

RankCSE : Unsupervised Sentence Representation Learning via Learning to Rank

Jiduan Liu^{1,2*}, Jiahao Liu³, Qifan Wang⁴, Jingang Wang³, Wei Wu³
Yunsen Xian³, Dongyan Zhao^{1,2,5,6†}, Kai Chen⁷, Rui Yan^{8,9†}

¹Wangxuan Institute of Computer Technology, Peking University

²Center for Data Science, AAIS, Peking University; ³Meituan; ⁴Meta AI

⁵National Key Laboratory of General Artificial Intelligence

⁶BIGAI, Beijing, China; ⁷School of Economics, Peking University

⁸Gaoling School of Artificial Intelligence, Renmin University of China

⁹Engineering Research Center of

Next-Generation Intelligent Search and Recommendation, Ministry of Education

{liujiduan, chen.kai, zhaody}@pku.edu.cn, ruiyan@ruc.edu.cn, wqfcr@fb.com

{liujiahao12, wangjingang02, xianyunsen}@meituan.com, wuwei19850318@gmail.com

Abstract

Unsupervised sentence representation learning is one of the fundamental problems in natural language processing with various downstream applications. Recently, contrastive learning has been widely adopted which derives high-quality sentence representations by pulling similar semantics closer and pushing dissimilar ones away. However, these methods fail to capture the fine-grained ranking information among the sentences, where each sentence is only treated as either positive or negative. In many real-world scenarios, one needs to distinguish and rank the sentences based on their similarities to a query sentence, e.g., very relevant, moderate relevant, less relevant, irrelevant, etc. In this paper, we propose a novel approach, RankCSE, for unsupervised sentence representation learning, which incorporates ranking consistency and ranking distillation with contrastive learning into a unified framework. In particular, we learn semantically discriminative sentence representations by simultaneously ensuring ranking consistency between two representations with different dropout masks, and distilling listwise ranking knowledge from the teacher. An extensive set of experiments are conducted on both semantic textual similarity (STS) and transfer (TR) tasks. Experimental results demonstrate the superior performance of our approach over several state-of-the-art baselines.

1 Introduction

Sentence representation learning refers to the task of encoding sentences into fixed-dimensional em-

* Work done during internship at Meituan.

† Corresponding authors: Dongyan Zhao (zhaody@pku.edu.cn) and Rui Yan (ruiyan@ruc.edu.cn).

Target Sentences	Label	SimCSE	RankCSE
• A woman is breaking eggs	4.80 (1)	0.93 (2)	0.97 (1)
• A man is cracking eggs	3.60 (2)	0.94 (1)	0.91 (2)
• A woman is talking to a man	1.60 (3)	0.45 (5)	0.65 (3)
• A man and a woman are speaking	1.40 (4)	0.47 (3)	0.61 (4)
• A man is talking to a boy	1.00 (5)	0.46 (4)	0.56 (5)
Query Sentence: A woman is cracking eggs			
• Broccoli are being cut by a woman	4.80 (1)	0.82 (2)	0.95 (1)
• A woman is slicing vegetables	4.20 (2)	0.83 (1)	0.91 (2)
• A woman is cutting some plants	3.50 (3)	0.74 (5)	0.87 (3)
• There is no woman cutting broccoli	3.40 (4)	0.76 (3)	0.85 (4)
• A woman is cutting some flowers	2.87 (5)	0.71 (7)	0.81 (5)
• A man is slicing tomatoes	2.60 (6)	0.75 (4)	0.79 (6)
• A man is cutting tomatoes	2.40 (7)	0.73 (6)	0.76 (7)
Query Sentence: A woman is cutting broccoli			

Table 1: Two examples of a query sentence and several target sentences from the STS datasets, with their similarity scores and rankings. The label scores are from human annotations. The SimCSE (Gao et al., 2021) and RankCSE similarity scores are from the model predictions respectively, with the corresponding ranking positions. It can be seen that sentence rankings based on SimCSE are incorrect, while RankCSE generates more effective scores with accurate rankings.

beddings. The sentence embeddings can be leveraged in various applications, including information retrieval (Le and Mikolov, 2014), text clustering (Ma et al., 2016) and semantic textual similarity comparison (Agirre et al., 2012). With the recent success of pre-trained language models (PLMs), such as BERT/RoBERTa (Devlin et al., 2019; Liu et al., 2019), a straightforward way to generate sentence representations is to directly use the [CLS] token embedding or the average token embeddings from the last layer of PLMs (Reimers and Gurevych, 2019). However, several studies (Ethayarajh, 2019; Li et al., 2020) have found that the native sentence representations derived by PLMs

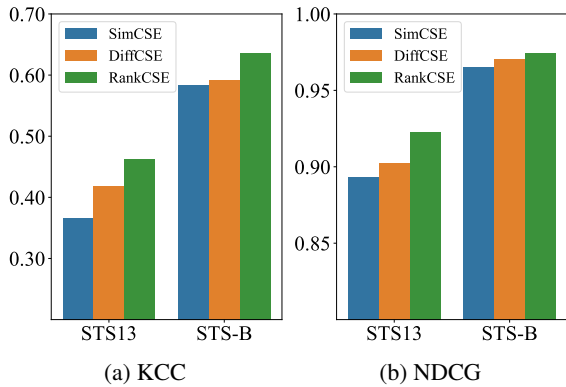


Figure 1: Sentence representation performances on ranking metrics KCC and NDCG (detailed in Appendix G). It can be seen that RankCSE captures more fine-grained ranking information than SimCSE (Gao et al., 2021) and DiffCSE (Chuang et al., 2022).

occupy a narrow cone in the vector space, and thus severely limits their representation capabilities, which is known as the *anisotropy* problem.

Supervised methods like SBERT (Reimers and Gurevych, 2019) usually generate better sentence representations, but require fine-tuning on a large amount of labeled data. Recent unsupervised models (Carlsson et al., 2021; Zhang et al., 2021; Giorgi et al., 2021) adopt contrastive learning framework without any labels, which pulls similar semantics closer and pushes dissimilar ones away. These methods usually design different augmentation algorithms for generating positive examples, such as back-translation (Zhang et al., 2021), dropout (Gao et al., 2021) and token shuffling or cutoff (Yan et al., 2021). In-batch negatives are further combined with the positives. Despite achieving promising results, they treat positives/negatives equally without capturing the fine-grained semantic ranking information, resulting in less effective sentence representations which fail to distinguish between very similar and less similar sentences. For example, Table 1 shows two examples of a query sentence and several target sentences from semantic textual similarity datasets. It is clear that the similarity scores produced by the contrastive learning method SimCSE are not optimized, where the sentence rankings are not preserved in the learned representations. On the other hand, our RankCSE generates effective sentence representations with consistent rankings to the ground-truth labels. Figure 1 further shows the advantage of RankCSE in terms of two ranking metrics. The fine-grained ranking information is crucial in various real-world

applications including search and recommendation. The ability to differentiate between subtle distinctions in sentence meaning can help these systems provide more relevant and accurate results, leading to a better user experience. Therefore, it is an important problem to learn ranking preserving sentence representations from unsupervised data.

To obtain semantically discriminative sentence representations, we propose a novel approach, RankCSE, which incorporates ranking consistency and ranking distillation with contrastive learning into a unified framework. Specifically, our model ensures ranking consistency between two representations with different dropout masks and minimizes the Jensen-Shannon (JS) divergence as the learning objective. In the meanwhile, our model also distills listwise ranking knowledge from the teacher model to the learned sentence representations. In our work, we explore two listwise ranking methods, ListNet (Cao et al., 2007) and ListMLE (Xia et al., 2008), and utilize the pre-trained SimCSE (Gao et al., 2021) models with coarse-grained semantic ranking information as the teachers to provide pseudo ranking labels. Our RankCSE is able to generalize fine-grained ranking information from the weak ranking knowledge learned by SimCSE. We conduct an extensive set of experiments on semantic textual similarity (STS) and transfer (TR) tasks. Experimental results show that RankCSE outperforms the existing state-of-the-art baselines.

2 Related Work

Unsupervised Sentence Representation Learning Early works typically augment the idea of word2vec (Mikolov et al., 2013) to learn sentence representations, including Skip-Thought (Kiros et al., 2015), FastSent (Hill et al., 2016) and Quick-Thought (Logeswaran and Lee, 2018). With the great success of PLMs, various attempts focus on generating sentence representations by leveraging the embedding of [CLS] token or applying mean pooling on the last layer of BERT (Reimers and Gurevych, 2019). However, Ethayarajh (2019) identifies the *anisotropy* problem in language representations, which means the native learned embeddings from PLMs occupy a narrow cone in the vector space. BERT-flow (Li et al., 2020) and BERT-whitening (Su et al., 2021) propose to resolve the *anisotropy* problem through post-processing.

Recently, contrastive learning has been adopted to learn sentence representations by designing dif-

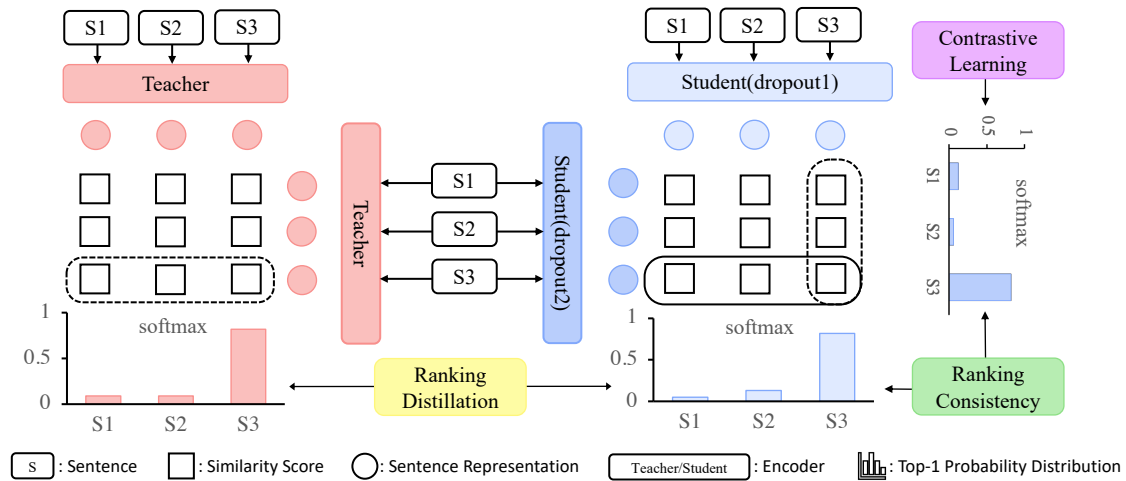


Figure 2: The framework of RankCSE which consists of three components: (1) contrastive learning object; (2) ranking consistency loss which ensures ranking consistency between two representations with different dropout masks; (3) ranking distillation loss which distills listwise ranking knowledge from the teacher.

ferent augmentation methods (Zhang et al., 2020; Carlsson et al., 2021; Giorgi et al., 2021; Yan et al., 2021; Kim et al., 2021; Gao et al., 2021). A typical example SimCSE uses dropout as a data augmentation strategy and is also the foundation of many following works. ArcCSE (Zhang et al., 2022) enhances the pairwise discriminative power and models the entailment relation among triplet sentences. DCLR (Zhou et al., 2022) alleviates the influence of improper negatives. DiffCSE (Chuang et al., 2022) introduces equivariant contrastive learning to SimCSE. PCL (Wu et al., 2022a) proposes contrastive representation learning with diverse augmentation strategies for an inherent anti-bias ability. InfoCSE (Wu et al., 2022b) learns sentence representations with the ability to reconstruct the original sentence fragments. Generative learning techniques (Wang et al., 2021; Wu and Zhao, 2022) have also been proposed to enhance the linguistic interpretability of sentence representations. Although achieving promising results, these methods fail to capture the fine-grained ranking knowledge among the sentences.

Learning to Rank Given a query example, learning to rank aims to rank a list of examples according to their similarities with the query. Learning to rank methods can be divided into three categories: pointwise (Li et al., 2007), pairwise (Burges et al., 2005, 2006) and listwise (Cao et al., 2007; Xia et al., 2008; Volkovs and Zemel, 2009; Pobrotyn and Bialobrzeski, 2021). Pointwise methods optimize the similarity between the query and each example, while pairwise approaches learn to cor-

rectly model the preference between two examples. Listwise methods directly evaluate the ranking of a list of examples based on the ground truth. In our framework, we leverage listwise ranking objectives for learning effective sentence representations, which have shown better performance compared to pointwise and pairwise methods.

3 Preliminary

We provide some conceptual explanations and definitions in learning to rank.

Top One Probability Given the scores of all objects $S = \{s_i\}_{i=1}^n$, the top one probability of an object is the probability of its being ranked at top-1: $\tilde{s}_i = \frac{e^{s_i/\tau}}{\sum_{j=1}^n e^{s_j/\tau}}$ where τ is a temperature hyperparameter, usually utilized to smooth the distribution. We simply denote the formulation for calculating the top one distribution based on the scores S as: $\tilde{S}_\tau = \sigma(S/\tau)$.

Permutation Probability Let $\pi = \{\pi(i)\}_{i=1}^n$ denote a permutation of the object indexes, which represents that the $\pi(i)$ -th sample is ranked i -th. The probability of a specific permutation π is given as: $P(\pi|S, \tau) = \prod_{i=1}^n \frac{e^{s_{\pi(i)}/\tau}}{\sum_{j=i}^n e^{s_{\pi(j)}/\tau}}$.

4 Methodology

4.1 Problem Formulation

Our goal is to learn sentence representations such that semantic similar sentences stay close while dissimilar ones should be far away in an unsupervised manner. Specifically, We aim to find an optimal

function f that maps a sentence $s \in p_s$ to a d -dimensional vector $f(s) \in p_e \subseteq \mathcal{R}^d$, where p_s and p_e denote the distributions of sentences and sentence representations, respectively. Supposing s_1 and s_2 are more semantic similar than s_1 and s_3 ($s_1, s_2, s_3 \in p_s$), a good mapping function f should satisfy that the distance between $f(s_1)$ and $f(s_2)$ is smaller than that between $f(s_1)$ and $f(s_3)$, i.e., $d(f(s_1), f(s_2)) < d(f(s_1), f(s_3))$, where d is the distance metric such as Euclidean distance and cosine distance. In this way, the similarities among the sentences are preserved in the learned sentence representations.

The general idea of RankCSE is to learn semantically discriminative sentence representations by capturing the ranking information among the sentences. As shown in Figure 2, our model consists of three components: (1) contrastive learning objective (Section 4.2); (2) ranking consistency loss which ensures ranking consistency between two representations with different dropout masks (Section 4.3); (3) ranking distillation loss which distills listwise ranking knowledge from the teacher (Section 4.4).

4.2 Contrastive Learning

Contrastive learning aims to learn effective representations by pulling similar semantics closer and pushing away dissimilar ones. SimCSE (Gao et al., 2021) creates positive examples by applying different dropout masks and takes a cross-entropy object with in-batch negatives (Chen et al., 2017). More specifically, for any sentence x_i in a min-batch, we send it to the encoder $f(\cdot)$ twice and obtain two representations with different dropout masks $f(x_i)$, $f(x_i)'$. SimCSE use the InfoNCE loss (van den Oord et al., 2018) as the training objective:

$$\mathcal{L}_{\text{InfoNCE}} = - \sum_{i=1}^N \log \frac{e^{\phi(f(x_i), f(x_i)')/\tau_1}}{\sum_{j=1}^N e^{\phi(f(x_i), f(x_j)')/\tau_1}}, \quad (1)$$

where N is the batch size, τ_1 is a temperature hyperparameter and $\phi(f(x_i), f(x_j)') = \frac{f(x_i)^\top f(x_j)'}{\|f(x_i)\| \cdot \|f(x_j)'\|}$ is the cosine similarity used in this work. Essentially, the contrastive learning objective is equivalent to maximizing the top one probability of the positive sample.

Although contrastive learning is effective in separating positive sentences with negative ones, it ignores the continuity modeling of the similarity. In other words, it is not effective in distinguishing highly similar sentences with moderate similar ones. To address this issue, we propose to

directly model the ranking information among the sentences, which could enhance the discrimination of semantic similarity in the learned sentence representations.

4.3 Ranking Consistency

The main drawback of contrastive learning is that the distinction between the in-batch negatives is not modeled, resulting in less effective sentence representations in capturing the fine-grained sentence similarity. Therefore, instead of treating the negatives equivalently, we propose to explicitly model the ranking information within the sentences by ensuring the ranking consistency between the two similarity sets (circled by the solid and dashed curves respectively in the right part of Figure 2).

Concretely, by taking a close look at the contrastive modeling in Section 4.2, there are two sets of sentence representations, $f(x_i)$ and $f(x_i)'$, derived from different dropout masks. For each sentence x_i , two lists of similarities with other sentences can be naturally obtained from the two representations, i.e., $S(x_i) = \{\phi(f(x_i), f(x_j)')\}_{j=1}^N$ and $S(x_i)' = \{\phi(f(x_i)', f(x_j))\}_{j=1}^N$. We then enforce the ranking consistency between these two similarity lists in our modeling. Intuitively, all corresponding elements in $S(x_i)$ and $S(x_i)'$ should have the same ranking positions.

Given two similarity lists $S(x_i)$ and $S(x_i)'$, we can obtain their top one probability distributions $\tilde{S}_{\tau_1}(x_i) = \sigma(S(x_i)/\tau_1)$, $\tilde{S}_{\tau_1}(x_i)' = \sigma(S(x_i)'/\tau_1)$. The ranking consistency can be ensured by minimizing the Jensen-Shannon (JS) divergence between the two top one probability distributions:

$$\begin{aligned} \mathcal{L}_{\text{consistency}} &= \sum_{i=1}^N \text{JS}(P_i \| Q_i) \\ &= \frac{1}{2} \sum_{i=1}^N (\text{KL}(P_i \| \frac{P_i + Q_i}{2}) + \text{KL}(Q_i \| \frac{P_i + Q_i}{2})) \quad (2) \\ &= \frac{1}{2} \sum_{i=1}^N (P_i \log(\frac{2P_i}{P_i + Q_i}) + Q_i \log(\frac{2Q_i}{P_i + Q_i})), \end{aligned}$$

where P_i and Q_i represents $\tilde{S}_{\tau_1}(x_i)$ and $\tilde{S}_{\tau_1}(x_i)'$ respectively. The reason we choose JS divergence instead of Kullback-Leibler (KL) divergence is that the two distributions are symmetric rather than one side being the ground truth.

4.4 Ranking Distillation

Contrastive learning based methods like SimCSE learn effective sentence representations with coarse-grained semantic ranking information (shown in

Appendix F and G), which have demonstrated their effectiveness in various downstream tasks. Orthogonal to ranking consistency, we further introduce ranking distillation by distilling the ranking knowledge from pre-trained teacher models into our learned sentence representations, to generalize effective ranking information from the weak ranking knowledge learned in the teachers. More specifically, for each sentence in a min-batch, we obtain the similarity score list from the teacher model, which is then served as pseudo ranking labels in the ranking distillation. The intuitive idea is to transfer the ranking knowledge from the teacher to the student as guidance for learning ranking preserved sentence representations. In the ranking distillation, ListNet (Cao et al., 2007) and ListMLE (Xia et al., 2008) methods are utilized. Formally they are defined as:

$$\mathcal{L}_{\text{rank}} = \sum_{i=1}^N \text{rank}(S(x_i), S^T(x_i)), \quad (3)$$

where $S(x_i)$ and $S^T(x_i)$ are the similarity score lists obtained from the student and the teacher, respectively, $\text{rank}(\cdot, \cdot)$ is the listwise method.

ListNet The original ListNet minimizes the cross entropy between the permutation probability distribution and the ground truth as the training objective. However, the computations will be intractable when the number of examples n is large, since the number of permutations is $n!$. To reduce the computation complexity, the top one probability distribution is usually adopted as a substitute:

$$\mathcal{L}_{\text{ListNet}} = - \sum_{i=1}^N \sigma(S^T(x_i)/\tau_3) \cdot \log(\sigma(S(x_i)/\tau_2)), \quad (4)$$

where τ_2 and τ_3 are temperature hyperparameters.¹

ListMLE Different from ListNet, ListMLE aims to maximize the likelihood of the ground truth permutation π_i^T which represents the sorted indexes of the similarity scores calculated by the teacher model. The objective of ListMLE is defined as:

$$\mathcal{L}_{\text{ListMLE}} = - \sum_{i=1}^N \log P(\pi_i^T | S(x_i), \tau_2). \quad (5)$$

In this work, we propose to use a multi-teacher from which more listwise ranking knowledge can

¹In practice, we exclude the score of the positive pair from the list to calculate the top one distribution used in Eq.(4), to enhance the ranking information of negatives, because the score of the positive pair occupies most in the full top one distribution calculated by the teacher SimCSE.

be transferred and preserved. In our experiments, we utilize the weighted average similarity scores of two teachers as pseudo ranking labels: $S^T(x_i) = \alpha S_1^T(x_i) + (1 - \alpha) S_2^T(x_i)$ where α is a hyperparameter to balance the weight of the teachers.

The contrastive learning loss $\mathcal{L}_{\text{infoNCE}}$ pushes apart the representations of different sentences to maximize the representation space, while the ranking consistency loss $\mathcal{L}_{\text{consistency}}$ and the ranking distillation loss $\mathcal{L}_{\text{rank}}$ pull similar negatives closer, thus capturing fine-grained semantic ranking information. Combining the above three loss functions, we can obtain the overall objective:

$$\mathcal{L}_{\text{final}} = \mathcal{L}_{\text{infoNCE}} + \beta \mathcal{L}_{\text{consistency}} + \gamma \mathcal{L}_{\text{rank}}, \quad (6)$$

where β and γ are hyperparameters to balance different losses.

5 Experiment

5.1 Setup

We evaluate our approach on two sentence related tasks, Semantic Textual Similarity (STS) and Transfer (TR). The SentEval toolkit (Conneau and Kiela, 2018) is used in our experiments. For STS tasks, we evaluate on seven datasets: STS12-16 (Agirre et al., 2012, 2013, 2014, 2015, 2016), STS Benchmark (Cer et al., 2017) and SICK-Relatedness (Marelli et al., 2014). These datasets contain pairs of sentences with similarity score labels from 0 to 5. Following SimCSE, we directly compute the cosine similarity between the sentence representations which means all the STS experiments are fully unsupervised, and report the Spearman’s correlation. For TR tasks, we evaluate on seven datasets with the default configurations from SentEval: 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). We use a logistic regression classifier trained on top of the frozen sentence representations, and report the classification accuracy.

For fair comparison, we use the same 10^6 randomly sampled sentences from English Wikipedia provided by SimCSE. Following previous works, we start from pre-trained checkpoints of BERT (Devlin et al., 2019) and RoBERTa (Liu et al., 2019), and utilize the embedding corresponding to [CLS] token as the representation of the input sentence. First we train SimCSE models including four variants: SimCSE-BERT_{base}, SimCSE-

PLMs	Methods	STS12	STS13	STS14	STS15	STS16	STS-B	SICK-R	avg.
Non-BERT	GloVe(avg.)	55.14	70.66	59.73	68.25	63.66	58.02	53.76	61.32
	USE	64.49	67.80	64.61	76.83	73.18	74.92	76.69	71.22
BERT _{base}	first-last avg.	39.70	59.38	49.67	66.03	66.19	53.87	62.06	56.70
	+flow	58.40	67.10	60.85	75.16	71.22	68.66	64.47	66.55
	+whitening	57.83	66.90	60.90	75.08	71.31	68.24	63.73	66.28
	+IS	56.77	69.24	61.21	75.23	70.16	69.21	64.25	66.58
	+ConSERT	64.64	78.49	69.07	79.72	75.95	73.97	67.31	72.74
	+SimCSE	68.40	82.41	74.38	80.91	78.56	76.85	72.23	76.25
	+DCLR	70.81	83.73	75.11	82.56	78.44	78.31	71.59	77.22
	+ArcCSE	72.08	84.27	76.25	82.32	79.54	79.92	72.39	78.11
	+DiffCSE	72.28	84.43	76.47	83.90	80.54	80.59	71.23	78.49
	+PaSeR	70.21	83.88	73.06	83.87	77.60	79.19	65.31	76.16
	+PCL	72.84	83.81	76.52	83.06	79.32	80.01	73.38	78.42
	+RankCSE _{listNet}	<u>74.38</u>	<u>85.97</u>	<u>77.51</u>	<u>84.46</u>	81.31	<u>81.46</u>	75.26	<u>80.05</u>
+RankCSE _{listMLE}	75.66	86.27	77.81	84.74	<u>81.10</u>	81.80	<u>75.13</u>	80.36	
BERT _{large}	+SimCSE	70.88	84.16	76.43	84.50	79.76	79.26	73.88	78.41
	+DCLR	71.87	84.83	77.37	84.70	79.81	79.55	74.19	78.90
	+ArcCSE	73.17	86.19	77.90	84.97	79.43	80.45	73.50	79.37
	+PCL	<u>74.87</u>	86.11	78.29	85.65	80.52	81.62	73.94	80.14
	+RankCSE _{listNet}	74.75	86.46	<u>78.52</u>	85.41	<u>80.62</u>	81.40	76.12	<u>80.47</u>
	+RankCSE _{listMLE}	75.48	86.50	78.60	<u>85.45</u>	81.09	<u>81.58</u>	<u>75.53</u>	80.60
RoBERTa _{base}	+SimCSE	70.16	81.77	73.24	81.36	80.65	80.22	68.56	76.57
	+DCLR	70.01	83.08	75.09	83.66	81.06	81.86	70.33	77.87
	+DiffCSE	70.05	83.43	75.49	82.81	82.12	82.38	71.19	78.21
	+PCL	71.13	82.38	75.40	83.07	81.98	81.63	69.72	77.90
	+RankCSE _{listNet}	<u>72.91</u>	<u>85.72</u>	<u>76.94</u>	<u>84.52</u>	82.59	83.46	71.94	<u>79.73</u>
	+RankCSE _{listMLE}	73.20	85.95	77.17	84.82	<u>82.58</u>	<u>83.08</u>	<u>71.88</u>	79.81
RoBERTa _{large}	+SimCSE	72.86	83.99	75.62	84.77	81.80	81.98	71.26	78.90
	+DCLR	73.09	84.57	76.13	85.15	81.99	82.35	71.80	79.30
	+PCL	74.08	84.36	76.42	85.49	81.76	82.79	71.51	79.49
	+RankCSE _{listNet}	<u>73.47</u>	<u>85.77</u>	78.07	85.65	<u>82.51</u>	<u>84.12</u>	73.73	80.47
	+RankCSE _{listMLE}	73.20	85.83	<u>78.00</u>	<u>85.63</u>	82.67	84.19	<u>73.64</u>	80.45

Table 2: Sentence representations performance on STS tasks (Spearman’s correlation). We directly import the results from the original papers and mark the best (bold) and second-best (underlined) results among models with the same PLMs. Results are statistically significant with respect to all baselines on each PLM (all p-value < 0.005).

BERT_{large}, SimCSE-RoBERTa_{base} and SimCSE-RoBERTa_{large}. We utilize the first two as a multi-teacher for RankCSE-BERT_{base} and RankCSE-BERT_{large}, while the last two for RankCSE-RoBERTa_{base} and RankCSE-RoBERTa_{large}. We evaluate our model every 125 training steps on the dev set of STS-B and keep the best checkpoint for the evaluation on test sets of all STS and TR tasks. More training details can be found in Appendix A.

We compare RankCSE with several unsupervised sentence representation learning methods, including average GloVe embeddings (Pennington et al., 2014), USE (Cer et al., 2018) and Skip-thought (Kiros et al., 2015), average BERT embeddings from the last layer, post-processing methods such as BERT-flow (Li et al., 2020) and BERT-whitening (Su et al., 2021), and contrastive learning methods such as IS-BERT (Zhang et al., 2020) and ConSERT (Yan et al., 2021). We also include recent strong unsupervised sentence representation baselines, including SimCSE (Gao et al., 2021),

DCLR (Zhou et al., 2022), ArcCSE (Zhang et al., 2022), DiffCSE (Chuang et al., 2022), PaSER (Wu and Zhao, 2022) and PCL (Wu et al., 2022a). Since RankCSE and the teacher model SimCSE are using the same unsupervised training data, the comparison between RankCSE and baselines is fair.

5.2 Main Results

Results on STS Tasks As shown in Table 2, it is clear that RankCSE significantly outperforms the previous methods on all PLMs, which demonstrates the effectiveness of our approach. For example, compared with SimCSE, RankCSE has brought noticeable improvements: 4.11% on BERT_{base}, 2.19% on BERT_{large}, 3.24% on RoBERTa_{base} and 1.57% on RoBERTa_{large}. RankCSE-BERT_{base} even outperforms SimCSE-BERT_{large} by nearly 2%. Compared with the previous state-of-the-art methods, RankCSE still achieves consistent improvements, which validates that RankCSE is able to obtain more semantically discriminative

PLMs	Methods	MR	CR	SUBJ	MPQA	SST	TREC	MRPC	avg.
Non-BERT	GloVe(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
BERT _{base}	last avg.	78.66	86.25	94.37	88.66	84.40	92.80	69.54	84.94
	+IS	81.09	87.18	94.96	88.75	85.96	88.64	74.24	85.83
	+SimCSE	81.18	86.46	94.45	88.88	85.50	89.80	74.43	85.81
	+ArcCSE	79.91	85.25	99.58	89.21	84.90	89.20	74.78	86.12
	+DiffCSE [†]	81.76	86.20	94.76	89.21	86.00	87.60	75.54	85.87
	+PCL	80.11	85.25	94.22	89.15	85.12	87.40	76.12	85.34
	+RankCSE _{listNet}	83.21	<u>88.08</u>	<u>95.25</u>	90.00	88.58	<u>90.00</u>	<u>76.17</u>	87.33
+RankCSE _{listMLE}	<u>83.07</u>	88.27	<u>95.06</u>	<u>89.90</u>	<u>87.70</u>	89.40	76.23	<u>87.09</u>	
BERT _{large}	+SimCSE	85.36	89.38	95.39	89.63	90.44	91.80	76.41	88.34
	+ArcCSE	84.34	88.82	99.58	89.79	90.50	92.00	74.78	88.54
	+PCL	82.47	87.87	95.04	89.59	87.75	93.00	76.00	87.39
	+RankCSE _{listNet}	<u>85.11</u>	89.56	95.39	90.30	90.77	93.20	77.16	88.78
+RankCSE _{listMLE}	84.63	<u>89.51</u>	<u>95.50</u>	<u>90.08</u>	<u>90.61</u>	93.20	<u>76.99</u>	<u>88.65</u>	
RoBERTa _{base}	+SimCSE	81.04	87.74	93.28	86.94	86.60	84.60	73.68	84.84
	+DiffCSE [†]	82.42	88.34	93.51	87.28	87.70	86.60	76.35	86.03
	+PCL	81.83	87.55	92.92	87.21	87.26	85.20	<u>76.46</u>	85.49
	+RankCSE _{listNet}	83.53	89.22	94.07	88.97	89.95	89.20	76.52	87.35
+RankCSE _{listMLE}	<u>83.32</u>	<u>88.61</u>	<u>94.03</u>	<u>88.88</u>	<u>89.07</u>	90.80	<u>76.46</u>	<u>87.31</u>	
RoBERTa _{large}	+SimCSE	82.74	87.87	93.66	88.22	88.58	92.00	69.68	86.11
	+PCL	84.47	89.06	94.60	89.26	89.02	94.20	74.96	87.94
	+RankCSE _{listNet}	<u>84.47</u>	89.51	94.65	<u>89.87</u>	<u>89.46</u>	<u>93.00</u>	75.88	88.12
	+RankCSE _{listMLE}	84.61	<u>89.27</u>	94.47	89.99	89.73	92.60	74.43	87.87

Table 3: Sentence representations performance on transfer tasks (accuracy). The results of DiffCSE[†] are obtained from the publicly available code and checkpoints, while others are imported from the original papers. We mark the best (bold) and second-best (underlined) results among models with the same PLMs. Results are statistically significant with respect to all baselines on each PLM (all p-value < 0.005).

Models	STS(avg.)	TR(avg.)
SimCSE	76.25	85.81
RankCSE _{listNet}	80.05	87.33
w/o $\mathcal{L}_{\text{consistency}}$	79.56	86.80
w/o $\mathcal{L}_{\text{infoNCE}}$	79.72	86.91
w/o $\mathcal{L}_{\text{consistency}}, \mathcal{L}_{\text{infoNCE}}$	79.41	86.76
RankCSE _{listMLE}	80.36	87.09
w/o $\mathcal{L}_{\text{consistency}}$	79.88	86.65
w/o $\mathcal{L}_{\text{infoNCE}}$	79.95	86.73
w/o $\mathcal{L}_{\text{consistency}}, \mathcal{L}_{\text{infoNCE}}$	79.73	86.24
RankCSE w/o $\mathcal{L}_{\text{rank}}$	76.93	85.97
RankCSE w/o $\mathcal{L}_{\text{infoNCE}}, \mathcal{L}_{\text{rank}}$	73.74	85.56

Table 4: Ablation studies of different loss functions based on BERT_{base}. Other PLMs yield similar patterns to BERT_{base}.

representations by incorporating ranking consistency and ranking distillation. We also observe that the performances of RankCSE_{listNet} and RankCSE_{listMLE} are very consistent across all datasets, which demonstrates the effectiveness of both listwise ranking methods.

Results on TR Tasks It can be seen in Table 3 that RankCSE achieves the best performance among all the compared baselines on all PLMs.

Teacher	RankCSE	
	ListNet	ListMLE
SimCSE _{base}	77.48	77.75
DiffCSE _{base}	78.87	79.06
SimCSE _{large}	79.66	79.81
SimCSE _{base} +DiffCSE _{base}	79.10	79.28
SimCSE _{base} +SimCSE _{large}	80.05	80.36
DiffCSE _{base} +SimCSE _{large}	80.20	80.47

Table 5: Comparisons of different teachers based on BERT. Results of RankCSE are average STS performance using BERT_{base}.

Note that for DiffCSE, we obtain the results from the publicly available code and checkpoints, because DiffCSE uses different dev sets to find the best hyperparameters for TR tasks than other baselines. More detailed explanation and comprehensive comparison are provided in Appendix B. Another observation is that the performance of the RankCSE_{listNet} is slightly better than that of the RankCSE_{listMLE}. Our hypothesis is that the inaccurate pseudo ranking labels introduce more errors in the calculation of the permutation probability than the top one probability. Nevertheless, both listwise methods achieve better results than the baselines, which is consistent with the results in Table 2.

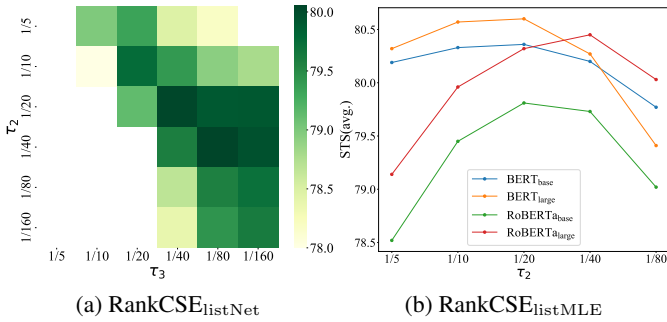


Figure 3: Effect of the temperatures τ_2 and τ_3 . Results are average STS performance, and RankCSE_{listNet} is based on BERT_{base} while RankCSE_{listMLE} is based on different PLMs. We do not demonstrate results below 78 to make the variation obvious.

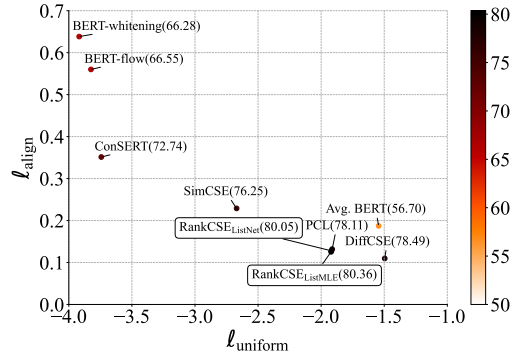


Figure 4: $\ell_{\text{align}} - \ell_{\text{uniform}}$ plot for different sentence representation methods based on BERT_{base} measured on the STS-B dev set. Color of points represents average STS performance.

PLMs	RankCSE _{listNet}		RankCSE _{listMLE}		SimCSE	
	STS(avg.)	TR(avg.)	STS(avg.)	TR(avg.)	STS(avg.)	TR(avg.)
BERT _{base}	80.00±0.13	87.28±0.19	80.39±0.04	87.05±0.06	75.52±0.70	85.44±0.47
BERT _{large}	80.41±0.10	88.74±0.14	80.59±0.05	88.63±0.06	77.79±0.64	88.10±0.36
RoBERTa _{base}	79.67±0.09	87.46±0.13	79.78±0.05	87.30±0.07	76.45±0.56	84.74±0.38
RoBERTa _{large}	80.46±0.11	87.97±0.14	80.34±0.08	87.82±0.08	78.53±0.49	86.29±0.33

Table 6: Mean and standard deviation across five different runs of RankCSE and SimCSE.

5.3 Analysis and Discussion

Ablation Study To investigate the impact of different losses in our approach, we conduct a set of ablation studies by removing $\mathcal{L}_{\text{infoNCE}}$, $\mathcal{L}_{\text{consistency}}$ and $\mathcal{L}_{\text{rank}}$ from Eq.(6). The average results on STS and TR tasks are reported in Table 4. There are several observations from the results. First, when $\mathcal{L}_{\text{rank}}$ is removed, the performance significantly drops in both STS and TR tasks, which indicates the effectiveness of $\mathcal{L}_{\text{rank}}$ in our modeling. Second, it is also clear that without $\mathcal{L}_{\text{infoNCE}}$ or $\mathcal{L}_{\text{consistency}}$, the model performance also decreases, especially on TR tasks. Thirdly, it is worth mentioning that RankCSE with only $\mathcal{L}_{\text{rank}}$ can also outperform the teachers on STS tasks. The reason is that RankCSE is able to preserve ranking knowledge from multiple teachers, and generalize fine-grained ranking information from multiple coarse-grained representations. Fourthly, since $\mathcal{L}_{\text{consistency}}$ does not explicitly distinguish the positives from negatives, RankCSE with only $\mathcal{L}_{\text{consistency}}$ will preserve inaccurate rankings leading to significant performance drop. Finally, RankCSE with all components achieves the best performance on both STS and TR tasks.

Comparisons of Different Teachers We conduct experiments to explore the impact of dif-

ferent teachers on the performance of RankCSE. As shown in Table 5, RankCSE outperforms the teacher model which indicates that incorporating ranking consistency and ranking distillation leads to more semantically discriminative sentence representations. Comparing the performance of RankCSE using different teachers, we observe that better teacher leads to better RankCSE, which is consistent with our expectation since accurate ranking labels yield more effective ranking knowledge transfer. Another observation is that the performance of RankCSE with a multi-teacher is better than that with a single teacher, which verifies that RankCSE is able to preserve listwise ranking knowledge from more than one teacher. It is also interesting to see that using DiffCSE-BERT_{base} and SimCSE-BERT_{large} as multi-teacher leads to even higher performance than the results in Table 2. We plan to conduct more investigation along this direction to explore the upper bound of improvements.

Effect of Hyperparameters To study the effect of temperature hyperparameters, we conduct experiments by setting different τ_2 and τ_3 . As shown in Figure 3a, we find that large discrepancy between τ_2 and τ_3 leads to significant drop in the performance of RankCSE_{listNet}. The best temperature setting for RankCSE_{listNet} is $\tau_2 : \tau_3 =$

2 : 1. The performance of RankCSE_{ListMLE} has similar trends based on different PLMs, as shown in Figure 3b. For both RankCSE_{ListNet} and RankCSE_{ListMLE}, the temperature should be set moderate.

Robustness of RankCSE We conduct 5 runs of model training with the hyperparameter settings which can be referred to Appendix A with different random seeds, and then calculate the mean and standard deviation values. The results provided in Table 6 demonstrate both the superior performance and the robustness of our model. It can also be seen that RankCSE_{ListMLE} achieves similar performance but more stable results compared with RankCSE_{ListNet}.

Alignment and Uniformity Following previous works (Wang and Isola, 2020), we use alignment and uniformity to measure the quality of representation space. Alignment measures the distance between similar instances, while uniformity measures how well the representations are uniformly distributed (detailed in Appendix H). For both measures, the smaller value indicates the better result. We plot the distribution of $\ell_{\text{align}} - \ell_{\text{uniform}}$ for different models using BERT_{base} which are measured on the STS-B dev set. As shown in Figure 4, RankCSE effectively improves both alignment and uniformity compared with average BERT embeddings, while SimCSE and DiffCSE only improve uniformity and alignment respectively. Since RankCSE pulls similar negatives closer during incorporating ranking consistency and ranking distillation, RankCSE has smaller alignment and bigger uniformity than SimCSE. We consider that RankCSE achieves a better trade-off than SimCSE. When compared with DiffCSE, RankCSE has smaller uniformity whereas similar alignment. We can also observe that RankCSE outperforms PCL on both metrics.

6 Conclusion

In this work, we propose RankCSE, an unsupervised approach to learn more semantically discriminative sentence representations. The core idea of RankCSE is incorporating ranking consistency and ranking distillation with contrastive learning into a unified framework. When simultaneously ensuring ranking consistency and distilling listwise ranking knowledge from the teacher, RankCSE can learn how to make fine-grained distinctions in semantics,

leading to more semantically discriminative sentence representations. Experimental results on STS and TR tasks demonstrate that RankCSE outperforms previous state-of-the-art methods. We also conduct thorough ablation study and analysis to demonstrate the effectiveness of each component and justify the inner workings of our approach. We leave what is the upper bound of improvements of the teacher for future work.

Limitations

In this section, we discuss the limitations of our work as follows. First, despite achieving promising results, our model needs to calculate pseudo ranking labels of the teacher which requires additional training time per epoch than the teacher. The training efficiency of RankCSE and SimCSE can be seen in Appendix D. Second, we directly use SimCSE_{base} and SimCSE_{large} as a multi-teacher in our implementation and experiments. However, how to choose the best combination of the teacher models is worth further exploration. It could help researchers to better understand the upper bound of improvements. We plan to investigate more along this direction in the future.

Acknowledgements

This work is supported by Ministry of Science and Technology Key R&D Program (2030 Artificial Intelligence) (No. 2020AAA0106600) and National Natural Science Foundation of China (NSFC Grant No. 62122089). We sincerely thank all reviewers for their valuable comments and suggestions, which are crucial for improving our work. We would also like to acknowledge Angela Li for her contributions in creating the figures used in this work.

References

- Hervé Abdi. 2007. [The kendall rank correlation coefficient](#). *Encyclopedia of measurement and statistics*, 2:508–510.
- Eneko Agirre, Carmen Banea, Claire Cardie, Daniel M. Cer, Mona T. Diab, Aitor Gonzalez-Agirre, Weiwei Guo, Iñigo Lopez-Gazpio, Montse Maritxalar, Rada Mihalcea, German Rigau, Larratiz Uriia, 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@NAACL-HLT 2015, Denver, Colorado, USA, June 4-5, 2015*,

- pages 252–263. The Association for Computer Linguistics.
- Eneko Agirre, Carmen Banea, Claire Cardie, Daniel M. Cer, Mona T. 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@COLING 2014, Dublin, Ireland, August 23-24, 2014*, pages 81–91. The Association for Computer Linguistics.
- Eneko Agirre, Carmen Banea, Daniel M. Cer, Mona T. 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@NAACL-HLT 2016, San Diego, CA, USA, June 16-17, 2016*, pages 497–511. The Association for Computer Linguistics.
- Eneko Agirre, Daniel M. Cer, Mona T. Diab, and Aitor Gonzalez-Agirre. 2012. [Semeval-2012 task 6: A pilot on semantic textual similarity](#). In *Proceedings of the 6th International Workshop on Semantic Evaluation, SemEval@NAACL-HLT 2012, Montréal, Canada, June 7-8, 2012*, pages 385–393. The Association for Computer Linguistics.
- Eneko Agirre, Daniel M. Cer, Mona T. Diab, Aitor Gonzalez-Agirre, and Weiwei Guo. 2013. [*sem 2013 shared task: Semantic textual similarity](#). In *Proceedings of the Second Joint Conference on Lexical and Computational Semantics, *SEM 2013, June 13-14, 2013, Atlanta, Georgia, USA*, pages 32–43. Association for Computational Linguistics.
- Christopher J. C. Burges, Robert Ragno, and Quoc Viet Le. 2006. [Learning to rank with nonsmooth cost functions](#). In *Advances in Neural Information Processing Systems 19, Proceedings of the Twentieth Annual Conference on Neural Information Processing Systems, Vancouver, British Columbia, Canada, December 4-7, 2006*, pages 193–200. MIT Press.
- Christopher J. C. Burges, Tal Shaked, Erin Renshaw, Ari Lazier, Matt Deeds, Nicole Hamilton, and Gregory N. Hullender. 2005. [Learning to rank using gradient descent](#). In *Machine Learning, Proceedings of the Twenty-Second International Conference (ICML 2005), Bonn, Germany, August 7-11, 2005*, volume 119 of *ACM International Conference Proceeding Series*, pages 89–96. ACM.
- Zhe Cao, Tao Qin, Tie-Yan Liu, Ming-Feng Tsai, and Hang Li. 2007. [Learning to rank: from pairwise approach to listwise approach](#). In *Machine Learning, Proceedings of the Twenty-Fourth International Conference (ICML 2007), Corvallis, Oregon, USA, June 20-24, 2007*, volume 227 of *ACM International Conference Proceeding Series*, pages 129–136. ACM.
- Fredrik Carlsson, Amaru Cuba Gyllensten, Evangelia Gogoulou, Erik Ylipää Hellqvist, and Magnus Sahlgren. 2021. [Semantic re-tuning with contrastive tension](#). In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net.
- 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, EMNLP 2018: System Demonstrations, Brussels, Belgium, October 31 - November 4, 2018*, pages 169–174. Association for Computational Linguistics.
- Daniel M. Cer, Mona T. 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@ACL 2017, Vancouver, Canada, August 3-4, 2017*, pages 1–14. Association for Computational Linguistics.
- Ting Chen, Yizhou Sun, Yue Shi, and Liangjie Hong. 2017. [On sampling strategies for neural network-based collaborative filtering](#). In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Halifax, NS, Canada, August 13 - 17, 2017*, pages 767–776. ACM.
- Yung-Sung Chuang, Rumen Dangovski, Hongyin Luo, Yang Zhang, Shiyu Chang, Marin Soljagic, Shang-Wen Li, Scott Yih, Yoon Kim, and James R. Glass. 2022. [Diffcse: Difference-based contrastive learning for sentence embeddings](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL 2022, Seattle, WA, United States, July 10-15, 2022*, pages 4207–4218. 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, May 7-12, 2018*. European Language Resources Association (ELRA).
- 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, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 4171–4186. Association for Computational Linguistics.
- William B. Dolan and Chris Brockett. 2005. [Automatically constructing a corpus of sentential paraphrases](#). In *Proceedings of the Third International Workshop on Paraphrasing, IWP@IJCNLP 2005, Jeju Island,*

- Korea, October 2005, 2005. Asian Federation of Natural Language Processing.
- 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 2019, Hong Kong, China, November 3-7, 2019*, pages 55–65. 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, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, pages 6894–6910. Association for Computational Linguistics.
- John M. Giorgi, Osvald Nitski, Bo Wang, and Gary D. Bader. 2021. [Declutr: Deep contrastive learning for unsupervised textual representations](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, pages 879–895. Association for Computational Linguistics.
- Felix Hill, Kyunghyun Cho, and Anna Korhonen. 2016. [Learning distributed representations of sentences from unlabelled data](#). In *NAACL HLT 2016, The 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, San Diego California, USA, June 12-17, 2016*, pages 1367–1377. The Association for Computational Linguistics.
- 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, Seattle, Washington, USA, August 22-25, 2004*, pages 168–177. ACM.
- Kalervo Järvelin and Jaana Kekäläinen. 2002. [Cumulated gain-based evaluation of IR techniques](#). *ACM Trans. Inf. Syst.*, 20(4):422–446.
- Taeuk Kim, Kang Min Yoo, and Sang-goo Lee. 2021. [Self-guided contrastive learning for BERT sentence representations](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, pages 2528–2540. Association for Computational Linguistics.
- Ryan Kiros, Yukun Zhu, Ruslan Salakhutdinov, Richard S. Zemel, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. [Skip-thought vectors](#). In *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, pages 3294–3302.
- Quoc V. Le and Tomáš Mikolov. 2014. [Distributed representations of sentences and documents](#). In *Proceedings of the 31th International Conference on Machine Learning, ICML 2014, Beijing, China, 21-26 June 2014*, volume 32 of *JMLR Workshop and Conference Proceedings*, pages 1188–1196. JMLR.org.
- 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 2020, Online, November 16-20, 2020*, pages 9119–9130. Association for Computational Linguistics.
- Ping Li, Christopher J. C. Burges, and Qiang Wu. 2007. [Mcrank: Learning to rank using multiple classification and gradient boosting](#). In *Advances in Neural Information Processing Systems 20, Proceedings of the Twenty-First Annual Conference on Neural Information Processing Systems, Vancouver, British Columbia, Canada, December 3-6, 2007*, pages 897–904. Curran Associates, Inc.
- 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](#). *CoRR*, abs/1907.11692.
- Lajanugen Logeswaran and Honglak Lee. 2018. [An efficient framework for learning sentence representations](#). In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net.
- Shutian Ma, Chengzhi Zhang, and Daqing He. 2016. [Document representation methods for clustering bilingual documents](#). In *Creating Knowledge, Enhancing Lives through Information & Technology - Proceedings of the 2016 Annual Meeting of the Association for Information Science and Technology, ASIST 2016, Copenhagen, Denmark, October 14-18, 2016*, volume 53 of *Proc. Assoc. Inf. Sci. Technol.*, pages 1–10. Wiley.
- 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 2014, Reykjavik, Iceland, May 26-31, 2014*, pages 216–223. European Language Resources Association (ELRA).
- Tomáš Mikolov, Ilya Sutskever, Kai Chen, Gregory S. Corrado, and Jeffrey Dean. 2013. [Distributed representations of words and phrases and their compositionality](#). In *Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural*

- Information Processing Systems 2013. Proceedings of a meeting held December 5-8, 2013, Lake Tahoe, Nevada, United States*, pages 3111–3119.
- 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 of the Association for Computational Linguistics, 21-26 July, 2004, Barcelona, Spain*, pages 271–278. ACL.
- Bo Pang and Lillian Lee. 2005. [Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales](#). In *ACL 2005, 43rd Annual Meeting of the Association for Computational Linguistics, Proceedings of the Conference, 25-30 June 2005, University of Michigan, USA*, pages 115–124. The Association for Computer Linguistics.
- Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. [Glove: Global vectors for word representation](#). In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, October 25-29, 2014, Doha, Qatar, A meeting of SIGDAT, a Special Interest Group of the ACL*, pages 1532–1543. ACL.
- Przemyslaw Pobrotyn and Radoslaw Bialobrzeski. 2021. [NeuralIndcg: Direct optimisation of a ranking metric via differentiable relaxation of sorting](#). *CoRR*, abs/2102.07831.
- 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 2019, Hong Kong, China, November 3-7, 2019*, pages 3980–3990. Association for Computational Linguistics.
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Y. 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, EMNLP 2013, 18-21 October 2013, Grand Hyatt Seattle, Seattle, Washington, USA, A meeting of SIGDAT, a Special Interest Group of the ACL*, pages 1631–1642. ACL.
- Jianlin Su, Jiarun Cao, Weijie Liu, and Yangyiwen Ou. 2021. [Whitening sentence representations for better semantics and faster retrieval](#). *CoRR*, abs/2103.15316.
- Aäron van den Oord, Yazhe Li, and Oriol Vinyals. 2018. [Representation learning with contrastive predictive coding](#). *CoRR*, abs/1807.03748.
- Maksims Volkovs and Richard S. Zemel. 2009. [Boltzrank: learning to maximize expected ranking gain](#). In *Proceedings of the 26th Annual International Conference on Machine Learning, ICML 2009, Montreal, Quebec, Canada, June 14-18, 2009*, volume 382 of *ACM International Conference Proceeding Series*, pages 1089–1096. ACM.
- Ellen M. Voorhees and Dawn M. Tice. 2000. [Building a question answering test collection](#). In *SIGIR 2000: Proceedings of the 23rd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, July 24-28, 2000, Athens, Greece*, pages 200–207. ACM.
- Kexin Wang, Nils Reimers, and Iryna Gurevych. 2021. [TSDAE: using transformer-based sequential denoising auto-encoder for unsupervised sentence embedding learning](#). In *Findings of the Association for Computational Linguistics: EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 16-20 November, 2021*, pages 671–688. Association for Computational Linguistics.
- Tongzhou Wang and Phillip Isola. 2020. [Understanding contrastive representation learning through alignment and uniformity on the hypersphere](#). In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 9929–9939. PMLR.
- Janyce Wiebe, Theresa Wilson, and Claire Cardie. 2005. [Annotating expressions of opinions and emotions in language](#). *Lang. Resour. Evaluation*, 39(2-3):165–210.
- Bohong Wu and Hai Zhao. 2022. [Sentence representation learning with generative objective rather than contrastive objective](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022*, pages 3356–3368. Association for Computational Linguistics.
- Qiyu Wu, Chongyang Tao, Tao Shen, Can Xu, Xiubo Geng, and Daxin Jiang. 2022a. [PCL: peer-contrastive learning with diverse augmentations for unsupervised sentence embeddings](#). *CoRR*, abs/2201.12093.
- Xing Wu, Chaochen Gao, Zijia Lin, Jizhong Han, Zhongyuan Wang, and Songlin Hu. 2022b. [Infocse: Information-aggregated contrastive learning of sentence embeddings](#). *CoRR*, abs/2210.06432.
- Fen Xia, Tie-Yan Liu, Jue Wang, Wensheng Zhang, and Hang Li. 2008. [Listwise approach to learning to rank: theory and algorithm](#). In *Machine Learning, Proceedings of the Twenty-Fifth International Conference (ICML 2008), Helsinki, Finland, June 5-9, 2008*, volume 307 of *ACM International Conference Proceeding Series*, pages 1192–1199. ACM.
- 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, ACL/IJCNLP*

2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021, pages 5065–5075. Association for Computational Linguistics.

Yan Zhang, Ruidan He, Zuozhu Liu, Lidong Bing, and Haizhou Li. 2021. [Bootstrapped unsupervised sentence representation learning](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, pages 5168–5180. Association for Computational Linguistics.

Yan Zhang, Ruidan He, Zuozhu Liu, Kwan Hui Lim, and Lidong Bing. 2020. [An unsupervised sentence embedding method by mutual information maximization](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*, pages 1601–1610. 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 Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, pages 4892–4903. Association for Computational Linguistics.

Kun Zhou, Beichen Zhang, Xin Zhao, and Ji-Rong Wen. 2022. [Debiased contrastive learning of unsupervised sentence representations](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, pages 6120–6130. Association for Computational Linguistics.

A Training Details

We implement all experiments with the deep learning framework PyTorch on a single NVIDIA Tesla A100 GPU (40GB memory). We carry out grid-search of learning rate $\in \{2e-5, 3e-5\}$ and temperatures $\tau_2, \tau_3 \in \{0.0125, 0.025, 0.05\}$, while setting batch size to 128, temperature τ_1 to 0.05, α to 1/3, β to 1, γ to 1 and the rate of linear scheduling warm-up to 0.05 for all the experiments. We train our models for 4 epochs, and evaluate the model every 125 steps on the dev set of STS-B and keep the best checkpoint for the final evaluation on test sets of all STS and TR tasks. The hyperparameter settings we adopt are shown in Table 9. Following SimCSE, we utilize the embedding corresponding to [CLS] token as the representation of the input sentence. We utilize SimCSE-BERT_{base} and SimCSE-BERT_{large} as a multi-teacher for RankCSE-BERT_{base} and RankCSE-BERT_{large}, while SimCSE-RoBERTa_{base} and SimCSE-RoBERTa_{large} as a multi-teacher for RankCSE-RoBERTa_{base} and RankCSE-RoBERTa_{large}.

B DiffCSE Settings for Transfer Tasks

DiffCSE uses different dev sets to find the best hyperparameters for the two tasks (STS-B dev set for STS tasks, dev sets of 7 TR tasks for TR tasks), while other methods only use the STS-B dev set for both tasks, which is not fair. Therefore we obtain the results in Table 3 from its publicly available code and checkpoints for STS tasks² instead of directly importing the results from its original paper. For a more comprehensive comparison with DiffCSE on TR tasks, we also use dev sets of 7 TR tasks to find the best hyperparameters and checkpoints. As shown in Table 10, RankCSE still outperforms DiffCSE in this setting.

C Data Statistics

The complete listings of train/dev/test stats of STS and TR datasets can be found in Table 7 and 8, respectively. Note that for STS tasks, we only use test sets for the final evaluation and dev set of STS-B to find best hyperparameters and checkpoints. The train sets of all STS datasets are not used in our experiments. For TR tasks, we follow the default settings of SentEval toolkit (Conneau and Kiela, 2018) to use 10-fold evaluation for all TR datasets

²<https://github.com/voidism/DiffCSE>

Dataset	Train	Dev	Test
STS12	-	-	3108
STS13	-	-	1500
STS14	-	-	3750
STS15	-	-	3000
STS16	-	-	1186
STS-B	5749	1500	1379
SICK-R	4500	500	4927

Table 7: A listing of train/dev/test stats of STS datasets.

Dataset	Train	Dev	Test
MR	10662	-	-
CR	3775	-	-
SUBJ	10000	-	-
MPQA	10606	-	-
SST	67349	872	1821
TREC	5452	-	500
MRPC	4076	-	1725

Table 8: A listing of train/dev/test stats of TR datasets.

except SST. We can directly use the already split datasets to evaluate on SST.

D Training Efficiency

We compare the training efficiency of SimCSE and RankCSE, which are tested on a single NVIDIA Tesla A100 GPU (40GB memory). We set batch size to 128 for both SimCSE and RankCSE, and training epoch to their original settings (1 for SimCSE, 4 for RankCSE). RankCSE utilizes SimCSE_{base} and SimCSE_{large} as a multi-teacher to provide pseudo ranking labels. As shown in Table 11, RankCSE_{base} and RankCSE_{large} can be trained within 2 hours and 3.7 hours respectively. Since RankCSE needs to calculate pseudo ranking labels of the teacher, it requires additional training time per epoch than SimCSE.

E Cosine Similarity Distribution

We demonstrate the distribution of cosine similarities for sentence pairs of STS-B dev set in Figure 5. We can observe that cosine similarity distributions from all models are consistent with human ratings. However, the cosine similarities of RankCSE are slightly higher than that of SimCSE under the same human rating, as RankCSE pulls similar negatives closer during incorporating ranking consistency and ranking distillation, and shows lower variance. Compared with DiffCSE, RankCSE shows a more scattered distribution. This observation further validates that RankCSE can achieve a better alignment-uniformity balance.

	RankCSE-BERT				RankCSE-RoBERTa			
	base		large		base		large	
	listNet	listMLE	listNet	listMLE	listNet	listMLE	listNet	listMLE
Batch size	128	128	128	128	128	128	128	128
Learning rate	3e-5	2e-5	3e-5	2e-5	2e-5	3e-5	3e-5	3e-5
τ_1	0.05	0.05	0.05	0.05	0.05	0.05	0.05	0.05
τ_2	0.025	0.05	0.05	0.05	0.05	0.05	0.025	0.025
τ_3	0.0125	-	0.025	-	0.025	-	0.0125	-

Table 9: The hyperparameter values for RankCSE training.

PLMs	Methods	MR	CR	SUBJ	MPQA	SST	TREC	MRPC	avg.
BERT _{base}	+DiffCSE	82.69	87.23	<u>95.23</u>	89.28	86.60	90.40	76.58	86.86
	+RankCSE _{listNet}	83.64	88.32	95.26	<u>89.99</u>	89.02	90.80	77.10	87.73
	+RankCSE _{listMLE}	<u>83.05</u>	<u>88.03</u>	95.13	90.00	<u>88.41</u>	<u>90.60</u>	<u>76.81</u>	<u>87.43</u>
RoBERTa _{base}	+DiffCSE	82.82	88.61	94.32	87.71	88.63	90.40	<u>76.81</u>	87.04
	+RankCSE _{listNet}	83.84	<u>88.93</u>	<u>94.21</u>	<u>89.17</u>	90.23	91.60	77.28	87.89
	+RankCSE _{listMLE}	<u>83.38</u>	89.04	94.17	89.23	<u>89.51</u>	<u>91.40</u>	76.58	<u>87.62</u>

Table 10: Sentence representations performance on TR tasks (accuracy) using the dev sets of 7 TR tasks to find the best hyperparameters. The results of DiffCSE are from its original paper. We mark the best (bold) and second-best (underlined) results among models with the same PLMs.

	SimCSE		RankCSE	
	base	large	base	large
Batch size	128	128	128	128
Epoch	1	1	4	4
Time	20min	45min	120min	220min
Time per epoch	20min	45min	30min	55min

Table 11: Training efficiency of SimCSE and RankCSE. SimCSE_{base} and SimCSE_{large} provide pseudo ranking labels for every RankCSE model.

F Case Study

We present another two examples of a query sentence and several target sentences from the STS datasets, with their similarity scores and rankings in Table 12. It is obvious that the similarity scores produced by RankCSE are more effective than SimCSE, with consistent rankings to the ground-truth labels. It further demonstrates that SimCSE only captures coarse-grained semantic ranking information via contrastive learning, while RankCSE can capture fine-grained semantic ranking information. For example, SimCSE can distinguish between similar and dissimilar sentences, however, it can not distinguish between very similar and less similar sentences as RankCSE.

G Ranking Tasks

We build the ranking task based on each STS dataset to verify that RankCSE can capture fine-grained semantic ranking information. For one sentence x_i , if there are more than three sentence pairs (x_i, x_i^j) containing x_i with similarity score label y_i^j in the dataset, we view $\{x_i, x_i^j, y_i^j\}_{j=1}^k$ ($k > 3$)

as a sample of the ranking task, as shown in Table 12. We adopt KCC (Kendall’s correlation coefficient (Abdi, 2007)) and NDCG (normalized discounted cumulative gain (Järvelin and Kekäläinen, 2002)) as evaluation metrics for ranking tasks, and demonstrate the results in Table 13. RankCSE outperforms SimCSE and DiffCSE on both KCC and NDCG, which validates that RankCSE can capture fine-grained semantic ranking information by incorporating ranking consistency and ranking distillation. Another observation is that SimCSE and DiffCSE also achieve moderate results, which shows they can distinguish coarse-grained semantic differences via contrastive learning.

H Alignment and Uniformity

Wang and Isola (2020) use two properties related to contrastive learning, alignment and uniformity, to measure the quality of representation space. Alignment calculates expected distance between normalized representations of positive pairs p_{pos} :

$$\ell_{\text{align}} \triangleq \mathbb{E}_{(x, x^+) \sim p_{\text{pos}}} \|f(x) - f(x^+)\|^2, \quad (7)$$

while uniformity measures how well the normalized representations are uniformly distributed:

$$\ell_{\text{uniform}} \triangleq \log \mathbb{E}_{x, y \stackrel{i.i.d.}{\sim} p_{\text{data}}} e^{-2\|f(x) - f(y)\|^2}, \quad (8)$$

where p_{data} denotes the distribution of sentence pairs. Smaller alignment means positive instances have been pulled closer, while smaller uniformity means random instances scatter on the hypersphere. These two measures are smaller the better, and well aligned with the object of contrastive learning.

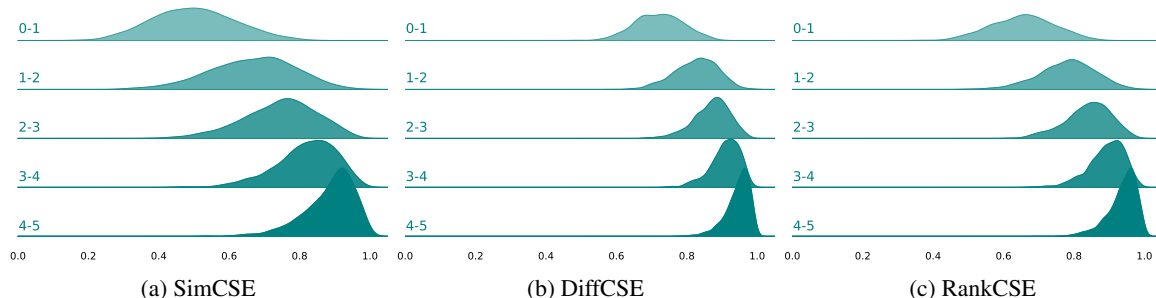


Figure 5: The distribution of cosine similarity for sentence pairs of STS-B dev set. Along the y-axis are 5 groups of pairs split based on ground truth ratings, and x-axis is the cosine similarity.

Target Sentences	Label	SimCSE	RankCSE
• a and c are on the same closed path with the battery	3.60 (1)	0.81 (1)	0.90 (1)
• bulb a and bulb c affect each other.	2.80 (2)	0.58 (3)	0.75 (2)
• they are on the same wire	1.60 (3)	0.60 (2)	0.68 (3)
• because breaking one bulb then affects the ability of the others to light up.	1.20 (4)	0.37 (5)	0.59 (4)
• if one bulb is removed, the others stop working	0.60 (5)	0.38 (4)	0.54 (5)
Query Sentence: a and c are in the same closed path			
• because by measuring voltage, you find the gap where there's a difference in electrical states.	3.80 (1)	0.86 (1)	0.90 (1)
• it allows you to measure electrical states between terminals	3.20 (2)	0.64 (3)	0.84 (2)
• it checks the electrical state between two terminals.	2.60 (3)	0.65 (2)	0.78 (3)
• find where there are different electrical states	2.60 (3)	0.55 (5)	0.78 (3)
• you can see where the gap is	2.20 (5)	0.62 (4)	0.69 (5)
Query Sentence: measuring voltage indicates the place where the electrical state changes due to a gap.			

Table 12: Two examples of a query sentence and several target sentences from the STS datasets, with their similarity scores and rankings. The label scores are from human annotations. The SimCSE and RankCSE similarity scores are from the model predictions respectively, with the corresponding ranking positions. It can be seen that sentence rankings based on SimCSE are incorrect, while RankCSE generates more effective scores with accurate rankings.

Metrics	Methods	STS12	STS13	STS14	STS15	STS16	STS-B	SICK-R	avg.
KCC	+SimCSE	36.08	36.60	44.14	49.02	54.66	58.44	54.65	47.66
	+DiffCSE	<u>38.59</u>	<u>41.89</u>	42.37	<u>51.19</u>	58.90	<u>59.21</u>	53.42	<u>49.37</u>
	+RankCSE	42.79	46.26	44.53	52.00	<u>57.21</u>	63.64	57.40	51.98
NDCG	+SimCSE	97.80	89.33	92.71	<u>96.93</u>	94.28	96.49	<u>98.44</u>	95.14
	+DiffCSE	98.35	90.22	<u>93.05</u>	96.91	94.79	97.05	98.34	<u>95.53</u>
	+RankCSE	<u>98.20</u>	92.27	93.46	97.21	95.24	97.45	98.67	96.07

Table 13: Sentence representations performance on ranking tasks (KCC and NDCG) using BERT_{base}. The results of SimCSE and DiffCSE are obtained from their publicly available codes and checkpoints. We mark the best (bold) and second-best (underlined) results.

ACL 2023 Responsible NLP Checklist

A For every submission:

- A1. Did you describe the limitations of your work?
Limitations
- A2. Did you discuss any potential risks of your work?
Sections 1 and 6
- A3. Do the abstract and introduction summarize the paper’s main claims?
Abstract and Section 1
- A4. Have you used AI writing assistants when working on this paper?
Left blank.

B Did you use or create scientific artifacts?

Section 5

- B1. Did you cite the creators of artifacts you used?
Section 5
- B2. Did you discuss the license or terms for use and / or distribution of any artifacts?
We were unable to find the license for the dataset we used.
- B3. Did you discuss if your use of existing artifact(s) was consistent with their intended use, provided that it was specified? For the artifacts you create, do you specify intended use and whether that is compatible with the original access conditions (in particular, derivatives of data accessed for research purposes should not be used outside of research contexts)?
Section 5
- B4. Did you discuss the steps taken to check whether the data that was collected / used contains any information that names or uniquely identifies individual people or offensive content, and the steps taken to protect / anonymize it?
The datasets we use are the commonly-used benchmarks.
- B5. Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.?
Section 5
- B6. Did you report relevant statistics like the number of examples, details of train / test / dev splits, etc. for the data that you used / created? Even for commonly-used benchmark datasets, include the number of examples in train / validation / test splits, as these provide necessary context for a reader to understand experimental results. For example, small differences in accuracy on large test sets may be significant, while on small test sets they may not be.
Appendix D

C Did you run computational experiments?

Appendix E

- C1. Did you report the number of parameters in the models used, the total computational budget (e.g., GPU hours), and computing infrastructure used?
Section 5 and Appendix E

The Responsible NLP Checklist used at ACL 2023 is adopted from NAACL 2022, with the addition of a question on AI writing assistance.

- C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values?

Section 5 and Appendix A

- C3. Did you report descriptive statistics about your results (e.g., error bars around results, summary statistics from sets of experiments), and is it transparent whether you are reporting the max, mean, etc. or just a single run?

Section 5

- C4. If you used existing packages (e.g., for preprocessing, for normalization, or for evaluation), did you report the implementation, model, and parameter settings used (e.g., NLTK, Spacy, ROUGE, etc.)?

Section 5

D Did you use human annotators (e.g., crowdworkers) or research with human participants?

Left blank.

- D1. Did you report the full text of instructions given to participants, including e.g., screenshots, disclaimers of any risks to participants or annotators, etc.?

No response.

- D2. Did you report information about how you recruited (e.g., crowdsourcing platform, students) and paid participants, and discuss if such payment is adequate given the participants' demographic (e.g., country of residence)?

No response.

- D3. Did you discuss whether and how consent was obtained from people whose data you're using/curating? For example, if you collected data via crowdsourcing, did your instructions to crowdworkers explain how the data would be used?

No response.

- D4. Was the data collection protocol approved (or determined exempt) by an ethics review board?

No response.

- D5. Did you report the basic demographic and geographic characteristics of the annotator population that is the source of the data?

No response.