

# ENHANCING UNSUPERVISED SENTENCE EMBEDDINGS VIA KNOWLEDGE-DRIVEN DATA AUGMENTATION AND GAUSSIAN-DECAYED CONTRASTIVE LEARNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Recently, using large language models (LLMs) for data augmentation has led to considerable improvements in unsupervised sentence embedding models. However, existing methods encounter two primary challenges: limited data diversity and high data noise. Current approaches often neglect fine-grained knowledge, such as entities and quantities, leading to insufficient diversity. Additionally, unsupervised data frequently lacks discriminative information, and the generated synthetic samples may introduce noise. In this paper, we propose a pipeline-based data augmentation method via LLMs and introduce the Gaussian-decayed gradient-assisted Contrastive Sentence Embedding (GCSE) model to enhance unsupervised sentence embeddings. To tackle the issue of low data diversity, our pipeline utilizes knowledge graphs (KGs) to extract entities and quantities, enabling LLMs to generate more diverse, knowledge-enriched samples. To address high data noise, the GCSE model uses a Gaussian-decayed function to limit the impact of false hard negative samples, enhancing the model’s discriminative capability. Experimental results show that our approach achieves state-of-the-art performance in semantic textual similarity (STS) tasks, using fewer data samples and smaller LLMs, demonstrating its efficiency and robustness across various models.

## 1 INTRODUCTION

Sentence representation learning, a fundamental task in natural language processing (NLP), aims to produce accurate sentence embeddings, thereby improving performance in downstream tasks such as semantic inference (Reimers & Gurevych, 2019), retrieval (Thakur et al., 2021; Wang et al., 2022a), and question answering (Sen et al., 2020). To enhance computational efficiency and reduce labor costs, unsupervised sentence embedding methods based on contrastive learning, such as SimCSE (Gao et al., 2021) and ESIMCSE (Wu et al., 2022c), have emerged as highly effective paradigms. In general, contrastive learning methods operate on the principle that effective sentence embeddings should pull similar sentences closer while pushing dissimilar ones further apart. The performance of unsupervised contrastive learning methods largely depend on the quantity and quality of the samples (Chen et al., 2022), making it crucial to develop strategies that effectively improve both.

Previous studies mainly focused on increasing the number of samples using rule-based word modifications (Wang & Dou, 2023; Wu et al., 2022c) or feature sampling and perturbation techniques (Xu et al., 2023; Chuang et al., 2022a). Recent studies (Zhang et al., 2023; Wang et al., 2024a) use either few-shot manually constructed samples or zero-shot generalized refactoring instructions to create prompts that guide large language models (LLMs) in generating new samples from original sentences, increasing both the quantity and quality of the data. Although these methods have achieved commendable performance, two limitations remain:

**Low Data Diversity.** Diverse data samples in sentence representation learning should contain varied expressions of the same knowledge. However, existing approaches often struggle to distinguish fine-grained semantic knowledge like entities and quantities in the context. Traditional methods modify sentences using limited patterns without considering fine-grained knowledge, restricting their effectiveness in enhancing sample diversity. Recent LLM-based methods like Wang et al. (2024b),

SynCSE (Zhang et al., 2023) and MultiCSR (Wang et al., 2024a), adjust topic and entailment categories in prompts to guide the model in generating varied samples. These methods focus on the global context but lack precise control over the knowledge in the samples. Consequently, the diversity of generated samples is constrained by the probability distributions of LLMs, resulting in unpredictable data quality.

**High Data Noise.** Unsupervised sentence representation learning often suffers from data noise caused by confusing negative samples, which mainly arise from two sources. First, traditional methods generate datasets by duplicating samples to create positive instances, leading to negatives with similar surface-level semantics that affect the model’s performance (Miao et al., 2023; Zhou et al., 2022). Second, in data synthesis, differences in semantic distributions can cause the LLM’s criteria for distinguishing between positive and negative samples to misalign with the target domain, introducing additional noise (Huang et al., 2023; Poerner & Schütze, 2019). The existing MultiCSR method attempts to remove noisy samples using linear programming, but this can eliminate potentially valuable samples and reduce data diversity. Figure 1 compares various baselines on the STS-Benchmark development set. The results show that the prediction of false positives outnumber false negatives, and data synthesis in SynCSE increases false negatives, further supporting the above analysis.

In this paper, we propose a pipeline-based data augmentation method using LLMs and introduce the Gaussian-decayed gradient-assisted Contrastive Sentence Embedding (GCSE) model to improve the performance of unsupervised sentence embedding methods. To address the issue of *low data diversity*, we begin by extracting entities and quantities from the data samples and constructing a knowledge graph (KG) with the extracted data. Next, we create a sentence construction prompt using the extracted knowledge to guide LLM in generating more diverse positive samples. To tackle *high data noise*, we employ an evaluation model to annotate the synthesized data and initially filter out false samples. However, this procedure is ineffective in filtering out false negatives with similar surface-level semantics. To balance sample diversity while minimizing the impact of noise from false negatives, we aim to align all hard negatives with the distribution of the evaluation model in the initial training step. Then, we leverage other in-batch negative samples to optimize the semantic space. Therefore, we propose the GCSE model that employs a Gaussian-decayed function to calculate the prediction distinctions between GCSE and the evaluation model. It first declines the gradients of hard negatives. As training progresses, the gradient weights for hard negatives that diverge farther from the evaluation model’s distribution progressively recover. This function helps prevent false negatives from being pushed further away in the semantic space, leading to a more uniform distribution. We highlight the key innovations of our approach in Table 1: (i) We are the first to incorporate fine-grained knowledge for sample synthesis in LLM-based methods. (ii) Unlike MultiCSR’s denoising approach, our method retains more false samples for training rather than discarding them. (iii) Our data selection strategy focuses on domain-specific samples, using a local LLM with fewer samples for synthesis, leading to improved performance. Experimental results demonstrate the efficiency of our model, outperforming previous best methods in average scores for semantic textual similarity (STS) tasks by 1.05% with BERT-base, 1.62% with BERT-large, 0.49% with RoBERTa-base, and 1.50% with RoBERTa-large.

In summary, our contributions are as follows: (1) *New method.* We introduce a pipeline-based data augmentation method using LLM for few-shot domain data and propose a Gaussian-decayed

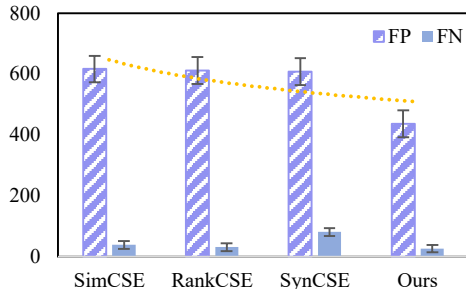


Figure 1: Comparison of false positives (FP) and negatives (FN). Both the predicted scores and labels are normalized (see details in Appendix I), where positives have a score greater than the label, while negatives lower than the label. False samples are identified when the root mean square error (RMSE) between the prediction and the label exceeds 0.2.

Methods	Synthesis Approach	Use Knowledge	Denoise
SynCSE	Few-shot Synthesis	No	No
MultiCSR	Zero-shot Synthesis	No	Yes
Ours	Zero-shot Synthesis	Yes	Yes

Table 1: Comparison of our methods and related LLM-based methods.

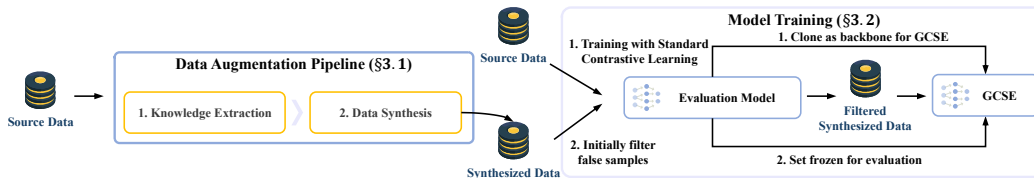


Figure 2: The overall workflow of our method.

gradient-assisted Contrastive Sentence Embedding (GCSE) model to reduce data noise. (2) *New perspective*. To the best of our knowledge, we are the first to explore combining knowledge graphs with LLM to synthesize data, enhancing fine-grained sentence representation learning by generating diverse positive and negative samples. (3) *State-of-the-art performance*. Experimental results demonstrate that our method achieves superior performance on STS tasks while using fewer samples for data synthesis with smaller LLM parameters.

## 2 RELATED WORK

Early work on sentence embeddings builds on the distributional hypothesis, predicting surrounding sentences (Kiros et al., 2015; Logeswaran & Lee, 2018; Hill et al., 2016) or extending the word2vec framework (Mikolov et al., 2013) with n-gram embeddings (Pagliardini et al., 2018). Post-processing techniques like BERT-flow (Li et al., 2020) and BERT-whitening (Su et al., 2021) address the anisotropy issue in pre-trained language models (PLMs), and more recent methods focus on generative approaches (Wang et al., 2021; Wu & Zhao, 2022) and regularizing embeddings to prevent representation degeneration (Huang et al., 2021). Recently, contrastive learning approaches have become prominent, using various augmentation methods to derive different views of the same sentence (Zhang et al., 2020; Giorgi et al., 2021; Kim et al., 2021; Gao et al., 2021). Among these, SimCSE uses dropout as a simple augmentation and achieves strong results in unsupervised STS tasks, inspiring further approaches like ArcCSE (Zhang et al., 2022), DiffCSE (Chuang et al., 2022a), GS-InfoNCE (Wu et al., 2022b), and RankCSE (Liu et al., 2023).

With the advent of LLM (OpenAI, 2023; Bai et al., 2023; Touvron et al., 2023), some works attempt to utilize LLM for sentence representation learning. For example, Ni et al. (2022) uses T5 with mean pooling to obtain a sentence embedding model by fine-tuning on a large-scale NLI corpus; Cheng et al. (2023) uses prompt learning to measure the semantic similarity of sentence pairs; Springer et al. (2024) employs sentence repetition to enhance the capacity for sentence representation; AoE (Li & Li, 2024a) optimize angle differences for improving supervised text embedding; and BeLLM (Li & Li, 2024b) designs a Siamese structure for learning sentence embeddings.

## 3 METHODOLOGY

In this section, we present the data augmentation pipeline via LLM and the specific structure of the GCSE. As shown in Figure 2, we start by using a data augmentation pipeline to synthesize new samples from the source data, and then train our model with the filtered synthetic data.

### 3.1 DATA AUGMENTATION

In the data augmentation pipeline, we utilize both domain data and partial general data to balance domain-specific relevance and general-domain applicability. We start by extracting knowledge from the source data and then synthesize new data for our model training. The detailed structure of the pipeline is shown in Figure 3.

**Knowledge Extraction and Integration.** The variety and relationships between samples directly impact model performance in sentence representation learning. A major challenge with existing LLM-based data synthesis methods is the limited diversity they generate for each short text. To trade off the low diversity of the generated samples with their relevance to the domain semantic space, we first design an extraction prompt to obtain entities and quantities from the given data.

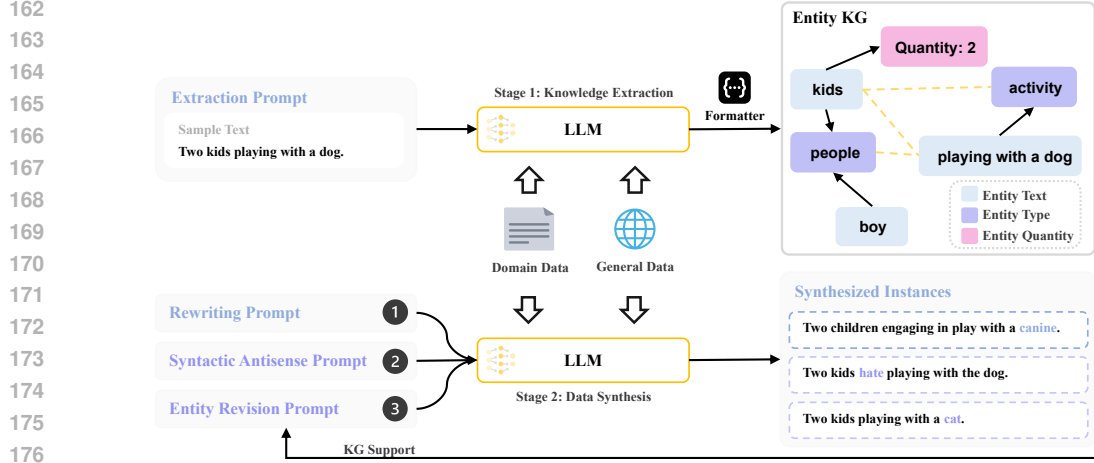


Figure 3: The pipeline of knowledge extraction and data synthesis, where the solid black arrows in the Entity KG are hard edges, and dotted yellow lines are soft edges.

Formally, we denote the extraction prompt as  $\mathcal{P}_e$ , and LLM  $\mathcal{L}$ , suppose we finally extract instances with  $d$  sample number, the knowledge set  $\mathcal{K}_i = \{k_{i1}, \dots, k_{in}\}$  of each instance  $x_i$  is computed in Equation 1, where  $t_j$ ,  $c_j$  and  $q_j$  represent the entity text, entity type and quantity of  $k_i$ .  $n_i$  is the size of  $\mathcal{K}_i$ , and  $\mathcal{F}(\cdot)$  is the formatting function that convert text to triplet. Next, we integrate all knowledge by establishing an entity knowledge graph  $\mathcal{G} = \langle V, E \rangle$ , where the node set  $V$  contains all the  $\langle t, c, q \rangle$  from  $\mathcal{K}$ :

$$\mathcal{K} = \bigcup_{i=1}^d \mathcal{F}([\mathcal{P}_e; x_i], \mathcal{L}) = \bigcup_{i=1}^d \{\langle t_{ij}, c_{ij}, q_{ij} \rangle \mid j \in [1, n_i]\}, \quad (1)$$

$$V = \{t_{ij}, c_{ij}, q_{ij} \mid i \in [1, d]; j \in [1, n_i]\}. \quad (2)$$

The edges  $E$  consist of hard edges  $E_r$  and soft edges  $E_s$ . As shown in Equations 3 and 4,  $E_r$  represents the relationship between the entity text, type and quantity of each  $k \in \mathcal{K}$ , and  $E_s$  indicates the relationship between entity text in  $k_{ij}$  and other entity text or type in the same instance  $x_i$ .

$$E_r = \{(t_{ij}, c_{ij}) \cup (t_{ij}, q_{ij}) \mid i \in [1, d]; j \in [1, n_i]\}, \quad (3)$$

$$E_s = \bigcup_{i=1}^d \{(t_{ij}, t_{ik}), (t_{ij}, c_{il}) \mid k, l \neq j; j, k, l \in [1, n_i]\}. \quad (4)$$

By defining hard and soft edges, we can more efficiently identify and replace entity nodes near the current node, improving the correlation between the synthesized instance and the source instance.

**Data Synthesis via LLM.** Empirical evidence and model performance on standard datasets show that sentence embedding models struggle more with accurately identifying negative samples than positives (Chuang et al., 2022a; Miao et al., 2023). In the contrastive learning methods, the model acquires sentence embedding representation by calculating the distance between sentence-pairs. It aims to minimize the spatial distance between positive pairs and increase the spatial distance between negative pairs. Thus, it is essential to obtain negative samples that closely resemble the source instance in surface-level features, while positive samples should have diverse representations but still convey the same meaning as the source instance.

In this study, we use LLM to generate positive samples through a rewrite prompt. We also focus on the impact of variations in entities and quantities within the samples. Negative samples are generated by the LLM at both the syntactic and fine-grained knowledge levels. The data synthesis prompts are divided into three main types: (1) Rewriting prompt, (2) Syntactic antisense prompt, and (3) Entity revision prompt. The first type is used to create positive samples, while the second and third types are used to create negative samples at the syntactic and knowledge levels, respectively.

The “rewriting prompt” can be classified into three forms: directly requesting LLM to generate a new sentence instance using the “rewrite” instruction, creating the preceding part of the sentence instance, and generating based on the knowledge set of the instance. As the diversity of synthetic samples increases, the likelihood of generating false positives also rises. To address this, the next section involves scoring the generated samples using an evaluation model. The “syntactic antisense prompt” aims to modify the semantics to create a contradiction at the syntactic level. Such as transforming it into a positive/negative statement using explicit positive/negative words, or by expressing a contrary sentiment. This is an initial approach to synthesizing negative samples that preserves a strong coherence with the source instance in terms of sequence structure. However, it is deficient in generation diversity. To alleviate the issue, the “entity revision prompt” aims to enhance text diversity by replacing the entity text and quantity compared to the source instance. Simultaneously, to ensure the semantic relevance between the synthetic samples and the source instance, replacement entities are selected by searching for neighboring nodes on entity KG. We define  $\mathcal{T}(\cdot)$  as the search function, and the replacement entity of  $t_{ij}$  are computed as:

$$\mathcal{T}_r(t_{ij}) = \{t_{ip} \mid (t_{ij}, c_{ik}) \in E_r \wedge (t_{ip}, c_{ik}) \in E_r\}, \quad (5)$$

$$\mathcal{T}_s(t_{ij}) = \{t_{ip} \mid (t_{ij}, t_{ip}) \in E_s\}, \quad (6)$$

$$\mathcal{T}_p(t_{ij}) = \{t_{ip} \mid t_{ik} \in \mathcal{T}_s(t_{ij}) \cap \mathcal{T}_s(t_{ip}) \wedge t_{ip} \in \mathcal{T}_r(t_{ij})\}, \quad (7)$$

$$\mathcal{T}(t_{ij}) = \mathcal{T}_r(t_{ij}) \cup \mathcal{T}_p(t_{ij}), \quad (8)$$

where the function  $\mathcal{T}_r(\cdot)$  is used to search for entities that have a hard edge with the current entity, and  $\mathcal{T}_s(\cdot)$  is used to search for entities that have a soft edge with the current entity.  $\mathcal{T}_p(\cdot)$  aims to search for  $t_{ip}$ , that is of the same type as  $t_{ij}$ , and they both have soft edges with another in-context entity  $t_{ik}$ . Finally, the replacement entity can be randomly selected from the result of the search function  $\mathcal{T}(t_{ij})$ . Compared to randomly replacing entities, our strategy enhances the semantic relevance between the generated sample and the source instance.

### 3.2 MODEL TRAINING

The training process of our model consists of two stages. First, we combine general and domain-specific data to train an evaluation model using standard unsupervised contrastive learning. This improves the uniformity of sentence embeddings in general scenarios and reduces the impact of semantic distribution limitations in the synthesized data, enhancing model robustness. Then, we freeze the evaluation model to filter synthetic data and help the GCSE model eliminate false hard negative sample noise.

**General Contrastive Learning.** In the first stage, we follow the formulation of SimCSE (Gao et al., 2021) to train the evaluation model. Formally, we define the encoder of the evaluation model as  $E'$ , each unlabeled sentence instance as  $x_i$ , and its positive sample as  $x_i^+ = x_i$ . The representation of each instance is denoted as  $\mathbf{h}' = \mathcal{F}_{E'}(x)$ , the representations of  $x_i$  and  $x_i^+$  are computed as  $\mathbf{h}'_i$  and  $\mathbf{h}'_i^+$ , respectively. Since the dropout mask in  $E'$  is random,  $\mathbf{h}'_i$  and  $\mathbf{h}'_i^+$  are computed with the same input but with slightly different results. Then, the loss of evaluation model is defined as:

$$-\log \frac{e^{\text{sim}(\mathbf{h}'_i, \mathbf{h}'_i^+)/\tau}}{\sum_{j=1}^N e^{\text{sim}(\mathbf{h}'_i, \mathbf{h}'_j^+)/\tau}}, \quad (9)$$

where  $N$  represents the size of each mini-batch,  $\tau$  is a temperature hyperparameter, and  $\text{sim}(\cdot)$  is the cosine similarity function.

**Denosing Training.** In the second stage, we adopt a copy of the evaluation model as the backbone of GCSE and continue training on synthesized data. In this stage, each input is set as a triplet  $(x_i, x_i^+, x_i^-)$ , where  $x_i^+$  and  $x_i^-$  stand for the positive and negative samples of  $x_i$ , respectively. Nevertheless, the synthesized data contains many potential false positive and false negative samples, necessitating the implementation of a filtering process. We use the frozen evaluation model to initially correct these inaccurate samples and build the ultimate triplet dataset. Let  $\mathcal{S}(x_i) = \{\hat{x}_{i1}, \dots, \hat{x}_{im}\}$  denotes the synthetic data set of  $x_i$ , where  $m$  is the size of the set, and  $x_i^+, x_i^-$  are calculated as:

$$x_i^+ = \begin{cases} \hat{x}_{ij}, & \text{sim}(\mathbf{h}'_i, \hat{\mathbf{h}}'_{ij}) \geq \alpha, j \in [1, m] \\ x_i, & \text{else} \end{cases}, \quad (10)$$

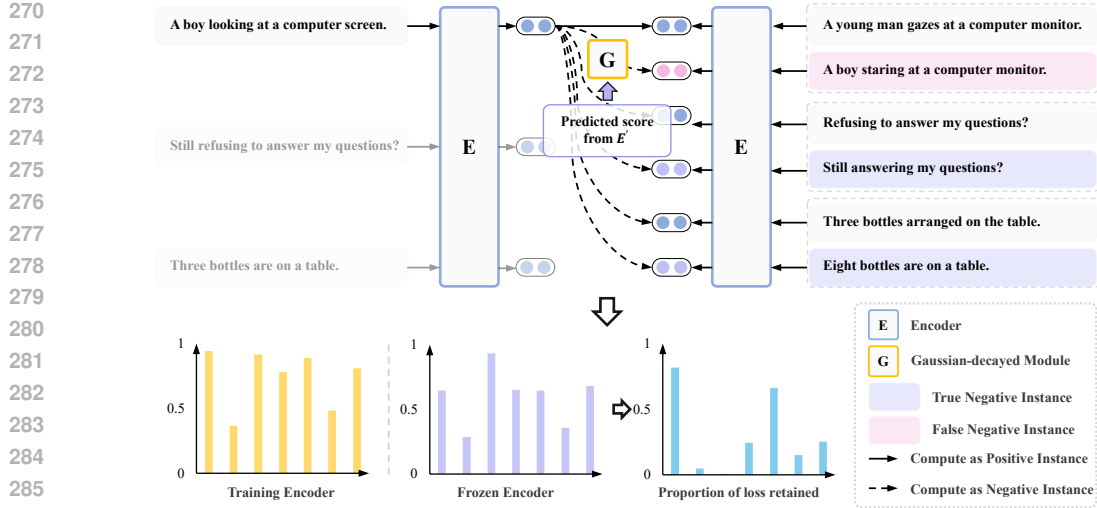


Figure 4: In-batch training with Gaussian-decayed on GCSE.

$$x_i^- = \begin{cases} \hat{x}_{ij}, & \text{sim}(\mathbf{h}'_i, \hat{\mathbf{h}}'_{ij}) \leq \beta, j \in [1, m] \\ x_k, & k \in [1, N], k \neq i \end{cases}, \quad (11)$$

where  $\alpha, \beta$  are the threshold for positives and negatives, respectively.  $x_k$  denotes a randomly selected instance from in-batch data. We can set a high value for  $\alpha$  to reduce false positive samples. However, filtering out false negatives in synthetic data is more challenging. In theory, smaller  $\beta$  can reduce more false negatives, but samples with low similarity to the source instance are easy to distinguish due to significant surface-level differences. As a result, training on these samples does not effectively improve the model’s ability to distinguish fine-grained false positives. Therefore, we opt for a higher value of  $\beta$ . During training, we use a Gaussian-decayed function to align the distances of hard negative samples between the GCSE encoder  $E$  and the frozen encoder  $E'$ . As shown in Figure 4, for each mini-batch of triplet inputs, both  $E$  and  $E'$  compute similarity scores for the negative samples and their corresponding source instances. The loss for each instance in GCSE is defined as:

$$-\log \frac{e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_i^+)/\tau}}{\sum_{j=1}^N e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_j^+)/\tau} + \sum_{\substack{j=1 \\ j \neq i}}^N e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_j^-)/\tau} + G(s_i, s'_i, \tau, \sigma)}, \quad (12)$$

$$G(s_i, s'_i, \tau, \sigma) = s_i \left( 1 - e^{-\frac{(s_i - s'_i)^2 \tau^2}{2\sigma^2}} \right), \quad (13)$$

where  $s_i = \text{sim}(\mathbf{h}_i, \mathbf{h}_i^-)$ ,  $s'_i = \text{sim}(\mathbf{h}'_i, \mathbf{h}'_i^-)$ .  $G(\cdot)$  is the Gaussian-decayed function, where the loss attenuation of the hard negative sample grows as the distance between  $s_i$  and  $s'_i$  decreases, and  $\sigma$  is a hyperparameter that controls the width of  $G(\cdot)$ . This implies that when  $E$  initially calculates the hard negative sample, it follows the spatial distribution of  $E'$  as the “established guidelines” and uses other in-batch negative samples to further increase the spatial distance between negatives, effectively reducing the influence of false negatives. As training progresses, the spatial distribution of true hard negatives between  $E$  and  $E'$  will progressively increase, and its gradient will be restored.

## 4 EXPERIMENT

### 4.1 EXPERIMENT SETUP

**Training:** We utilize the subset of NLI dataset from Gao et al. (2021) as the general data, and use the training sets from STS-Benchmark (STS-B) (Cer et al., 2017) with 5.7k samples and SICK (Marelli et al., 2014) with 4.5k samples as the domain data for a fair comparison with related approaches.

To simulate the unsupervised scenario, we exclusively include unlabeled samples from the dataset. In this experiment, the ratio of sample numbers between domain data and general data was 1:3. We adopt ChatGLM3-6B (GLM et al., 2024), GLM4-9B (GLM et al., 2024) and ChatGPT (OpenAI, 2022) as LLMs for data synthesis, respectively. We choose BERT (Devlin et al., 2019) and RoBERTa (Liu et al., 2019) as the backbone models of GCSE. In the stage of Gaussian-decayed training on synthesized data, the filtering thresholds of  $\alpha$  and  $\beta$  are set as 0.9 and 0.75, respectively. The temperature of  $\tau$  is set as 0.05, and the  $\sigma$  of  $G(\cdot)$  is set as 0.01. In the first stage training, the evaluation model is firstly trained on the unlabeled dataset of all general data and domain data. One copy instance of the evaluation model is then utilized as the pre-trained model for GCSE, while the original instance is set to be frozen to filter synthesized data and provide guidance for GCSE. In the second stage, GCSE is trained on the filtered synthesized data, and the sentence embedding is obtained from the last output hidden states of the first token.

**Evaluation:** To validate our method for sentence embeddings, we evaluated the model’s performance on semantic textual similarity (STS) tasks, we use the standard evaluation method, measuring model performance with Spearman’s correlation, and we adopt SentEval<sup>1</sup> (Conneau & Kiela, 2018) as the evaluation tool, which contains seven STS subsets: STS 2012-2016 (Agirre et al., 2012; 2013; 2014; 2015; 2016), the STS-Benchmark (Cer et al., 2017) and the SICK Relatedness (Marelli et al., 2014). To compare the ranking performance of our method on retrieval tasks, we evaluated the model using the MTEB benchmark (Muennighoff et al., 2023) with four reranking datasets: AskUbuntuDupQuestions (Lei et al., 2016), MindSmallReranking (Wu et al., 2020), SciDocsRR (Cohan et al., 2020) and StackOverflowDupQuestions (Liu et al., 2018), and follow the same settings of Zhang et al. (2023) by using Mean Average Precision (MAP) as the metric. Additionally, we compared the performance of our model with other methods on transfer tasks in SentEval to evaluate its applicability in Appendix C.

**Baselines:** We compare our method with mainstream unsupervised sentence embedding baselines: BERT-whitening (Su et al., 2021), SimCSE (Gao et al., 2021), DiffCSE (Chuang et al., 2022b), PromptBERT (Jiang et al., 2022), PCL (Wu et al., 2022a), CARDS (Wang et al., 2022b), DebCSE (Miao et al., 2023) and RankCSE (Liu et al., 2023). In addition, we further compare two baselines: SynCSE (Zhang et al., 2023) and MultiCSR (Wang et al., 2024a), which use LLM for data synthesizing in whole NLI datasets. To verify the effectiveness of our data synthesis method, we choose their results of using ChatGPT for comparison.

## 4.2 MAIN RESULTS

**STS Tasks:** The overall results of the STS tasks are shown in Table 2. Our approach, utilizing synthetic samples from ChatGPT, achieves state-of-the-art performance across all backbones when compared to other unsupervised baselines. Even with synthetic samples from ChatGLM3-6B, our method still outperforms previous approaches on BERT-base, BERT-large, and RoBERTa-large. This highlights the applicability of our method, as it can be effectively applied to multiple models. Compared to the standard unsupervised SimCSE, Spearman’s correlation of GCSE (ChatGLM3-6B) is improved by an average of 17.24% on the base models and 3.44% on the large models. On the strong baseline RankCSE, GCSE (ChatGLM3-6B) achieved a 1.36% improvement over its average performance, demonstrating the effectiveness of the LLM data synthesis process. Furthermore, we compare two baseline models: SynCSE and MultiCSR, both of which utilize LLM as the data synthesis model. We specifically analyze the results of using ChatGPT for both models and the results show that our approach outperforms both models in most cases. It should be noted that our method only utilizes 14% of the sample size compared to the other two methods that employ the entire NLI datasets. This demonstrates the effectiveness of our data synthesis strategy and domain-oriented sample selection strategy.

**Reranking Tasks:** Table 3 presents the MAP results of our approach and related baselines on the reranking benchmark, and all models are evaluated on the test sets of the reranking benchmark without using the training sets. The results indicate that various approaches exhibit varying performance on different datasets, which can be attributed to the distinct semantic distribution and evaluation scale of each dataset. Our GCSE outperforms SynCSE by 0.39% in average MAP score and achieves the

<sup>1</sup><https://github.com/facebookresearch/SentEval>

Model	Method	STS-12	STS-13	STS-14	STS-15	STS-16	STS-B	SICK-R	Avg.
BERT-base	whitening <sup>†</sup>	57.83	66.90	60.90	75.08	71.31	68.24	63.73	66.28
	SimCSE <sup>†</sup>	68.40	82.41	74.38	80.91	78.56	76.85	72.23	76.25
	DiffCSE <sup>†</sup>	72.28	84.43	76.47	83.90	80.54	80.59	71.23	78.49
	PromptBERT <sup>♣</sup>	71.56	84.58	76.98	84.47	80.60	81.60	69.87	78.54
	PCL <sup>♠</sup>	72.84	83.81	76.52	83.06	79.32	80.01	73.38	78.42
	DebCSE <sup>†</sup>	76.15	84.67	78.91	85.41	80.55	82.99	73.60	80.33
	RankCSE <sup>♠</sup>	75.66	<b>86.27</b>	77.81	84.74	81.10	81.80	75.13	80.36
	SynCSE (ChatGPT)*	75.86	82.19	78.71	<b>85.63</b>	81.11	82.35	<u>78.79</u>	80.66
	MultiCSR (ChatGPT) <sup>♣</sup>	74.86	84.19	79.46	84.70	80.34	<u>83.59</u>	<b>79.37</b>	80.93
	<b>GCSE (ChatGLM3-6B)</b>	76.91	<u>86.23</u>	<u>80.49</u>	<u>85.16</u>	<u>81.45</u>	<u>82.54</u>	<u>75.71</u>	<u>81.21</u>
	<b>GCSE (GLM4-9B)</b>	78.19	84.88	<u>80.28</u>	<u>84.39</u>	<b>81.81</b>	<b>83.89</b>	<u>77.74</u>	<u>81.60</u>
<b>GCSE (ChatGPT)</b>	<b>78.20</b>	85.90	<b>81.17</b>	84.88	81.44	83.56	78.69	<b>81.98</b>	
BERT-large	SimCSE <sup>†</sup>	70.88	84.16	76.43	84.50	79.76	79.26	73.88	78.41
	PCL <sup>♠</sup>	74.87	86.11	78.29	85.65	80.52	81.62	73.94	80.14
	DebCSE <sup>†</sup>	76.82	86.36	79.81	<b>85.80</b>	80.83	83.45	74.67	81.11
	RankCSE <sup>♠</sup>	75.48	86.50	78.60	85.45	81.09	81.58	75.53	80.60
	SynCSE (ChatGPT)*	74.24	85.31	79.41	85.71	<b>81.76</b>	82.61	<u>79.25</u>	81.18
	<b>GCSE (ChatGLM3-6B)</b>	76.99	<b>87.34</b>	80.88	85.47	<b>80.55</b>	<b>82.97</b>	<u>75.68</u>	81.41
	<b>GCSE (GLM4-9B)</b>	76.94	86.69	81.16	85.53	81.44	<b>84.47</b>	78.88	<u>82.16</u>
<b>GCSE (ChatGPT)</b>	<b>78.70</b>	<u>87.30</u>	<b>81.94</b>	<b>86.10</b>	<u>81.60</u>	<u>84.08</u>	<b>79.86</b>	<b>82.80</b>	
RoBERTa-base	whitening <sup>†</sup>	46.99	63.24	57.23	71.36	68.99	61.36	62.91	61.73
	SimCSE <sup>†</sup>	70.16	81.77	73.24	81.36	80.65	80.22	68.56	76.57
	DiffCSE <sup>†</sup>	70.05	83.43	75.49	82.81	82.12	82.38	71.19	78.21
	PromptRoBERTa <sup>♣</sup>	73.94	84.74	77.28	84.99	81.74	81.88	69.50	79.15
	PCL <sup>♠</sup>	71.13	82.38	75.40	83.07	81.98	81.63	69.72	77.90
	DebCSE <sup>†</sup>	74.29	85.54	79.46	85.68	81.20	83.96	74.04	80.60
	RankCSE <sup>♠</sup>	73.20	<b>85.95</b>	77.17	84.82	82.58	83.08	71.88	79.81
	SynCSE (ChatGPT) <sup>††</sup>	74.61	83.76	77.89	85.09	82.28	82.71	78.88	80.75
	MultiCSR (ChatGPT) <sup>♣</sup>	75.61	84.33	80.10	84.98	82.13	84.54	<b>79.67</b>	81.62
	<b>GCSE (ChatGLM3-6B)</b>	76.06	85.30	<u>80.38</u>	85.28	<b>83.26</b>	<u>84.07</u>	<u>74.55</u>	<u>81.27</u>
<b>GCSE (GLM4-9B)</b>	77.13	85.05	<u>80.25</u>	84.89	83.08	<b>84.78</b>	76.63	<u>81.69</u>	
<b>GCSE (ChatGPT)</b>	<b>78.03</b>	83.79	<b>80.61</b>	<b>86.28</b>	82.76	84.31	<u>79.01</u>	<b>82.11</b>	
RoBERTa-large	SimCSE <sup>†</sup>	72.86	83.99	75.62	84.77	81.80	81.98	71.26	78.90
	PCL <sup>♠</sup>	74.08	84.36	76.42	85.49	81.76	82.79	71.51	79.49
	DebCSE <sup>†</sup>	77.68	87.17	80.53	85.90	83.57	85.36	73.89	82.01
	RankCSE <sup>♠</sup>	73.20	85.83	78.00	85.63	82.67	84.19	73.64	80.45
	SynCSE (ChatGPT) <sup>††</sup>	75.45	85.01	80.28	86.55	83.95	84.49	<b>80.61</b>	82.33
	<b>GCSE (ChatGLM3-6B)</b>	<b>78.24</b>	<u>87.24</u>	<u>81.93</u>	<u>86.80</u>	83.52	85.08	<u>76.70</u>	<u>82.79</u>
	<b>GCSE (GLM4-9B)</b>	77.18	<u>86.72</u>	<b>82.62</b>	85.89	83.97	85.75	77.97	82.87
<b>GCSE (ChatGPT)</b>	77.76	<b>87.45</b>	<b>82.62</b>	<b>88.38</b>	<b>84.43</b>	<b>86.08</b>	80.09	<b>83.83</b>	

Table 2: Comparison of Spearman’s correlation results on STS tasks, where the value highlighted in bold is the best value, and the value underlined is the second-best value. “<sup>†</sup>”: results from Miao et al. (2023), “<sup>♣</sup>”: results from Wang et al. (2024a), “<sup>♠</sup>”: results from Liu et al. (2023), “<sup>††</sup>”: results from Zhang et al. (2023). “\*”: we reproduce the results with the officially released corpus from Zhang et al. (2023). **GCSE has significant differences with all comparable baselines on the t-test ( $p < 0.5\%$ ).**

best results in all backbone models, demonstrating the efficacy of our approach in enhancing the precision of unsupervised ranking tasks.

### 4.3 ANALYSIS

**Ablation Studies:** We analyze the impact of each module or strategy in GCSE and report the results in Table 4. First, “w/o stage-2” refers to the results obtained without training in the second stage. This leads to a significant decrease in performance compared to the default model, which is the performance of the evaluation model and is similar to the conventional unsupervised SimCSE. Then, “w randomly” refers to the direct use of the instance itself as a positive sample in the combination dataset of domain and general data, while randomly selecting a negative instance from the dataset. We can observe that its performance in this case is even worse than the evaluation model. This demonstrates that the diversity of positive samples and the quality of negative samples significantly impact the performance of the model. “w/o filtering” indicates the results of training by skipping evaluation model filtering and directly using the data synthesized by LLM. The results show that the performance of the model is significantly affected when false positive and negative samples are introduced without filtering. We investigate the impact of the Gaussian-decayed function by removing it, and the results are shown in “w/o decay”. We can observe that the default model performs better overall than when the Gaussian-decayed function is removed, indicating that it can filter out potential false negative sample noise. Finally, we analyze the necessity of including general data and domain data in “w/o general” and “w/o domain” respectively. It can be observed that



removing either of them results in a decline in performance, which indicates the significance of domain data and the essentiality of general data in our method.

#### Analysis of entities and quantities awareness:

We analyze GCSE awareness of entities and quantities by constructing a dataset using the data synthesis method in Section 3.1 on the STS-Benchmark development set. Then, the similarity scores of each triplet in the dataset are annotated by two supervised pre-trained models: “sup-simcse-bert-large” and “sup-simcse-roberta-large”. The final label is the average score of the similarity calculated by both models. We evaluate Spearman’s correlation scores of GCSE and the other three strong baselines on the backbone of the BERT-base model, and the results are shown in Table 5. Our GCSE achieves the best result and outperforms RankCSE by 14.03%. In this case, both SynCSE and GCSE achieve significant improvements over methods without LLM. This might be due to the similarity of the semantic representation space between the training set and the development set, both of which are synthesized via LLM. Nevertheless, GCSE shows a notable enhancement in performance of 2.19% compared to SynCSE, demonstrating that its understanding of the entities and quantities in sentences has enhanced to a certain degree.

Method	Spearman’s
unsup-SimCSE	75.59
RankCSE	79.74
SynCSE (ChatGPT)	91.58
GCSE (ChatGLM3-6B)	<b>93.77</b>

Table 5: Comparison of Spearman’s correlation results on the synthetic data of the STS-Benchmark development set.

**Impact on the ratio between domain and general data:** Figure 5 presents the trend of the GCSE Spearman’s correlation result as the proportion of general data introduced increases, where “d” represents that only using the domain data. The results show that adding a certain amount of general data improves performance on STS tasks. However, when the size of general data exceeds three times that of domain data, performance starts to decline. This suggests that incorporating a moderate amount of external data enhances the uniformity of sentence embeddings. But as the out-of-domain data grows, the influence of domain-specific data on training weakens. Overall, the results indicate that domain data improves the model’s ability to represent target domain sentences, while general data helps with sentence embedding uniformity.

**Impact of the Gaussian-decayed:** To further investigate the effectiveness of the Gaussian-decayed function, we analyze the GCSE performance against the weight of  $\sigma$  on the synthesized data, both with and without filtering. As shown in Figure 6, we use the synthesized data without filtering to evaluate the efficacy of the Gaussian-decayed function in eliminating false negative samples, and

Model	Method	AskU.	Mindsmall	SciDocsRR	StackO.	Avg.
BERT-base	SimCSE	51.89	28.68	67.88	<u>39.60</u>	47.01
	PCL	52.46	28.72	68.03	<b>41.30</b>	47.63
	SynCSE (ChatGPT)*	<u>52.61</u>	<b>29.17</b>	<u>68.46</u>	38.60	47.21
	GCSE (ChatGLM3-6B)	<b>52.62</b>	28.79	<b>70.67</b>	39.53	<b>47.90</b>
BERT-large	SimCSE	53.10	29.59	71.94	40.68	48.83
	PCL	52.03	29.11	70.30	<b>42.33</b>	48.44
	SynCSE (ChatGPT)*	<u>53.24</u>	<b>30.09</b>	71.45	39.24	48.50
	GCSE (ChatGLM3-6B)	<b>53.40</b>	29.43	<b>73.04</b>	39.68	<b>48.89</b>
RoBERTa-base	SimCSE††	52.78	29.91	65.96	39.25	46.95
	CARDS††	52.94	27.92	64.62	<b>41.51</b>	46.75
	PCL††	51.85	27.92	64.70	41.18	46.41
	SynCSE (ChatGPT)††	53.27	<b>30.29</b>	67.55	39.39	47.63
	GCSE (ChatGLM3-6B)	<b>53.44</b>	29.35	<b>67.89</b>	41.13	<b>47.95</b>
RoBERTa-large	SimCSE††	<u>55.10</u>	29.23	68.54	<u>42.56</u>	48.86
	CARDS††	53.83	29.07	68.26	<b>43.24</b>	48.60
	PCL††	53.43	28.56	66.06	41.54	47.40
	SynCSE (ChatGPT)††	<b>55.48</b>	30.27	70.85	40.00	49.15
	GCSE (ChatGLM3-6B)	54.05	<b>30.30</b>	<b>71.23</b>	41.65	<b>49.31</b>

Table 3: Comparison of Mean Average Precision (MAP) results on reranking tasks, where the value highlighted in bold is the best value, and the value underlined is the second-best value. “††”: results from Zhang et al. (2023). “\*”: we reproduce the results with the officially released corpus from Zhang et al. (2023).

Method	STS-12	STS-13	STS-14	STS-15	STS-16	STS-B	SICK-R	Avg.
GCSE (ChatGLM3-6B)	<b>76.91</b>	<b>86.23</b>	<b>80.49</b>	85.16	<b>81.45</b>	82.54	75.71	<b>81.21</b>
w/o stage-2	71.85	83.65	76.84	83.37	78.74	79.10	71.69	77.89
w randomly	71.94	84.03	76.99	83.65	79.11	78.66	69.28	77.67
w/o filtering	74.65	83.54	77.39	83.27	79.97	79.66	74.27	78.96
w/o decay	76.15	85.83	79.77	<b>85.19</b>	80.72	<b>82.59</b>	75.55	80.83
w/o general	75.44	85.55	79.19	84.91	80.23	81.57	74.14	80.15
w/o domain	75.59	85.66	78.93	84.09	80.87	82.29	<b>76.00</b>	80.49

Table 4: Ablation studies of STS tasks on BERT-base. Other PLMs yield similar patterns to BERT-base.

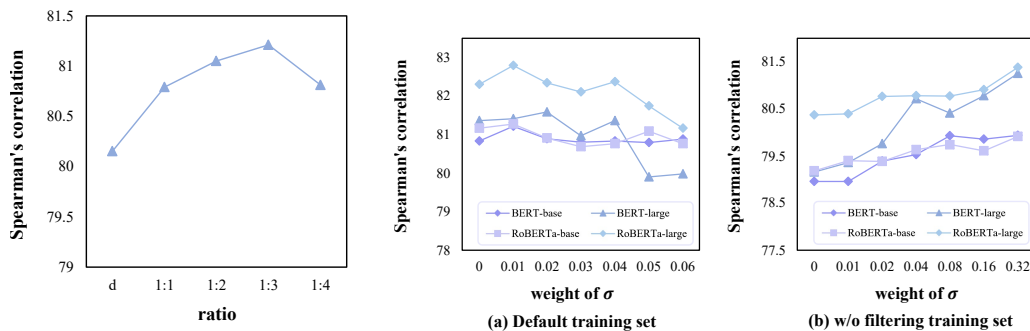


Figure 5: Spearman's correlation against the ratio of domain data to general data on the STS tasks.

Figure 6: Spearman's correlation against the weight of Gaussian-decay on the STS tasks.

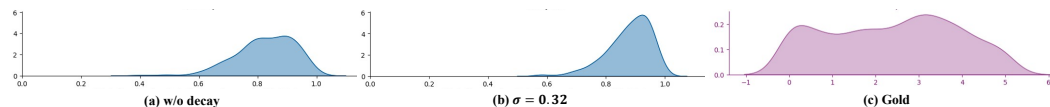


Figure 7: Density plots of the STS-Benchmark development set with labels  $\geq 4$ , which is evaluated by GCSE with different  $\sigma$  weights. (c) is the density plot of gold labels.

results are presented in Figure 6 (b). It is clear that the model's performance improves as the weight of  $\sigma$  grows. This suggests that a greater  $\sigma$  weight enhances the model's effectiveness in mitigating the impact of false negative samples. It is important to acknowledge that a higher  $\sigma$  does not necessarily indicate better performance. As shown in Figure 6 (a), an increase in  $\sigma$  at the initial stage contributes to enhancing the model's performance. Nevertheless, as the weight of  $\sigma$  increases, the performance of backbones generally declines, resulting in the model adhering too strictly to the "established guidelines". Consequently, it impacts the efficacy of learning from the hard negative samples. We further use the density plots to visualize the prediction on the STS-Benchmark development set in Figure 7. These models are trained on the synthesized data without filtering. We can observe that in Figure 7 (a), the distribution of prediction results for labels  $\geq 4$  is significantly shifted to the left. Compared with the results in Figure 7 (b), this issue is effectively alleviated, demonstrating the effectiveness of the Gaussian-decayed function in reducing the influence of false negative samples. To further verify the applicability of the Gaussian-decayed function, we applied it to SynCSE and verified the performance in Appendix E.

## 5 CONCLUSION

In this paper, we propose a pipeline-based data augmentation method using LLM to enhance data diversity in sentence representation learning. By leveraging knowledge of entities and quantities, our approach improves the model's ability to capture fine-grained semantic distinctions. The Gaussian-decayed function in our GCSE model further reduces noise in the generated data. Extensive experiments on STS and reranking tasks show that our method achieves state-of-the-art results with fewer synthesized samples and a more lightweight LLM, demonstrating its effectiveness and efficiency.

## REFERENCES

- 540  
541  
542 Eneko Agirre, Daniel Cer, Mona Diab, and Aitor Gonzalez-Agirre. SemEval-2012 task 6: A pilot on  
543 semantic textual similarity. In Eneko Agirre, Johan Bos, Mona Diab, Suresh Manandhar, Yuval  
544 Marton, and Deniz Yuret (eds.), *\*SEM 2012: The First Joint Conference on Lexical and Com-*  
545 *putational Semantics – Volume 1: Proceedings of the main conference and the shared task, and*  
546 *Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation (SemEval*  
547 *2012)*, pp. 385–393, Montréal, Canada, 7-8 June 2012. Association for Computational Linguis-  
548 tics. URL <https://aclanthology.org/S12-1051>.
- 549 Eneko Agirre, Daniel Cer, Mona Diab, Aitor Gonzalez-Agirre, and Weiwei Guo. \*SEM 2013  
550 shared task: Semantic textual similarity. In Mona Diab, Tim Baldwin, and Marco Baroni  
551 (eds.), *Second Joint Conference on Lexical and Computational Semantics (\*SEM), Volume 1:*  
552 *Proceedings of the Main Conference and the Shared Task: Semantic Textual Similarity*, pp.  
553 32–43, Atlanta, Georgia, USA, June 2013. Association for Computational Linguistics. URL  
554 <https://aclanthology.org/S13-1004>.
- 555 Eneko Agirre, Carmen Banea, Claire Cardie, Daniel Cer, Mona Diab, Aitor Gonzalez-Agirre, Wei-  
556 wei Guo, Rada Mihalcea, German Rigau, and Janyce Wiebe. SemEval-2014 task 10: Multi-  
557 lingual semantic textual similarity. In Preslav Nakov and Torsten Zesch (eds.), *Proceedings of*  
558 *the 8th International Workshop on Semantic Evaluation (SemEval 2014)*, pp. 81–91, Dublin, Ire-  
559 land, August 2014. Association for Computational Linguistics. doi: 10.3115/v1/S14-2010. URL  
560 <https://aclanthology.org/S14-2010>.
- 561 Eneko Agirre, Carmen Banea, Claire Cardie, Daniel Cer, Mona Diab, Aitor Gonzalez-Agirre,  
562 Weiwei Guo, Iñigo Lopez-Gazpio, Montse Maritxalar, Rada Mihalcea, German Rigau, Larraitz  
563 Uria, and Janyce Wiebe. SemEval-2015 task 2: Semantic textual similarity, English, Spanish  
564 and pilot on interpretability. In Preslav Nakov, Torsten Zesch, Daniel Cer, and David Jurgens  
565 (eds.), *Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015)*,  
566 pp. 252–263, Denver, Colorado, June 2015. Association for Computational Linguistics. doi:  
567 10.18653/v1/S15-2045. URL <https://aclanthology.org/S15-2045>.
- 568 Eneko Agirre, Carmen Banea, Daniel Cer, Mona Diab, Aitor Gonzalez-Agirre, Rada Mihal-  
569 cea, German Rigau, and Janyce Wiebe. SemEval-2016 task 1: Semantic textual similarity,  
570 monolingual and cross-lingual evaluation. In Steven Bethard, Marine Carpuat, Daniel Cer,  
571 David Jurgens, Preslav Nakov, and Torsten Zesch (eds.), *Proceedings of the 10th Interna-*  
572 *tional Workshop on Semantic Evaluation (SemEval-2016)*, pp. 497–511, San Diego, Califor-  
573 nia, June 2016. Association for Computational Linguistics. doi: 10.18653/v1/S16-1081. URL  
574 <https://aclanthology.org/S16-1081>.
- 575 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge,  
576 Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu,  
577 Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng  
578 Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Sheng-  
579 guang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao,  
580 Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang,  
581 Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. Qwen tech-  
582 nical report. *CoRR*, abs/2309.16609:1–59, 2023. doi: 10.48550/ARXIV.2309.16609. URL  
583 <https://doi.org/10.48550/arXiv.2309.16609>.
- 584 Daniel M. Cer, Mona T. Diab, Eneko Agirre, Iñigo Lopez-Gazpio, and Lucia Specia. Semeval-2017  
585 task 1: Semantic textual similarity - multilingual and cross-lingual focused evaluation. *CoRR*,  
586 abs/1708.00055, 2017. URL <http://arxiv.org/abs/1708.00055>.
- 587 Yiming Chen, Yan Zhang, Bin Wang, Zuozhu Liu, and Haizhou Li. Generate, discriminate and  
588 contrast: A semi-supervised sentence representation learning framework. In Yoav Goldberg,  
589 Zornitsa Kozareva, and Yue Zhang (eds.), *Proceedings of the 2022 Conference on Empirical*  
590 *Methods in Natural Language Processing*, pp. 8150–8161, Abu Dhabi, United Arab Emirates,  
591 December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.  
592 558. URL <https://aclanthology.org/2022.emnlp-main.558>.

- 594 Qinyuan Cheng, Xiaogui Yang, Tianxiang Sun, Linyang Li, and Xipeng Qiu. Improving contrastive  
595 learning of sentence embeddings from AI feedback. In Anna Rogers, Jordan Boyd-Graber,  
596 and Naoaki Okazaki (eds.), *Findings of the Association for Computational Linguistics: ACL*  
597 *2023*, pp. 11122–11138, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-acl.707. URL <https://aclanthology.org/2023.findings-acl.707>.
- 600 Yung-Sung Chuang, Rumen Dangovski, Hongyin Luo, Yang Zhang, Shiyu Chang, Marin Soljacic,  
601 Shang-Wen Li, Scott Yih, Yoon Kim, and James Glass. DiffCSE: Difference-based contrastive  
602 learning for sentence embeddings. In Marine Carpuat, Marie-Catherine de Marneffe,  
603 and Ivan Vladimir Meza Ruiz (eds.), *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 4207–4218, Seattle, United States, July 2022a. Association for Computational Linguistics. doi: 10.18653/v1/2022.naacl-main.311. URL <https://aclanthology.org/2022.naacl-main.311>.
- 608 Yung-Sung Chuang, Rumen Dangovski, Hongyin Luo, Yang Zhang, Shiyu Chang, Marin Soljacic,  
609 Shang-Wen Li, Scott Yih, Yoon Kim, and James R. Glass. Diffcse: Difference-based contrastive  
610 learning for sentence embeddings. In Marine Carpuat, Marie-Catherine de Marneffe, and Iván  
611 Vladimir Meza Ruíz (eds.), *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*, pp. 4207–4218. Association for Computational  
612 Linguistics, 2022b. doi: 10.18653/v1/2022.NAACL-MAIN.311. URL <https://doi.org/10.18653/v1/2022.naacl-main.311>.
- 616 Arman Cohan, Sergey Feldman, Iz Beltagy, Doug Downey, and Daniel Weld. SPECTER:  
617 Document-level representation learning using citation-informed transformers. In Dan Jurafsky,  
618 Joyce Chai, Natalie Schluter, and Joel Tetreault (eds.), *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 2270–2282, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.207. URL <https://aclanthology.org/2020.acl-main.207>.
- 622 Alexis Conneau and Douwe Kiela. Senteval: An evaluation toolkit for universal sentence representations. In Nicoletta Calzolari, Khalid Choukri, Christopher Cieri, Thierry Declerck, Sara Goggi, Kôiti Hasida, Hitoshi Isahara, Bente Maegaard, Joseph Mariani, H  l  ne Mazo, Asunci  n Moreno, Jan Odi  k, Stelios Piperidis, and Takenobu Tokunaga (eds.), *Proceedings of the Eleventh International Conference on Language Resources and Evaluation, LREC 2018, Miyazaki, Japan, May 7-12, 2018*. European Language Resources Association (ELRA), 2018. URL <http://www.lrec-conf.org/proceedings/lrec2018/summaries/757.html>.
- 629 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of  
630 deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and  
631 Tamar Solorio (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics*, pp. 4171–4186. Association for Computational  
632 Linguistics, June 2019. doi: 10.18653/v1/N19-1423. URL <https://aclanthology.org/N19-1423>.
- 635 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha  
636 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony  
637 Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark,  
638 Arun Rao, Aston Zhang, Aur  lien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Rozi  re,  
639 Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris  
640 Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong,  
641 Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny  
642 Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino,  
643 Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael  
644 Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson,  
645 Graeme Nail, Gr  goire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar,  
646 Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel M. Kloumann, Ishan  
647 Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy

- 648 Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak,  
649 Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Al-  
650 wala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, and et al. The  
651 llama 3 herd of models. *CoRR*, abs/2407.21783, 2024. doi: 10.48550/ARXIV.2407.21783. URL  
652 <https://doi.org/10.48550/arXiv.2407.21783>.
- 653 Tianyu Gao, Xingcheng Yao, and Danqi Chen. SimCSE: Simple contrastive learning of sen-  
654 tence embeddings. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott  
655 Wen-tau Yih (eds.), *Proceedings of the 2021 Conference on Empirical Methods in Natural  
656 Language Processing*, pp. 6894–6910. Association for Computational Linguistics, November  
657 2021. doi: 10.18653/v1/2021.emnlp-main.552. URL [https://aclanthology.org/  
658 2021.emnlp-main.552](https://aclanthology.org/2021.emnlp-main.552).
- 660 John Giorgi, Osvald Nitski, Bo Wang, and Gary Bader. DeCLUTR: Deep contrastive learning for un-  
661 supervised textual representations. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli  
662 (eds.), *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics  
663 and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long  
664 Papers)*, pp. 879–895, Online, August 2021. Association for Computational Linguistics. doi: 10.  
665 18653/v1/2021.acl-long.72. URL <https://aclanthology.org/2021.acl-long.72>.
- 666 Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Diego Rojas, Guanyu  
667 Feng, Hanlin Zhao, Hanyu Lai, Hao Yu, Hongning Wang, Jiadai Sun, Jiajie Zhang, Jiale Cheng,  
668 Jiayi Gui, Jie Tang, Jing Zhang, Juanzi Li, Lei Zhao, Lindong Wu, Lucen Zhong, Mingdao Liu,  
669 Minlie Huang, Peng Zhang, Qinkai Zheng, Rui Lu, Shuaiqi Duan, Shudan Zhang, Shulin Cao,  
670 Shuxun Yang, Weng Lam Tam, Wenyi Zhao, Xiao Liu, Xiao Xia, Xiaohan Zhang, Xiaotao Gu,  
671 Xin Lv, Xinghan Liu, Xinyi Liu, Xinyue Yang, Xixuan Song, Xunkai Zhang, Yifan An, Yifan  
672 Xu, Yilin Niu, Yuantao Yang, Yueyan Li, Yushi Bai, Yuxiao Dong, Zehan Qi, Zhaoyu Wang,  
673 Zhen Yang, Zhengxiao Du, Zhenyu Hou, and Zihan Wang. Chatglm: A family of large language  
674 models from glm-130b to glm-4 all tools, 2024.
- 675 Felix Hill, Kyunghyun Cho, and Anna Korhonen. Learning distributed representations of sen-  
676 tences from unlabelled data. In Kevin Knight, Ani Nenkova, and Owen Rambow (eds.), *Pro-  
677 ceedings of the 2016 Conference of the North American Chapter of the Association for Com-  
678 putational Linguistics: Human Language Technologies*, pp. 1367–1377, San Diego, California,  
679 June 2016. Association for Computational Linguistics. doi: 10.18653/v1/N16-1162. URL  
680 <https://aclanthology.org/N16-1162>.
- 681 Mingqing Hu and Bing Liu. Mining and summarizing customer reviews. In Won Kim, Ron Ko-  
682 havi, Johannes Gehrke, and William DuMouchel (eds.), *Proceedings of the Tenth ACM SIGKDD  
683 International Conference on Knowledge Discovery and Data Mining, Seattle, Washington, USA,  
684 August 22-25, 2004*, pp. 168–177. ACM, 2004. doi: 10.1145/1014052.1014073. URL <https://doi.org/10.1145/1014052.1014073>.
- 686 Junjie Huang, Duyu Tang, Wanjun Zhong, Shuai Lu, Linjun Shou, Ming Gong, Daxin Jiang, and  
687 Nan Duan. Whiteningbert: An easy unsupervised sentence embedding approach. In Marie-  
688 Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (eds.), *Findings of the As-  
689 sociation for Computational Linguistics: EMNLP 2021, Virtual Event / Punta Cana, Dominican  
690 Republic, 16-20 November, 2021*, pp. 238–244. Association for Computational Linguistics, 2021.  
691 doi: 10.18653/V1/2021.FINDINGS-EMNLP.23. URL [https://doi.org/10.18653/v1/  
692 2021.findings-emnlp.23](https://doi.org/10.18653/v1/2021.findings-emnlp.23).
- 693 Yongxin Huang, Kexin Wang, Sourav Dutta, Raj Patel, Goran Glavaš, and Iryna Gurevych. AdaSent:  
694 Efficient domain-adapted sentence embeddings for few-shot classification. In Houda Bouamor,  
695 Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Meth-  
696 ods in Natural Language Processing*, pp. 3420–3434, Singapore, December 2023. Associa-  
697 tion for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.208. URL <https://aclanthology.org/2023.emnlp-main.208>.
- 700 Ting Jiang, Jian Jiao, Shaohan Huang, Zihan Zhang, Deqing Wang, Fuzhen Zhuang, Furu Wei,  
701 Haizhen Huang, Denvy Deng, and Qi Zhang. PromptBERT: Improving BERT sentence embed-  
dings with prompts. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Proceedings of*

- 702 *the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 8826–8837,  
703 Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.603. URL <https://aclanthology.org/2022.emnlp-main.603>.  
704  
705  
706
- 707 Taeuk Kim, Kang Min Yoo, and Sang-goo Lee. Self-guided contrastive learning for BERT sentence  
708 representations. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (eds.), *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pp. 2528–2540, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.197. URL <https://aclanthology.org/2021.acl-long.197>.  
709  
710  
711  
712
- 713 Ryan Kiros, Yukun Zhu, Russ R Salakhutdinov, Richard Zemel, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. Skip-thought vectors. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015. URL [https://proceedings.neurips.cc/paper\\_files/paper/2015/file/f442d33fa06832082290ad8544a8da27-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2015/file/f442d33fa06832082290ad8544a8da27-Paper.pdf).  
714  
715  
716  
717
- 718 Tao Lei, Hrishikesh Joshi, Regina Barzilay, Tommi Jaakkola, Kateryna Tymoshenko, Alessandro Moschitti, and Lluís Màrquez. Semi-supervised question retrieval with gated convolutions. In Kevin Knight, Ani Nenkova, and Owen Rambow (eds.), *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 1279–1289, San Diego, California, June 2016. Association for Computational Linguistics. doi: 10.18653/v1/N16-1153. URL <https://aclanthology.org/N16-1153>.  
719  
720  
721  
722  
723  
724  
725
- 726 Bohan Li, Hao Zhou, Junxian He, Mingxuan Wang, Yiming Yang, and Lei Li. On the sentence embeddings from pre-trained language models. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu (eds.), *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 9119–9130, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.733. URL <https://aclanthology.org/2020.emnlp-main.733>.  
727  
728  
729  
730  
731
- 732 Xianming Li and Jing Li. AoE: Angle-optimized embeddings for semantic textual similarity. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1825–1839, Bangkok, Thailand, August 2024a. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.101. URL <https://aclanthology.org/2024.acl-long.101>.  
733  
734  
735  
736  
737
- 738 Xianming Li and Jing Li. BeLLM: Backward dependency enhanced large language model for sentence embeddings. In Kevin Duh, Helena Gomez, and Steven Bethard (eds.), *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 792–804, Mexico City, Mexico, June 2024b. Association for Computational Linguistics. doi: 10.18653/v1/2024.naacl-long.45. URL <https://aclanthology.org/2024.naacl-long.45>.  
739  
740  
741  
742  
743
- 744 Jiduan Liu, Jiahao Liu, Qifan Wang, Jingang Wang, Wei Wu, Yunsen Xian, Dongyan Zhao, Kai Chen, and Rui Yan. Rankcse: Unsupervised sentence representations learning via learning to rank. In Anna Rogers, Jordan L. Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*, pp. 13785–13802. Association for Computational Linguistics, 2023. doi: 10.18653/V1/2023.ACL-LONG.771. URL <https://doi.org/10.18653/v1/2023.acl-long.771>.  
745  
746  
747  
748  
749  
750
- 751 Xueqing Liu, Chi Wang, Yue Leng, and ChengXiang Zhai. Linkso: a dataset for learning to retrieve similar question answer pairs on software development forums. In Yijun Yu, Erik M. Fredericks, and Premkumar T. Devanbu (eds.), *Proceedings of the 4th ACM SIGSOFT International Workshop on NLP for Software Engineering, NLASE@ESEC/SIGSOFT FSE 2018, Lake Buena Vista, FL, USA, November 4, 2018*, pp. 2–5. ACM, 2018. doi: 10.1145/3283812.3283815. URL <https://doi.org/10.1145/3283812.3283815>.  
752  
753  
754  
755

- 756 Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike  
757 Lewis, Luke Zettlemoyer, and Veselin Stoyanov. RoBERTa: A robustly optimized BERT pretrain-  
758 ing approach. *CoRR*, abs/1907.11692:1–13, 2019. URL <http://arxiv.org/abs/1907.11692>.  
759 11692.
- 760 Lajanugen Logeswaran and Honglak Lee. An efficient framework for learning sentence repre-  
761 sentations. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=rJvJXZb0W>.  
762 763
- 764 Marco Marelli, Stefano Menini, Marco Baroni, Luisa Bentivogli, Raffaella Bernardi, and Roberto  
765 Zamparelli. A SICK cure for the evaluation of compositional distributional semantic models. In  
766 Nicoletta Calzolari, Khalid Choukri, Thierry Declerck, Hrafn Loftsson, Bente Maegaard, Joseph  
767 Mariani, Asunción Moreno, Jan Odijk, and Stelios Piperidis (eds.), *Proceedings of the Ninth In-*  
768 *ternational Conference on Language Resources and Evaluation, LREC 2014, Reykjavik, Iceland,*  
769 *May 26-31, 2014*, pp. 216–223. European Language Resources Association (ELRA), 2014. URL  
770 <http://www.lrec-conf.org/proceedings/lrec2014/summaries/363.html>.
- 771 Pu Miao, Zeyao Du, and Junlin Zhang. Debcase: Rethinking unsupervised contrastive sentence  
772 embedding learning in the debiasing perspective. In Ingo Frommholz, Frank Hopfgartner, Mark  
773 Lee, Michael Oakes, Mounia Lalmas, Min Zhang, and Rodrygo L. T. Santos (eds.), *Proceedings*  
774 *of the 32nd ACM International Conference on Information and Knowledge Management, CIKM*  
775 *2023, Birmingham, United Kingdom, October 21-25, 2023*, pp. 1847–1856. ACM, 2023. doi:  
776 10.1145/3583780.3614833. URL <https://doi.org/10.1145/3583780.3614833>.
- 777 Tomas Mikolov, Ilya Sutskever, Kai Chen, Gregory S. Corrado, and Jeffrey Dean. Distributed  
778 representations of words and phrases and their compositionality. In Christopher J. C. Burges, Leon  
779 Bottou, Zoubin Ghahramani, and Kilian Q. Weinberger (eds.), *Advances in Neural Information*  
780 *Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems 2013.*  
781 *Proceedings of a meeting held December 5-8, 2013, Lake Tahoe, Nevada, United States*, pp.  
782 3111–3119, 2013. URL <https://proceedings.neurips.cc/paper/2013/hash/9aa42b31882ec039965f3c4923ce901b-Abstract.html>.  
783
- 784 Niklas Muennighoff, Nouamane Tazi, Loic Magne, and Nils Reimers. MTEB: Massive text em-  
785 bedding benchmark. In Andreas Vlachos and Isabelle Augenstein (eds.), *Proceedings of the 17th*  
786 *Conference of the European Chapter of the Association for Computational Linguistics*, pp. 2014–  
787 2037, Dubrovnik, Croatia, May 2023. Association for Computational Linguistics. doi: 10.18653/  
788 v1/2023.eacl-main.148. URL <https://aclanthology.org/2023.eacl-main.148>.
- 789 Jianmo Ni, Gustavo Hernandez Abrego, Noah Constant, Ji Ma, Keith B. Hall, Daniel Cer, and  
790 Yinfei Yang. Sentence-t5: Scalable sentence encoders from pre-trained text-to-text models. In  
791 Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.), *Findings of the Association*  
792 *for Computational Linguistics: ACL 2022, Dublin, Ireland, May 22-27, 2022*, pp. 1864–1874.  
793 Association for Computational Linguistics, 2022. doi: 10.18653/V1/2022.FINDINGS-ACL.146.  
794 URL <https://doi.org/10.18653/v1/2022.findings-acl.146>.
- 795 OpenAI. Chatgpt: Optimizing language models for dialogue, 2022. URL <https://openai.com/blog/chatgpt/>. Accessed: 2024-11-19.  
796 797
- 798 OpenAI. GPT-4 technical report. *CoRR*, abs/2303.08774:1–100, 2023. doi: 10.48550/ARXIV.2303.  
799 08774. URL <https://doi.org/10.48550/arXiv.2303.08774>.
- 800 Matteo Pagliardini, Prakhar Gupta, and Martin Jaggi. Unsupervised learning of sentence embed-  
801 dings using compositional n-gram features. In Marilyn Walker, Heng Ji, and Amanda Stent  
802 (eds.), *Proceedings of the 2018 Conference of the North American Chapter of the Association*  
803 *for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pp.  
804 528–540, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi:  
805 10.18653/v1/N18-1049. URL <https://aclanthology.org/N18-1049>.
- 806 Bo Pang and Lillian Lee. A sentimental education: sentiment analysis using subjectivity summa-  
807 rization based on minimum cuts. In *Proceedings of the 42nd Annual Meeting on Association for*  
808 *Computational Linguistics, ACL ’04*, pp. 271–es, USA, 2004. Association for Computational Lin-  
809 guistics. doi: 10.3115/1218955.1218990. URL <https://doi.org/10.3115/1218955.1218990>.

- 810 Bo Pang and Lillian Lee. Seeing stars: exploiting class relationships for sentiment categorization  
811 with respect to rating scales. In *Proceedings of the 43rd Annual Meeting on Association for Com-*  
812 *putational Linguistics*, ACL '05, pp. 115–124, USA, 2005. Association for Computational Lin-  
813 guistics. doi: 10.3115/1219840.1219855. URL [https://doi.org/10.3115/1219840.](https://doi.org/10.3115/1219840.1219855)  
814 1219855.
- 815 Nina Poerner and Hinrich Schütze. Multi-view domain adapted sentence embeddings for low-  
816 resource unsupervised duplicate question detection. In Kentaro Inui, Jing Jiang, Vincent Ng, and  
817 Xiaojun Wan (eds.), *Proceedings of the 2019 Conference on Empirical Methods in Natural Lan-*  
818 *guage Processing and the 9th International Joint Conference on Natural Language Processing*  
819 *(EMNLP-IJCNLP)*, pp. 1630–1641, Hong Kong, China, November 2019. Association for Com-  
820 putational Linguistics. doi: 10.18653/v1/D19-1173. URL [https://aclanthology.org/](https://aclanthology.org/D19-1173)  
821 D19-1173.
- 822 Nils Reimers and Iryna Gurevych. Sentence-BERT: Sentence embeddings using siamese bert-  
823 networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language*  
824 *Processing and the 9th International Joint Conference on Natural Language Processing*, pp.  
825 3973–3983, 2019.
- 826 Jaydeep Sen, Chuan Lei, Abdul Quamar, Fatma Özcan, Vasilis Efthymiou, Ayushi Dalmia, Greg  
827 Stager, Ashish R. Mittal, Diptikalyan Saha, and Karthik Sankaranarayanan. ATHENA++: natural  
828 language querying for complex nested SQL queries. *Proc. VLDB Endow.*, 13(11):2747–2759,  
829 2020. URL <http://www.vldb.org/pvldb/vol13/p2747-sen.pdf>.
- 830 Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and  
831 Christopher Potts. Recursive deep models for semantic compositionality over a sentiment tree-  
832 bank. In David Yarowsky, Timothy Baldwin, Anna Korhonen, Karen Livescu, and Steven Bethard  
833 (eds.), *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Process-*  
834 *ing*, pp. 1631–1642, Seattle, Washington, USA, October 2013. Association for Computational  
835 Linguistics. URL <https://aclanthology.org/D13-1170>.
- 836 Jacob Mitchell Springer, Suhas Kotha, Daniel Fried, Graham Neubig, and Aditi Raghunathan. Rep-  
837 etition improves language model embeddings. *CoRR*, abs/2402.15449, 2024. doi: 10.48550/  
838 ARXIV.2402.15449. URL <https://doi.org/10.48550/arXiv.2402.15449>.
- 839 Jianlin Su, Jiarun Cao, Weijie Liu, and Yangyiwen Ou. Whitening sentence representations for  
840 better semantics and faster retrieval. *CoRR*, abs/2103.15316, 2021. URL [https://arxiv.](https://arxiv.org/abs/2103.15316)  
841 [org/abs/2103.15316](https://arxiv.org/abs/2103.15316).
- 842 Qwen Team. Qwen2.5: A party of foundation models, September 2024. URL [https://qwenlm.](https://qwenlm.github.io/blog/qwen2.5/)  
843 [github.io/blog/qwen2.5/](https://qwenlm.github.io/blog/qwen2.5/).
- 844 Nandan Thakur, Nils Reimers, Andreas Rücklé, Abhishek Srivastava, and Iryna Gurevych. BEIR: A  
845 heterogeneous benchmark for zero-shot evaluation of information retrieval models. In *Thirty-fifth*  
846 *Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round*  
847 *2)*, 2021. URL <https://openreview.net/forum?id=wCu6T5xFjeJ>.
- 848 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
849 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Ar-  
850 mand Joulin, Edouard Grave, and Guillaume Lample. LLaMA: Open and efficient foundation  
851 language models. *CoRR*, abs/2302.13971:1–27, 2023. doi: 10.48550/ARXIV.2302.13971. URL  
852 <https://doi.org/10.48550/arXiv.2302.13971>.
- 853 Ellen M. Voorhees and Dawn M. Tice. Building a question answering test collection. In *Pro-*  
854 *ceedings of the 23rd Annual International ACM SIGIR Conference on Research and Develop-*  
855 *ment in Information Retrieval*, SIGIR '00, pp. 200–207, New York, NY, USA, 2000. Asso-  
856 ciation for Computing Machinery. ISBN 1581132263. doi: 10.1145/345508.345577. URL  
857 <https://doi.org/10.1145/345508.345577>.
- 858 Bin Wang, C.-C. Jay Kuo, and Haizhou Li. Just rank: Rethinking evaluation with word and sentence  
859 similarities. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.), *Proceedings*  
860 *of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long*  
861 *and Short Papers)*, pp. 100–109, 2022. Association for Computational Linguistics. doi: 10.18653/v1/W22-1001. URL <https://doi.org/10.18653/v1/W22-1001>.



- 864 *Papers*), pp. 6060–6077, Dublin, Ireland, May 2022a. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.419. URL <https://aclanthology.org/2022.acl-long.419>.
- 865  
866  
867
- 868 Hao Wang and Yong Dou. Sncse: Contrastive learning for unsupervised sentence embedding  
869 with soft negative samples. In *International Conference on Intelligent Computing*, pp. 419–431.  
870 Springer, 2023.
- 871 Huiming Wang, Zhaodonghui Li, Liying Cheng, De Wen Soh, and Lidong Bing. Large lan-  
872 guage models can contrastively refine their generation for better sentence representation learn-  
873 ing. In Kevin Duh, Helena Gómez-Adorno, and Steven Bethard (eds.), *Proceedings of the 2024*  
874 *Conference of the North American Chapter of the Association for Computational Linguistics:*  
875 *Human Language Technologies (Volume 1: Long Papers), NAACL 2024, Mexico City, Mex-*  
876 *ico, June 16-21, 2024*, pp. 7874–7891. Association for Computational Linguistics, 2024a. doi:  
877 10.18653/v1/2024.NAACL-LONG.436. URL <https://doi.org/10.18653/v1/2024.naacl-long.436>.
- 878
- 879 Kexin Wang, Nils Reimers, and Iryna Gurevych. TSDAE: Using transformer-based sequen-  
880 tial denoising auto-encoder for unsupervised sentence embedding learning. In Marie-Francine  
881 Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (eds.), *Findings of the Associ-*  
882 *ation for Computational Linguistics: EMNLP 2021*, pp. 671–688, Punta Cana, Dominican Re-  
883 public, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.  
884 findings-emnlp.59. URL <https://aclanthology.org/2021.findings-emnlp.59>.
- 885
- 886 Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang, Rangan Majumder, and Furu Wei. Im-  
887 proving text embeddings with large language models. In Lun-Wei Ku, Andre Martins, and  
888 Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Com-*  
889 *putational Linguistics (Volume 1: Long Papers)*, pp. 11897–11916, Bangkok, Thailand, August  
890 2024b. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.642. URL  
891 <https://aclanthology.org/2024.acl-long.642>.
- 892
- 893 Wei Wang, Liangzhu Ge, Jingqiao Zhang, and Cheng Yang. Improving contrastive learning of  
894 sentence embeddings with case-augmented positives and retrieved negatives. In *Proceedings of*  
895 *the 45th International ACM SIGIR Conference on Research and Development in Information*  
896 *Retrieval, SIGIR ’22*, pp. 2159–2165, New York, NY, USA, 2022b. Association for Computing  
897 Machinery. ISBN 9781450387323. doi: 10.1145/3477495.3531823. URL <https://doi.org/10.1145/3477495.3531823>.
- 898
- 899 Janyce Wiebe, Theresa Wilson, and Claire Cardie. Annotating expressions of opinions and  
900 emotions in language. *Lang. Resour. Evaluation*, 39(2-3):165–210, 2005. doi: 10.1007/  
901 S10579-005-7880-9. URL <https://doi.org/10.1007/s10579-005-7880-9>.
- 902
- 903 Bohong Wu and Hai Zhao. Sentence representation learning with generative objective rather than  
904 contrastive objective. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Proceed-*  
905 *ings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 3356–  
906 3368, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Lin-  
907 guistics. doi: 10.18653/v1/2022.emnlp-main.221. URL <https://aclanthology.org/2022.emnlp-main.221>.
- 908
- 909 Fangzhao Wu, Ying Qiao, Jiun-Hung Chen, Chuhan Wu, Tao Qi, Jianxun Lian, Danyang Liu, Xing  
910 Xie, Jianfeng Gao, Winnie Wu, and Ming Zhou. MIND: A large-scale dataset for news recom-  
911 mendation. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (eds.), *Proceedings*  
912 *of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 3597–3606, On-  
913 line, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.331.  
914 URL <https://aclanthology.org/2020.acl-main.331>.
- 915
- 916 Qiyu Wu, Chongyang Tao, Tao Shen, Can Xu, Xiubo Geng, and Daxin Jiang. PCL: Peer-contrastive  
917 learning with diverse augmentations for unsupervised sentence embeddings. In Yoav Goldberg,  
918 Zornitsa Kozareva, and Yue Zhang (eds.), *Proceedings of the 2022 Conference on Empirical*  
919 *Methods in Natural Language Processing*, pp. 12052–12066, Abu Dhabi, United Arab Emi-  
920 rates, December 2022a. Association for Computational Linguistics. doi: 10.18653/v1/2022.  
921 emnlp-main.826. URL <https://aclanthology.org/2022.emnlp-main.826>.

- 918 Xing Wu, Chaochen Gao, Yipeng Su, Jizhong Han, Zhongyuan Wang, and Songlin Hu. Smoothed  
919 contrastive learning for unsupervised sentence embedding. In Nicoletta Calzolari, Chu-Ren  
920 Huang, Hansaem Kim, James Pustejovsky, Leo Wanner, Key-Sun Choi, Pum-Mo Ryu, Hsin-  
921 Hsi Chen, Lucia Donatelli, Heng Ji, Sadao Kurohashi, Patrizia Paggio, Nianwen Xue, Seokhwan  
922 Kim, Younggyun Hahm, Zhong He, Tony Kyungil Lee, Enrico Santus, Francis Bond, and Seung-  
923 Hoon Na (eds.), *Proceedings of the 29th International Conference on Computational Linguistics*,  
924 pp. 4902–4906, Gyeongju, Republic of Korea, October 2022b. International Committee on Com-  
925 putational Linguistics. URL <https://aclanthology.org/2022.coling-1.434>.
- 926 Xing Wu, Chaochen Gao, Liangjun Zang, Jizhong Han, Zhongyuan Wang, and Songlin Hu. ES-  
927 imCSE: Enhanced sample building method for contrastive learning of unsupervised sentence em-  
928 bedding. In Nicoletta Calzolari, Chu-Ren Huang, Hansaem Kim, James Pustejovsky, Leo Wan-  
929 ner, Key-Sun Choi, Pum-Mo Ryu, Hsin-Hsi Chen, Lucia Donatelli, Heng Ji, Sadao Kurohashi,  
930 Patrizia Paggio, Nianwen Xue, Seokhwan Kim, Younggyun Hahm, Zhong He, Tony Kyungil  
931 Lee, Enrico Santus, Francis Bond, and Seung-Hoon Na (eds.), *Proceedings of the 29th In-  
932 ternational Conference on Computational Linguistics*, pp. 3898–3907, Gyeongju, Republic of  
933 Korea, October 2022c. International Committee on Computational Linguistics. URL <https://aclanthology.org/2022.coling-1.342>.
- 934  
935 Bo Xu, Shouang Wei, Luyi Cheng, Shizhou Huang, Hui Song, Ming Du, and Hongya Wang. Hsim-  
936 cse: Improving contrastive learning of unsupervised sentence representation with adversarial hard  
937 positives and dual hard negatives. In *2023 International Joint Conference on Neural Networks  
938 (IJCNN)*, pp. 1–8, 2023. doi: 10.1109/IJCNN54540.2023.10191335.
- 939  
940 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li,  
941 Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang,  
942 Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai,  
943 Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng  
944 Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai  
945 Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan  
946 Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang  
947 Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zhihao Fan. Qwen2  
948 technical report. *arXiv preprint arXiv:2407.10671*, 2024.
- 949 Bowen Zhang, Kehua Chang, and Chunping Li. Simple techniques for enhancing sentence em-  
950 beddings in generative language models. In De-Shuang Huang, Zhanjun Si, and Qinhu Zhang  
951 (eds.), *Advanced Intelligent Computing Technology and Applications*, pp. 52–64, Singapore,  
952 2024. Springer Nature Singapore. ISBN 978-981-97-5669-8.
- 953  
954 Junlei Zhang, Zhenzhong Lan, and Junxian He. Contrastive learning of sentence embeddings  
955 from scratch. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023  
956 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore,  
957 December 6-10, 2023*, pp. 3916–3932. Association for Computational Linguistics, 2023. doi:  
958 10.18653/v1/2023.EMNLP-MAIN.238. URL [https://doi.org/10.18653/v1/2023.  
959 emnlp-main.238](https://doi.org/10.18653/v1/2023.emnlp-main.238).
- 960 Yan Zhang, Ruidan He, Zuozhu Liu, Kwan Hui Lim, and Lidong Bing. An unsupervised sen-  
961 tence embedding method by mutual information maximization. In Bonnie Webber, Trevor Cohn,  
962 Yulan He, and Yang Liu (eds.), *Proceedings of the 2020 Conference on Empirical Methods  
963 in Natural Language Processing (EMNLP)*, pp. 1601–1610, Online, November 2020. Associ-  
964 ation for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.124. URL <https://aclanthology.org/2020.emnlp-main.124>.
- 965  
966 Yuhao Zhang, Hongji Zhu, Yongliang Wang, Nan Xu, Xiaobo Li, and Binqiang Zhao. A contrastive  
967 framework for learning sentence representations from pairwise and triple-wise perspective in an-  
968 gular space. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.), *Proceedings  
969 of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long  
970 Papers)*, pp. 4892–4903, Dublin, Ireland, May 2022. Association for Computational Linguis-  
971 tics. doi: 10.18653/v1/2022.acl-long.336. URL [https://aclanthology.org/2022.  
acl-long.336](https://aclanthology.org/2022.acl-long.336).

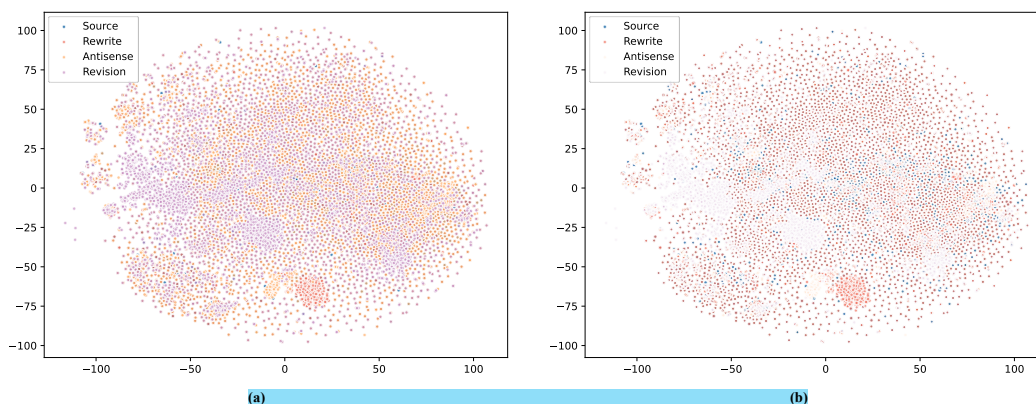
972 Kun Zhou, Beichen Zhang, Xin Zhao, and Ji-Rong Wen. Debaised contrastive learning of unsuper-  
 973 vised sentence representations. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio  
 974 (eds.), *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics*  
 975 (*Volume 1: Long Papers*), pp. 6120–6130, Dublin, Ireland, May 2022. Association for Compu-  
 976 tational Linguistics. doi: 10.18653/v1/2022.acl-long.423. URL <https://aclanthology.org/2022.acl-long.423>.

## 978 APPENDIX

### 981 A DATA SYNTHESIS PROMPTS

982 In this section, we provide the specifics of our prompts for knowledge extraction and integration,  
 983 and data synthesis. The particular prompts are presented in Table 6.

### 984 B VISUALIZATION OF SYNTHETIC SAMPLE DISTRIBUTION



989 Figure 8: t-SNE visualization of the synthetic sample generated by ChatGLM3-6B, where the trans-  
 990 parency of “Antisense” and “Revision” samples in subgraph (b) is reduced to 10% for better obser-  
 991 vation.

992 In this section, we use the supervised SimCSE model to generate sentence embeddings for the syn-  
 993 thesized samples and utilize t-SNE to project the vectors into two-dimensional space for a visual  
 994 analysis of the diversity. To facilitate observation, we group the synthesized samples into three cate-  
 995 gories: “Rewrite” refers to positive samples synthesized using “Rewriting Prompt 1” and “Rewriting  
 996 Prompt 2” from Table 6, while “Antisense” denotes the negative samples generated using “Syntac-  
 997 tic Antisense Prompt”. “Revision” denotes the negative samples generated using “Entity Revision  
 998 Prompt”, “Quantity Revision Prompt” and “Rewriting Prompt 3”, which are related to knowledge  
 999 modification. And “Source” indicates the original samples from the dataset. We randomly selected  
 1000 5k “Source” samples and corresponding synthetic samples from our dataset for visualization, and  
 1001 the results are illustrated in Figure 8. We observe that “Rewrite” samples basically cover the spa-  
 1002 tial distribution of “Source” samples while expanding into the neighborhood space to some extent.  
 1003 “Antisense” and “Revision” samples further enhance the information density within the target se-  
 1004 mantic space. Comparing Figure 8 (a) and (b), it can be observed that the “Revision” samples cover  
 1005 areas with sparse information, while their overall spatial distribution remains consistent with the  
 1006 semantic distribution of “Source” samples. This indicates that the sample synthesis with knowledge  
 1007 effectively increases sample diversity within the semantic space.

### 1008 C PERFORMANCE ON TRANSFER TASKS

1009 We also evaluate our GCSE following the same settings as SimCSE on seven transfer tasks: MR  
 1010 (Pang & Lee, 2005), CR (Hu & Liu, 2004), SUBJ (Pang & Lee, 2004), MPQA (Wiebe et al., 2005),  
 1011

1026	<b>Knowledge Extraction Prompt</b>	<b>Instruction:</b> Predicts the subject categories, contained entities, and quantified information of the following text <b>Rules:</b> The category is an item in $\{\{categories\_name\}, \dots\}$ , quantified information refers to information contained in the text with numerical values or units, such as '2GB', 'three cups', 'two dogs', etc Output format: json format data, the data format is: { cls: [], // category entities: $\{\{text: \text{""}, type: \text{""}\}\}$ , // entities, 'text' must be subsequences in the Input text quantities: $\{\{text: \text{""}, type: \text{""}, quantity: 0\}\}$ // To quantify the information, 'text' must be a subsequence in the Input text } <b>Input:</b> $\{x\}$
1035	<b>Rewriting Prompt 1</b>	<b>Instruction:</b> You are an excellent storyteller; rewrite the input sentence in a different way. Please try to recreate the sentence using different expressions, including varied tones, synonyms, and sentence patterns, while ensuring that the new sentence has the same meaning as the original sentence. <b>Input:</b> $\{x\}$
1039	<b>Rewriting Prompt 2</b>	<b>Instruction:</b> You are a great storyteller; I would be grateful if you could employ your creativity to devise an illustration of the preceding segment of the sentence. The preceding statement must not exceed $\{number\}$ words, and it follows the original text. <b>Input:</b> $\{x\}$
1043	<b>Rewriting Prompt 3</b>	<b>Instruction:</b> You are a great rewriter, and I want you to generate new sentence according to the classification, entities and quantities info provided by the json. <b>Rules:</b> You should aware that the new text in "quantities" should be rewrite follows the "quantity" value. e.g. "text": "A man", "quantity": 5 should rewrite as "five men". <b>Metadata:</b> { "cls": " $\{categories\_name\}$ ", "entities": [{  "text": " $\{entity\_text\}$ ", "type": " $\{entity\_type\}$ " }], "quantities": [{ "text": " $\{entity\_text\}$ ", "quantity": $\{entity\_quantity\}$ }], ... } <b>Input:</b> $\{x\}$
1056	<b>Syntactic Antisense Prompt</b>	<b>Instruction:</b> You are dishonest; you ought to reformulate the input sentence so that the NLI model perceives it as an opposing sample. <b>Rules:</b> 1. If the statement asserts negation, you should affirm; conversely, if the statement asserts affirmation, you should negate. 2. If an individual loves something, one should assert that it does not reciprocate that affection. 3. If an individual is engaged in one activity, state that they are performing a different activity. 4. If the statement is affirmative/negative, express it as negative/affirmative. <b>Input:</b> $\{x\}$
1063	<b>Entity Revision Prompt</b>	<b>Instruction:</b> You are a great story teller, rewrites the input sentence, and change the entity ' $\{original\_entity\_text\}$ ' to another $\{entity\_type\}$ ' $\{new\_entity\_text\}$ '. <b>Input:</b> $\{x\}$
1065	<b>Quantity Revision Prompt</b>	<b>Instruction:</b> You are a great story teller, rewrites the input sentence, and change the quantity $\{original\_quantity\_value\}$ of ' $\{original\_quantity\_text\}$ ' to $\{random\_quantity\_value\}$ . <b>Input:</b> $\{x\}$

Table 6: Examples of data synthesis prompts, where  $\{variable\ name\}$  refers to a variable.

1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076 SST2 (Socher et al., 2013), TREC (Voorhees & Tice, 2000), and MRPC (Voorhees & Tice, 2000).  
1077 The results are shown in Table 7, it can be observed that our GCSE (ChatGPT) achieves the best per-  
1078 formance on all backbone models, outperforming second-best methods in average scores of 0.89%  
1079 with BERT-base, 0.79% with BERT-large, 0.44% with RoBERTa-base, and 0.40% with RoBERTa-  
large, demonstrating the potential capability in downstream tasks.

Model	Method	MR	CR	SUBJ	MPQA	SST2	TREC	MRPC	Avg.
BERT-base	SimCSE♠	68.40	82.41	74.38	80.91	78.56	76.85	72.23	76.25
	DiffCSE♠	72.28	84.43	76.47	83.90	80.54	80.59	71.23	78.49
	PCL♠	72.84	83.81	76.52	83.06	79.32	80.01	73.38	78.42
	RankCSE♠	75.66	86.27	77.81	84.74	81.10	81.80	75.13	80.36
	MultiCSR (ChatGPT)♣	82.70	88.15	94.97	90.08	86.87	87.70	75.46	86.56
	SynCSE (ChatGPT)*	83.34	88.80	93.88	90.39	88.96	83.60	75.94	86.42
	<b>GCSE (ChatGLM3-6B)</b>	82.22	88.43	94.59	90.09	86.88	<b>89.40</b>	76.06	<b>86.81</b>
	<b>GCSE (GLM4-9B)</b>	<b>84.63</b>	89.78	<b>95.01</b>	<b>90.54</b>	88.96	86.00	76.12	87.29
<b>GCSE (ChatGPT)</b>	84.59	<b>90.15</b>	94.97	90.39	<b>89.68</b>	86.00	<b>76.35</b>	<b>87.45</b>	
BERT-large	SimCSE♠	70.88	84.16	76.43	84.50	79.76	79.26	73.88	78.41
	PCL♠	74.87	86.11	78.29	85.65	80.52	81.62	73.94	80.14
	RankCSE♠	75.48	86.50	78.60	85.45	81.09	81.58	75.53	80.60
	SynCSE (ChatGPT)*	85.78	90.47	94.77	90.41	90.50	89.00	<b>75.77</b>	88.10
	<b>GCSE (ChatGLM3-6B)</b>	83.97	89.38	95.13	90.22	89.57	90.60	75.71	87.80
	<b>GCSE (GLM4-9B)</b>	<b>86.01</b>	<b>90.94</b>	<b>95.40</b>	90.24	<b>92.15</b>	<b>92.00</b>	75.48	<b>88.89</b>
<b>GCSE (ChatGPT)</b>	85.93	90.44	94.94	<b>90.52</b>	92.04	88.80	75.25	88.27	
RoBERTa-base	SimCSE♠	70.16	81.77	73.24	81.36	80.65	80.22	68.56	76.57
	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♠	73.20	85.95	77.17	84.82	82.58	83.08	71.88	79.81
	MultiCSR (ChatGPT)♣	84.70	90.69	94.40	89.38	89.42	<b>89.62</b>	<b>77.01</b>	87.89
	SynCSE (ChatGPT)††	85.47	91.44	92.53	89.67	90.94	81.60	76.06	86.82
	<b>GCSE (ChatGLM3-6B)</b>	84.39	90.81	94.02	88.90	91.05	89.40	76.12	87.81
<b>GCSE (GLM4-9B)</b>	<b>86.49</b>	<b>92.24</b>	<b>94.70</b>	89.63	92.37	86.60	76.29	<b>88.33</b>	
<b>GCSE (ChatGPT)</b>	86.32	91.58	94.37	<b>90.04</b>	<b>92.42</b>	84.00	76.12	87.84	
RoBERTa-large	SimCSE♠	72.86	83.99	75.62	84.77	81.80	81.98	71.26	78.90
	PCL♠	74.08	84.36	76.42	85.49	81.76	82.79	71.51	79.49
	RankCSE♠	73.20	85.83	78.00	85.63	82.67	84.19	73.64	80.45
	SynCSE (ChatGPT)††	87.24	<b>92.16</b>	93.75	<b>90.81</b>	91.87	84.00	<b>76.29</b>	88.02
	<b>GCSE (ChatGLM3-6B)</b>	85.65	90.78	94.16	90.08	90.44	<b>92.80</b>	73.74	88.24
	<b>GCSE (GLM4-9B)</b>	87.45	91.60	<b>94.62</b>	90.30	<b>92.42</b>	88.40	71.77	88.08
<b>GCSE (ChatGPT)</b>	<b>87.56</b>	91.76	94.56	90.69	92.26	88.80	74.84	<b>88.64</b>	

Table 7: Comparison of different sentence embedding models accuracy on transfer tasks. “♠”: results from Liu et al. (2023), “♣”: results from Wang et al. (2024a), “††”: results from Zhang et al. (2023). “\*”: we reproduce the results with the officially released corpus from Zhang et al. (2023).

Premise	Hypothesis	Gold	SimCSE	RankCSE	SynCSE	GCSE
A woman is cooking <b>eggs</b> .	A woman is cooking <b>something</b> .	3.00	4.37 (1.372)	4.23 (1.320)	<u>3.66 (0.662)</u>	<b>3.24 (0.236)</b>
<b>Two</b> little girls are talking on the phone.	<b>A</b> little girl is walking down the street.	0.50	3.38 (2.881)	3.64 (3.139)	<u>1.97 (1.468)</u>	<b>1.85 (1.351)</b>
A chef is preparing <b>some</b> food.	<b>A</b> chef prepared <b>a</b> meal.	4.00	<b>4.27 (0.270)</b>	4.59 (0.588)	4.56 (0.561)	<u>4.41 (0.408)</u>
<b>Five</b> kittens are eating out of <b>five</b> dishes.	Kittens are eating <b>food</b> on trays.	2.75	3.81 (1.056)	3.71 (0.957)	<u>3.28 (0.535)</u>	<b>3.12 (0.373)</b>
A woman is cutting <b>some</b> herbs.	A woman is chopping <b>cilantro</b> .	2.80	3.58 (0.777)	3.58 (0.967)	<u>3.11 (0.313)</u>	<b>2.61 (0.185)</b>

Table 8: Case studies on model prediction similarity with gold labels in the STS-Benchmark development set, where Gold represents the label score of the sentence pair (ranging from zero to five). The similarity scores of all models are multiplied by a coefficient of five for better comparison, and the value in parentheses denotes the RMS error between the predicted score and the label. Words highlighted in blue denote the entity alteration in the sentence-pair, whereas words in yellow indicate the quantities that change inside the sentence-pair.

## D CASE STUDIES

To further verify the improvement in our method’s awareness of entity and quantity, we selected five sample sets from the STS-Benchmark development set that explicitly contained alterations in entity or quantity within the sentence-pair, and presented the prediction cosine-similarity scores of GCSE and related methodologies with the backbone of BERT-base in Table 8. We can observe from the results that the prediction score of our model achieves the minimum root-mean-square error compared to the label in most cases, which indicates that our model has a stronger capacity to distinguish information.

## E ABLATION STUDIES OF GAUSSIAN-DECAYED AND FEW-SHOT SAMPLES

Method	STS-12	STS-13	STS-14	STS-15	STS-16	STS-B	SICK-R	Avg.
SynCSE (ChatGPT)*	75.86	82.19	78.71	85.63	<b>81.11</b>	82.35	<b>78.79</b>	80.66
w sampled	75.48	85.60	78.76	84.78	80.38	82.12	76.46	80.51
w sampled & G.D.	75.71	85.24	79.09	85.15	80.82	82.68	77.54	80.89
w G.D.	<b>75.89</b>	85.26	79.24	85.67	80.79	82.63	78.19	<b>81.10</b>
w sampled & domain & G.D.	75.88	<b>86.02</b>	<b>79.46</b>	<b>86.10</b>	80.27	<b>82.87</b>	76.91	81.07

Table 9: Ablation studies of sample size and the Gaussian-decayed function by utilizing SynCSE. “\*”: we reproduce the results with the officially released corpus from Zhang et al. (2023).

We employ the Gaussian-decayed function on SynCSE and sample SynCSE training data with a sample size the same as our synthetic data to evaluate the efficacy of the proposed Gaussian-decayed function and our domain-oriented selection strategy in the ablation experiment. The data sample size is 64k, and the weight of  $\sigma$  in  $G(\cdot)$  is assigned the same value as specified in Section 4.1. The results of various policies implemented in SynCSE are presented in Table 9. “w sampled” denotes the utilization of purely the sampled data in SynCSE, and a performance decrease can be observed when training on a reduced number of samples without extra configurations. “w sampled & G.D.” denotes the additional incorporation of  $G(\cdot)$  based on “w sampled”. “w G.D.” indicates the results by training on the full dataset utilizing  $G(\cdot)$ . In both configurations, the average performance outperforms the vanilla model, illustrating the module’s efficacy. “w sampled & domain & G.D.” denotes the concurrent utilization of sample data, domain data, and  $G(\cdot)$ , with a sample size of 48k for the SynCSE dataset and 16k for the synthesized domain dataset. The results reveal that “w sampled & domain & G.D.” attains the second-best performance, suggesting that incorporating domain data can decrease the required training samples while enhancing model efficacy.

## F UNSUPERVISED SENTENCE EMBEDDING ON LLM

Model	Avg.	Model	Avg.
<i>Unsupervised</i>		<i>Data Augmentation</i>	
Llama3.2-3B LoRA	71.34	Llama3.2-3B LoRA	78.26
Llama-3-8B LoRA	<b>72.73</b>	Llama-3-8B LoRA	78.24
ChatGLM3-6B LoRA	69.38	ChatGLM3-6B LoRA	79.04
GLM4-9B LoRA	71.77	GLM4-9B LoRA	<b>79.52</b>
Qwen2.5-14B LoRA	68.49	Qwen2.5-14B LoRA	78.02

Table 10: Performance comparison of different LLMs on STS tasks, where results of “Unsupervised” refers to models trained on the same unsupervised settings as Gao et al. (2021), and “Data Augmentation” refers to models trained with the synthetic data generated by ChatGLM3-6B.

In this section, we utilize contrastive learning on multiple LLMs to evaluate the alignment of LLM-generated similarities with the gold labels and the effectiveness of our data augmentation strategy. We use Llama3.2-3B (Dubey et al., 2024), Llama3-8B (Dubey et al., 2024), ChatGLM3-6B (GLM et al., 2024), GLM4-9B (GLM et al., 2024) and Qwen2.5-14B (Team, 2024; Yang et al., 2024) with a low-rank adapter (LoRA) layer for training. The sentence embedding vectors are obtained from the output hidden states of the last position, which is followed by the method of pretended chain of thought (Pretended CoT) (Zhang et al., 2024). We may derive two major conclusion from the results in Table 10: (1) In conventional unsupervised settings, decoder-based LLMs have no significant performance advantage over encoder-based PLMs for sentence representation learning tasks. The model performance does not increase significantly with the increase of the number of model parameters. To reduce expenses, we assert that fully leveraging the capabilities of LLMs for distilling smaller models is the better option. (2) The application of our data augmentation technique to sentence representation learning tasks in LLMs significantly enhances performance relative to the “Unsupervised” settings, which further proves the applicability and efficacy of our strategy.

## G VISUALIZATION OF PREDICTION SCORES AND GRADIENT COMPARISONS

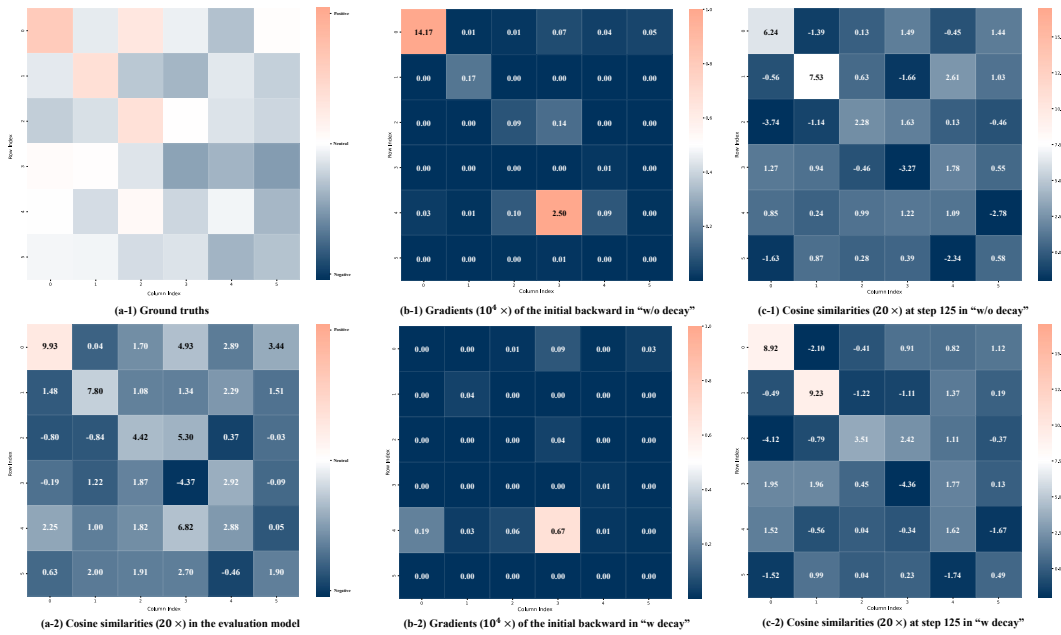


Figure 9: Heatmap visualization of the prediction scores and gradients.

To further analyze the effectiveness of the Gaussian-decayed function in mitigating the impact of false negative noise, we visualized the changes in predicted scores and gradients during the training process using heatmaps. In the training procedure of GCSE, each input consists of a source sample, its corresponding positive sample, and a hard negative sample. We visualize the cosine similarity scores and gradient heatmaps for negative samples within a batch in Figure 9. Each cell of a heatmap represents the relationship between the source sample and the negative sample, and the diagonal cells highlight the relationships between source samples and their hard negatives. Since synthetic samples lack manual annotations, we use supervised SimCSE models (Gao et al., 2021) based on different backbones to compute their similarity scores as the ground truth. We normalized the output scores of each model with min-max scaling and averaged them as the final scores to address distributional differences across models, and the results are shown in Figure 9 (a-1). It can be observed that several hard negatives on the diagonal display scores biased towards positive similarity, indicating the presence of false negative noise. In the framework of contrastive learning, when optimized using standard contrastive loss, these hard negatives are positioned further from the source samples in the semantic space, negatively impacting the model’s representational capacity. Figure 9 (a-2) displays the normalized cosine similarity scores of hard negatives in the initial step as calculated by the evaluation model in GCSE. The initial score distribution of hard negatives shows a strong correlation with the ground truth, suggesting that these scores could efficiently guide GCSE in gradient correction.

Figures 9 (b-1) and (b-2) present the backward gradient values of the model trained without and with the Gaussian-decayed function, respectively. For better visualization, all gradient values are amplified by  $10^4$ , and all similarities are amplified by 20 by the temperature. By comparing the gradients of hard negative samples in these two figures, it can be observed that the gradient values on false hard negatives are significantly smaller when the Gaussian-decayed function is applied. Additionally, Figures 9 (c-1) and (c-2) present a comparison of cosine similarity scores after 125 training steps with and without the Gaussian-decayed function. The scores for false hard negatives are significantly higher when the Gaussian-decayed function is employed, while the true hard negatives had lower scores. The overall score distribution aligns more accurately with the ground truth, and these results demonstrate that the Gaussian-decayed function effectively prevents false negatives from being pushed farther away from source samples in the semantic space, thereby validating its effectiveness in mitigating noise and improving model performance.

## H ABLATION ANALYSIS OF FILTERING THRESHOLDS

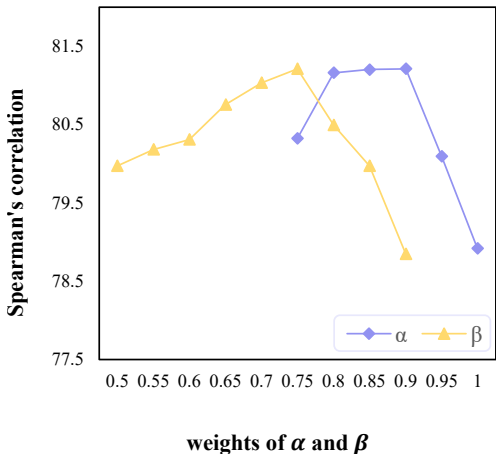


Figure 10: Spearman’s correlation against the weight of  $\alpha$  and  $\beta$  on the STS tasks. When adjusting the weight of one parameter, the other parameter is fixed at its default value as specified in the experimental settings.

To study the impact of different filtering thresholds, we evaluate the performance on the backbone of the BERT-base, and the results are shown in Figure 10. When  $\alpha > 0.9$ , the model’s performance declines significantly, primarily because the high threshold filters out too many samples, heavily reducing the number of positive samples. In the range  $\alpha \in [0.8, 0.9]$ , performance degradation is observed due to noise introduced by false positive samples. Similarly, when  $\alpha < 0.8$ , the model suffers from a performance drop caused by an excessive number of false positives being included in the training process. The threshold for  $\beta$  demonstrates a noticeable impact on model performance when it deviates from 0.75. Specifically, when  $\beta > 0.75$ , the model’s performance declines significantly due to the inclusion of excessive false negative noise, which severely affects the model performance. Conversely, when  $\beta < 0.75$ , the selected negative samples become easier for the model to distinguish, providing limited benefit for enhancing its representation learning capacity. The results highlight the influence of filtering thresholds on sample quality and distribution.

## I SCORE NORMALIZATION METHODOLOGY

In this work, the labels in datasets are normalized with standard min-max normalization. To address the discrepancy in score distributions among different models, we applied a variant min-max normalization method to align their predicted scores. For each label  $l \in [0, \text{MAX}]$ , we collect all predicted scores with  $l = 0$  as list  $C_0$ , and all predicted scores with  $l = \text{MAX}$  as list  $C_1$ . Specifically, we computed the median prediction scores for  $C_0$  and  $C_1$  as  $\text{min}_p = \text{median}(C_0)$  and  $\text{max}_p = \text{median}(C_1)$ , respectively. The use of medians, rather than the minimum predicted score for  $C_0$  or the maximum predicted score for  $C_1$ , avoids reliance on outlier values that may disproportionately skew the normalization, ensuring a more balanced score distribution. For a given score  $s$ , the normalized score  $s'$  is calculated as:

$$s' = \text{clip} \left( \frac{s - \text{min}_p}{\text{max}_p - \text{min}_p}, 0, 1 \right), \tag{14}$$

where the function  $\text{clip}(x, 0, 1)$  ensures the normalized score is bounded within  $[0, 1]$ . This method adjusts the score range to maintain consistency across models while preserving relative score differences.