

# A Sentence is Worth 128 Pseudo Tokens: A Semantic-Aware Contrastive Learning Framework for Sentence Embeddings

Anonymous ACL submission

## Abstract

Contrastive learning has shown great potential in unsupervised sentence embedding tasks, e.g., SimCSE (Gao et al., 2021). However, these existing solutions are heavily affected by superficial features like the length of sentences or syntactic structures. In this paper, we propose a semantic-aware contrastive learning framework for sentence embeddings, termed Pseudo-Token BERT (PT-BERT), which is able to explore the pseudo-token space (i.e., latent semantic space) representation of a sentence while eliminating the impact of superficial features such as sentence length and syntax. Specifically, we introduce an additional pseudo token embedding layer independent of the BERT encoder to map each sentence into a sequence of pseudo tokens in a fixed length. Leveraging these pseudo sequences, we are able to construct same-length positive and negative pairs based on the attention mechanism to perform contrastive learning. In addition, we utilize both the gradient-updating and momentum-updating encoders to encode instances while dynamically maintaining an additional queue to store the representation of sentence embeddings, enhancing the encoder’s learning performance for negative examples. Experiments show that our model outperforms the state-of-the-art baselines on six standard semantic textual similarity (STS) tasks. Furthermore, experiments on alignments and uniformity losses, as well as hard examples with different sentence lengths and syntax, consistently verify the effectiveness of our method.

## 1 Introduction

Sentence embedding serves as an essential technique in a wide range of applications, including semantic search, text clustering, text classification, etc. (Kiros et al., 2015; Logeswaran and Lee, 2018; Conneau et al., 2017; Cer et al., 2018; Reimers and Gurevych, 2019; Gao et al., 2021). Contrastive learning works on learning representations such

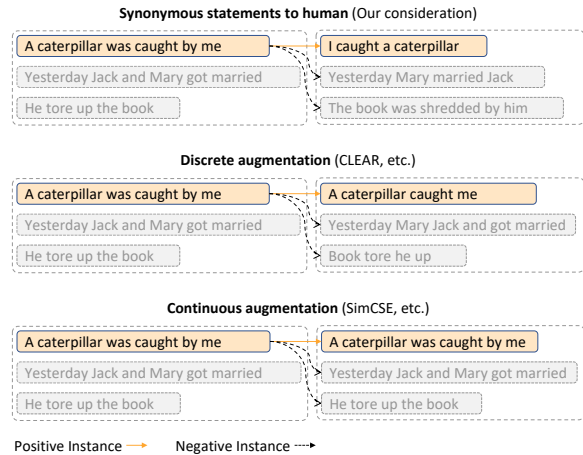


Figure 1: In a realistic scenario, negative examples have the same length and structure, while positive examples act in the opposite way. In comparison, discrete augmentation obtains positive instances with word deletion or reordering (Wu et al., 2020; Meng et al., 2021), which may misinterpret the meaning. The continuous method treats embeddings of the same original sentence as positive examples and augments sentences with the different encoding functions (Carlsson et al., 2021; Gao et al., 2021).

that similar examples stay close whereas dissimilar ones are far apart, and thus is suitable for sentence embeddings due to its natural availability of similar examples. Incorporating contrastive learning in sentence embeddings improves the efficiency of semantic information learning in an unsupervised manner (He et al., 2020; Chen et al., 2020) and has been shown to be effective on a variety of tasks (Reimers and Gurevych, 2019; Gao et al., 2021; Zhang et al., 2020).

In contrastive learning for sentence embeddings, a key challenge is how to construct positive instances. Both discrete and continuous augmentation methods have been studied recently. Methods in (Wu et al., 2018; Meng et al., 2021) (e.g., CLEAR) perform discrete operations directly on the original sentences, such as word deletion and

060 sentence shuffling, to get positive samples. How- 111  
061 ever, these methods may lead to unacceptable se- 112  
062 mantic distortions or even complete misinter- 113  
063 pretations of the original statement. In contrast, the 114  
064 SimCSE method (Gao et al., 2021) obtains two dif- 115  
065 ferent embeddings in the continuous embedding 116  
066 space as a positive pair for one sentence through 117  
067 different *dropout masks* (Srivastava et al., 2014) 118  
068 in the neural network for representation learning. 119  
069 Nonetheless, this method overly relies on superfi- 120  
070 cial features existing in the dataset like sentence 121  
071 lengths and syntactic structures and may pay less 122  
072 reflection on meaningful semantic information. As 123  
073 an illustrative example, the sentence-pair in Fig. 1 124  
074 “A caterpillar was caught by me.” and “I caught a 125  
075 caterpillar.” appear to organize differently in ex- 126  
076 pression but convey exactly the same semantics. 127

077 To overcome these drawbacks, in this paper, 128  
078 we propose a semantic-aware contrastive learn- 129  
079 ing framework for sentence embeddings, termed 130  
080 Pseudo-Token BERT (PT-BERT), that is able to 131  
081 capture the pseudo-token space (i.e., latent seman- 132  
082 tic space) representation while ignoring effects of 133  
083 superficial features like sentence lengths and syn- 134  
084 tactic structures. Inspired by previous work on 135  
085 prompt learning and sentence selection (Li and 136  
086 Liang, 2021; Liu et al., 2021; Humeau et al., 2020), 137  
087 which create a pseudo-sequence and have it serve 138  
088 the downstream tasks, we present PT-BERT to train 139  
089 pseudo token representations and then to map sen- 140  
090 tences into pseudo token spaces based on an atten- 141  
091 tion mechanism. 142

092 In particular, we train additional 128 pseudo 143  
093 token embeddings, together with sentence em- 144  
094 beddings extracted from the BERT model (i.e., 145  
095 gradient-encoder), and then use the attention mech- 146  
096 anism (Devlin et al., 2019) to map the sentence 147  
097 embedding to the pseudo token space (i.e., se- 148  
098 mantic space). We use another BERT model (i.e., 149  
099 momentum-encoder) to encode the original sen- 150  
100 tence, adopt a similar attention mechanism with 151  
101 the pseudo token embeddings, and finally output 152  
102 a continuously augmented version of the sentence 153  
103 embedding. We treat the representations of original 154  
104 sentence encoded by the gradient-encoder and the 155  
105 momentum-encoder as a positive pair. In addition, 156  
106 the momentum-encoder also generates negative ex- 157  
107 amples, dynamically maintains a queue to store 158  
108 these negative examples, and updates them over- 159  
109 time. By projecting all sentences onto the same 160  
110 pseudo sentence, the model greatly reduces the

dependence on sentence length and syntax when 111  
making judgments and makes the model more fo- 112  
cused on the semantic level information. 113

114 In our experiments, we compare our results with 115  
116 the previous state-of-the-art work. We train PT- 117  
118 BERT on  $10^6$  randomly sampled sentences from 119  
120 English Wikipedia and evaluate on seven standard 121  
122 semantic textual similarity (STS) tasks (Agirre 123  
124 et al., 2012, 2013, 2014, 2015, 2016) (Marelli et al., 125  
126 2014). Besides, we also compare our approach 127  
128 with a framework based on an advanced discrete 129  
130 augmentation we proposed. We obtain a new state-  
of-the-art on standard semantic textual similarity  
tasks with our PT-BERT, which achieves 77.74% of  
Spearman’s correlation. To show the effectiveness  
of pseudo tokens, we calculate the align-loss and  
uniformity loss (Wang and Isola, 2020) and verify  
our approach on a sub-dataset with hard examples  
sampled from STS-(2012-2016).

## 2 Related Work 130

131 In this section, we discuss related studies with 132  
133 respect to the contrastive learning framework and  
sentence embedding.

### 2.1 Contrastive Learning for Sentence Embedding 134

135 **Contrastive learning and MoCo.** Contrastive 136  
137 learning (Hadsell et al., 2006) has been used with 138  
139 much success in both natural language processing  
and computer vision (Yang et al., 2019; Klein and  
Nabi, 2020; Chen et al., 2020; He et al., 2020; Gao  
et al., 2021). In contrast to generative learning,  
contrastive learning requires learning to distinguish  
and match data at the abstract semantic level of the  
feature space. It focuses on learning common fea-  
tures between similar examples and distinguishing  
differences between non-similar examples. In order  
to compare the instances with more negative exam-  
ples and less computation, memory bank (Wu et al.,  
2018) is proposed to enhance the performance un-  
der the contrastive learning framework. While with  
a large capacity to store more samples, the mem-  
ory bank is not consistent enough, which could not  
update the "key" during comparison. Momentum-  
Contrast (MoCo) (He et al., 2020) uses a queue to  
maintain the dictionary of samples which allows  
the model to compare the query with more keys for  
each step and ensure the consistency of the frame-  
work. It updates the parameter of the dictionary in  
a momentum way. 159

**Discrete and continuous augmentation.** By equipping discrete augmentation that modifies sentences directly on token level with contrastive learning, significant success has been achieved in obtaining sentence embeddings. Such methods include word omission (Yang et al., 2019), entity replacement (Xiong et al., 2020), trigger words (Klein and Nabi, 2020) and traditional augmentations such as deletion, reorder and substitution (Wu et al., 2020; Meng et al., 2021). Examples with diverse expressions can be learned during training, making the model more robust to expressions of different sentence lengths and styles. However, these approaches are limited because there are huge difficulties in augmenting sentences precisely since a few changes can make the meaning completely different or even opposite.

Researchers have also explored the possibility of building sentences continuously, which instead applies operation in embedding space. CT-BERT (Carlsson et al., 2021) encodes same sentence with two different encoders. Unsup-SimCSE (Gao et al., 2021) compares the representations of the same sentence with different dropout masks among the mini-batch. These approaches continuously augment sentences while retaining the original meaning. However, positive pairs seen by SimCSE always have the same length and structure, whereas negative samples are likely to act oppositely. As a result, sentence length and structure are highly correlated to the similarity score of examples. During training, the model has never seen positive samples with diverse expressions, so that in real test scenarios, the model would be more inclined to classify the synonymous pairs with different expressions as negatives, and those sentences with the same length and structures are more likely to be grouped as positive pairs. This may cause a biased encoder.

## 2.2 Pseudo Tokens

In the domain of prompt learning (Liu et al., 2021; Jiang et al., 2020; Li and Liang, 2021; Gao et al., 2020), the way to create prompt can be divided into two types, namely discrete and continuous ways. Discrete methods usually search the natural language template as the prompt (Davison et al., 2019; Petroni et al., 2019), while the continuous way always directly works on the embedding space with "pseudo tokens" (Liu et al., 2021; Li and Liang, 2021). In retrieval and dialogue tasks, the current

	Sub-dataset	original
STS12	66.54	68.40
STS13	78.50	82.41
STS14	68.76	74.38
STS15	70.27	80.91
STS16	71.31	78.56

Table 1: SimCSE’s results on sub-dataset from STS12-16, comparing with original results.

	SimCSE <sub>32</sub>	SimCSE <sub>64</sub>	SimCSE <sub>128</sub>
Avg.	76.25	75.20	75.29

Table 2: Different acceptable sequence length of SimCSE would affect the result on STS tasks.

approach adopts "pseudo tokens", namely "poly codes" (Humeau et al., 2020), to jointly encode the query and response precisely and ensure the inference time when compared with the Cross-Encoders and Bi-Encoders (Wolf et al., 2019; Mazaré et al., 2018; Dinan et al., 2019). The essence of these methods is to create a pseudo-sequence and have it serve the downstream tasks without the need for humans to understand the exact meaning. The parameters of these pseudo tokens are independent of the natural language embeddings, and can be tuned based on a specific downstream task. In the following sections, we will show the idea to weaken the model’s consideration of sentence length and structures by introducing additional pseudo token embeddings on top of the BERT encoder.

## 3 Methods

In this section, we introduce PT-BERT, which provides novel contributions on combining advantages of both discrete and continuous augmentations to advance the state-of-art of sentence embeddings. We first present the setup of problems with a thorough analysis on the bias introduced by the textual similarity theoretically and experimentally. Then we show the details of Pseudo-Token representation and our model’s architecture.

### 3.1 Preliminary

Consider a sentence  $s$ , we say that the augmentation is continuous if  $s$  is augmented by different encoding functions,  $f(\cdot)$  and  $f'(\cdot)$ . Sentence embeddings  $\mathbf{h} = f(s)$  and  $\mathbf{h}' = f'(s)$  are obtained by these two functions. With a slight change of the encoding function (e.g., encoders with different *dropout* masks),  $\mathbf{h}'$  can be seen as a more precisely

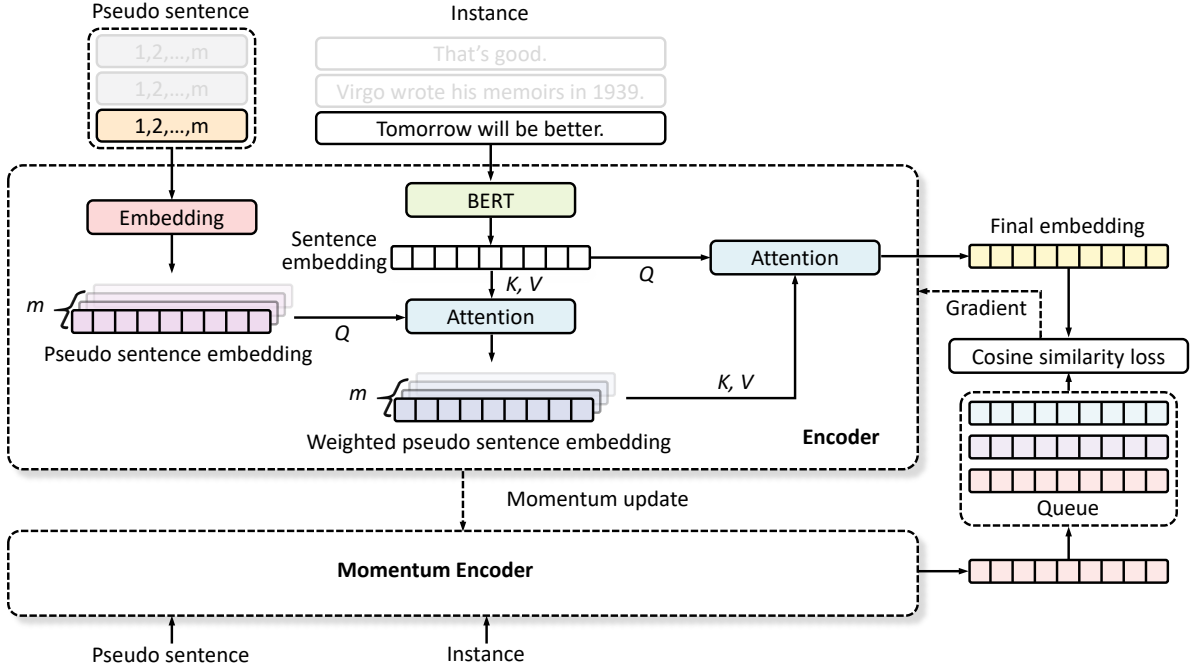


Figure 2: The model is divided into two parts, the upper part (Encoder) updates the learnable parameters with gradient, while the bottom (Momentum Encoder) inherits parameters from the upper part with momentum-updating. We repeatedly input the same sequence of pseudo tokens while processing the original sentences. An additional BERT attention mapping the pooler-output of BERT to pseudo sequence representation, extending the sentence embedding to a fixed length and mapping the syntactic structure to the style of the pseudo sentence. The two attentions in the figure are the same and with identical parameters.

244 augmented version of  $\mathbf{h}$  compared with the discrete  
 245 augmentation. Semantic information of  $\mathbf{h}'$  could  
 246 be the same as  $\mathbf{h}$ . Therefore,  $\mathbf{h}$  and  $\mathbf{h}'$  are a pair of  
 247 positive examples and we could randomly sample  
 248 a sentence to construct negative example pairs.

249 Previous state-of-the-art models (Gao et al.,  
 250 2021) adopt the continuous strategy that augments  
 251 sentences with *dropout* (Srivastava et al., 2014).  
 252 Through careful observation, we find that all the  
 253 positive examples in SimCSE have the same length  
 254 and structure while negative examples act oppo-  
 255 sitely. In this way, SimCSE will inevitably take  
 256 these two factors as hints during test. To further  
 257 verify this conjecture, we sort out the positive pairs  
 258 with a length difference of more than five words  
 259 and negative pairs of less than two words from  
 260 STS-(2012-2016).

261 Table 1 shows that the performance of SimCSE  
 262 plummets on this dataset. Besides, we also find  
 263 that SimCSE truncates all training corpus into 32  
 264 tokens, which shortens the discrepancy of the sen-  
 265 tence’s length. After we scale the max length that  
 266 SimCSE could accept from 32 to 64 and 128, the  
 267 performance degrades significantly during the test  
 268 even though the model is supposed to learn more

269 from the complete version of sentences. The reason  
 270 for this result may lie in the fact that, without  
 271 truncation, all positive pairs still have the same  
 272 length, whereas the difference in length between  
 273 the negative and positive ones is enlarged. There-  
 274 fore, the encoder will rely more on sentence length  
 275 and make the wrong decision.

### 3.2 Pseudo-Token BERT

276 We realize it is vital to train an unbiased encoder  
 277 that captures the semantics and also would not in-  
 278 troduce intermediate errors. This motivates us to  
 279 propose the PT-BERT, as evidence shows that the  
 280 encoder may fail to make predictions when trained  
 281 on a biased dataset with same-length positive pairs,  
 282 by learning the spurious correlations that work only  
 283 well on the training dataset (Arjovsky et al., 2019;  
 284 Nam et al., 2020).  
 285

286 **Pseudo-Token representations.** The idea of PT-  
 287 BERT is to reduce the model’s excessive depen-  
 288 dence on textual similarity when making predic-  
 289 tions. Discrete augmentation achieves this goal by  
 290 providing both positive and negative examples with  
 291 diverse expressions. Therefore the model does not



jump to conclusions based on sentence length and syntactic structure during the test.

Note that we achieve this same purpose in a seemingly opposite way: *mapping the representations of both positive and negative examples to a pseudo sentence with the same length and structure*. We take an additional embedding layer outside the BERT encoder to represent a pseudo sentence  $\{0, 1, \dots, m\}$  with fixed length  $m$  and unchangeable syntax. This embedding layer is fully independent of the BERT encoder, including the parameters and corresponding vocabulary. Random initialization is applied to this layer, and the parameter will be updated during training. The size of this layer depends on the length of pseudo sentences. Besides, adopting the attention mechanism (Vaswani et al., 2017; Bahdanau et al., 2015; Gehring et al., 2017), we take the pseudo sentence embeddings as the query while key and value are the sentence embeddings obtained from the BERT encoder. This allows the pseudo sentence to attend to the core part and ignore the redundant part while keeping the fixed length and pseudo syntactic structure.

Fig. 2 illustrates the framework of PT-BERT. Denoting the pseudo sentence embedding as  $\mathbf{P}$  and the sentence embedding encoded by BERT as  $\mathbf{Y}$ , we obtain the weighted pseudo sentence embedding of each sentence by mapping the sentence embedding to the pseudo tokens with attention:

$$\mathbf{Z}'_i = \text{Attention}(\mathbf{P}\mathbf{W}^{\mathbf{Q}}, \mathbf{Y}_i\mathbf{W}^{\mathbf{K}}, \mathbf{Y}_i\mathbf{W}^{\mathbf{V}}) \quad (1)$$

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^{\mathbf{T}}}{\sqrt{d_k}}\right)\mathbf{V}, \quad (2)$$

where  $d_k$  is the dimension of the model,  $\mathbf{W}^{\mathbf{Q}}$ ,  $\mathbf{W}^{\mathbf{K}}$ ,  $\mathbf{W}^{\mathbf{V}}$  are the learnable parameters with  $\mathbb{R}^{d_k \times d_k}$ ,  $i$  denotes the  $i$ -th sentence in the dataset. Then we obtain the final embedding  $\mathbf{h}_i$  with the same attention layer by mapping pseudo sentences back to original sentence embeddings:

$$\mathbf{h}_i = \text{Attention}(\mathbf{Y}_i\mathbf{W}^{\mathbf{Q}}, \mathbf{Z}'_i\mathbf{W}^{\mathbf{K}}, \mathbf{Z}'_i\mathbf{W}^{\mathbf{V}}). \quad (3)$$

Finally, we compare the cosine similarities between the obtained embeddings of  $\mathbf{h}$  and  $\mathbf{h}'$  using Eq. 4, where  $\mathbf{h}'$  are the samples encoded by the momentum-encoder and stored in a queue.

**Model architecture.** Instead of inputting the same sentence twice to the same encoder, we follow the architecture proposed in Momentum-Contrast

(MoCo) (He et al., 2020) such that PT-BERT can efficiently learn from more negative examples. Samples in PT-BERT are encoded into vectors with two encoders: gradient-update encoder (the upper encoder in Fig. 2) and momentum-update encoder (the momentum encoder in Fig. 2). We dynamically maintain a queue to store the sentence representations from momentum-update encoder.

This mechanism allows us to store as much negative samples as possible without re-computation. Once the queue is full, we replace the "oldest" negative sample with a "fresh" one encoded by the momentum-encoder.

Similar to the works based on continuous augmentation, at the very beginning of the framework, PT-BERT takes input sentence  $s$  and obtained  $\mathbf{h}_i$  and  $\mathbf{h}'_i$  with two different encoder functions. We measure the loss function with:

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

where  $\mathbf{h}_i$  denotes the representations extracted from the gradient-update encoder,  $\mathbf{h}'_i$  represents the sentence embedding in the queue, and  $M$  is the queue size. Our gradient-update and momentum-update encoder is based on the pre-trained language model with the same structure and dimensions as BERT-base-uncased (Devlin et al., 2019). The momentum encoder will update its parameters similar to MoCo:

$$\theta_k \leftarrow \lambda\theta_k + (1 - \lambda)\theta_q, \quad (5)$$

where  $\theta_k$  is the parameter of the momentum-contrast encoder that maintains the dictionary,  $\theta_q$  is the query encoder that updates the parameters with gradients, and  $\lambda$  is a hyperparameter used to control the updating process.

## 4 Experiments

In this section, we perform the standard semantic textual similarity (STS) (Agirre et al., 2012, 2013, 2014, 2015, 2016) tasks to test our model. For all tasks, we measure the Spearman's correlation to compare our performance with SimCSE (Gao et al., 2021). In the following, we will describe the training procedure in detail.

### 4.1 Training Data and Settings

**Datasets.** Following SimCSE, We train our model on 1-million sentences randomly sampled

Model	STS12	STS13	STS14	STS15	STS16	STS-B	SICK-R	Avg.
<i>Discrete Augmentation</i>								
CLEAR	49.00	48.90	57.40	63.60	65.60	72.50	<b>75.60</b>	61.80
MoCo	68.35	81.42	73.34	81.63	78.61	76.40	68.50	75.46
MoCo+reorder	66.14	80.06	73.14	81.35	76.01	73.99	65.76	73.78
MoCo+duplication	65.88	82.24	73.34	81.49	77.48	76.29	68.86	75.08
MoCo+deletion	67.86	81.43	72.8	81.48	77.84	76.91	69.46	75.40
MoCo+SRL	68.92	82.20	73.67	81.58	78.73	77.63	71.07	76.26
<i>Continuous Augmentation</i>								
CT-BERT	61.63	76.80	68.47	77.50	76.48	74.31	69.19	72.05
SimCSE-BERT <sub>base</sub>	68.40	82.41	74.38	80.91	78.56	76.85	72.23	76.25
PT-BERT <sub>base</sub>	<b>71.20</b>	<b>83.76</b>	<b>76.34</b>	<b>82.63</b>	<b>78.90</b>	<b>79.42</b>	71.94	<b>77.74</b>

Table 3: Sentence embedding performance on STS tasks with Spearman’s correlation measured. We highlight the highest number for each methods. CLEAR (Wu et al., 2020) is trained on both English Wikipedia and Book Corpus with 500k steps with their own version of pre-trained models. Result of CT-BERT (Carlsson et al., 2021) is based on the settings of SimCSE (Gao et al., 2021)

Models	STS-B dev
SimCSE-BERT <sub>base</sub> + None	82.50
SimCSE-BERT <sub>base</sub> + Crop	77.80
SimCSE-BERT <sub>base</sub> + Deletion	75.90
MoCo-BERT <sub>base</sub> + None	82.03
MoCo-BERT <sub>base</sub> + Reorder	81.89
MoCo-BERT <sub>base</sub> + Duplication	81.82
MoCo-BERT <sub>base</sub> + Deletion	82.97
MoCo-BERT <sub>base</sub> + SRL	82.40
PT-BERT <sub>base</sub>	<b>84.50</b>

Table 4: Results on STS-B development sets. Results of SimCSE (Gao et al., 2021) are reported from original paper.

from English Wikipedia, and evaluate the model every 125 steps to find the best checkpoints. Note that we do not fine-tune our model on any dataset, which indicates that our method is completely unsupervised.

**Hardware and schedule.** We train our model on the machine with one NVIDIA V100s GPU. Following the settings of SimCSE (Gao et al., 2021), it takes 50 minutes to run an epoch.

## 4.2 Implementations

In this subsection, we implement PT-BERT based on Huggingface transformers (Wolf et al., 2020) and initialize it with the released BERT<sub>base</sub> (Devlin et al., 2019). We initialize a new embedding for pseudo tokens with  $128 \times 768$ . During training, we create a pseudo sentence  $\{0, 1, 2, \dots, 127\}$  for every input and map the original sentence to this pseudo

sentence by attention. With batches of 64 sentences and an additional dynamically maintained queue of 256 sentences, each sentence has one positive sample and 255 negative samples. Adam (Kingma and Ba, 2014) optimizer is used to update the model parameters. We also take the original dropout strategy of BERT with rate  $p = 0.1$ . We set the momentum for the momentum-encoder with  $\lambda = 0.885$ .

## 4.3 Evaluation Setup

We evaluate the fine-tuned BERT encoder on STS-B development sets every 125 steps to select the best checkpoints. We report all the checkpoints based on the evaluation results reported in Table 4. The training process is fully unsupervised since no training corpus from STS is used. During the evaluation, we also calculate the trends of alignment-loss and uniformity-loss. Losses were compared with SimCSE (Gao et al., 2021) under the same experimental settings. After training and evaluation, we test models 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). We report the result of Spearman’s correlation for all the experiments.

## 4.4 Main Results and Analysis

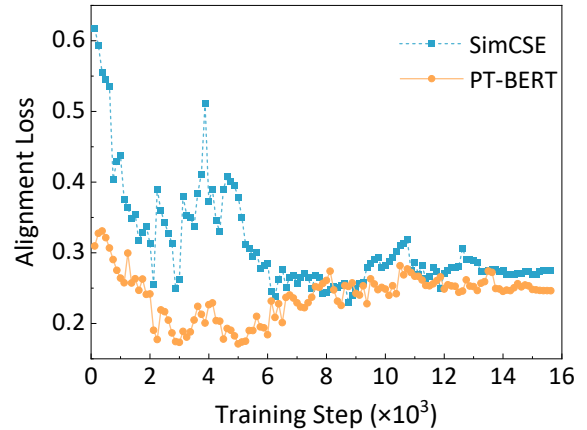
We first compare PT-BERT with our baseline: MoCo framework + BERT encoder (MoCo-BERT). MoCo-BERT could be seen as a version of PT-BERT without pseudo token embeddings. Then we apply traditional discrete augmentations such as re-

Method	STS12	STS13	STS14	STS15	STS16	STS-B	SICK-R	Avg.
<i>(a) Ablation studies on pseudo sequence length</i>								
L-64	67.04	82.04	73.65	81.12	78.64	77.35	71.33	75.88
L-90	68.94	82.08	74.53	81.22	79.06	78.01	71.49	76.48
<b>L-128(Ours)</b>	<b>71.20</b>	<b>83.76</b>	<b>76.34</b>	<b>82.63</b>	<b>78.90</b>	<b>79.42</b>	<b>71.94</b>	<b>77.74</b>
L-256	67.09	82.25	72.63	81.48	78.55	77.30	69.53	75.55
L-360	68.90	82.21	73.77	81.31	77.50	77.22	69.32	75.75
<i>(b) Ablation studies on queue size</i>								
Q-192	70.29	<b>83.78</b>	75.98	82.13	78.48	78.91	72.53	77.44
<b>Q-256(Ours)</b>	71.20	83.76	<b>76.34</b>	82.63	<b>78.90</b>	<b>79.42</b>	<b>71.94</b>	<b>77.74</b>
Q-320	<b>71.71</b>	83.36	75.00	<b>82.99</b>	78.76	79.17	<b>72.85</b>	77.69
<i>(c) Evaluations on hard sentence pairs with different length</i>								
SimCSE	66.54	78.50	68.76	70.27	71.31	-	-	71.08
PT-BERT	<b>72.02</b>	<b>80.24</b>	<b>72.92</b>	<b>74.50</b>	<b>72.50</b>	-	-	<b>74.44</b>

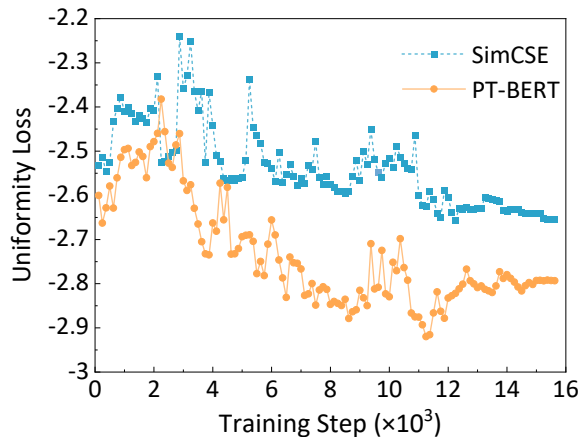
Table 5: Evaluation results of ablation studies and hard sentence pairs.

order, duplication, and deletion on this framework. We also compare our work with CLEAR (Wu et al., 2020) that substitutes and deletes the token spans. Besides, we argue that the performance of these methods is too weak. We additionally propose an advanced discrete augmentation approach that produces positive examples with the guidance of Semantic Role Labeling (SRL) (Gildea and Jurafsky, 2002) information, instead of random deletion and reordering. SRL-guided augmentation could compensate the errors caused by these factors, acting as a combination of deletion, duplication, and reordering with better accuracy. For the sentences with multiple predicates, we keep all the sets with order [ARG0, PRED, ARGM – NEG, ARG1] and concatenate them into a new sequence. For the sentences without recognized predicate-argument sets, we keep the original sentence as positive examples. In addition to the work based on discrete approaches, we also compare with SimCSE (Gao et al., 2021) which continuously augment sentences with *dropout*. In Table 3, PT-BERT with 128 pseudo tokens further pushed the state-of-the-art results to 77.74% and significantly outperformed SimCSE over six datasets.

In Fig 3, we observe that PT-BERT also achieves better alignment and uniformity against SimCSE, which indicates that pseudo tokens really help the learning of sentence representations. In detail, alignment and uniformity are proposed by (Wang and Isola, 2020) to evaluate the quality of representations in contrastive learning. The calculation of these two metrics are shown in the following



(a) Alignment loss comparison on STS-B



(b) Uniformity loss comparison on STS-B

Figure 3: Alignment and uniformity loss plot for PT-BERT and SimCSE. We visualize the checkpoints every 125 training steps. For both measurements, lower numbers are better.

formulas:

$$L_{alignment} = E_{(x,x^+) \sim p_{pos}} \|f(x) - f(x^+)\|^2, \quad (6)$$

$$L_{uniformity} = E_{(x,y) \sim p_{data}} e^{-2\|f(x)-f(y)\|^2}, \quad (7)$$

where  $(x, x^+)$  is the positive pair,  $(x, y)$  is the pair consisting of any two different sentences in the whole sentence set,  $f(x)$  is the normalized representation of  $x$ . We employ the final embedding  $\mathbf{h}$  to calculate these scores.

According to the above formulas, lower alignment loss means a shorter distance between the positive samples, and low uniformity loss implies the diversity of embeddings of all sentences. Both are our expectations for the representations based on contrastive learning. To evaluate our model’s performance on alignment and uniformity, we compare it with SimCSE on the STS-benchmark dataset (Cer et al., 2017), and the result is shown in Figure 3. The result demonstrates that PT-BERT outperforms SimCSE on these two metrics: our model has a lower alignment and uniformity than SimCSE in almost all the training steps, which indicates that the representations produced by our model are more in line with the goal of the contrastive learning.

## 5 Analysis

### 5.1 Ablation Studies

In this section, we first investigate the impact of different sizes of pseudo token embeddings. Then we would like to report the performance difference caused by queue size under the MoCo framework.

**Pseudo Sentence Length** Different lengths of pseudo tokens can affect the ability of the model to express the sentence representations. By mapping the original sentences to various lengths of pseudo tokens, the performance of PT-BERT could be different. In this section, we keep all the parts except the pseudo tokens and their embeddings unchanged. We scale the pseudo sequence length from 64 to 360. Table 5(a) shows a comparison between different lengths of pseudo sequence in PT-BERT. We find that during training, PT-BERT performs better when attending to pseudo sequences with 128 tokens. Too few pseudo tokens do not fully explain the semantics of the original sentence, while too many pseudo tokens increase the number of parameters and over-express the sentence.

**Queue Size** The introduction of more negative samples would make the model’s training more reliable. By training with different queue sizes, we report the result of PT-BERT with different performances due to the number of negative samples. In Table 5(b), queue size  $q = 4$  performs best. However, the difference in performance between the three sets of experiments is not large, suggesting that the model can learn well as long as it can see enough negative samples.

### 5.2 Exploration on Hard Examples with Different Length

To prove the effectiveness of PT-BERT that could weaken the hints caused by textual similarity, we further test PT-BERT on the sub-dataset introduced in Sec. 3.1. We sorted out the positive pairs with a length difference of more than five words and negative pairs of less than two words from STS-(2012-2016). PT-BERT significantly outperforms SimCSE with 3.36% Spearman’s correlation, indicating that PT-BERT could handle these hard examples better than SimCSE. This further proves that PT-BERT could debias the spurious correlation introduced by sentence length and syntax, and focus more on the semantics.

## 6 Conclusion

In this paper, we propose a semantic-aware contrastive learning framework for sentence embeddings, termed PT-BERT. Our proposed PT-BERT approach is able to weaken textual similarity information, such as sentence length and syntactic structures, by mapping the original sentence to a fixed pseudo sentence embedding. We provide analysis of these factors on methods based on continuous and discrete augmentation, showing that PT-BERT augments sentences more accurately than discrete methods while considering more semantics instead of textual similarity than continuous approaches. Lower uniformity loss and alignment loss prove the effectiveness of PT-BERT and further experiments also show that PT-BERT could handle hard examples better than existing approaches.

Providing a new perspective to the continuous data augmentation in sentence embeddings, we believe our proposed PT-BERT has great potential to be applied in broader downstream applications, such as text classification, text clustering, and sentiment analysis.



## References

- 558 Eneko Agirre, Carmen Banea, Claire Cardie, Daniel  
559 Cer, Mona Diab, Aitor Gonzalez-Agirre, Weiwei  
560 Guo, Iñigo Lopez-Gazpio, Montse Maritxalar, Rada  
561 Mihalcea, German Rigau, Larraitz Uria, and Janyce  
562 Wiebe. 2015. [SemEval-2015 task 2: Semantic textual similarity, English, Spanish and pilot on interpretability](#). In *Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015)*, pages 252–263, Denver, Colorado. Association for Computational Linguistics. 616
- 563  
564  
565  
566  
567
- 568 Eneko Agirre, Carmen Banea, Claire Cardie, Daniel Cer,  
569 Mona Diab, Aitor Gonzalez-Agirre, Weiwei Guo,  
570 Rada Mihalcea, German Rigau, and Janyce Wiebe.  
571 2014. [SemEval-2014 task 10: Multilingual semantic textual similarity](#). In *Proceedings of the 8th International Workshop on Semantic Evaluation (SemEval 2014)*, pages 81–91, Dublin, Ireland. Association for Computational Linguistics. 617
- 572  
573  
574  
575
- 576 Eneko Agirre, Carmen Banea, Daniel Cer, Mona Diab,  
577 Aitor Gonzalez-Agirre, Rada Mihalcea, German  
578 Rigau, and Janyce Wiebe. 2016. [SemEval-2016 task 1: Semantic textual similarity, monolingual and cross-lingual evaluation](#). In *Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016)*, pages 497–511, San Diego, California. Association for Computational Linguistics. 618
- 579  
580  
581  
582  
583
- 584 Eneko Agirre, Daniel Cer, Mona Diab, and Aitor  
585 Gonzalez-Agirre. 2012. [SemEval-2012 task 6: A pilot on semantic textual similarity](#). In *\*SEM 2012: The First Joint Conference on Lexical and Computational Semantics – Volume 1: Proceedings of the main conference and the shared task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation (SemEval 2012)*, pages 385–393, Montréal, Canada. Association for Computational Linguistics. 619
- 586  
587  
588  
589  
590  
591  
592  
593
- 594 Eneko Agirre, Daniel Cer, Mona Diab, Aitor Gonzalez-  
595 Agirre, and Weiwei Guo. 2013. [\\*SEM 2013 shared task: Semantic textual similarity](#). In *Second Joint Conference on Lexical and Computational Semantics (\*SEM), Volume 1: Proceedings of the Main Conference and the Shared Task: Semantic Textual Similarity*, pages 32–43, Atlanta, Georgia, USA. Association for Computational Linguistics. 620
- 596  
597  
598  
599  
600  
601
- 602 Martín Arjovsky, Léon Bottou, Ishaan Gulrajani, and  
603 David Lopez-Paz. 2019. Invariant risk minimization.  
604 *ArXiv*, abs/1907.02893. 621
- 605 Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Ben-  
606 gio. 2015. Neural machine translation by jointly  
607 learning to align and translate. *CoRR*, abs/1409.0473. 622
- 608 Fredrik Carlsson, Amaru Cuba Gyllensten, Evan-  
609 gelia Gogoulou, Erik Ylipää Hellqvist, and Magnus  
610 Sahlgren. 2021. Semantic re-tuning with contrastive  
611 tension. In *ICLR*. 623
- 612 Daniel Cer, Mona Diab, Eneko Agirre, Iñigo Lopez-  
613 Gazpio, and Lucia Specia. 2017. [SemEval-2017 task 1: Semantic textual similarity multilingual and crosslingual focused evaluation](#). In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages 1–14, Vancouver, Canada. Association for Computational Linguistics. 624
- 614  
615  
616  
617  
618
- 619 Daniel Cer, Yinfei Yang, Sheng-yi Kong, Nan Hua,  
620 Nicole Limtiaco, Rhomni St. John, Noah Constant,  
621 Mario Guajardo-Cespedes, Steve Yuan, Chris Tar,  
622 Brian Strope, and Ray Kurzweil. 2018. [Universal sentence encoder for English](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 169–174, Brussels, Belgium. Association for Computational Linguistics. 625
- 626  
627
- 628 Ting Chen, Simon Kornblith, Mohammad Norouzi, and  
629 Geoffrey Hinton. 2020. [A simple framework for contrastive learning of visual representations](#). In *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 1597–1607. PMLR. 628
- 630  
631  
632  
633  
634
- 635 Alexis Conneau, Douwe Kiela, Holger Schwenk, Loïc  
636 Barrault, and Antoine Bordes. 2017. [Supervised learning of universal sentence representations from natural language inference data](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 670–680, Copenhagen, Denmark. Association for Computational Linguistics. 629
- 637  
638  
639  
640  
641  
642
- 643 Joe Davison, Joshua Feldman, and Alexander Rush.  
644 2019. [Commonsense knowledge mining from pre-trained models](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1173–1178, Hong Kong, China. Association for Computational Linguistics. 630
- 645  
646  
647  
648  
649  
650
- 651 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and  
652 Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics. 631
- 653  
654  
655  
656  
657  
658  
659
- 660 Emily Dinan, Stephen Roller, Kurt Shuster, Angela  
661 Fan, Michael Auli, and Jason Weston. 2019. Wizard  
662 of Wikipedia: Knowledge-powered conversational  
663 agents. In *Proceedings of the International Conference on Learning Representations (ICLR)*. 632
- 664  
665  
666  
667
- 668 Tianyu Gao, Adam Fisch, and Danqi Chen. 2020.  
669 [Making pre-trained language models better few-shot learners](#). *CoRR*, abs/2012.15723. 633
- 670  
671
- 672 Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021.  
673 [SimCSE: Simple contrastive learning of sentence embeddings](#). In *Empirical Methods in Natural Language Processing (EMNLP)*. 634

672	Jonas Gehring, Michael Auli, David Grangier, Denis Yarats, and Yann N. Dauphin. 2017. <a href="#">Convolutional sequence to sequence learning</a> . In <i>Proceedings of the 34th International Conference on Machine Learning</i> , volume 70 of <i>Proceedings of Machine Learning Research</i> , pages 1243–1252. PMLR.	725
673		726
674		727
675		728
676		729
677		730
678	Daniel Gildea and Daniel Jurafsky. 2002. <a href="#">Automatic labeling of semantic roles</a> . <i>Comput. Linguist.</i> , 28(3):245–288.	731
679		732
680		
681	R. Hadsell, S. Chopra, and Y. LeCun. 2006. <a href="#">Dimensionality reduction by learning an invariant mapping</a> . In <i>2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)</i> , volume 2, pages 1735–1742.	
682		
683		
684		
685		
686	Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. 2020. <a href="#">Momentum contrast for unsupervised visual representation learning</a> . In <i>2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 9726–9735.	
687		
688		
689		
690		
691	Samuel Humeau, Kurt Shuster, Marie-Anne Lachaux, and Jason Weston. 2020. <a href="#">Poly-encoders: Architectures and pre-training strategies for fast and accurate multi-sentence scoring</a> . In <i>International Conference on Learning Representations</i> .	
692		
693		
694		
695		
696	Zhengbao Jiang, Frank F. Xu, Jun Araki, and Graham Neubig. 2020. <a href="#">How Can We Know What Language Models Know?</a> <i>Transactions of the Association for Computational Linguistics</i> , 8:423–438.	
697		
698		
699		
700	Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. <i>arXiv preprint arXiv:1412.6980</i> .	
701		
702		
703	Ryan Kiros, Yukun Zhu, Ruslan Salakhutdinov, Richard S. Zemel, Antonio Torralba, Raquel Urtasun, and Sanja Fidler. 2015. Skip-thought vectors. In <i>Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 2, NIPS'15</i> , page 3294–3302, Cambridge, MA, USA. MIT Press.	
704		
705		
706		
707		
708		
709		
710	Tassilo Klein and Moin Nabi. 2020. <a href="#">Contrastive self-supervised learning for commonsense reasoning</a> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7517–7523, Online. Association for Computational Linguistics.	
711		
712		
713		
714		
715		
716	Xiang Lisa Li and Percy Liang. 2021. <a href="#">Prefix-tuning: Optimizing continuous prompts for generation</a> .	
717		
718	Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. 2021. Gpt understands, too. <i>arXiv:2103.10385</i> .	
719		
720		
721	Lajanugen Logeswaran and Honglak Lee. 2018. <a href="#">An efficient framework for learning sentence representations</a> . In <i>International Conference on Learning Representations</i> .	
722		
723		
724		
	Marco Marelli, Stefano Menini, Marco Baroni, Luisa Bentivogli, Raffaella Bernardi, and Roberto Zamparelli. 2014. <a href="#">A SICK cure for the evaluation of compositional distributional semantic models</a> . In <i>Proceedings of the Ninth International Conference on Language Resources and Evaluation, LREC 2014, Reykjavik, Iceland, May 26-31, 2014</i> , pages 216–223. European Language Resources Association (ELRA).	733
		734
		735
		736
		737
		738
	Pierre-Emmanuel Mazaré, Samuel Humeau, Martin Raison, and Antoine Bordes. 2018. <a href="#">Training millions of personalized dialogue agents</a> . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 2775–2779, Brussels, Belgium. Association for Computational Linguistics.	739
		740
		741
		742
	Yu Meng, Chenyan Xiong, Payal Bajaj, Saurabh Tiwary, Paul Bennett, Jiawei Han, and Xia Song. 2021. <a href="#">Cocolm: Correcting and contrasting text sequences for language model pretraining</a> .	743
		744
		745
		746
	Junhyun Nam, Hyuntak Cha, Sungsoo Ahn, Jaeho Lee, and Jinwoo Shin. 2020. Learning from failure: Training debiased classifier from biased classifier. In <i>Advances in Neural Information Processing Systems</i> .	747
		748
		749
		750
		751
		752
		753
		754
		755
	Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019. <a href="#">Language models as knowledge bases?</a> In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 2463–2473, Hong Kong, China. Association for Computational Linguistics.	756
		757
		758
		759
		760
	Nils Reimers and Iryna Gurevych. 2019. <a href="#">Sentence-bert: Sentence embeddings using siamese bert-networks</a> . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing</i> . Association for Computational Linguistics.	761
		762
		763
		764
		765
	Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: A simple way to prevent neural networks from overfitting. <i>J. Mach. Learn. Res.</i> , 15(1):1929–1958.	766
		767
		768
		769
		770
	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. <a href="#">Attention is all you need</a> . In <i>Advances in Neural Information Processing Systems</i> , volume 30. Curran Associates, Inc.	771
		772
		773
		774
		775
	Tongzhou Wang and Phillip Isola. 2020. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In <i>International Conference on Machine Learning</i> , pages 9929–9939. PMLR.	776
		777
		778
		779
		780
	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le	

781 Scao, Sylvain Gugger, Mariama Drame, Quentin  
782 Lhoest, and Alexander M. Rush. 2020. [Transformers: State-of-the-art natural language processing](#). In  
783 *Proceedings of the 2020 Conference on Empirical*  
784 *Methods in Natural Language Processing: System*  
785 *Demonstrations*, pages 38–45, Online. Association  
786 for Computational Linguistics.  
787

788 Thomas Wolf, Victor Sanh, Julien Chaumond, and  
789 Clement Delangue. 2019. [Transfertransfo: A transfer](#)  
790 [learning approach for neural network based conver-](#)  
791 [sational agents](#).

792 Zhirong Wu, Yuanjun Xiong, Stella Yu, and Dahua  
793 Lin. 2018. [Unsupervised feature learning via non-](#)  
794 [parametric instance-level discrimination](#).

795 Zhuofeng Wu, Sinong Wang, Jiatao Gu, Madian Khabsa,  
796 Fei Sun, and Hao Ma. 2020. [Clear: Contrastive](#)  
797 [learning for sentence representation](#).

798 Wenhan Xiong, Jingfei Du, William Yang Wang, and  
799 Veselin Stoyanov. 2020. [Pretrained encyclopedia:](#)  
800 [Weakly supervised knowledge-pretrained language](#)  
801 [model](#). In *8th International Conference on Learning*  
802 *Representations, ICLR 2020, Addis Ababa, Ethiopia,*  
803 *April 26-30, 2020*. OpenReview.net.

804 Zonghan Yang, Yong Cheng, Yang Liu, and Maosong  
805 Sun. 2019. [Reducing word omission errors in neural](#)  
806 [machine translation: A contrastive learning approach](#).  
807 In *Proceedings of the 57th Annual Meeting of the As-*  
808 *sociation for Computational Linguistics*, pages 6191–  
809 6196, Florence, Italy. Association for Computational  
810 Linguistics.

811 Yan Zhang, Ruidan He, Zuozhu Liu, Kwan Hui Lim,  
812 and Lidong Bing. 2020. [An unsupervised sentence](#)  
813 [embedding method by mutual information maximiza-](#)  
814 [tion](#). In *Proceedings of the 2020 Conference on*  
815 *Empirical Methods in Natural Language Processing*  
816 *(EMNLP)*, pages 1601–1610, Online. Association for  
817 Computational Linguistics.