

SimCSE: Simple Contrastive Learning of Sentence Embeddings¹

Tianyu Gao^{†*} Xingcheng Yao^{‡*} Danqi Chen[†]

[†]Department of Computer Science, Princeton University

[‡]Institute for Interdisciplinary Information Sciences, Tsinghua University

{tianyug, dangic}@cs.princeton.edu

yxc18@mails.tsinghua.edu.cn

Abstract³

This paper presents SimCSE, a simple contrastive learning framework that greatly advances state-of-the-art sentence embeddings. We first describe an unsupervised approach, which takes an input sentence and predicts *itself* in a contrastive objective, with only standard dropout used as noise. This simple method works surprisingly well, performing on par with previous supervised counterparts. We find that dropout acts as minimal data augmentation, and removing it leads to a representation collapse. Then, we propose a supervised approach, which incorporates annotated pairs from natural language inference datasets into our contrastive learning framework by using “entailment” pairs as positives and “contradiction” pairs as hard negatives. We evaluate SimCSE on standard semantic textual similarity (STS) tasks, and our unsupervised and supervised models using BERT_{base} achieve an average of 76.3% and 81.6% Spearman’s correlation respectively, a 4.2% and 2.2% improvement compared to the previous best results. We also show—both theoretically and empirically—that the contrastive learning objective regularizes pre-trained embeddings’ anisotropic space to be more uniform, and it better aligns positive pairs when supervised signals are available.¹

embedding methods and demonstrate that a contrastive objective can be extremely effective when coupled with pre-trained language models such as BERT (Devlin et al., 2019) or RoBERTa (Liu et al., 2019). We present SimCSE, a simple contrastive sentence embedding framework, which can produce superior sentence embeddings, from either unlabeled or labeled data.

Our *unsupervised* SimCSE simply predicts the input sentence itself with only *dropout* (Srivastava et al., 2014) used as noise (Figure 1(a)). In other words, we pass the same sentence to the pre-trained encoder *twice*: by applying the standard dropout twice, we can obtain two different embeddings as “positive pairs”. Then we take other sentences in the same mini-batch as “negatives”, and the model predicts the positive one among the negatives. Although it may appear strikingly simple, this approach outperforms training objectives such as predicting next sentences (Logeswaran and Lee, 2018) and discrete data augmentation (e.g., word deletion and replacement) by a large margin, and even matches previous supervised methods. Through careful analysis, we find that dropout acts as minimal “data augmentation” of hidden representations while removing it leads to a representation collapse.

Our *supervised* SimCSE builds upon the recent success of using natural language inference (NLI) datasets for sentence embeddings (Conneau et al., 2017; Reimers and Gurevych, 2019) and incorporates annotated sentence pairs in contrastive learning (Figure 1(b)). Unlike previous work that casts it as a 3-way classification task (entailment, neutral, and contradiction), we leverage the fact that entailment pairs can be naturally used as positive instances. We also find that adding corresponding contradiction pairs as hard negatives further improves performance. This simple use of NLI datasets achieves a substantial improvement compared to prior methods using the same datasets. We also compare to other labeled sentence-pair

1 Introduction⁵

Learning universal sentence embeddings is a fundamental problem in natural language processing and has been studied extensively in the literature (Kiros et al., 2015; Hill et al., 2016; Conneau et al., 2017; Logeswaran and Lee, 2018; Cer et al., 2018; Reimers and Gurevych, 2019, *inter alia*). In this work, we advance state-of-the-art sentence

^{*}The first two authors contributed equally (listed in alphabetical order). This work was done when Xingcheng visited the Princeton NLP group remotely.

¹Our code and pre-trained models are publicly available at <https://github.com/princeton-nlp/SimCSE>.

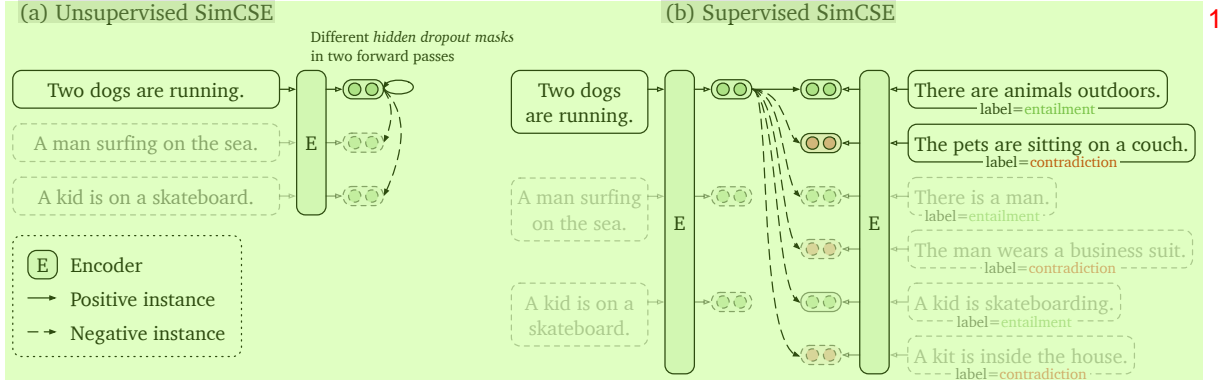


Figure 1: (a) Unsupervised SimCSE predicts the input sentence itself from in-batch negatives, with different hidden dropout masks applied. (b) Supervised SimCSE leverages the NLI datasets and takes the entailment (premise-hypothesis) pairs as positives, and contradiction pairs as well as other in-batch instances as negatives.

datasets and find that NLI datasets are especially effective for learning sentence embeddings.

To better understand the strong performance of SimCSE, we borrow the analysis tool from Wang and Isola (2020), which takes *alignment* between semantically-related positive pairs and *uniformity* of the whole representation space to measure the quality of learned embeddings. Through empirical analysis, we find that our unsupervised SimCSE essentially improves uniformity while avoiding degenerated alignment via dropout noise, thus improving the expressiveness of the representations. The same analysis shows that the NLI training signal can further improve alignment between positive pairs and produce better sentence embeddings. We also draw a connection to the recent findings that pre-trained word embeddings suffer from anisotropy (Ethayarajh, 2019; Li et al., 2020) and prove that—through a spectrum perspective—the contrastive learning objective “flattens” the singular value distribution of the sentence embedding space, hence improving uniformity.

We conduct a comprehensive evaluation of SimCSE on seven standard semantic textual similarity (STS) tasks (Agirre et al., 2012, 2013, 2014, 2015, 2016; Cer et al., 2017; Marelli et al., 2014) and seven transfer tasks (Conneau and Kiela, 2018). On the STS tasks, our unsupervised and supervised models achieve a 76.3% and 81.6% averaged Spearman’s correlation respectively using BERT_{base}, a 4.2% and 2.2% improvement compared to previous best results. We also achieve competitive performance on the transfer tasks. Finally, we identify an incoherent evaluation issue in the literature and consolidate the results of different settings for future work in evaluation of sentence embeddings.

2 Background: Contrastive Learning

Contrastive learning aims to learn effective representation by pulling semantically close neighbors together and pushing apart non-neighbors (Hadsell et al., 2006). It assumes a set of paired examples $\mathcal{D} = \{(x_i, x_i^+)\}_{i=1}^m$, where x_i and x_i^+ are semantically related. We follow the contrastive framework in Chen et al. (2020) and take a cross-entropy objective with in-batch negatives (Chen et al., 2017; Henderson et al., 2017): let \mathbf{h}_i and \mathbf{h}_i^+ denote the representations of x_i and x_i^+ , the training objective for (x_i, x_i^+) with a mini-batch of N pairs is:

$$\ell_i = -\log \frac{e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_i^+)/\tau}}{\sum_{j=1}^N e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_j^+)/\tau}}, \quad (1)$$

where τ is a temperature hyperparameter and $\text{sim}(\mathbf{h}_1, \mathbf{h}_2)$ is the cosine similarity $\frac{\mathbf{h}_1^\top \mathbf{h}_2}{\|\mathbf{h}_1\| \cdot \|\mathbf{h}_2\|}$. In this work, we encode input sentences using a pre-trained language model such as BERT (Devlin et al., 2019) or RoBERTa (Liu et al., 2019): $\mathbf{h} = f_\theta(x)$, and then fine-tune all the parameters using the contrastive learning objective (Eq. 1).

Positive instances. One critical question in contrastive learning is how to construct (x_i, x_i^+) pairs. In visual representations, an effective solution is to take two random transformations of the *same* image (e.g., cropping, flipping, distortion and rotation) as x_i and x_i^+ (Dosovitskiy et al., 2014). A similar approach has been recently adopted in language representations (Wu et al., 2020; Meng et al., 2021) by applying augmentation techniques such as word deletion, reordering, and substitution. However, data augmentation in NLP is inherently difficult because of its discrete nature. As we will see in §3,

simply using standard dropout on intermediate representations outperforms these discrete operators.

In NLP, a similar contrastive learning objective has been explored in different contexts (Henderson et al., 2017; Gillick et al., 2019; Karpukhin et al., 2020). In these cases, (x_i, x_i^+) are collected from supervised datasets such as question-passage pairs. Because of the distinct nature of x_i and x_i^+ , these approaches always use a *dual*-encoder framework, i.e., using two independent encoders f_{θ_1} and f_{θ_2} for x_i and x_i^+ . For sentence embeddings, Logeswaran and Lee (2018) also use contrastive learning with a dual-encoder approach, by forming current sentence and next sentence as (x_i, x_i^+) .

Alignment and uniformity. Recently, Wang and Isola (2020) identify two key properties related to contrastive learning—*alignment* and *uniformity*—and propose to use them to measure the quality of representations. Given a distribution of positive pairs p_{pos} , alignment calculates expected distance between embeddings of the paired instances (assuming representations are already normalized):

$$\ell_{\text{align}} \triangleq \mathbb{E}_{(x, x^+) \sim p_{\text{pos}}} \|f(x) - f(x^+)\|^2. \quad (2)$$

On the other hand, uniformity measures how well the embeddings are uniformly distributed:

$$\ell_{\text{uniform}} \triangleq \log \mathbb{E}_{x, y \stackrel{i.i.d.}{\sim} p_{\text{data}}} e^{-2\|f(x) - f(y)\|^2}, \quad (3)$$

where p_{data} denotes the data distribution. These two metrics are well aligned with the objective of contrastive learning: positive instances should stay close and embeddings for random instances should scatter on the hypersphere. In the following sections, we will also use the two metrics to justify the inner workings of our approaches.

3 Unsupervised SimCSE

The idea of unsupervised SimCSE is extremely simple: we take a collection of sentences $\{x_i\}_{i=1}^m$ and use $x_i^+ = x_i$. The key ingredient to get this to work with identical positive pairs is through the use of independently sampled *dropout masks* for x_i and x_i^+ . In standard training of Transformers (Vaswani et al., 2017), there are dropout masks placed on fully-connected layers as well as attention probabilities (default $p = 0.1$). We denote $\mathbf{h}_i^z = f_{\theta}(x_i, z)$ where z is a random mask for dropout. We simply feed the same input to the encoder *twice* and get

Data augmentation	STS-B		
None (unsup. SimCSE)	82.5		
Crop	10%	20%	30%
	77.8	71.4	63.6
Word deletion	10%	20%	30%
	75.9	72.2	68.2
Delete one word	75.9		
w/o dropout	74.2		
Synonym replacement	77.4		
MLM 15%	62.2		

Table 1: Comparison of data augmentations on STS-B development set (Spearman’s correlation). *Crop* $k\%$: keep $100-k\%$ of the length; *word deletion* $k\%$: delete $k\%$ words; *Synonym replacement*: use `nlpaug` (Ma, 2019) to randomly replace one word with its synonym; *MLM* $k\%$: use `BERTbase` to replace $k\%$ of words.

Training objective	f_{θ}	$(f_{\theta_1}, f_{\theta_2})$
Next sentence	67.1	68.9
Next 3 sentences	67.4	68.8
Delete one word	75.9	73.1
Unsupervised SimCSE	82.5	80.7

Table 2: Comparison of different unsupervised objectives (STS-B development set, Spearman’s correlation). The two columns denote whether we use one encoder or two independent encoders. *Next 3 sentences*: randomly sample one from the next 3 sentences. *Delete one word*: delete one word randomly (see Table 1).

two embeddings with different dropout masks z, z' , and the training objective of SimCSE becomes:

$$\ell_i = -\log \frac{e^{\text{sim}(\mathbf{h}_i^{z_i}, \mathbf{h}_i^{z'_i})/\tau}}{\sum_{j=1}^N e^{\text{sim}(\mathbf{h}_i^{z_i}, \mathbf{h}_j^{z'_j})/\tau}}, \quad (4)$$

for a mini-batch of N sentences. Note that z is just the *standard* dropout mask in Transformers and we do not add any additional dropout.

Dropout noise as data augmentation. We view it as a minimal form of data augmentation: the positive pair takes exactly the same sentence, and their embeddings only differ in dropout masks. We compare this approach to other training objectives on the STS-B development set (Cer et al., 2017)². Table 1 compares our approach to common data augmentation techniques such as crop, word deletion and replacement, which can be viewed as

²We randomly sample 10^6 sentences from English Wikipedia and fine-tune `BERTbase` with learning rate = $3e-5$, $N = 64$. In all our experiments, no STS training sets are used.

p	0.0	0.01	0.05	0.1
STS-B	71.1	72.6	81.1	82.5

p	0.15	0.2	0.5	Fixed 0.1
STS-B	81.4	80.5	71.0	43.6

Table 3: Effects of different dropout probabilities p on the STS-B development set (Spearman’s correlation, $\text{BERT}_{\text{base}}$). *Fixed 0.1*: default 0.1 dropout rate but apply the same dropout mask on both x_i and x_i^+ .

$\mathbf{h} = f_\theta(g(x), z)$ and g is a (random) discrete operator on x . We note that even deleting one word would hurt performance and none of the discrete augmentations outperforms dropout noise.

We also compare this self-prediction training objective to the next-sentence objective used in Logeswaran and Lee (2018), taking either one encoder or two independent encoders. As shown in Table 2, we find that SimCSE performs much better than the next-sentence objectives (82.5 vs 67.4 on STS-B) and using one encoder instead of two makes a significant difference in our approach.

Why does it work? To further understand the role of dropout noise in unsupervised SimCSE, we try out different dropout rates in Table 3 and observe that all the variants underperform the default dropout probability $p = 0.1$ from Transformers. We find two extreme cases particularly interesting: “no dropout” ($p = 0$) and “fixed 0.1” (using default dropout $p = 0.1$ but the same dropout masks for the pair). In both cases, the resulting embeddings for the pair are exactly the same, and it leads to a dramatic performance degradation. We take the checkpoints of these models every 10 steps during training and visualize the alignment and uniformity metrics³ in Figure 2, along with a simple data augmentation model “delete one word”. As clearly shown, starting from pre-trained checkpoints, all models greatly improve uniformity. However, the alignment of the two special variants also degrades drastically, while our unsupervised SimCSE keeps a steady alignment, thanks to the use of dropout noise. It also demonstrates that starting from a pre-trained checkpoint is crucial, for it provides good initial alignment. At last, “delete one word” improves the alignment yet achieves a smaller gain on the uniformity metric, and eventually underperforms unsupervised SimCSE.

³We take STS-B pairs with a score higher than 4 as p_{pos} and all STS-B sentences as p_{data} .

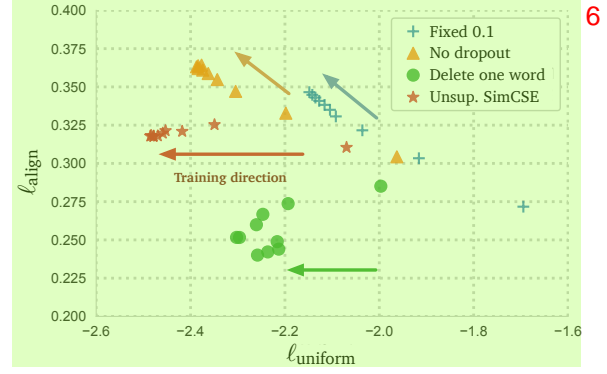


Figure 2: $\ell_{\text{align}}-\ell_{\text{uniform}}$ plot for unsupervised SimCSE, “no dropout”, “fixed 0.1”, and “delete one word”. We visualize checkpoints every 10 training steps and the arrows indicate the training direction. For both ℓ_{align} and ℓ_{uniform} , lower numbers are better.

4 Supervised SimCSE

We have demonstrated that adding dropout noise is able to keep a good alignment for positive pairs $(x, x^+) \sim p_{\text{pos}}$. In this section, we study whether we can leverage supervised datasets to provide better training signals for improving alignment of our approach. Prior work (Conneau et al., 2017; Reimers and Gurevych, 2019) has demonstrated that supervised natural language inference (NLI) datasets (Bowman et al., 2015; Williams et al., 2018) are effective for learning sentence embeddings, by predicting whether the relationship between two sentences is *entailment*, *neutral* or *contradiction*. In our contrastive learning framework, we instead directly take (x_i, x_i^+) pairs from supervised datasets and use them to optimize Eq. 1.

Choices of labeled data. We first explore which supervised datasets are especially suitable for constructing positive pairs (x_i, x_i^+) . We experiment with a number of datasets with sentence-pair examples, including 1) QQP⁴: Quora question pairs; 2) Flickr30k (Young et al., 2014): each image is annotated with 5 human-written captions and we consider any two captions of the same image as a positive pair; 3) ParaNMT (Wieting and Gimpel, 2018): a large-scale back-translation paraphrase dataset⁵; and finally 4) NLI datasets: SNLI (Bowman et al., 2015) and MNLI (Williams et al., 2018).

We train the contrastive learning model (Eq. 1) with different datasets and compare the results in

⁴<https://www.quora.com/q/quoradata/>

⁵ParaNMT is automatically constructed by machine translation systems. Strictly speaking, we should not call it “supervised”. It underperforms our unsupervised SimCSE though.

Table 4. For a fair comparison, we also run experiments with the same # of training pairs. Among all the options, using entailment pairs from the NLI (SNLI + MNLI) datasets performs the best. We think this is reasonable, as the NLI datasets consist of high-quality and crowd-sourced pairs. Also, human annotators are expected to write the hypotheses manually based on the premises and two sentences tend to have less lexical overlap. For instance, we find that the lexical overlap (F1 measured between two bags of words) for the entailment pairs (SNLI + MNLI) is 39%, while they are 60% and 55% for QQP and ParaNMT.

Contradiction as hard negatives. Finally, we further take the advantage of the NLI datasets by using its contradiction pairs as hard negatives⁶. In NLI datasets, given one premise, annotators are required to manually write one sentence that is absolutely true (*entailment*), one that might be true (*neutral*), and one that is definitely false (*contradiction*). Therefore, for each premise and its entailment hypothesis, there is an accompanying contradiction hypothesis⁷ (see Figure 1 for an example).

Formally, we extend (x_i, x_i^+) to (x_i, x_i^+, x_i^-) , where x_i is the premise, x_i^+ and x_i^- are entailment and contradiction hypotheses. The training objective ℓ_i is then defined by (N is mini-batch size):

$$-\log \frac{e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_i^+)/\tau}}{\sum_{j=1}^N \left(e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_j^+)/\tau} + e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_j^-)/\tau} \right)}.$$

As shown in Table 4, adding hard negatives can further improve performance ($84.9 \rightarrow 86.2$) and this is our final supervised SimCSE. We also tried to add the ANLI dataset (Nie et al., 2020) or combine it with our unsupervised SimCSE approach, but didn’t find a meaningful improvement. We also considered a dual encoder framework in supervised SimCSE and it hurt performance ($86.2 \rightarrow 84.2$).

5 Connection to Anisotropy

Recent work identifies an *anisotropy* problem in language representations (Ethayarajh, 2019; Li et al., 2020), i.e., the learned embeddings occupy a narrow cone in the vector space, which severely limits their expressiveness. Gao et al. (2019)

⁶We also experimented with adding neutral hypotheses as hard negatives. See Section 6.3 for more discussion.

⁷In fact, one premise can have multiple contradiction hypotheses. In our implementation, we only sample one as the hard negative and we did not find a difference by using more.

Dataset	sample	full
Unsup. SimCSE (1m)	-	82.5
QQP (134k)	81.8	81.8
Flickr30k (318k)	81.5	81.4
ParaNMT (5m)	79.7	78.7
SNLI+MNLI		
entailment (314k)	84.1	84.9
neutral (314k) ⁸	82.6	82.9
contradiction (314k)	77.5	77.6
all (942k)	81.7	81.9
SNLI+MNLI		
entailment + hard neg.	-	86.2
+ ANLI (52k)	-	85.0

Table 4: Comparisons of different supervised datasets as positive pairs. Results are Spearman’s correlations on the STS-B development set using BERT_{base} (we use the same hyperparameters as the final SimCSE model). Numbers in brackets denote the # of pairs. *Sample*: subsampling 134k positive pairs for a fair comparison among datasets; *full*: using the full dataset. In the last block, we use entailment pairs as positives and contradiction pairs as hard negatives (our final model).

demonstrate that language models trained with tied input/output embeddings lead to anisotropic word embeddings, and this is further observed by Ethayarajh (2019) in pre-trained contextual representations. Wang et al. (2020) show that singular values of the word embedding matrix in a language model decay drastically: except for a few dominating singular values, all others are close to zero.

A simple way to alleviate the problem is post-processing, either to eliminate the dominant principal components (Arora et al., 2017; Mu and Viswanath, 2018), or to map embeddings to an isotropic distribution (Li et al., 2020; Su et al., 2021). Another common solution is to add regularization during training (Gao et al., 2019; Wang et al., 2020). In this work, we show that—both theoretically and empirically—the contrastive objective can also alleviate the anisotropy problem.

The anisotropy problem is naturally connected to *uniformity* (Wang and Isola, 2020), both highlighting that embeddings should be evenly distributed in the space. Intuitively, optimizing the contrastive learning objective can improve uniformity (or ease the anisotropy problem), as the objective pushes negative instances apart. Here, we take a singular spectrum perspective—which is a common practice

⁸Though our final model only takes entailment pairs as positive instances, here we also try taking neutral and contradiction pairs from the NLI datasets as positive pairs.

in analyzing word embeddings (Mu and Viswanath, 2018; Gao et al., 2019; Wang et al., 2020), and show that the contrastive objective can “flatten” the singular value distribution of sentence embeddings and make the representations more isotropic.

Following Wang and Isola (2020), the asymptotics of the contrastive learning objective (Eq. 1) can be expressed by the following equation when the number of negative instances approaches infinity (assuming $f(x)$ is normalized):

$$-\frac{1}{\tau} \mathbb{E}_{(x, x^+) \sim p_{\text{pos}}} [f(x)^\top f(x^+)] + \mathbb{E}_{x \sim p_{\text{data}}} \left[\log \mathbb{E}_{x^- \sim p_{\text{data}}} \left[e^{f(x)^\top f(x^-)/\tau} \right] \right], \quad (6)$$

where the first term keeps positive instances similar and the second pushes negative pairs apart. When p_{data} is uniform over finite samples $\{x_i\}_{i=1}^m$, with $\mathbf{h}_i = f(x_i)$, we can derive the following formula from the second term with Jensen’s inequality:

$$\begin{aligned} & \mathbb{E}_{x \sim p_{\text{data}}} \left[\log \mathbb{E}_{x^- \sim p_{\text{data}}} \left[e^{f(x)^\top f(x^-)/\tau} \right] \right] \\ &= \frac{1}{m} \sum_{i=1}^m \log \left(\frac{1}{m} \sum_{j=1}^m e^{\mathbf{h}_i^\top \mathbf{h}_j / \tau} \right) \\ &\geq \frac{1}{\tau m^2} \sum_{i=1}^m \sum_{j=1}^m \mathbf{h}_i^\top \mathbf{h}_j. \end{aligned} \quad (7)$$

Let \mathbf{W} be the sentence embedding matrix corresponding to $\{x_i\}_{i=1}^m$, i.e., the i -th row of \mathbf{W} is \mathbf{h}_i . Optimizing the second term in Eq. 6 essentially minimizes an upper bound of the summation of all elements in $\mathbf{W}\mathbf{W}^\top$, i.e., $\text{Sum}(\mathbf{W}\mathbf{W}^\top) = \sum_{i=1}^m \sum_{j=1}^m \mathbf{h}_i^\top \mathbf{h}_j$.

Since we normalize \mathbf{h}_i , all elements on the diagonal of $\mathbf{W}\mathbf{W}^\top$ are 1 and then $\text{tr}(\mathbf{W}\mathbf{W}^\top)$ (the sum of all eigenvalues) is a constant. According to Merikoski (1984), if all elements in $\mathbf{W}\mathbf{W}^\top$ are positive, which is the case in most times according to Figure G.1, then $\text{Sum}(\mathbf{W}\mathbf{W}^\top)$ is an upper bound for the largest eigenvalue of $\mathbf{W}\mathbf{W}^\top$. When minimizing the second term in Eq. 6, we reduce the top eigenvalue of $\mathbf{W}\mathbf{W}^\top$ and inherently “flatten” the singular spectrum of the embedding space. Therefore, contrastive learning is expected to alleviate the representation degeneration problem and improve uniformity of sentence embeddings.

Compared to post-processing methods in Li et al. (2020); Su et al. (2021), which only aim to encourage isotropic representations, contrastive learning

also optimizes for aligning positive pairs by the first term in Eq. 6, which is the key to the success of SimCSE. A quantitative analysis is given in §7.

6 Experiment

6.1 Evaluation Setup

We conduct our experiments on 7 semantic textual similarity (STS) tasks. Note that all our STS experiments are fully **unsupervised** and no STS training sets are used. Even for supervised SimCSE, we simply mean that we take extra labeled datasets for training, following previous work (Conneau et al., 2017). We also evaluate 7 transfer learning tasks and provide detailed results in Appendix E. We share a similar sentiment with Reimers and Gurevych (2019) that the main goal of sentence embeddings is to cluster semantically similar sentences and hence take STS as the main result.

Semantic textual similarity tasks. We evaluate on 7 STS tasks: STS 2012–2016 (Agirre et al., 2012, 2013, 2014, 2015, 2016), STS Benchmark (Cer et al., 2017) and SICK-Relatedness (Marelli et al., 2014). When comparing to previous work, we identify invalid comparison patterns in published papers in the evaluation settings, including (a) whether to use an additional regressor, (b) Spearman’s vs Pearson’s correlation, and (c) how the results are aggregated (Table B.1). We discuss the detailed differences in Appendix B and choose to follow the setting of Reimers and Gurevych (2019) in our evaluation (no additional regressor, Spearman’s correlation, and “all” aggregation). We also report our replicated study of previous work as well as our results evaluated in a different setting in Table B.2 and Table B.3. We call for unifying the setting in evaluating sentence embeddings for future research.

Training details. We start from pre-trained checkpoints of BERT (Devlin et al., 2019) (uncased) or RoBERTa (Liu et al., 2019) (cased) and take the [CLS] representation as the sentence embedding⁹ (see §6.3 for comparison between different pooling methods). We train unsupervised SimCSE on 10^6 randomly sampled sentences from English Wikipedia, and train supervised SimCSE on the combination of MNLI and SNLI datasets (314k). More training details can be found in Appendix A.

⁹There is an MLP layer over [CLS] in BERT’s original implementation and we keep it with random initialization.

Model	STS12	STS13	STS14	STS15	STS16	STS-B	SICK-R	Avg.
<i>Unsupervised models</i>								
GloVe embeddings (avg.) [♣]	55.14	70.66	59.73	68.25	63.66	58.02	53.76	61.32
BERT _{base} (first-last avg.)	39.70	59.38	49.67	66.03	66.19	53.87	62.06	56.70
BERT _{base} -flow	58.40	67.10	60.85	75.16	71.22	68.66	64.47	66.55
BERT _{base} -whitening	57.83	66.90	60.90	75.08	71.31	68.24	63.73	66.28
IS-BERT _{base} [♡]	56.77	69.24	61.21	75.23	70.16	69.21	64.25	66.58
CT-BERT _{base}	61.63	76.80	68.47	77.50	76.48	74.31	69.19	72.05
* SimCSE-BERT _{base}	68.40	82.41	74.38	80.91	78.56	76.85	72.23	76.25
RoBERTa _{base} (first-last avg.)	40.88	58.74	49.07	65.63	61.48	58.55	61.63	56.57
RoBERTa _{base} -whitening	46.99	63.24	57.23	71.36	68.99	61.36	62.91	61.73
DeCLUTR-RoBERTa _{base}	52.41	75.19	65.52	77.12	78.63	72.41	68.62	69.99
* SimCSE-RoBERTa _{base}	70.16	81.77	73.24	81.36	80.65	80.22	68.56	76.57
* SimCSE-RoBERTa _{large}	72.86	83.99	75.62	84.77	81.80	81.98	71.26	78.90
<i>Supervised models</i>								
InferSent-GloVe [♣]	52.86	66.75	62.15	72.77	66.87	68.03	65.65	65.01
Universal Sentence Encoder [♣]	64.49	67.80	64.61	76.83	73.18	74.92	76.69	71.22
SBERT _{base} [♣]	70.97	76.53	73.19	79.09	74.30	77.03	72.91	74.89
SBERT _{base} -flow	69.78	77.27	74.35	82.01	77.46	79.12	76.21	76.60
SBERT _{base} -whitening	69.65	77.57	74.66	82.27	78.39	79.52	76.91	77.00
CT-SBERT _{base}	74.84	83.20	78.07	83.84	77.93	81.46	76.42	79.39
* SimCSE-BERT _{base}	75.30	84.67	80.19	85.40	80.82	84.25	80.39	81.57
SRoBERTa _{base} [♣]	71.54	72.49	70.80	78.74	73.69	77.77	74.46	74.21
SRoBERTa _{base} -whitening	70.46	77.07	74.46	81.64	76.43	79.49	76.65	76.60
* SimCSE-RoBERTa _{base}	76.53	85.21	80.95	86.03	82.57	85.83	80.50	82.52
* SimCSE-RoBERTa _{large}	77.46	87.27	82.36	86.66	83.93	86.70	81.95	83.76

Table 5: Sentence embedding performance on STS tasks (Spearman’s correlation, “all” setting). We highlight the highest numbers among models with the same pre-trained encoder. ♣: results from Reimers and Gurevych (2019); ♡: results from Zhang et al. (2020); all other results are reproduced or reevaluated by ourselves. For BERT-flow (Li et al., 2020) and whitening (Su et al., 2021), we only report the “NLI” setting (see Table C.1).

6.2 Main Results

We compare unsupervised and supervised SimCSE to previous state-of-the-art sentence embedding methods on STS tasks. Unsupervised baselines include average GloVe embeddings (Pennington et al., 2014), average BERT or RoBERTa embeddings¹⁰, and post-processing methods such as BERT-flow (Li et al., 2020) and BERT-whitening (Su et al., 2021). We also compare to several recent methods using a contrastive objective, including 1) IS-BERT (Zhang et al., 2020), which maximizes the agreement between global and local features; 2) DeCLUTR (Giorgi et al., 2021), which takes different spans from the same document as positive pairs; 3) CT (Carlsson et al., 2021), which aligns embeddings of the same sentence from two different encoders.¹¹ Other supervised

methods include InferSent (Conneau et al., 2017), Universal Sentence Encoder (Cer et al., 2018), and SBERT/SRoBERTa (Reimers and Gurevych, 2019) with post-processing methods (BERT-flow, whitening, and CT). We provide more details of these baselines in Appendix C.

Table 5 shows the evaluation results on 7 STS tasks. SimCSE can substantially improve results on all the datasets with or without extra NLI supervision, greatly outperforming the previous state-of-the-art models. Specifically, our unsupervised SimCSE-BERT_{base} improves the previous best averaged Spearman’s correlation from 72.05% to 76.25%, even comparable to supervised baselines. When using NLI datasets, SimCSE-BERT_{base} further pushes the state-of-the-art results to 81.57%. The gains are more pronounced on RoBERTa encoders, and our supervised SimCSE achieves 83.76% with RoBERTa_{large}.

In Appendix E, we show that SimCSE also achieves on par or better transfer task performance compared to existing work, and an auxiliary MLM objective can further boost performance.

¹⁰Following Su et al. (2021), we take the average of the first and the last layers, which is better than only taking the last.

¹¹We do not compare to CLEAR (Wu et al., 2020), because they use their own version of pre-trained models, and the numbers appear to be much lower. Also note that CT is a concurrent work to ours.

Pooler	Unsup.	Sup.
[CLS]		
w/ MLP	81.7	86.2
w/ MLP (train)	82.5	85.8
w/o MLP	80.9	86.2
First-last avg.	81.2	86.1

Table 6: Ablation studies of different pooling methods in unsupervised and supervised SimCSE. *[CLS] w/ MLP (train)*: using MLP on [CLS] during training but removing it during testing. The results are based on the development set of STS-B using BERT_{base}.

Hard neg	N/A	Contradiction	Contra.+ Neutral
α	-	0.5 1.0 2.0	1.0
STS-B	84.9	86.1 86.2 86.2	85.3

Table 7: STS-B development results with different hard negative policies. “N/A”: no hard negative.

6.3 Ablation Studies

We investigate the impact of different pooling methods and hard negatives. All reported results in this section are based on the STS-B development set. We provide more ablation studies (normalization, temperature, and MLM objectives) in Appendix D.

Pooling methods. Reimers and Gurevych (2019); Li et al. (2020) show that taking the average embeddings of pre-trained models (especially from both the first and last layers) leads to better performance than [CLS]. Table 6 shows the comparison between different pooling methods in both unsupervised and supervised SimCSE. For [CLS] representation, the original BERT implementation takes an extra MLP layer on top of it. Here, we consider three different settings for [CLS]: 1) keeping the MLP layer; 2) no MLP layer; 3) keeping MLP during training but removing it at testing time. We find that for unsupervised SimCSE, taking [CLS] representation with MLP only during training works the best; for supervised SimCSE, different pooling methods do not matter much. By default, we take [CLS] with MLP (train) for unsupervised SimCSE and [CLS] with MLP for supervised SimCSE.

Hard negatives. Intuitively, it may be beneficial to differentiate hard negatives (contradiction examples) from other in-batch negatives. Therefore, we extend our training objective defined in Eq. 5 to

incorporate weighting of different negatives:

$$-\log \frac{e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_i^+)/\tau}}{\sum_{j=1}^N (e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_j^+)/\tau} + \alpha \mathbb{1}_i^j e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_j^-)/\tau})}, \quad (8)$$

where $\mathbb{1}_i^j \in \{0, 1\}$ is an indicator that equals 1 if and only if $i = j$. We train SimCSE with different values of α and evaluate the trained models on the development set of STS-B. We also consider taking neutral hypotheses as hard negatives. As shown in Table 7, $\alpha = 1$ performs the best, and neutral hypotheses do not bring further gains.

7 Analysis

In this section, we conduct further analyses to understand the inner workings of SimCSE.

Uniformity and alignment. Figure 3 shows uniformity and alignment of different sentence embedding models along with their averaged STS results. In general, models which have both better alignment and uniformity achieve better performance, confirming the findings in Wang and Isola (2020). We also observe that (1) though pre-trained embeddings have good alignment, their uniformity is poor (i.e., the embeddings are highly anisotropic); (2) post-processing methods like BERT-flow and BERT-whitening greatly improve uniformity but also suffer a degeneration in alignment; (3) unsupervised SimCSE effectively improves uniformity of pre-trained embeddings whereas keeping a good alignment; (4) incorporating supervised data in SimCSE further amends alignment. In Appendix F, we further show that SimCSE can effectively flatten singular value distribution of pre-trained embeddings. In Appendix G, we demonstrate that SimCSE provides more distinguishable cosine similarities between different sentence pairs.

Qualitative comparison. We conduct a small-scale retrieval experiment using SBERT_{base} and SimCSE-BERT_{base}. We use 150k captions from Flickr30k dataset and take any random sentence as query to retrieve similar sentences (based on cosine similarity). As several examples shown in Table 8, the retrieved sentences by SimCSE have a higher quality compared to those retrieved by SBERT.

8 Related Work

Early work in sentence embeddings builds upon the distributional hypothesis by predicting surrounding sentences of a given one (Kiros et al., 2015; Hill

	SBERT _{base}	Supervised SimCSE-BERT _{base}
Query: A man riding a small boat in a harbor.		
#1	A group of men traveling over the ocean in a small boat.	A man on a moored blue and white boat.
#2	Two men sit on the bow of a colorful boat.	A man is riding in a boat on the water.
#3	A man wearing a life jacket is in a small boat on a lake.	A man in a blue boat on the water.
Query: A dog runs on the green grass near a wooden fence.		
#1	A dog runs on the green grass near a grove of trees.	The dog by the fence is running on the grass.
#2	A brown and white dog runs through the green grass.	Dog running through grass in fenced area.
#3	The dogs run in the green field.	A dog runs on the green grass near a grove of trees.

Table 8: Retrieved top-3 examples by SBERT and supervised SimCSE from Flickr30k (150k sentences). 1

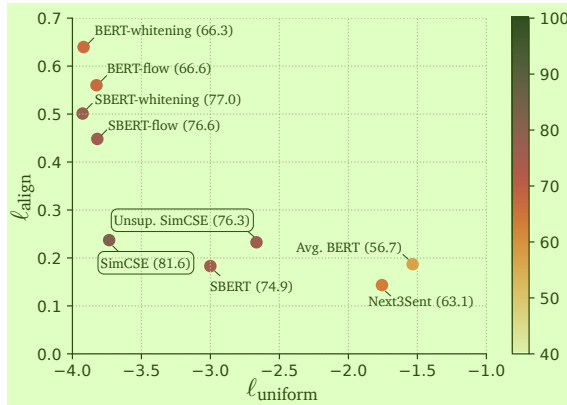


Figure 3: $\ell_{\text{align}}-\ell_{\text{uniform}}$ plot of models based on BERT_{base}. Color of points and numbers in brackets represent average STS performance (Spearman’s correlation). Next3Sent: “next 3 sentences” from Table 2. 4

et al., 2016; Logeswaran and Lee, 2018). Pagliarini et al. (2018) show that simply augmenting the idea of word2vec (Mikolov et al., 2013) with n-gram embeddings leads to strong results. Several recent (and concurrent) approaches adopt contrastive objectives (Zhang et al., 2020; Giorgi et al., 2021; Wu et al., 2020; Meng et al., 2021; Carlsson et al., 2021; Kim et al., 2021; Yan et al., 2021) by taking different views—from data augmentation or different copies of models—of the same sentence or document. Compared to these work, SimCSE uses the simplest idea by taking different outputs of the same sentence from standard dropout, and performs the best on STS tasks. 5

Supervised sentence embeddings are promised to have stronger performance compared to unsupervised counterparts. Conneau et al. (2017) propose to fine-tune a Siamese model on NLI datasets, which is further extended to other encoders or pre-trained models (Cer et al., 2018; Reimers and Gurevych, 2019). Furthermore, Wieting and Gimpel (2018); Wieting et al. (2020) demonstrate that 6

bilingual and back-translation corpora provide useful supervision for learning semantic similarity. Another line of work focuses on regularizing embeddings (Li et al., 2020; Su et al., 2021; Huang et al., 2021) to alleviate the representation degeneration problem (as discussed in §5), and yields substantial improvement over pre-trained language models. 7

9 Conclusion 8

In this work, we propose SimCSE, a simple contrastive learning framework, which greatly improves state-of-the-art sentence embeddings on semantic textual similarity tasks. We present an unsupervised approach which predicts input sentence itself with dropout noise and a supervised approach utilizing NLI datasets. We further justify the inner workings of our approach by analyzing alignment and uniformity of SimCSE along with other baseline models. We believe that our contrastive objective, especially the unsupervised one, may have a broader application in NLP. It provides a new perspective on data augmentation with text input, and can be extended to other continuous representations and integrated in language model pre-training. 9

Acknowledgements 10

We thank Tao Lei, Jason Lee, Zhengyan Zhang, Jinhyuk Lee, Alexander Wettig, Zexuan Zhong, and the members of the Princeton NLP group for helpful discussion and valuable feedback. This research is supported by a Graduate Fellowship at Princeton University and a gift award from Apple. 11

References¹

- 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.
- 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.
- 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. 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.
- 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.
- Sanjeev Arora, Yingyu Liang, and Tengyu Ma. 2017. [A simple but tough-to-beat baseline for sentence embeddings](#). In *International Conference on Learning Representations (ICLR)*.
- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. [A large annotated corpus for learning natural language inference](#). In *Empirical Methods in Natural Language Processing (EMNLP)*, pages 632–642.
- Fredrik Carlsson, Amaru Cuba Gyllensten, Evangelia Gogoulou, Erik Ylipää Hellqvist, and Magnus Sahlgren. 2021. [Semantic re-tuning with contrastive tension](#). In *International Conference on Learning Representations (ICLR)*.
- 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.
- 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 *Empirical Methods in Natural Language Processing (EMNLP): System Demonstrations*, pages 169–174.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. [A simple framework for contrastive learning of visual representations](#). In *International Conference on Machine Learning (ICML)*, pages 1597–1607.
- Ting Chen, Yizhou Sun, Yue Shi, and Liangjie Hong. 2017. [On sampling strategies for neural network-based collaborative filtering](#). In *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 767–776.
- Alexis Conneau and Douwe Kiela. 2018. [SentEval: An evaluation toolkit for universal sentence representations](#). In *International Conference on Language Resources and Evaluation (LREC)*.
- Alexis Conneau, Douwe Kiela, Holger Schwenk, Loïc Barrault, and Antoine Bordes. 2017. [Supervised learning of universal sentence representations from natural language inference data](#). In *Empirical Methods in Natural Language Processing (EMNLP)*, pages 670–680.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, pages 4171–4186.
- William B. Dolan and Chris Brockett. 2005. [Automatically constructing a corpus of sentential paraphrases](#). In *Proceedings of the Third International Workshop on Paraphrasing (IWP2005)*.
- Alexey Dosovitskiy, Jost Tobias Springenberg, Martin Riedmiller, and Thomas Brox. 2014. [Discriminative unsupervised feature learning with convolutional neural networks](#). In *Advances in Neural Information Processing Systems (NIPS)*, volume 27.
- Kawin Ethayarajh. 2019. [How contextual are contextualized word representations? comparing the geometry of BERT, ELMo, and GPT-2 embeddings](#). In *Empirical Methods in Natural Language Processing and International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 55–65.
- Jun Gao, Di He, Xu Tan, Tao Qin, Liwei Wang, and Tieyan Liu. 2019. [Representation degeneration problem in training natural language generation](#)

- models. In *International Conference on Learning Representations (ICLR)*.
- Dan Gillick, Sayali Kulkarni, Larry Lansing, Alessandro Presta, Jason Baldridge, Eugene Ie, and Diego Garcia-Olano. 2019. [Learning dense representations for entity retrieval](#). In *Computational Natural Language Learning (CoNLL)*, pages 528–537.
- John Giorgi, Osvald Nitski, Bo Wang, and Gary Bader. 2021. [DeCLUTR: Deep contrastive learning for unsupervised textual representations](#). In *Association for Computational Linguistics and International Joint Conference on Natural Language Processing (ACL-IJCNLP)*, pages 879–895.
- Raia Hadsell, Sumit Chopra, and Yann LeCun. 2006. [Dimensionality reduction by learning an invariant mapping](#). In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 2, pages 1735–1742. IEEE.
- Matthew Henderson, Rami Al-Rfou, Brian Strope, Yun-Hsuan Sung, László Lukács, Ruiqi Guo, Sanjiv Kumar, Balint Miklos, and Ray Kurzweil. 2017. [Efficient natural language response suggestion for smart reply](#). *arXiv preprint arXiv:1705.00652*.
- Felix Hill, Kyunghyun Cho, and Anna Korhonen. 2016. [Learning distributed representations of sentences from unlabelled data](#). In *North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, pages 1367–1377.
- Minqing Hu and Bing Liu. 2004. [Mining and summarizing customer reviews](#). In *ACM SIGKDD international conference on Knowledge discovery and data mining*.
- Junjie Huang, Duyu Tang, Wanjun Zhong, Shuai Lu, Linjun Shou, Ming Gong, Daxin Jiang, and Nan Duan. 2021. [Whiteningbert: An easy unsupervised sentence embedding approach](#). *arXiv preprint arXiv:2104.01767*.
- Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. [Dense passage retrieval for open-domain question answering](#). In *Empirical Methods in Natural Language Processing (EMNLP)*, pages 6769–6781.
- Taeuk Kim, Kang Min Yoo, and Sang-goo Lee. 2021. [Self-guided contrastive learning for BERT sentence representations](#). In *Association for Computational Linguistics and International Joint Conference on Natural Language Processing (ACL-IJCNLP)*, pages 2528–2540.
- Ryan Kiros, Yukun Zhu, Ruslan Salakhutdinov, Richard S Zemel, Antonio Torralba, Raquel Urtasun, and Sanja Fidler. 2015. [Skip-thought vectors](#). In *Advances in Neural Information Processing Systems (NIPS)*, pages 3294–3302.
- Bohan Li, Hao Zhou, Junxian He, Mingxuan Wang, Yiming Yang, and Lei Li. 2020. [On the sentence embeddings from pre-trained language models](#). In *Empirical Methods in Natural Language Processing (EMNLP)*, pages 9119–9130.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. [Roberta: A robustly optimized bert pretraining approach](#). *arXiv preprint arXiv:1907.11692*.
- Lajanugen Logeswaran and Honglak Lee. 2018. [An efficient framework for learning sentence representations](#). In *International Conference on Learning Representations (ICLR)*.
- Edward Ma. 2019. [Nlp augmentation](#). <https://github.com/makcedward/nlpaug>.
- 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 *International Conference on Language Resources and Evaluation (LREC)*, pages 216–223.
- Yu Meng, Chenyan Xiong, Payal Bajaj, Saurabh Tiwary, Paul Bennett, Jiawei Han, and Xia Song. 2021. [COCO-LM: Correcting and contrasting text sequences for language model pretraining](#). *arXiv preprint arXiv:2102.08473*.
- Jorma Kaarlo Merikoski. 1984. [On the trace and the sum of elements of a matrix](#). *Linear Algebra and its Applications*, 60:177–185.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, G. Corrado, and J. Dean. 2013. [Distributed representations of words and phrases and their compositionality](#). In *Advances in Neural Information Processing Systems (NIPS)*.
- Jiaqi Mu and Pramod Viswanath. 2018. [All-but-the-top: Simple and effective postprocessing for word representations](#). In *International Conference on Learning Representations (ICLR)*.
- Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. 2020. [Adversarial NLI: A new benchmark for natural language understanding](#). In *Association for Computational Linguistics (ACL)*, pages 4885–4901.
- Matteo Pagliardini, Prakhar Gupta, and Martin Jaggi. 2018. [Unsupervised learning of sentence embeddings using compositional n-gram features](#). In *North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, pages 528–540.
- Bo Pang and Lillian Lee. 2004. [A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts](#). In *Association for Computational Linguistics (ACL)*, pages 271–278.

- Bo Pang and Lillian Lee. 2005. Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales. In *Association for Computational Linguistics (ACL)*, pages 115–124. 1
- 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. 2
- Nils Reimers, Philip Beyer, and Iryna Gurevych. 2016. Task-oriented intrinsic evaluation of semantic textual similarity. In *International Conference on Computational Linguistics (COLING)*, pages 87–96. 3
- Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence embeddings using Siamese BERT-networks. In *Empirical Methods in Natural Language Processing and International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3982–3992. 4
- 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 *Empirical Methods in Natural Language Processing (EMNLP)*, pages 1631–1642. 5
- Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: a simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research (JMLR)*, 15(1):1929–1958. 6
- Jianlin Su, Jiarun Cao, Weijie Liu, and Yangyiwen Ou. 2021. Whitening sentence representations for better semantics and faster retrieval. *arXiv preprint arXiv:2103.15316*. 7
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems (NIPS)*, pages 6000–6010. 8
- Ellen M Voorhees and Dawn M Tice. 2000. Building a question answering test collection. In *the 23rd annual international ACM SIGIR conference on Research and development in information retrieval*, pages 200–207. 9
- Lingxiao Wang, Jing Huang, Kevin Huang, Ziniu Hu, Guangtao Wang, and Quanquan Gu. 2020. Improving neural language generation with spectrum control. In *International Conference on Learning Representations (ICLR)*. 10
- Tongzhou Wang and Phillip Isola. 2020. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In *International Conference on Machine Learning (ICML)*, pages 9929–9939. 11
- Janyce Wiebe, Theresa Wilson, and Claire Cardie. 2005. Annotating expressions of opinions and emotions in language. *Language resources and evaluation*, 39(2-3):165–210. 12
- John Wieting and Kevin Gimpel. 2018. ParaNMT-50M: Pushing the limits of paraphrastic sentence embeddings with millions of machine translations. In *Association for Computational Linguistics (ACL)*, pages 451–462. 13
- John Wieting, Graham Neubig, and Taylor Berg-Kirkpatrick. 2020. A bilingual generative transformer for semantic sentence embedding. In *Empirical Methods in Natural Language Processing (EMNLP)*, pages 1581–1594. 14
- Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference. In *North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, pages 1112–1122. 15
- 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 *Empirical Methods in Natural Language Processing (EMNLP): System Demonstrations*, pages 38–45. 16
- Zhuofeng Wu, Sinong Wang, Jiatao Gu, Madian Khabsa, Fei Sun, and Hao Ma. 2020. Clear: Contrastive learning for sentence representation. *arXiv preprint arXiv:2012.15466*. 17
- 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 *Association for Computational Linguistics and International Joint Conference on Natural Language Processing (ACL-IJCNLP)*, pages 5065–5075. 18
- Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. 2014. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *Transactions of the Association for Computational Linguistics*, 2:67–78. 19
- Yan Zhang, Ruidan He, Zuozhu Liu, Kwan Hui Lim, and Lidong Bing. 2020. An unsupervised sentence embedding method by mutual information maximization. In *Empirical Methods in Natural Language Processing (EMNLP)*, pages 1601–1610. 20

A Training Details¹

We implement SimCSE with transformers² package (Wolf et al., 2020). For supervised SimCSE, we train our models for 3 epochs, evaluate the model every 250 training steps on the development set of STS-B and keep the best checkpoint for the final evaluation on test sets. We do the same for the unsupervised SimCSE, except that we train the model for one epoch. We carry out grid-search of batch size $\in \{64, 128, 256, 512\}$ and learning rate $\in \{1e-5, 3e-5, 5e-5\}$ on STS-B development set and adopt the hyperparameter settings in Table A.1. We find that SimCSE is not sensitive to batch sizes as long as tuning the learning rates accordingly, which contradicts the finding that contrastive learning requires large batch sizes (Chen et al., 2020). It is probably due to that all SimCSE models start from pre-trained checkpoints, which already provide us a good set of initial parameters.

	Unsupervised				Supervised	
	BERT		RoBERTa		base	large
Batch size	64	64	512	512	512	512
Learning rate	3e-5	1e-5	1e-5	3e-5	5e-5	1e-5

Table A.1: Batch sizes and learning rates for SimCSE.³

For both unsupervised and supervised SimCSE,⁵ we take the [CLS] representation with an MLP layer on top of it as the sentence representation. Specially, for unsupervised SimCSE, we discard the MLP layer and only use the [CLS] output during test, since we find that it leads to better performance (ablation study in §6.3).

Finally, we introduce one more optional variant⁶ which adds a masked language modeling (MLM) objective (Devlin et al., 2019) as an auxiliary loss to Eq. 1: $\ell + \lambda \cdot \ell^{\text{mlm}}$ (λ is a hyperparameter). This helps SimCSE avoid catastrophic forgetting of token-level knowledge. As we will show in Table D.2, we find that adding this term can help improve performance on transfer tasks (not on sentence-level STS tasks).

B Different Settings for STS Evaluation⁷

We elaborate the differences in STS evaluation settings in previous work in terms of (a) whether to use additional regressors; (b) reported metrics; (c) different ways to aggregate results.

Additional regressors. The default SentEval⁹ implementation applies a linear regressor on top of

Paper	Reg.	Metric	Aggr.
Hill et al. (2016)		Both	all
Conneau et al. (2017)	✓	Pearson	mean
Conneau and Kiela (2018)	✓	Pearson	mean
Reimers and Gurevych (2019)		Spearman	all
Zhang et al. (2020)		Spearman	all
Li et al. (2020)		Spearman	wmean
Su et al. (2021)		Spearman	wmean
Wieting et al. (2020)		Pearson	mean
Giorgi et al. (2021)		Spearman	mean
Ours		Spearman	all

Table B.1: STS evaluation protocols used in different papers. “Reg.”: whether an additional regressor is used; “aggr.”: methods to aggregate different subset results.¹¹

frozen sentence embeddings for STS-B and SICK-R,¹² and train the regressor on the training sets of the two tasks, while most sentence representation papers take the raw embeddings and evaluate in an unsupervised way. In our experiments, *we do not apply any additional regressors and directly take cosine similarities for all STS tasks.*

Metrics. Both Pearson’s and Spearman’s correlation coefficients are used in the literature.¹³ Reimers et al. (2016) argue that Spearman correlation, which measures the rankings instead of the actual scores, better suits the need of evaluating sentence embeddings. *For all of our experiments, we report Spearman’s rank correlation.*

Aggregation methods. Given that each year’s STS challenge contains several subsets, there are different choices to gather results from them: one way is to concatenate all the topics and report the overall Spearman’s correlation (denoted as “all”), and the other is to calculate results for different subsets separately and average them (denoted as “mean” if it is simple average or “wmean” if weighted by the subset sizes). However, most papers do not claim the method they take, making it challenging for a fair comparison. We take some of the most recent work: SBERT (Reimers and Gurevych, 2019), BERT-flow (Li et al., 2020) and BERT-whitening (Su et al., 2021)¹² as an example: In Table B.2, we compare our reproduced results to reported results of SBERT and BERT-whitening, and find that Reimers and Gurevych (2019) take the “all” setting but Li et al. (2020); Su et al. (2021) take the “wmean” setting, even though Li et al. (2020) claim that they take the same setting as Reimers

¹²Li et al. (2020) and Su et al. (2021) have consistent results, so we assume that they take the same evaluation and just take BERT-whitening in experiments here.

Model	STS12	STS13	STS14	STS15	STS16	STS-B	SICK-R	Avg.
SBERT (all)	70.97	76.53	73.19	79.09	74.30	76.98	72.91	74.85
SBERT (wmean)	66.35	73.76	73.88	77.33	73.62	76.98	72.91	73.55
SBERT♣	70.97	76.53	73.19	79.09	74.30	77.03	72.91	74.89
BERT-whitening (NLI, all)	57.83	66.90	60.89	75.08	71.30	68.23	63.73	66.28
BERT-whitening (NLI, wmean)	61.43	65.90	65.96	74.80	73.10	68.23	63.73	67.59
BERT-whitening (NLI)♠	61.69	65.70	66.02	75.11	73.11	68.19	63.60	67.63
BERT-whitening (target, all)	42.88	77.77	66.27	63.60	67.58	71.34	60.40	64.26
BERT-whitening (target, wmean)	63.38	73.01	69.13	74.48	72.56	71.34	60.40	69.19
BERT-whitening (target)♠	63.62	73.02	69.23	74.52	72.15	71.34	60.60	69.21

Table B.2: Comparisons of our reproduced results using different evaluation protocols and the original numbers. ♣: results from Reimers and Gurevych (2019); ♠: results from Su et al. (2021); Other results are reproduced by us. From the table we see that SBERT takes the “all” evaluation and BERT-whitening takes the “wmean” evaluation.

Model	STS12	STS13	STS14	STS15	STS16	STS-B	SICK-R	Avg.
BERT _{base} (first-last avg.)♠	57.86	61.97	62.49	70.96	69.76	59.04	63.75	63.69
+ flow (NLI)♠	59.54	64.69	64.66	72.92	71.84	58.56	65.44	65.38
+ flow (target)♠	63.48	72.14	68.42	73.77	75.37	70.72	63.11	69.57
+ whitening (NLI)♠	61.69	65.70	66.02	75.11	73.11	68.19	63.60	67.63
+ whitening (target)♠	63.62	73.02	69.23	74.52	72.15	71.34	60.60	69.21
* Unsup. SimCSE-BERT _{base}	70.14	79.56	75.91	81.46	79.07	76.85	72.23	76.46
SBERT _{base} (first-last avg.)♠	68.70	74.37	74.73	79.65	75.21	77.63	74.84	75.02
+ flow (NLI)♠	67.75	76.73	75.53	80.63	77.58	79.10	78.03	76.48
+ flow (target)♠	68.95	78.48	77.62	81.95	78.94	81.03	74.97	77.42
+ whitening (NLI)♠	69.11	75.79	75.76	82.31	79.61	78.66	76.33	76.80
+ whitening (target)♠	69.01	78.10	77.04	80.83	77.93	80.50	72.54	76.56
* Sup. SimCSE-BERT _{base}	70.90	81.49	80.19	83.79	81.89	84.25	80.39	80.41

Table B.3: STS results with “wmean” setting (Spearman). ♠: from Li et al. (2020); Su et al. (2021).

and Gurevych (2019). Since the “all” setting fuses data from different topics together, it makes the evaluation closer to real-world scenarios, and unless specified, we take the “all” setting.

We list evaluation settings for a number of previous work in Table B.1. Some of the settings are reported by the paper and some of them are inferred by comparing the results and checking their code. As we can see, the evaluation protocols are very incoherent across different papers. We call for unifying the setting in evaluating sentence embeddings for future research. We will also release our evaluation code for better reproducibility. Since previous work uses different evaluation protocols from ours, we further evaluate our models in these settings to make a direct comparison to the published numbers. We evaluate SimCSE with “wmean” and Spearman’s correlation to directly compare to Li et al. (2020) and Su et al. (2021) in Table B.3.

C Baseline Models

We elaborate on how we obtain different baselines for comparison in our experiments:

- For average GloVe embedding (Pennington et al., 2014), InferSent (Conneau et al., 2017) and Universal Sentence Encoder (Cer et al., 2018), we directly report the results from Reimers and Gurevych (2019), since our evaluation setting is the same as theirs.

- For BERT (Devlin et al., 2019) and RoBERTa (Liu et al., 2019), we download the pre-trained model weights from HuggingFace’s Transformers¹³, and evaluate the models with our own scripts.

- For SBERT and SRoBERTa (Reimers and Gurevych, 2019), we reuse the results from the original paper. For results not reported by Reimers and Gurevych (2019), such as the performance of SRoBERTa on transfer tasks, we download the model weights from SentenceTransformers¹⁴ and evaluate them.

¹³<https://github.com/huggingface/transformers>

¹⁴<https://www.sbert.net/>

Model	STS12	STS13	STS14	STS15	STS16	STS-B	SICK-R	Avg.
BERT-flow (NLI)	58.40	67.10	60.85	75.16	71.22	68.66	64.47	66.55
BERT-flow (target)	53.15	78.38	66.02	62.09	70.84	71.70	61.97	66.31
BERT-whitening (NLI)	57.83	66.90	60.90	75.08	71.31	68.24	63.73	66.28
BERT-whitening (target)	42.88	77.77	66.28	63.60	67.58	71.34	60.40	64.26
SBERT-flow (NLI)	69.78	77.27	74.35	82.01	77.46	79.12	76.21	76.60
SBERT-flow (target)	66.18	82.69	76.22	73.72	75.71	79.99	73.82	75.48
SBERT-whitening (NLI)	69.65	77.57	74.66	82.27	78.39	79.52	76.91	77.00
SBERT-whitening (target)	52.91	81.91	75.44	72.24	72.93	80.50	72.54	72.64

Table C.1: Comparison of using NLI or target data for postprocessing methods (“all”, Spearman’s correlation). 2

τ	N/A	0.001	0.01	0.05	0.1	1
STS-B	85.9	84.9	85.4	86.2	82.0	64.0

Table D.1: STS-B development results (Spearman’s correlation) with different temperatures. “N/A”: Dot product instead of cosine similarity. 4

Model	STS-B	Avg. transfer
w/o MLM	86.2	85.8
w/ MLM		
$\lambda = 0.01$	85.7	86.1
$\lambda = 0.1$	85.7	86.2
$\lambda = 1$	85.1	85.8

Table D.2: Ablation studies of the MLM objective based on the development sets using BERT_{base}. 12

- For DeCLUTR (Giorgi et al., 2021) and contrastive tension (Carlsson et al., 2021), we reevaluate their checkpoints in our setting. 5
- For BERT-flow (Li et al., 2020), since their original numbers take a different setting, we retrain their models using their code¹⁵, and evaluate the models using our own script. 6
- For BERT-whitening (Su et al., 2021), we implemented our own version of whitening script following the same pooling method in Su et al. (2021), i.e. first-last average pooling. Our implementation can reproduce the results from the original paper (see Table B.2). 7

For both BERT-flow and BERT-whitening, they have two variants of postprocessing: one takes the NLI data (“NLI”) and one directly learns the embedding distribution on the target sets (“target”). We find that in our evaluation setting, “target” is generally worse than “NLI” (Table C.1), so we only report the NLI variant in the main results. 8

D Ablation Studies 9

Normalization and temperature. We train SimCSE using both dot product and cosine similarity with different temperatures and evaluate them on the STS-B development set. As shown in Table D.1, with a carefully tuned temperature $\tau = 0.05$, cosine similarity is better than dot product. 10

MLM auxiliary task. Finally, we study the impact of the MLM auxiliary objective with different λ . As shown in Table D.2, the token-level MLM objective improves the averaged performance on transfer tasks modestly, yet it brings a consistent drop in semantic textual similarity tasks. 13

E Transfer Tasks 14

We evaluate our models on the following transfer tasks: MR (Pang and Lee, 2005), CR (Hu and Liu, 2004), SUBJ (Pang and Lee, 2004), MPQA (Wiebe et al., 2005), SST-2 (Socher et al., 2013), TREC (Voorhees and Tice, 2000) and MRPC (Dolan and Brockett, 2005). A logistic regression classifier is trained on top of (frozen) sentence embeddings produced by different methods. We follow default configurations from SentEval¹⁶. 15

Table E.1 shows the evaluation results on transfer tasks. We find that supervised SimCSE performs on par or better than previous approaches, although the trend of unsupervised models remains unclear. We find that adding this MLM term consistently improves performance on transfer tasks, confirming our intuition that sentence-level objective may not directly benefit transfer tasks. We also experiment with post-processing methods (BERT- 16

¹⁵<https://github.com/bohanli/BERT-flow>

¹⁶<https://github.com/facebookresearch/SentEval>

Model	MR	CR	SUBJ	MPQA	SST	TREC	MRPC	Avg.
<i>Unsupervised models</i>								
GloVe embeddings (avg.)♣	77.25	78.30	91.17	87.85	80.18	83.00	72.87	81.52
Skip-thought♡	76.50	80.10	93.60	87.10	82.00	92.20	73.00	83.50
Avg. BERT embeddings♣	78.66	86.25	94.37	88.66	84.40	92.80	69.54	84.94
BERT-[CLS] embedding♣	78.68	84.85	94.21	88.23	84.13	91.40	71.13	84.66
IS-BERT _{base} ♡	81.09	87.18	94.96	88.75	85.96	88.64	74.24	85.83
* SimCSE-BERT _{base}	81.18	86.46	94.45	88.88	85.50	89.80	74.43	85.81
w/ MLM	82.92	87.23	95.71	88.73	86.81	87.01	78.07	86.64
* SimCSE-RoBERTa _{base}	81.04	87.74	93.28	86.94	86.60	84.60	73.68	84.84
w/ MLM	83.37	87.76	95.05	87.16	89.02	90.80	75.13	86.90
* SimCSE-RoBERTa _{large}	82.74	87.87	93.66	88.22	88.58	92.00	69.68	86.11
w/ MLM	84.66	88.56	95.43	87.50	89.46	95.00	72.41	87.57
<i>Supervised models</i>								
InferSent-GloVe♣	81.57	86.54	92.50	90.38	84.18	88.20	75.77	85.59
Universal Sentence Encoder♣	80.09	85.19	93.98	86.70	86.38	93.20	70.14	85.10
SBERT _{base} ♣	83.64	89.43	94.39	89.86	88.96	89.60	76.00	87.41
* SimCSE-BERT _{base}	82.69	89.25	94.81	89.59	87.31	88.40	73.51	86.51
w/ MLM	82.68	88.88	94.52	89.82	88.41	87.60	76.12	86.86
SRoBERTa _{base}	84.91	90.83	92.56	88.75	90.50	88.60	78.14	87.76
* SimCSE-RoBERTa _{base}	84.92	92.00	94.11	89.82	91.27	88.80	75.65	88.08
w/ MLM	85.08	91.76	94.02	89.72	92.31	91.20	76.52	88.66
* SimCSE-RoBERTa _{large}	88.12	92.37	95.11	90.49	92.75	91.80	76.64	89.61
w/ MLM	88.45	92.53	95.19	90.58	93.30	93.80	77.74	90.23

Table E.1: Transfer task results of different sentence embedding models (measured as accuracy). ♣: results from Reimers and Gurevych (2019); ♡: results from Zhang et al. (2020). We highlight the highest numbers among models with the same pre-trained encoder. MLM: adding MLM as an auxiliary task with $\lambda = 0.1$.

flow/whitening) and find that they both hurt performance compared to their base models, showing that good uniformity of representations does not lead to better embeddings for transfer learning. As we argued earlier, we think that transfer tasks are not a major goal for sentence embeddings, and thus we take the STS results for main comparison.

F Distribution of Singular Values

Figure F.1 shows the singular value distribution of SimCSE together with other baselines. For both unsupervised and supervised cases, singular value drops the fastest for vanilla BERT or SBERT embeddings, while SimCSE helps flatten the spectrum distribution. Postprocessing-based methods such as BERT-flow or BERT-whitening flatten the curve even more since they directly aim for the goal of mapping embeddings to an isotropic distribution.

G Cosine-similarity Distribution

To directly show the strengths of our approaches on STS tasks, we illustrate the cosine similarity distributions of STS-B pairs with different groups of

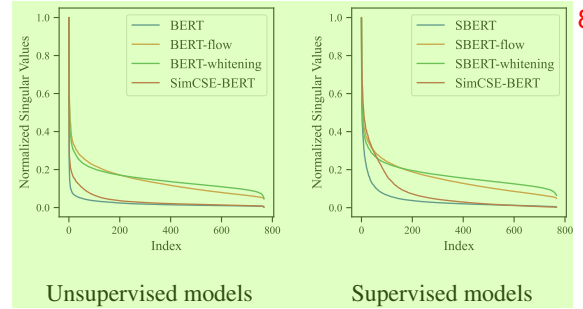


Figure F.1: Singular value distributions of sentence embedding matrix from sentences in STS-B. We normalize the singular values so that the largest one is 1.

human ratings in Figure G.1. Compared to all the baseline models, both unsupervised and supervised SimCSE better distinguish sentence pairs with different levels of similarities, thus leading to a better performance on STS tasks. In addition, we observe that SimCSE generally shows a more scattered distribution than BERT or SBERT, but also preserves a lower variance on semantically similar sentence pairs compared to whitened distribution. This observation further validates that SimCSE can achieve a better alignment-uniformity balance.

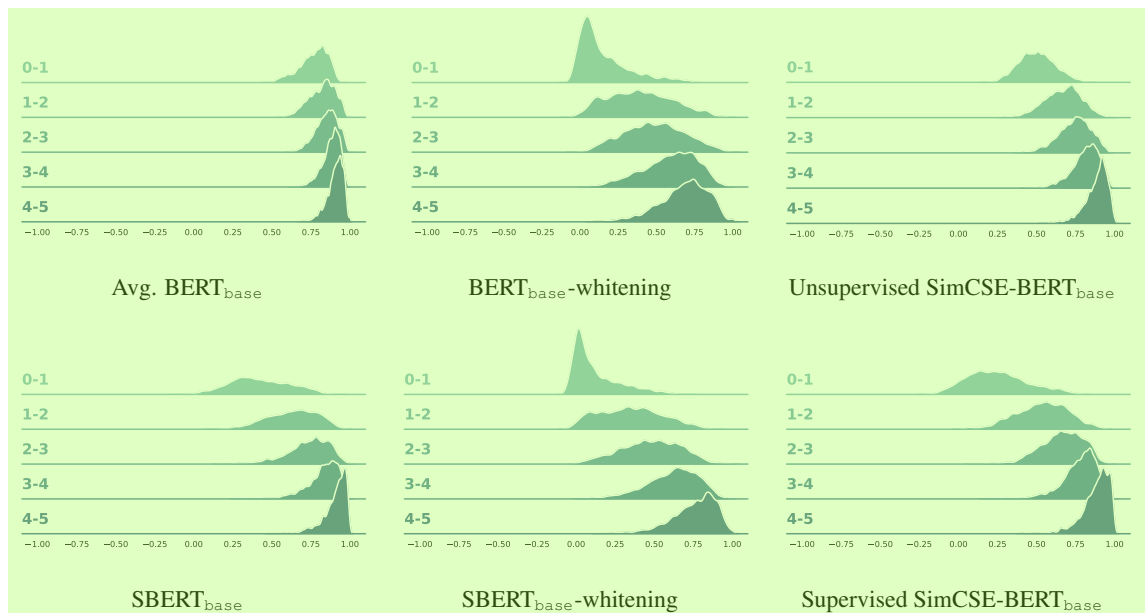


Figure G.1: Density plots of cosine similarities between sentence pairs in STS-B. Pairs are divided into 5 groups based on ground truth ratings (higher means more similar) along the y-axis, and x-axis is the cosine similarity.