavy lifting, the rest can be pruned](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5797–5808, Florence, Italy. Association for Computational Linguistics.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. [Transformers: State-of-the-art natural language processing](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.
- Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. 2019. [Xlnet: Generalized autoregressive pretraining for language understanding](#). In *Advances in Neural Information Processing Systems*, volume 32, pages 5753–5763. Curran Associates, Inc.
- Mengjie Zhao, Philipp Dufter, Yadollah Yaghoobzadeh, and Hinrich Schütze. 2020. [Quantifying the contextualization of word representations with semantic class probing](#).

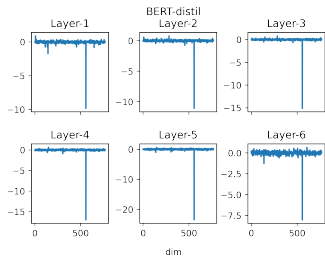


Figure 7: Average vectors for each layer of BERT-distil.

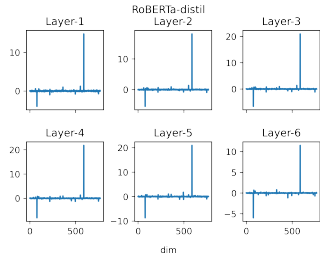


Figure 8: Average vectors for each layer of RoBERTa-distil.

A Outliers of distilled and large models

For BERT-distil, Figure 7 shows the patterns of BERT-distil across all layers. The 557th element is an outlier. For RoBERTa-distil, Figure 8 shows the patterns of RoBERTa-distil across all layers. the 77th and 588th elements are two outliers. For BERT-large, Figure 9 shows the patterns of BERT-large across all layers. From the first layer to the tenth layer, the 896th element is an outlier. From the tenth layer to the seventeenth layer, the 678th element is an outlier. From the sixteenth layer to the nineteenth layer, the 122nd element is an outlier. From the nineteenth layer to the twenty-third layer, the 928th element is an outlier. The final layer does not have outliers. For RoBERTa-large, Figure

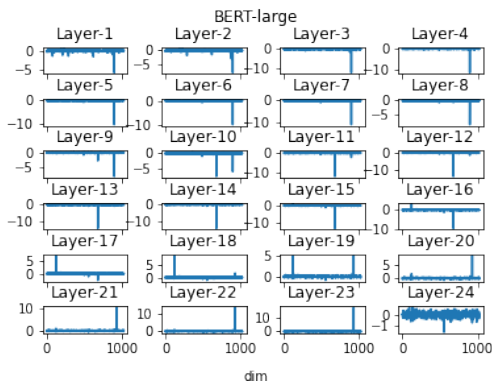


Figure 9: Average vectors for each layer of BERT-large.

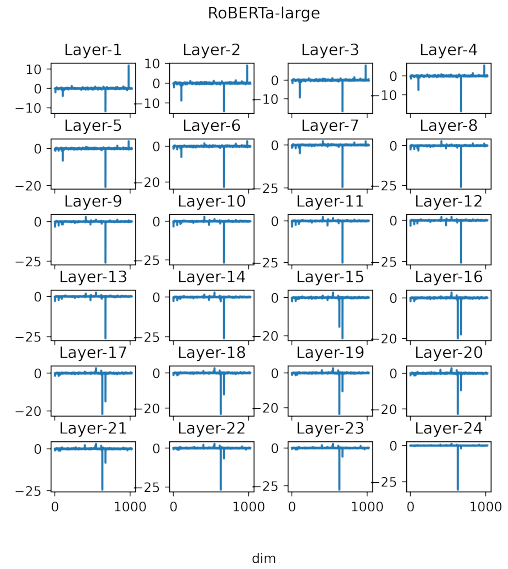


Figure 10: Average vectors for each layer of RoBERTa-large.

10 shows the patterns of RoBERTa-large across all layers. From the first layer to the twenty-third layer, the 673rd element is an outlier. From the fifteenth layer to the final layer, the 631st element is an outlier. From the first layer to the sixth layer, the 981st element is an outlier.

B Neuron-level analysis

B.1 Heatmaps of base models

Figure 11 and 12 show the heatmaps of the outlier neurons and the highest non-outlier contribution values.

B.2 Distilled and large models

Figure 13 show the accuracy scores of position prediction of distilled and large models.

Distil-models Figure 14 shows the contribution value of distilled models' outlier neurons on position prediction.

Large-models Figure 15 shows the contribution value of large models' outlier neurons on position prediction.

C Our Pre-training Models

C.1 Hyper-parameters

Table 4 shows the hyper-parameters of pre-training our RoBERTa-base models.

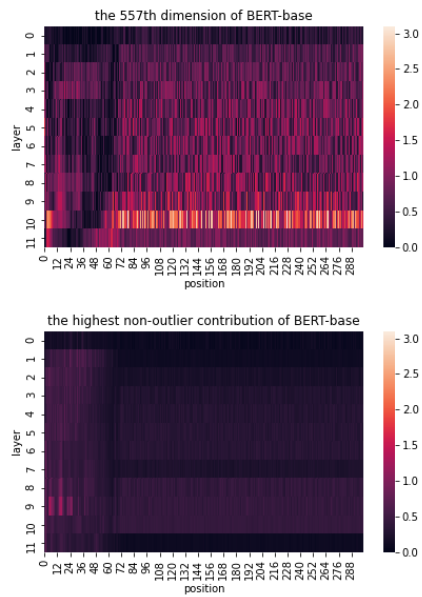


Figure 11: Up: contribution values heatmap of the outlier neuron of BERT-base. Down: the highest non-outlier contribution value of BERT-base.

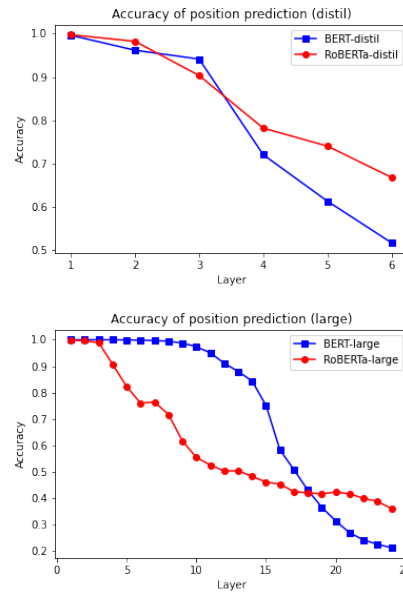


Figure 13: Up: accuracy of position prediction of distilled models. Down: accuracy of position prediction of large models.

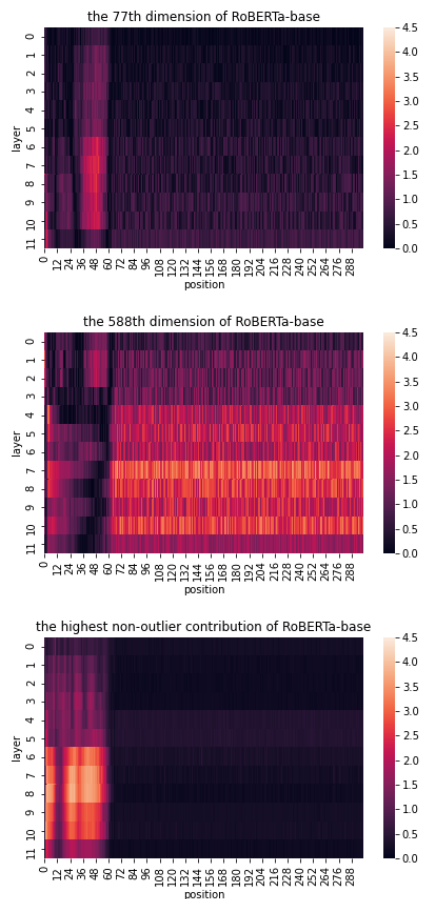


Figure 12: Up: contribution values heatmap of the 77th dimension of RoBERTa-base. Mid: contribution values heatmap of the 588th dimension of RoBERTa-base. Down: the highest non-outlier contribution value of RoBERTa-base.

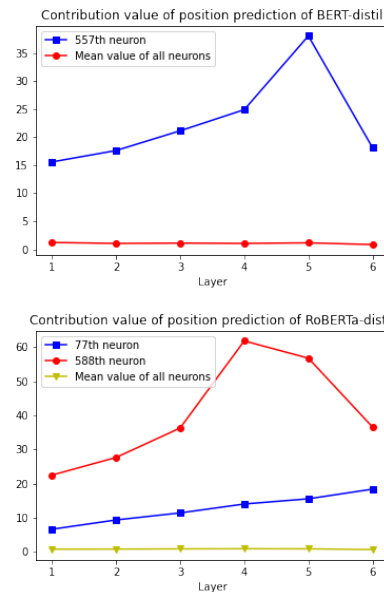


Figure 14: The contribution value of distilled models' outlier neurons on position prediction.

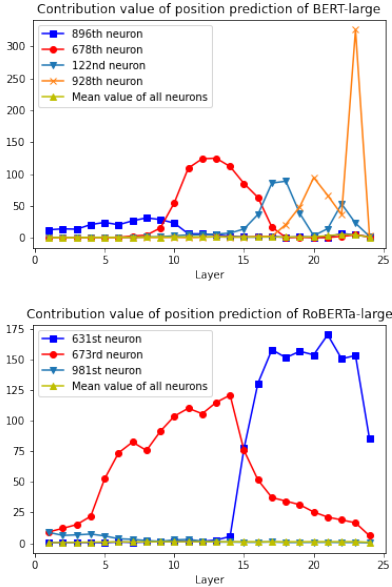


Figure 15: The contribution value of large models outlier neurons on position prediction.

Hyper-parameter	Our RoBERTa-base
Number of Layers	12
Hidden size	768
FNN inner hidden size	3072
Attention Heads	12
Attention Head size	64
Dropout	0.1
Warmup Steps	10k
Max Steps	200k
Learning Rates	1e-4
Batch Size	256
Weight Decay	0.01
Learning Rate Decay	Polynomial
Adam (ϵ , β_1 , β_2)	(1e-6, 0.9, 0.98)
Gradient Clipping	0.5

Table 4: Hyper-parameters for pre-training our RoBERTa-base models.

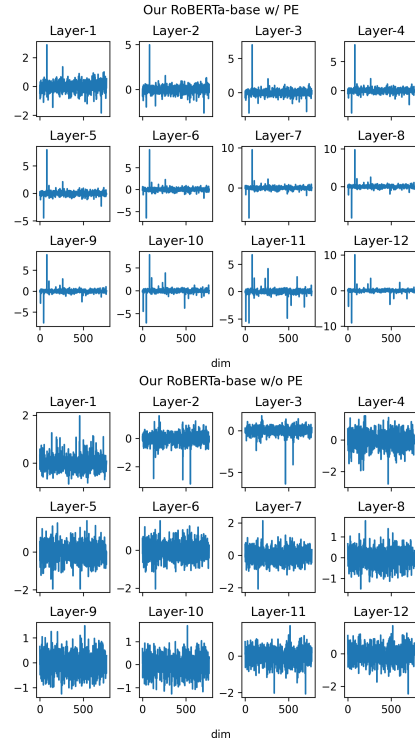


Figure 16: Average vectors for each layer of our RoBERTa-base w/ or w/o PE.

C.2 Average subword vectors

Figure 16 show the average vectors for each of our models.

D Clipping the outliers

D.1 Geometry of vector space

Distil-models Figure 17 shows the anisotropic measurement of distilled models and the self-similarity measurement of distilled models.

Large-models Figure 18 shows the anisotropic measurement of large models and Figure 19 shows the self-similarity measurement of large models. We “clip” different outlier neurons in different layers. For BERT-large, we zero-out the 896th neuron from the first layer to the tenth layer, the 678th neuron from the tenth layer to the seventeenth layer, the 122nd neuron from the sixteenth layer to the nineteenth layer and the 928th neuron from the nineteenth layer to the twenty-third layer. For RoBERTa-large, we zero-out the 673rd neuron for all non-input layers, the 981st neuron for the first 9 layers and the 631st neuron for the last 10 layers.

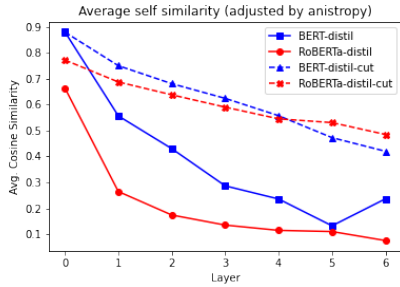
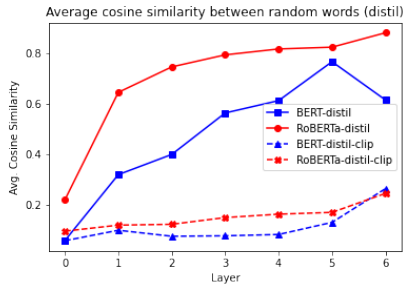


Figure 17: Up: average cosine similarity between random words of distil-models. Down: self-similarity measurement of BERT-distil and RoBERTa-distil (adjusted by anisotropy) before and after “clipping the outliers”.

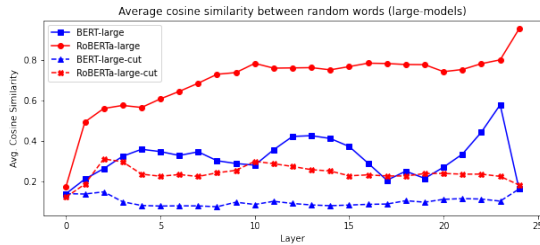


Figure 18: Average cosine similarity between random words of large-models.

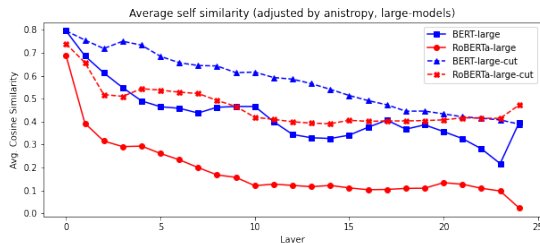


Figure 19: Self-similarity measurement of BERT-large and RoBERTa-large (adjusted by anisotropy) before and after “clipping the outliers”.

Model	Layer	Threshold	Acc.
Baseline	-	-	50.0%
Before clipping			
BERT-distil	5	0.9	66.5%
RoBERTa-distil	5	0.9	63.7%
BERT-large	12	0.7	70.2%
RoBERTa-large	10	0.9	70.4%
After clipping			
BERT-distil-clip	6	0.6	67.3%
RoBERTa-distil-clip	5	0.6	66.7%
BERT-large-clip	12	0.6	70.3%
RoBERTa-large-clip	16	0.6	71.3%

Table 5: The best accuracy scores on WiC dataset for distilled and large models. **Bold** indicates that the best result increases after clipping.

Dataset	BERT distil	RoBERTa distil	BERT distil clip	RoBERTa distil clip
STS-B	59.65(6)	56.06(5)	56.62(6)	58.47(5)
SICK-R	62.64(6)	62.63(5)	62.42(6)	62.73(6)
STS-12	42.96(1)	40.19(1)	46.47(1)	42.36(1)
STS-13	59.33(1)	56.42(5)	55.74(1)	60.64(6)
STS-14	53.81(6)	49.59(6)	50.57(1)	52.51(2)
STS-15	61.40(6)	65.10(5)	61.48(1)	65.93(2)
STS-16	61.43(6)	62.90(5)	60.75(6)	64.49(5)

Table 6: Experimental results on semantic textual similarity of distilled models. The number in the parenthesis denotes that this result belongs to the specific layer. **Bold** indicates that the best result increases after clipping.

D.2 Word sense

Table 5 shows the accuracy scores of distill-models and large-models on WiC dataset before and after “clipping the outliers”.

D.3 Sentence embedding

Table 6 shows the results on semantic textual similarity tasks of distilled models before and after “clipping the outliers”.

Table 7 shows the results on semantic textual similarity tasks of large models before and after “clipping the outliers”.

Dataset	BERT large	RoBERTa large	BERT large clip	RoBERTa large clip
STS-B	62.56(1)	59.71(19)	66.43(3)	62.01(23)
SICK-R	64.47(24)	63.08(14)	65.72(23)	63.50(16)
STS-12	54.05(1)	44.72(1)	56.44(3)	49.69(1)
STS-13	68.80(2)	61.68(8)	71.07(2)	62.82(10)
STS-14	60.46(1)	51.39(8)	63.35(1)	57.33(1)
STS-15	73.91(1)	65.98(7)	76.51(1)	69.71(1)
STS-16	66.35(17)	66.50(14)	71.41(3)	68.25(11)

Table 7: Experimental results on semantic textual similarity of large models. The number in the parenthesis denotes that this result belongs to the specific layer. **Bold** indicates that the best result increases after clipping.