

# BERT Has More to Offer: BERT Layers Combination Yields Better Sentence Embeddings

MohammadSaleh Hosseini\*, Munawara Saiyara Munia†, Latifur Khan\*

\*Department of Computer Science

†Department of Electrical and Computer Engineering

The University of Texas at Dallas

{seyyedmohammadsaleh.hosseini, munawarasaiyara.munia, lkhan}@utdallas.edu

## Abstract

Obtaining sentence representations from BERT-based models as feature extractors is invaluable as it takes much less time to pre-compute a one-time representation of the data and then use it for the downstream tasks, rather than fine-tune the whole BERT. Most previous works acquire a sentence’s representation by passing it to BERT and averaging its last layer. In this paper, we propose that the combination of certain layers of a BERT-based model rested on the data set and model can achieve substantially better results. We empirically show the effectiveness of our method for different BERT-based models on different tasks and data sets. Specifically, on seven standard semantic textual similarity data sets, we outperform the baseline BERT by improving the Spearman’s correlation by up to 25.75% and on average 16.32% without any further training. We also achieved state-of-the-art results on eight transfer data sets by reducing the relative error by up to 37.41% and on average 17.92%.<sup>1</sup>

## 1 Introduction

Learning sentence vector representations is a crucial problem in natural language processing (NLP) and has been widely studied in the literature (Conneau et al., 2017; Cer et al., 2018; Li et al., 2020). Given a sentence, the goal is to acquire a vector that semantically and/or syntactically represents it. BERT (Devlin et al., 2019) has set new state-of-the-art records on many NLP tasks (Madabushi et al., 2020; Hu et al., 2022; Skorupa Parolin et al., 2022; Wang and Kuo, 2020). However, this is achieved by fine-tuning all of BERT’s layers. The disadvantage of fine-tuning is that it is computationally expensive as even BERT-base has 110M parameters; hence, pre-computing a representation of the data and using it for the downstream task is much less computationally expensive (Devlin et al., 2019).

<sup>1</sup>Our code is available at: <https://github.com/DiamondRock/BERT-Layers-Combination>

Moreover, for sentence-pair tasks, BERT uses a cross-encoder; nonetheless, this setup is inappropriate for certain pair regression tasks, such as finding the most similar sentence in a data set to a specific sentence due to the large number of possible combinations (Reimers and Gurevych, 2019).

Considering the aforementioned drawbacks, researchers have tried to derive fixed-sized sentence embeddings from BERT or proposed new BERTs with the exact same architecture but different ways of training (Reimers and Gurevych, 2019; Li et al., 2020). After training, the resultant BERT is used in a feature-based manner by passing the sentence to it and obtaining its embedding vector in different ways, such as averaging the last layer of BERT.

In this paper, we present a simple, yet effective and novel method called **BERT-LC** (BERT Layers Combination). BERT-LC combines certain layers of BERT in order to obtain the representation of a sentence. As we will show, this model significantly outperforms its correspondent BERT baseline with no need of any further training. Our work was inspired by Jawahar et al. (2019), who show that different layers of BERT carry different features, such as surface, syntactic, and semantic. We argue that each data set with its unique distribution might need a different set of features for its sentences, which can only be fully exploited by combining different layers of BERT in an unsupervised way.

Our contributions are as follows: (1) We propose a new method called BERT-LC that is capable of acquiring superior results by combining certain layers of BERT instead of just the last layer, in an unsupervised manner. We also include the embedding layer, which was to our knowledge ignored in previous works. (2) We additionally show that our method improves SBERT (Reimers and Gurevych, 2019) and SimCSE (Gao et al., 2021), which were specifically designed for obtaining sentence representations (opposite to BERT and RoBERTa (Liu et al., 2019)). (3) We developed an algorithm that

speeds up the process of finding the best layer combination (among  $2^{13}$  layer combinations in base cases) by a factor of 189 times. (4) We propose an innovative method that integrates the layer combination method with the CLS pooling head, improving the performance metrics for certain models. (5) We achieve state-of-the-art performance on the transfer tasks using layer combination.

We demonstrate the superiority of our approach through conducting extensive experiments on seven standard semantic textual similarity (STS) data sets and eight transfer tasks. On the STS data sets, our method is able to outperform its corresponding baseline by up to 25.75% and on average 16.32% for BERT-large-uncased. We also achieve the state-of-the-art performances on transfer tasks, reducing the previous best model’s relative error rate by an average of 17.92% and up to 37.41%.

## 2 Related Work

Learning sentence embeddings is a well-studied realm in NLP. There are mainly two methods used for this purpose: methods that use unlabeled data or labeled data. Although the latter use labeled data, the target data sets and tasks on which they are tested are different from the training data set and task. Early work on sentence embedding utilized the distributional hypothesis by predicting surrounding sentences of a sentence (Kiros et al., 2015; Hill et al., 2016; Logeswaran and Lee, 2018). Pagliardini et al. 2018 build on the idea of word2vec (Mikolov et al., 2013) using n-gram embeddings. Some researchers simply use the average of BERT’s last layer embeddings as the sentence embedding (Reimers and Gurevych, 2019). Recently, contrastive learning has proven to be very powerful in many domains (Khorram et al., 2022; Munia et al., 2021; Hu et al., 2021); thus, some recent methods exploit contrastive learning (Gao et al., 2021; Zhang et al., 2022; Yan et al., 2021) and utilize the same sentence by looking at it from multiple angles. For instance, the popular method SimCSE leverages different outputs of the same sentence from BERT’s standard dropout.

There are some previous works (Ethayarajh, 2019; Jawahar et al., 2019; Bommasani et al., 2020) that investigate the impact of different BERT layers. However, they either investigate each layer individually or are restricted to considering a set of consecutive layers strictly starting from layer 1.

We apply our method to different variations of

BERT, RoBERTa, SBERT, and SimCSE. To the best of our knowledge, no similar work has been done that whether combines different and arbitrary layers of these models or further integrates layer combination and the CLS pooling head.

## 3 Approach

Given an input sentence  $S = (s_1, s_2, \dots, s_N)$ , the goal of a sentence embedding model is to output a vector  $E_S \in \mathbb{R}^d$  which carries the semantic and/or syntactic information of the sentence. In order to obtain a sentence embedding, we first pass the sentence to a BERT-like model, which outputs the tensor  $H \in \mathbb{R}^{L' \times N \times d}$ , in which  $d$  is the dimension of the token vectors in each layer, and  $L' = L + 1$  is the number of layers (including Layer 0) in BERT. We then apply to this tensor a pooling function  $p$ , which can be *max* or *mean*, but we choose mean as it yielded better results in our experiments. The pooling is done across all the tokens in all the desired layers set,  $D$ . For example, for mean pooling,  $p$  is defined as:

$$p(H, D) = \frac{1}{|D|} \sum_{l \in D} \frac{1}{N} \sum_{n=1}^N H_{l,n,:} \quad (1)$$

where  $H_{l,n,:} \in \mathbb{R}^d$  is the Transformer vector at layer number  $l$  corresponding to the  $n$ th token.

Previous works usually set  $D$  to  $\{L\}$  or use CLS pooler, the output of the MLP layer attached to the last layer’s first token ([CLS]) to obtain the sentence embedding. Using CLS pooler was shown to underperform last layer averaging (Reimers and Gurevych, 2019). Further, as we will empirically show, choosing  $D = \{L\}$  leads to significant underperformance as well; hence, in this work, we iterate through all possible  $D$ s ( $D \in \mathcal{P}(A) - \emptyset$ , where  $A = \{0, \dots, L\}$ ), and choose the best-performing  $D$  as our layers to pool from.

As iterating through all possible  $D$ s is very time-consuming, we propose an algorithm that speeds up the process of finding the best layer combination ( $2^{13}$  layer combinations in base cases) by a factor of 189 times, which can be found in Appendix A.

We further propose an extension to our method, and that is exploiting the MLP head and layer combination simultaneously. The idea is to pass  $p(H, D)$  in Eq. 1 to the MLP head and use the output of the MLP head as the new  $p$ . We spot that this method works better than merely using the layer combination on SimCSE. We conjecture that this is because SimCSE’s MLP head was trained for learn-

ing sentence embeddings as opposed to the other methods, such as BERT, SBERT, and RoBERTa; hence, it carries important information to be utilized. Consequently, we use the two-step pipeline of layer combination and MLP for SimCSE models, and we propose that the two-step pipeline is utilized for any other BERT-based sentence embedding model whose MLP head has been trained for sentence embedding.

## 4 Experiments and Results

We carry out our experiments on two different tasks: transfer and STS tasks. In this section, we discuss the transfer tasks, while the STS tasks are discussed in Appendix B.

### 4.1 Data Sets and Evaluation Setup

For transfer tasks, we use eight data sets from the popular SentEval toolkit (Conneau and Kiela, 2018), which is used for assessing the quality of sentence embeddings: MR (Pang and Lee, 2005), CR (Hu and Liu, 2004), SUBJ (Pang and Lee, 2004), MPQA (Wiebe et al., 2005), SST and SSTM (binary and six-class Stanford Sentiment Treebanks)(Socher et al., 2013), TREC (Li and Roth, 2002), and MRPC (Dolan et al., 2004).

For each data set, we combine all of its data subsets and randomly split them into training-development (train-dev) and test data with ratios 85% and 15%. This *outer* cross-validation is done randomly for 10 times, and the average accuracy results on the test data are reported. We further utilize a 10-fold *inner* cross-validation on train-dev data set, with ratios 82% (train) and 18% (dev).

The sentence embeddings by our method or the baselines are used as feature vectors for a logistic regression classifier. Note that even though there exists a training data set, it is only utilized by the logistic regression classifier and not in our method.

### 4.2 Baselines and Hyper-parameters

We use the following baselines and apply our layer combination method to them: BERT base/large (un)cased (**BERT-(B/L)<sub>(un)cased</sub>**), RoBERTa base/large (**RoBERTa-B/L**), SBERT and SRoBERTa base/large (**SBERT/SRoBERTa-B/L**), (un)supervised SimCSE BERT/RoBERTa base/large ((Un)SupSimCSE-(B/R)(B/L)), and **SupSimCSE-RB/RL<sub>M</sub>**, which are SimCSE models trained with an additional MLM head.

For the base and large models, we consider layer

combinations with up to four and three layers, respectively, as we did not spot much difference when combining more layers.

For the logistic regression, we use SAGA (Defazio et al., 2014) as the optimizer, 0.01 as the tolerance level, 200 as the maximum number of iterations, and 10 as the regularization parameter.

### 4.3 Results and Discussion

The results for our proposed method and the baselines on transfer tasks are shown in Table 1. The results are statistically significant (p-value < 0.01). Note that our methods are denoted by an **LC** at the end of them (e.g., RoBERTa-L-LC means RoBERTa-L when layer combination is used).

Since all models have accuracies close to 100% on most of the transfer data sets, we believe it is more reasonable to see how much improvement we are acquiring by considering the relative error reduction. Relative error for a model and its baseline with respective accuracies of  $A_M$  and  $A_{MB}$  is defined as  $\frac{A_M - A_{MB}}{1 - A_{MB}}$ . From Table 1, we observe that our method significantly outperforms its corresponding baselines for all the models on all the data sets. For example, RoBERTa-B-LC, SBERT-L-LC, and UnSupSimCSE-RL-LC improve on their baselines by reducing the relative errors by up to 36.19%, 45.93%, and 40.80% and on average 19.89%, 14.95%, and 27.31%, respectively.

Furthermore, we spot that our method is more effective on unsupervised versions of SimCSE models. For instance, the relative error improvement over UnSupSimCSE-RL is 27.31%, whereas it is 14.61% for SupSimCSE-RL. However, there is still important information to be extracted from all the layers of the supervised SimCSE models as well.

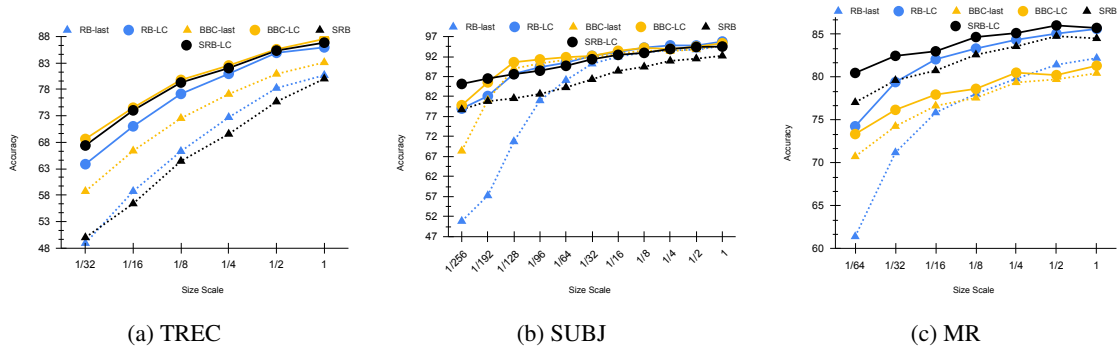
We see that even an unsupervised model such as RoBERTa-L when upgraded with LC can achieve 84.58%, beating SupSimCSE-RL, a supervised model which held the previous best average result (83.73%). RoBERTa-L-LC’s average accuracy does not fall very short of the new state-of-the-art model, SupSimCSE-RL-LC (84.68% vs. 85.89%). Finally, we have improved the previous state-of-the-art models on all of the transfer data sets by reducing the relative (absolute) errors by up to 37.41% (3.67%) and on average 17.92% (1.90%), respectively. We achieve a state-of-the-art average accuracy of 86.36% on the transfer data sets. We have provided more baselines, which can be found in Appendix D.

	MR	CR	SUBJ	MPQA	SSTM	TREC	MRPC	SST	Avg.		MR	CR	SUBJ	MPQA	SSTM	TREC	MRPC	SST	Avg.
BERT-B <sub>uncased</sub> -CLS	80.16	83.17	93.97	84.35	46.66	74.74	70.87	85.01	77.37	SBERT-B	82.90	88.96	93.93	89.58	47.53	80.04	74.34	89.19	80.81
BERT-B <sub>uncased</sub> -last avg	81.19	86.17	95.07	88.10	47.53	85.76	73.66	87.08	80.57	SBERT-B-LC	<b>83.61</b>	<b>89.95</b>	<b>95.17</b>	<b>91.06</b>	<b>49.01</b>	<b>88.56</b>	<b>78.71</b>	<b>89.65</b>	<b>83.22</b>
BERT-B <sub>uncased</sub> -LC	<b>82.24</b>	<b>86.60</b>	<b>95.40</b>	<b>90.42</b>	<b>49.22</b>	<b>88.76</b>	<b>77.17</b>	<b>88.17</b>	<b>82.25</b>	SRoBERTa-B	84.67	90.12	92.57	89.21	50.59	81.79	77.13	90.12	82.03
BERT-B <sub>cased</sub> -CLS	77.88	83.07	92.33	85.29	45.19	69.97	70.64	83.45	75.98	SRoBERTa-B-LC	<b>85.76</b>	<b>91.75</b>	<b>94.80</b>	<b>90.51</b>	<b>53.65</b>	<b>87.95</b>	<b>78.92</b>	<b>90.78</b>	<b>84.26</b>
BERT-B <sub>cased</sub> -last avg	80.55	85.19	94.64	87.63	46.48	83.65	74.30	85.83	79.78	SupSimCSE-BB	<b>81.85</b>	<b>89.31</b>	94.60	89.73	<b>50.16</b>	82.75	74.60	88.81	81.48
BERT-B <sub>cased</sub> -LC	<b>81.42</b>	<b>86.56</b>	<b>95.11</b>	<b>89.88</b>	<b>47.76</b>	<b>88.00</b>	<b>75.36</b>	<b>86.86</b>	<b>81.37</b>	SupSimCSE-BB-LC	<b>81.85</b>	<b>89.31</b>	<b>95.68</b>	<b>90.87</b>	<b>50.16</b>	<b>89.09</b>	<b>77.68</b>	<b>89.67</b>	<b>83.04</b>
RoBERTa-B-last avg	82.58	84.69	94.55	85.83	50.26	81.90	72.51	87.28	79.95	SupSimCSE-RB	84.29	91.50	93.12	90.19	52.69	81.12	76.14	90.14	82.40
RoBERTa-B-LC	<b>85.89</b>	<b>90.23</b>	<b>95.68</b>	<b>89.24</b>	<b>52.39</b>	<b>86.92</b>	<b>74.85</b>	<b>89.65</b>	<b>83.11</b>	SupSimCSE-RB-LC	<b>86.00</b>	<b>92.38</b>	<b>95.21</b>	<b>90.89</b>	<b>54.00</b>	<b>87.97</b>	<b>78.76</b>	<b>90.97</b>	<b>84.52</b>
UnSupSimCSE-BB	80.67	85.29	94.29	88.75	45.14	83.96	72.94	85.86	79.61	SupSimCSE-RB <sub>M</sub>	84.56	91.89	93.29	89.49	52.28	81.88	76.23	90.13	82.47
UnSupSimCSE-BB-LC	<b>81.61</b>	<b>86.77</b>	<b>95.16</b>	<b>90.03</b>	<b>47.07</b>	<b>89.05</b>	<b>77.29</b>	<b>87.49</b>	<b>81.81</b>	SupSimCSE-RB <sub>M</sub> -LC	<b>86.11</b>	<b>92.38</b>	<b>95.81</b>	<b>90.63</b>	<b>54.10</b>	<b>88.40</b>	<b>78.14</b>	<b>91.00</b>	<b>84.57</b>
UnSupSimCSE-RB	80.85	85.36	92.40	86.70	47.25	76.13	72.74	85.67	78.39	SBERT-L	84.69	90.62	94.40	90.25	49.24	79.71	74.99	90.92	81.85
UnSupSimCSE-RB-LC	<b>84.00</b>	<b>88.99</b>	<b>94.69</b>	<b>88.84</b>	<b>50.50</b>	<b>86.02</b>	<b>77.59</b>	<b>88.33</b>	<b>82.37</b>	SBERT-L-LC	<b>85.41</b>	<b>91.64</b>	<b>95.61</b>	<b>91.15</b>	<b>51.84</b>	<b>89.03</b>	<b>79.03</b>	<b>91.46</b>	<b>84.40</b>
BERT-L <sub>uncased</sub> -CLS	83.22	85.64	93.03	80.75	46.73	69.45	69.10	84.01	76.49	SRoBERTa-L	86.85	90.83	93.16	90.69	50.82	83.14	77.36	92.44	83.16
BERT-L <sub>uncased</sub> -last avg	83.88	88.29	95.52	86.51	49.88	83.81	71.49	88.48	80.96	SRoBERTa-L-LC	<b>87.95</b>	<b>92.49</b>	<b>95.35</b>	<b>92.06</b>	<b>54.23</b>	<b>88.51</b>	<b>79.56</b>	<b>92.94</b>	<b>85.39</b>
BERT-L <sub>uncased</sub> -LC	<b>85.21</b>	<b>89.88</b>	<b>96.08</b>	<b>90.17</b>	<b>51.06</b>	<b>88.89</b>	<b>76.57</b>	<b>89.96</b>	<b>83.48</b>	SupSimCSE-BL	85.47	<b>90.69</b>	95.01	90.38	51.16	84.28	73.70	90.83	82.69
BERT-L <sub>cased</sub> -CLS	81.65	82.79	91.43	83.49	45.76	64.23	70.02	84.83	75.53	SupSimCSE-BL-LC	<b>85.60</b>	<b>90.69</b>	<b>95.88</b>	<b>91.30</b>	<b>52.20</b>	<b>89.47</b>	<b>77.01</b>	<b>91.39</b>	<b>84.19</b>
BERT-L <sub>cased</sub> -last avg	84.28	88.68	94.95	88.03	49.23	84.10	72.62	88.75	81.33	SupSimCSE-RL	88.00	<b>90.69</b>	94.80	90.86	51.89	86.52	74.21	92.84	83.73
BERT-L <sub>cased</sub> -LC	<b>85.34</b>	<b>90.19</b>	<b>95.41</b>	<b>90.50</b>	<b>50.92</b>	<b>88.62</b>	<b>77.17</b>	<b>90.00</b>	<b>83.52</b>	SupSimCSE-RL-LC	<b>89.24</b>	<b>90.69</b>	<b>96.29</b>	<b>92.01</b>	<b>55.39</b>	<b>90.19</b>	<b>79.82</b>	<b>93.49</b>	<b>85.89</b>
RoBERTa-L-last avg	84.30	85.22	94.93	87.25	50.59	81.75	67.24	89.41	80.09	SupSimCSE-RL <sub>M</sub>	87.78	92.10	94.72	90.51	52.43	83.85	74.39	92.60	83.55
RoBERTa-L-LC	<b>88.04</b>	<b>91.68</b>	<b>96.51</b>	<b>91.11</b>	<b>54.11</b>	<b>88.06</b>	<b>75.01</b>	<b>92.08</b>	<b>84.58</b>	SupSimCSE-RL <sub>M</sub> -LC	<b>89.09</b>	<b>93.44</b>	<b>96.41</b>	<b>91.43</b>	<b>56.10</b>	<b>88.06</b>	<b>78.97</b>	<b>93.44</b>	<b>85.87</b>
UnSupSimCSE-BL	<b>84.85</b>	88.15	95.09	89.13	48.84	83.63	74.09	89.02	81.60										
UnSupSimCSE-BL-LC	<b>84.85</b>	<b>89.84</b>	<b>95.77</b>	<b>90.62</b>	<b>50.93</b>	<b>89.07</b>	<b>77.10</b>	<b>90.00</b>	<b>83.53</b>										
UnSupSimCSE-RL	82.29	86.24	92.77	88.12	45.98	82.02	73.91	88.08	79.93										
UnSupSimCSE-RL-LC	<b>86.75</b>	<b>91.01</b>	<b>95.72</b>	<b>90.85</b>	<b>52.83</b>	<b>87.93</b>	<b>79.13</b>	<b>91.57</b>	<b>84.47</b>										

(a) Unsupervised Models

(b) Supervised Models

Table 1: Transfer task results for different sentence embedding models (measured as accuracy \* 100)



(a) TREC

(b) SUBJ

(c) MR

Figure 1: Varying the size of the training data for transfer tasks.

## 5 Ablation Studies

While more ablation studies, such as the effect of the number of layers in layer combination and the effect of excluding a particular layer from layer combination, can be found in Appendix E, in this section, we discuss the performance of layer combination in limited data settings. For this, we reduce the size of training data for MR, SUBJ, and TREC, while keeping the test data untouched. Fig. 1 shows the results for RoBERTa-B (RB), BERT-B<sub>cased</sub> (BBC) and SRoBERTa-B (SRB) on these data sets. We report the average results of 10 different seeds. As we can see, for all the models on all the data sets, the more reduction in size is done, the more the difference is between the LC version and the last layer average. For instance, for RB on SUBJ, on the original training data, the accuracies for RB-LC and RB are 94.77% and 93.73%, respectively, while when reducing the size to 1/256th, the respective accuracies are 78.99% and 50.83%. This suggests that it is even more beneficial to use LC when one has small training data.

## 6 Conclusion and Future Work

In this paper, we proposed a new method called BERT Layer Combination, a simple, yet effective framework, which when applied to various BERT-based models, significantly improves them for the downstream tasks of STS and transfer learning. Further, it achieves the state-of-the-art performances on eight transfer tasks. Our method combines certain layers of BERT-based models in an unsupervised manner, which shows that different layers of BERT hold important information which was previously ignored. We demonstrated the effectiveness of our approach by conducting comprehensive experiments on various BERT-based models and on a host of different tasks and data sets.

As future work, we would like to apply the technique of layer combination to other NLP tasks (e.g., punctuation insertion (Hosseini and Sameti, 2017) and question answering (Qu et al., 2019)) and also utilize it in other domains where deep learning models are used (e.g., biosignal analysis (Munia et al., 2020, Munia et al., 2023) and image captioning



(Huang et al., 2019)).

## 7 Limitation

One limitation of our work is that though we applied our layer combination technique to a raft of different models, all of them are BERT based; ergo, whether or not this technique works on other deep learning NLP models remains unsettled. We leave this experimentation as future work.

## 8 Acknowledgement

This research was supported in part by NSF awards DMS-1737978, DGE-2039542, OAC-1828467, OAC-1931541, and DGE-1906630. The material presented here is partially based on High Performance Computing (HPC) resources supported by the University of Arizona TRIF, UITS, and Research, Innovation, and Impact (RII) and maintained by the UArizona Research Technologies department.

## 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.
- Rishi Bommasani, Kelly Davis, and Claire Cardie. 2020. *Interpreting Pretrained Contextualized Representations via Reductions to Static Embeddings*. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4758–4781, Online. Association for Computational Linguistics.
- 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.
- Daniel Cer, Yinfei Yang, Sheng-yi Kong, Nan Hua, Nicole Limtiaco, Rhomni St. John, Noah Constant, Mario Guajardo-Cespedes, Steve Yuan, Chris Tar, Brian Strope, and Ray Kurzweil. 2018. *Universal sentence encoder for English*. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 169–174, Brussels, Belgium. Association for Computational Linguistics.
- 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).
- Alexis Conneau, Douwe Kiela, Holger Schwenk, Loïc Barrault, and Antoine Bordes. 2017. *Supervised learning of universal sentence representations from natural language inference data*. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 670–680, Copenhagen, Denmark. Association for Computational Linguistics.
- Aaron Defazio, Francis Bach, and Simon Lacoste-Julien. 2014. *Saga: A fast incremental gradient method with support for non-strongly convex composite objectives*. In *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc.
- 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.
- Bill Dolan, Chris Quirk, and Chris Brockett. 2004. [Un-supervised construction of large paraphrase corpora: Exploiting massively parallel news sources](#). In *COLING 2004: Proceedings of the 20th International Conference on Computational Linguistics*, pages 350–356, Geneva, Switzerland. COLING.
- 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.
- Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. [SimCSE: Simple contrastive learning of sentence embeddings](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6894–6910, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Felix Hill, Kyunghyun Cho, and Anna Korhonen. 2016. [Learning distributed representations of sentences from unlabelled data](#). In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1367–1377, San Diego, California. Association for Computational Linguistics.
- Seyyed MohammadSaleh Hosseini and Hossein Sameti. 2017. [Creating a corpus for automatic punctuation prediction in persian texts](#). In *2017 Iranian Conference on Electrical Engineering (ICEE)*, pages 1537–1542.
- Minqing Hu and Bing Liu. 2004. [Mining and summarizing customer reviews](#). In *Proceedings of the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '04*, page 168–177, New York, NY, USA. Association for Computing Machinery.
- Xiang Hu, Teng Li, Tong Zhou, Yu Liu, and Yuanxi Peng. 2021. [Contrastive learning based on transformer for hyperspectral image classification](#). *Applied Sciences*, 11(18):8670.
- Yibo Hu, MohammadSaleh Hosseini, Erick Sko-rupa Parolin, Javier Osorio, Latifur Khan, Patrick Brandt, and Vito D’Orazio. 2022. [ConflBERT: A pre-trained language model for political conflict and violence](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5469–5482, Seattle, United States. Association for Computational Linguistics.
- Lun Huang, Wenmin Wang, Jie Chen, and Xiao-Yong Wei. 2019. [Attention on attention for image captioning](#). In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4634–4643.
- Ganesh Jawahar, Benoît Sagot, and Djamé Seddah. 2019. [What does BERT learn about the structure of language?](#) In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3651–3657, Florence, Italy. Association for Computational Linguistics.
- Soheil Khorram, Jaeyoung Kim, Anshuman Tripathi, Han Lu, Qian Zhang, and Hasim Sak. 2022. [Contrastive siamese network for semi-supervised speech recognition](#). In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 7207–7211. IEEE.
- Ryan Kiros, Yukun Zhu, Russ R Salakhutdinov, Richard Zemel, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. [Skip-thought vectors](#). In *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc.
- Bohan Li, Hao Zhou, Junxian He, Mingxuan Wang, Yiming Yang, and Lei Li. 2020. [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.
- Xin Li and Dan Roth. 2002. [Learning question classifiers](#). In *Proceedings of the 19th International Conference on Computational Linguistics - Volume 1, COLING '02*, page 1–7, USA. 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](#).
- Lajanugen Logeswaran and Honglak Lee. 2018. [An efficient framework for learning sentence representations](#). *CoRR*, abs/1803.02893.
- Harish Tayyar Madabushi, Elena Kochkina, and Michael Castelle. 2020. [Cost-sensitive bert for generalisable sentence classification with imbalanced data](#). *arXiv preprint arXiv:2003.11563*.
- Marco Marelli, Stefano Menini, Marco Baroni, Luisa Bentivogli, Raffaella Bernardi, and Roberto Zamparelli. 2014. [A SICK cure for the evaluation of compositional distributional semantic models](#). In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, pages 216–223, Reykjavik, Iceland. European Language Resources Association (ELRA).

- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. [Distributed representations of words and phrases and their compositionality](#). In *Advances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc.
- Munawara Saiyara Munia, Seyyed MohammadSaleh Hosseini, Mehrdad Nourani, Jay Harvey, and Hina Dave. 2021. [Imbalanced EEG analysis using one-shot learning with siamese neural network](#). IEEE.
- Munawara Saiyara Munia, Mehrdad Nourani, Jay Harvey, and Hina Dave. 2023. Interictal epileptiform discharge detection using multi-head deep convolutional neural network. In *2023 45th Annual International Conference of the IEEE Engineering in Medicine Biology Society (EMBC)*.
- Munawara Saiyara Munia, Mehrdad Nourani, and Sammy Houari. 2020. [Biosignal oversampling using wasserstein generative adversarial network](#). In *2020 IEEE International Conference on Healthcare Informatics (ICHI)*, pages 1–7.
- Matteo Pagliardini, Prakhar Gupta, and Martin Jaggi. 2018. [Unsupervised learning of sentence embeddings using compositional n-gram features](#). 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 528–540, New Orleans, Louisiana. Association for Computational Linguistics.
- Bo Pang and Lillian Lee. 2004. [A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts](#). In *Proceedings of the 42nd Annual Meeting on Association for Computational Linguistics, ACL '04*, page 271–es, USA. Association for Computational Linguistics.
- Bo Pang and Lillian Lee. 2005. [Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales](#). In *Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL'05)*, pages 115–124, Ann Arbor, Michigan. Association for Computational Linguistics.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. [Pytorch: An imperative style, high-performance deep learning library](#). In *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc.
- Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. 2011. Scikit-learn: Machine learning in python. *Journal of machine learning research*, 12(Oct):2825–2830.
- Chen Qu, Liu Yang, Minghui Qiu, W. Bruce Croft, Yongfeng Zhang, and Mohit Iyyer. 2019. [Bert with history answer embedding for conversational question answering](#). In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR'19*, page 1133–1136, New York, NY, USA. Association for Computing Machinery.
- 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.
- Erick Skorupa Parolin, MohammadSaleh Hosseini, Yibo Hu, Latifur Khan, Patrick T. Brandt, Javier Osorio, and Vito D'Orazio. 2022. [Multi-coped: A multilingual multi-task approach for coding political event data on conflict and mediation domain](#). In *Proceedings of the 2022 AAAI/ACM Conference on AI, Ethics, and Society, AIES '22*, page 700–711, New York, NY, USA. Association for Computing Machinery.
- 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.
- Bin Wang and C-C Jay Kuo. 2020. Sbert-wk: A sentence embedding method by dissecting bert-based word models. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 28:2146–2157.
- Janyce Wiebe, Theresa Wilson, and Claire Cardie. 2005. Annotating expressions of opinions and emotions in language. *Language Resources and Evaluation*, 39:165–210.
- Yuanmeng Yan, Rumei Li, Sirui Wang, Fuzheng Zhang, Wei Wu, and Weiran Xu. 2021. [ConSERT: A contrastive framework for self-supervised sentence representation transfer](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 5065–5075, Online. Association for Computational Linguistics.
- Yuhao Zhang, Hongji Zhu, Yongliang Wang, Nan Xu, Xiaobo Li, and Binqiang Zhao. 2022. [A contrastive framework for learning sentence representations from pairwise and triple-wise perspective in angular space](#). In *Proceedings of the 60th Annual Meeting of the*



## A Optimized Algorithm for Finding the Best Layer Combination

For a data set  $X = \{(x_1, y_1), \dots, (x_{|X|}, y_{|X|})\}$  where each sample consists of a list of sentences  $x_i = \{S_i^1, \dots, S_i^n\}$  ( $n=2$  for sentence-pair tasks such as STS) and a label  $y_i$ , the best-performing layers set ( $D^*$ ) is calculated by the following equation:

$$D^* = \operatorname{argmax}_{D \in \mathcal{P}^+(A)} m((f_1, \dots, f_{|X|}), (y_1, \dots, y_{|X|})) \quad (2)$$

where  $m$  is the desired metric function, and  $f_i$  is  $f_i = f(p(H_i^1, D), \dots, p(H_i^n, D))$  (3) in which  $H_j^i$  is the tensor of the  $i$ th sample’s  $j$ th sentence, and  $f$  is a function, such as cosine similarity.

Eq. 3 requires calculating  $p(H, D)$  for each sentence in  $X$  and for every possible  $D$ , which can be very time-consuming. In this subsection, we discuss our proposed algorithm for efficiently calculating all  $p(H, D)$ s in a step-by-step manner. Each method improves over the previous one, with Method 1 being the naive algorithm. The algorithms discussed here utilize *mean* as the pooling function, yet they can be easily modified to use other functions such as *max*.

**Method 1:** For every  $D \in \mathcal{P}^+(A) = \mathcal{P}(A) - \emptyset$ , calculate  $p(H, D)$ .

**Method 2:** For every  $D \in \{\{i\} | 0 \leq i \leq L\}$  calculate  $p_l = p(H, D) = \frac{1}{N} \sum_{n=1}^N H_{l,n}$ . Then, for every other  $D \in \mathcal{P}^+(A)$ , calculate  $p(H, D)$  as  $p(H, D) = \frac{1}{|D|} \sum_{l \in D} p_l$ .

**Method 3:** This method is explained in Algos. 1 and 2 (*maxOptim=False*). The idea is that for calculating every  $p(H, D)$ , we can exploit other  $P(H, D')$ s that we have calculated before. For instance,  $p(H, \{1, 3, 6, 8, 10\}) = \frac{3 * p(H, \{1, 3, 6\}) + 2 * p(H, \{8, 10\})}{5}$ . This method uses bottom-up dynamic programming to store the previous  $p(H, D)$ s. *powerset(A, i)* in Algo. 1 is equal to  $\{D | D \in \mathcal{P}(A) \wedge |D| = i\}$ . *mem* is a hash map with key-value pairs as  $(D, p(H, D))$ . *greedyPart(D, mem)* returns an array of partitions of  $D$  such that all partitions are present in *mem*, and each partition has the highest length possible. This is done in a greedy manner, by first

dividing  $D$  into two partitions, and then three partitions, and so on. If no such partitions exist, the resultant array is empty.

**Method 4:** This method further improves over Method 3 by avoiding certain multiplications (Algo. 2, *maxOptim=True*). For instance, if  $\alpha = p(H, \{1, \dots, 4\})$  and  $\beta = p(H, \{5, \dots, 8\})$ , then  $p(H, \{1, \dots, 8\})$  can be calculated as  $\frac{\alpha + \beta}{2}$  instead of  $\frac{\alpha * 4 + \beta * 4}{8}$ , reducing two vector multiplications. *getPartLens(parts)* in Algo. 2 returns a set of the lengths of *parts*’ elements.

---

### Algorithm 1: Layer Combination Iterator

---

**input** : Tensor  $H$ , Set  $A$ , Bool  $maxOptim$ , Int  $maxMem$   
**output** : Set  $D$ , Array  $P$

- 1 from Algo. 2 use `pool`
- 2 **Function** `layerIterator(H, A, maxOptim, maxMem)` :
  - 3 `mem = hashMap()`
  - 4 **for**  $i = 1$  to  $A.length$  **do**
    - 5 **for**  $D$  in `powerset(A, i)` **do**
      - 6 `P = pool(H, D, mem, maxMem)`
      - 7 `yield D, P`

---

To show the effectiveness of our algorithm, we carried out an experiment. We randomly select 1000 samples from the SICK dataset and try to find the best layer combination among all possible 8192 layer combinations, once with our algorithm and once with the naive algorithm (Method 2). We do this experimentation five times and report the average results. Our algorithm takes 5.65 seconds to find the best layer combination after the tensor of all layer and token vectors are obtained by BERT-base-uncased, while Method 2 takes 1067 seconds. The one-time forward pass of BERT takes 10 seconds.

## B STS Tasks’ Experiments and Results

### B.1 Data Sets and Evaluation Setup

For STS tasks, we use the seven standard STS data sets: **STS2012-2016** (Agirre et al., 2012, 2013, 2014, 2015, 2016), **STS Benchmark** (Cer et al., 2017), and **SICK Relatedness** (Marelli et al., 2014). These data sets consist of sentence pairs and a similarity score from 1 to 5 assigned to each one of them. For each data set, we first combine all of its data subsets, shuffle it, and choose a small



---

**Algorithm 2:** Layer Combination Pooler

---

```
input : Tensor  $H$ , Set  $D$ , HashMap  $mem$ ,  
        Bool  $maxOptim$ , Int  $maxMem$   
output : Array  $P$   
1 Function pool ( $H, D, mem, maxOptim,$   
    $maxMem$ ) :  
2    $parts = greedyPart (D, mem)$   
3   if  $parts \neq []$  then  
4      $P = array (parts [0].length, 0)$  ;  
       //  $array(N, v)$  returns an array  
       of size  $N$  filled with  $vs$   
5      $denom = D.length$   
6     if  $maxOptim$  then  
7        $lens = getPartLens (parts)$   
8       foreach  $len$  in  $lens$  do  
9          $cur_P = array (P.length, 0)$   
10        foreach  $part$  of size  $len$  in  $parts$   
11          do  
12             $cur_P += mem [part]$   
13            if  $lens.length > 1$  then  
14               $cur_P *= len$   
15            else  
16               $denom = len$   
17             $P += cur_P$   
18          else  
19            foreach  $part$  in  $parts$  do  
20               $P += mem [part] * len$   
21             $P /= denom$   
22          else  
23             $P = p(H, D)$   
24          if  $mem.length < maxMem$  then  
25             $mem [part] = P$   
26        return  $P$ 
```

---

subset (the first 350 sentence pairs) as our development data set and the rest as our test data set. We do this splitting 5 times randomly and report the average test performance. Please note that our method does not require training data. Nonetheless, we note that our method still needs a development data set in order for it to be bootstrapped. However, the required data set can be as small as 350 samples. Furthermore, as opposed to the other models, this data is not leveraged for fine-tuning, which is very time-consuming.

We use the sentence embeddings obtained by

our method or the baselines to calculate the cosine similarity between two sentences and then use Spearman’s correlation ( $\rho$ ) for evaluating the models’ performances as suggested by (Reimers and Gurevych, 2019).

## B.2 Baselines and Hyper-parameters

We use the same baselines that we use for the transfer tasks except for the MLM versions of the SimCSE models as those were proposed in the SimCSE paper to strengthen their models for the transfer tasks.

For the base models with 13 layers, we iterate through all the possible layer combinations for choosing the best one, yet for the large models with 25 layers, we only iterate through all the layer combinations with up to eight layers since in our experiments, we observed that combining more than eight layers hurts the performance.

## B.3 Results and Discussion

The results for STS tasks are shown in Table 2. From this table, we observe that our method significantly outperforms its corresponding variations of BERT, RoBERTa, SBERT, and SRoBERTa baselines on all the STS data sets. For instance, for BERT- $L_{uncased}$ , RoBERTa-L, and SRoBERTa-B, we obtain  $\rho$  improvements of up to 25.75%, 13.94%, and 3.95% and on average 16.32%, 9.52%, and 2.75%, respectively, compared to their best baselines (last avg). SimCSE-LC versions also outperform most of their baselines by a decent margin. For example, UnSupSimCSE-BL-LC improves its baseline’s  $\rho$  by up to 2.30% and on average 0.70%. Finally, the best SimCSE (Sup-RL)’s performances are also improved by up to 1.14%.

We can also see that our method is more effective on unsupervised models such as BERT and RoBERTa, as we spot a higher improvement compared with the supervised models such as SBERT or SupSimCSE. This suggests that the unsupervised models’ layers hold more important information than the supervised models for STS tasks. However, there is still crucial information to be exploited from all the layers of the supervised models as well.

We also observe that we have higher improvements for uncased BERT versions over their baseline compared to the cased versions. Nonetheless, after applying our method to either one, the final  $\rho$ s are within a 1.25% range of each other. This suggests that a lot of BERT- $uncased$ ’s information

is carried in its layers, and when exploited, it can perform more closely to BERT<sub>cased</sub>-LC.

## C Computation Setup

To conduct the experiments discussed in this paper, we used a computer with 128GB of RAM and one Intel Core i9-9980XE CPU. Our code uses PyTorch (Paszke et al., 2019), huggingface<sup>2</sup>, and sci-kit learn (Pedregosa et al., 2011).

## D More Baselines

In this section, we compare the results of the layer combination technique with more baselines on both STS and transfer tasks. We use the same baseline models as shown in Table 2, but instead of only comparing with the last layer for BERT/RobERTa models and original SimCSE and SBERT/SRobERTa models (which use CLS pooler and last layer average, respectively), we also compare the results with these two new baselines: last four layers average and random layers. The results are shown in Table 3 and Table 4. The results show that our method performs better than the new baselines as well.

## E Ablation studies

### E.1 Effect of the Number of Layers

In this subsection, we show the effect of combining only  $N$  layers of a model,  $N$  varying from 1 to 13 for the base models. To this end, we show this effect on the data sets STS16, STSB, and SICK, and for the models BERT-B<sub>cased</sub>, RoBERTa-B, SBERT-B, and SRobERTa-B. From Figs. 2 and 3, we can see that (1) For all the data sets and all the models adding more layers shows an upward trend in performance up to some peak point (e.g.,  $N = 4$  for BERT-B<sub>cased</sub> on STS16), and adding more layers after that point shows a downward trend in performance; however, this point is different for different models and data sets, yet it is always at most 6. (2) For BERT-B<sub>cased</sub>, RoBERTa-B, and SBERT-B, moving from one layer to two layers leads to a huge increase in performances for all but one data set-model pairs, yet combining more than two layers always leads to a further substantial increase in performance. For instance, on STSB, an absolute  $\rho$  improvement of 1.22% can be obtained when moving from two layers to six layers. We, nonetheless,

do not see drastic changes for SRobERTa-B when moving from one layer to two layers.

### E.2 Effect of Excluding a Particular Layer from Layer Combination

In this section, we discuss how excluding one particular layer from layer combination affects the performances. For this, we show the results for the models BERT-B<sub>cased</sub>, RoBERTa-B, and SBERT-B on the data sets STS16, STSB, and SICK in Fig. 4. We can mention these observations:

1. Most of the time, excluding layer 11 hurts the performance, showing that this is an important layer.
2. Excluding either of the layers 4, 5, or 6 does not hurt the performance, suggesting that these layers are not important for these models and data sets.
3. Last layer (L12) is always important for SBERT-B, is important for BERT-B<sub>cased</sub> in two cases (STS16 and STSB), and is only important for RoBERTa-B in one case (STSB).
4. Layer 0, which is the embedding layer, proves to be important in more than half of the cases here, showing that it carries important information to be considered for STS tasks, and it cannot be ignored.
5. Excluding any particular layer in any of the nine model-data set pairs shown in Fig. 4, leads to a maximum decrease of 0.97% (L12 for SBERT-B on STS16) in the performance. However as shown in Table 2, in any of these 9 cases, we get at least an improvement of 1.86% by combining certain layers (considering all the layers). This suggests that layer combination still outperforms the baselines even if one (any) layer is not considered at all.

<sup>2</sup><https://huggingface.co/transformers/>







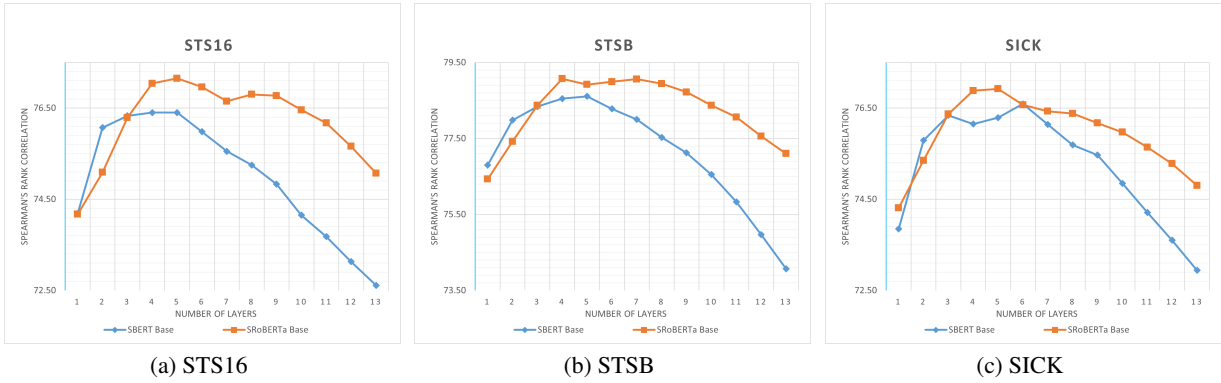


Figure 3: Effect of combining  $N$  Layers for SBERT-B and SRoBERTa



Figure 4: Effect of Excluding a Particular Layer