

# CLL-RetICL: Contrastive Linguistic Label Retrieval-based In-Context Learning for Text Classification via Large Language Models

Anonymous ACL submission

## Abstract

Recent research has delved into Retrieval-based In-Context Learning (RetICL), leveraging the power of large language models (LLMs) for text classification. Despite its promise, a persistent challenge lies in effectively retrieving relevant demonstrations from a support set. Many existing approaches have overlooked the essential role of linguistic label information in guiding this retrieval process. To bridge this gap, we present Contrastive Linguistic Label Retrieval-based In-Context Learning (CLL-RetICL), a novel framework designed to identify the most relevant and impactful sentences without altering the model parameters. Our approach uniquely integrates sentence-query similarity with sentence-label similarity, enabling a more nuanced and comprehensive evaluation of relevance. We tested CLL-RetICL across diverse text classification tasks and evaluated its performance on various LLMs. Experimental results demonstrate that CLL-RetICL consistently outperforms previous retrieval methods that do not incorporate linguistic label information. These findings highlight the critical importance of linguistic label-aware selection in enhancing text classification accuracy.<sup>1</sup>

## 1 Introduction

A linguistic label represents the semantics of a category and plays a vital role in text classification tasks. Human annotators rely on the meaning conveyed by these labels to accurately categorize text. Depending on the specific requirements of a custom classification task, a linguistic label can often be substituted with synonyms or more descriptive phrases to better align with the task's context.

Recently, researchers have begun exploring few-shot in-context learning (ICL) using LLMs for text classification tasks. (Luo et al., 2024; Yu et al., 2023; Chae and Davidson, 2023; Rouzegar and

<sup>1</sup>Our code is available: <http://acl-org.github.io/ACL-PUB/formatting.html>

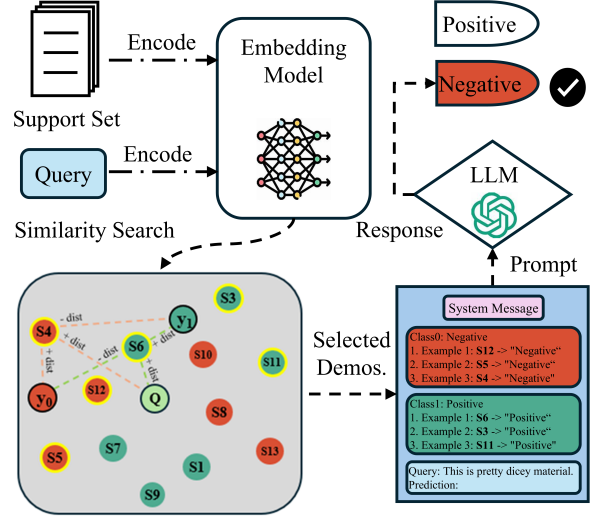


Figure 1: An illustration of CLL-RetICL with  $N = 2$  and  $k = 3$ , demonstrating a prediction between Positive and Negative classes. Here,  $y_0$  and  $y_1$  represent the vector representations of the linguistic labels "Negative" and "Positive", respectively, in a pre-trained sentence embedding model. Similarly,  $s_0, s_1, \dots$  represent the vector representations of the sentences in a support set within the same pre-trained sentence embedding model.

Makrehchi, 2024). Instead of selecting static, pre-defined demonstration sets for ICL, RetICL adopts a dynamic, context-sensitive approach. At its core, adaptive demonstration selection leverages a specialized retriever to intelligently curate tailored demonstrations for each task input. RetICL has gained popularity because prior research suggests that context-insensitive demonstrations can limit the full potential of LLMs (Luo et al., 2024; Wu et al., 2022). Despite RetICL consistently surpassing approaches based on random or static demonstrations, it still remains an open challenge to retrieve relevant demonstrations.

To address the problem, previous researchers have proposed various strategies, including  $k$ -nearest neighbors (KNN), NwayKshot, and clustering-based RetICL (Li et al., 2024; Pecher et al., 2024; Zhang et al., 2022a). However, these

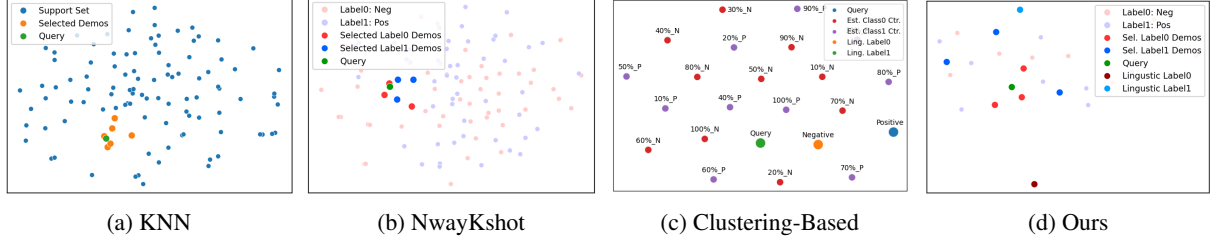


Figure 2: A comparison of four different approaches to RetICL strategies. (a) KNN suffers from two key weaknesses: the copying effect and misleading by similarity. (b) NwayKshot always ignores any linguistic cues conveyed through the labels. (c) Clustering-based approaches are hindered by the difficulty in estimating category centers and the neglect of query similarity. (d) Our method avoids the copying effect, prevents misleading similarity, incorporates linguistic label information, utilizes fixed label category centers, and integrates query similarity.

methods suffer from various challenges, as shown in Figure 2. To identify the most effective demonstrations, we analyzed failure cases. Our investigation revealed that there always exists a specific combination of demonstrations that enables LLMs to classify accurately. Additionally, our analysis uncovered that failure cases are error-prone: they often lie closer to the representation of an opposing linguistic label or near the representation of an incorrect label cluster center, despite their similarity to the query. In contrast, when the demonstrations are correctly combined, they align more closely with the representation of the intended linguistic label. A detailed discussion of these findings is presented in Section 3.

Building on these observations, we present a novel RetICL framework, CLL-RetICL (Contrastive Linguistic Label Retrieval-based In-Context Learning) as illustrated in Figure 1. Our approach introduces a trade-off method that computes a relevance score by integrating both sentence-query and sentence-label similarities, thereby effectively leveraging label information. Furthermore, to optimize the effectiveness of CLL-RetICL, we developed a universal  $N$ -way  $K$ -shot prompt structure applicable to all text classification tasks. This prompt design mitigates the copying effect and prevents LLMs from being misled by overly similar examples. Moreover, we demonstrate that the sentence embeddings of linguistic labels can serve as clustering centers—generated by a pre-trained sentence embedding model—to address the challenge of estimating clustering centers. Additionally, we initiate four variations for integrating the linguistic label style into RetICL and evaluate their effectiveness on four text classification datasets. Finally, to assess the generalizability of CLL-RetICL, we conduct experiments using Gemini (Team et al., 2024), Llama (Dubey et al.,

2024), and Mistral (Jiang et al., 2024). Empirical experiments show that CLL-RetICL consistently outperforms both previous RetICL baselines and other variants across multiple datasets and LLMs. Ablation studies further reveal several key findings: (1) Effectiveness across variations: CLL-RetICL maintains strong performance across different  $k$ -shot settings, various pre-trained sentence embedding models, and multiple similarity functions. (2) Component dependency: The proposed method relies on the original component responsible for calculating sentence-query similarity; omitting this component degrades performance. (3) Impact of hyperparameters: Trade-off hyperparameters have a minor influence on the final classification accuracy. The following summarizes our main contributions:

- We present a novel perspective in which sentence embeddings of linguistic labels serve as highly accurate clustering centers, free from the biases introduced by limited support data and independent of data-driven constraints.
- We propose an innovative method, CLL-RetICL, which employs a rigorous relevance scoring metric that leverages linguistic label information to select high-quality demonstrations for improving LLMs in text classification tasks. Our approach does not require fine-tuning the pre-trained weights of either the sentence embedding models or LLMs.
- We conduct extensive experiments to evaluate the proposed method, achieving better performance on most datasets compared to existing RetICL methods.

## 2 Related Work

**Text Classification via LLMs.** Text classification via LLMs has recently demonstrated excep-

tional generalizability and reasoning capabilities, attracting significant research interest in their application to text classification tasks (Zhang et al., 2024; Wang et al., 2024; Fields et al., 2024). Existing methods can be broadly divided into two groups, depending on whether they involve adapting the parameters of LLMs or not. The first group concentrates on fine-tuning the parameters of LLMs to excel in custom text classification tasks (Chae and Davidson, 2023; Zhang et al., 2024; Yu et al., 2023; Jin et al., 2023). However, this approach generally demands significant computational resources to load the full LLM model parameters, and fine-tuning these models can often diminish their generalizability. The other category is known as ICL, or prompt engineering (Guo et al., 2024; Luo et al., 2024; Fan et al., 2024). While this method avoids the need to update LLM model parameters, it heavily depends on well-designed prompts, making it challenging to guide LLMs to consistently meet human expectations (Shi et al., 2023; Mavromatis et al., 2023; Edwards and Camacho-Collados, 2024).

**RetICL.** RetICL can generally be divided into two categories: approaches that retrain or fine-tune a retriever for specific text classification tasks, and approaches that utilize pre-trained language models without additional fine-tuning. An intuitive strategy for RetICL involves directly selecting a few similar sentences, leveraging readily available demonstration retrievers like those based on sentence embeddings. Existing methods include KATE (Liu et al., 2021), Z-ICL (Lyu et al., 2022) and ICL-ML (Milios et al., 2023). However, recent research has shown that selecting the most similar demonstrations can lead to the copying effect and misleading by similarity, degrading performance in text classification tasks (Olsson et al., 2022; Zhang et al., 2022b). To mitigate the issue of homogeneity in retrieval, clustering retrieval approaches ensure the selection of a diverse and representative set of demonstrations, which is critical to its effectiveness (Luo et al., 2024). Several methods exist, including NwayKshot (Li et al., 2024), Votek (Su et al., 2022) and D-CALM (Hassan and Alikhani, 2023). While these approaches leverage label information and offer improvements, accurately estimating the clustering center for each category remains challenging. This difficulty arises because clustering center estimation is a data-driven process that depends on a support set.

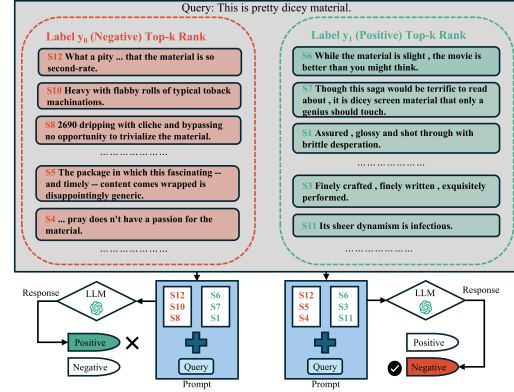


Figure 3: A comparison of the correct and incorrect demonstration combinations is presented. On the left, NwayKshot retrieves the top- $k$  sentences most similar to the query from each group; however, this approach fails to classify the query correctly. In contrast, on the right, CLL-RetICL does not rely solely on proximity to the query, resulting in an accurate classification.

The second category of RetICL involves fine-tuning or retraining a retriever model to rank relevant sentences using either in-domain or out-of-domain datasets for text classification tasks. There are established methods, such as PEFT (Tunstall et al., 2022), UDR (Li et al., 2023) and Ambig-ICL (Gao et al., 2023). These methods utilize label information and feedback to optimize model parameters, highlighting the essential role of labeled data in yielding valuable insights for text classification tasks. However, they often demand substantial computational resources and considerable time to construct a retriever.

### 3 Linguistic Label Retrieval Hypothesis

Previous studies have shown that retrieving sentences closest to the query and applying a clustering-based selection method can enhance the diversity of demonstrations while mitigating the risk of misleading results due to similarity (Li et al., 2020; Luo et al., 2024). Therefore, a question arises: are the clustering centers reliable? To explore this further, we analyze the distribution of clustering centers, as shown in Appendix C. Varying the proportion of fully supported data from 10% to 100% reveals that the distribution of clustering centers shifts according to the number of sentences in the support set. Notably, negative-labeled clustering centers tend to be less distinct within a certain range compared to positive-labeled ones. These findings suggest that clustering center estimation is inherently data-driven and prone to bias, making it difficult to accurately identify true clustering

centers. On the other hand, by analyzing failure cases, we find that, for a given query, there is an optimal combination of demonstrations that can effectively guide LLMs to classify the query correctly. However, relying solely on the top-ranked closest demonstrations retrieved does not always yield accurate results. An example of this limitation is illustrated in Figure 3. To further investigate, we compared cases where the top- $k$  closest demonstrations led to incorrect results versus cases where randomly selected demonstrations produced correct outcomes. We provide five examples of such instances in Appendix C. We found that incorrect nearest-neighbor demonstrations exhibit an error-prone tendency, being either closer to the linguistic representation of an opposite label, closer to the center of an incorrect label cluster, or both—despite being similar to the query. Conversely, in correct combinations, the selected demonstrations exhibit a stronger alignment with the correct tendency. For example, sentences with a Negative label tend to show higher similarity to the linguistic word "Negative" and the same holds for "Positive" label. Although correct demonstrations align closely with their respective cluster centers, we observe exceptions where a correct output contains sentences that are nearer to the center of an incorrect label cluster. Furthermore, even sentences closest to their correct cluster centers can still lead to classification errors due to inaccurate estimation of those centers.

Based on these observations, we hypothesize that the vector representations of linguistic labels should be explicitly incorporated into the retrieval process rather than relying on cluster center estimation. Compared to traditional clustering center estimation, this approach offers two advantages: (1) Independence from data Bias – The linguistic label clustering center is not data-driven, preventing bias introduced by the support set. (2) Leveraging linguistic information – Linguistic labels play a crucial role in zero-shot ICL, as LLMs rely entirely on these labels for text classification tasks.

## 4 Our Method: CLL-RetICL

**Preliminary.** Let the query set  $Q$  represent a task, where  $q \in Q$  denotes a sample query for which we aim to find an answer via an LLM. In the context of RetICL, multiple demonstrations ( $d_1, \dots, d_k$ ) are retrieved from a support set  $C$ . Each demonstration  $d_i$  consists of a sentence and its label,  $(s_i, y_i) \in C$ , where  $y_i$  belongs to the label set  $Y$ .

**Overview.** We present CLL-RetICL, a novel RetICL approach leveraging information extraction between demonstrations and linguistic labels to predict the correct label for a given query input  $q_i$  (Wang et al., 2023). Unlike earlier methods (Liu et al., 2021; Su et al., 2022; Li et al., 2022; Milios et al., 2023) that create input-label pairs by retrieving sentences closest to a given query, CLL-RetICL selects demonstrations that balance a trade-off by augmenting the corresponding label while penalizing others.

CLL-RetICL involves three key steps, as illustrated in Figure 1: (1) Retrieving more relevant sentences by integrating sentence-query similarity with sentence-label similarity (detailed in Section 4.1), (2) Forming demonstrations by organizing the retrieved demonstrations into an N-way K-shot format (discussed in Section 4.2), and (3) Making inferences through ICL (explained in Section 4.3).

### 4.1 Linguistic Label Retriever

RetICL employs a retrieval mechanism to identify  $k$  examples from  $C$  that are most relevant to a given query  $q$ . This process is guided by a similarity function,  $sim$ , which quantifies the relationship between a sentence  $s_i$  and a query  $q$ . The corresponding formula is as follows:

$$score_{RetICL} = sim(q, s_i) \quad (1)$$

To build on this hypothesis, CLL-RetICL incorporates sentence-query similarity with sentence-label similarity. Rather than solely considering the similarity distance between a sentence  $s_i$  and the query  $q$ , CLL-RetICL employs the following formula:

$$score_{c-RetICL} = sim(q, s_i) + w_1 * \log \frac{\exp^{sim(s_i, y_i)}}{\frac{1}{n-1} \sum_{y \in Y, y \neq y_i} \exp^{sim(s_i, y)}} \quad (2)$$

where  $w_1$  is a trade-off hyperparameter that balances the relative importance of the corresponding terms in the objective function.

CLL-RetICL considers the relationship between sentences and linguistic labels by utilizing a similarity function. It increases the score based on the similarity between a sentence and its assigned correct label (referred to as the positive label) while decreasing the score based on the similarity between the sentence and other labels (referred to as



negative labels). Additionally, we propose several variations and evaluate their performance through experiments. These include Positive Label Augment (PLA), Negative Label Penalty (NLP), and Contrastive Label (CTL). The corresponding formulas are provided below:

$$score_{PLA} = sim(q, x_i) + w_1 * sim(x_i, y_i) \quad (3)$$

$$score_{NLP} = sim(q, x_i) - w_1 * \frac{1}{n-1} \sum_{y \in Y, y \neq y_i} sim(x_i, y) \quad (4)$$

$$score_{CTL} = sim(q, x_i) + w_1 * sim(x_i, y_i) - w_2 * \frac{1}{n-1} \sum_{y \in Y, y \neq y_i} sim(x_i, y) \quad (5)$$

where  $w_1$  and  $w_2$  are trade-off hyperparameters.

Our methods ensure that the selected sentences (1) maintain a safe distance from  $q$  to prevent the copying effect (Olsson et al., 2022; Zhang et al., 2022b), (2) incorporate the information between sentences and linguistic labels and (3) align closely with the requirements of the custom text classification task.

## 4.2 *N*-way *K*-shot

We adopt a clustering-based retrieval method, as prior research suggests that *N*-way *K*-shot effectively addresses the issue of homogeneity (Li and Qiu, 2023). Here, we partition all sentences into  $N$  sub-groups, aiming to cluster sentences that share the same label. Our retriever selects top  $K$  high demonstrations according to above score formula from each sub-group, resulting in a final set of  $N \times K$  demonstrations.

## 4.3 Inference

Finally, CLL-RetICL constructs a prompt by concatenating *N*-way *K*-shot input-label pairs  $(x_1, y_1), (x_2, y_2), \dots, (x_k, y_k)$  for each *N*-way label, along with the query input  $q$ . This prompt is then fed into a LLM, which generates a prediction using  $argmax_{y \in Y} P(y|prompt)$ . The universal prompt template for each text classification task is outlined in Table 5 in Appendix B.

# 5 Experimental Analysis

## 5.1 Experimental Setup

We evaluate multiple LLMs to identify factors affecting classification accuracy across four tasks. Key results are summarized in the main text, with additional details presented in the Appendix D.

### 5.1.1 Datasets

We conduct experiments on four widely recognized text classification tasks: SST2 (Socher et al., 2013), CoLA (Warstadt et al., 2018), CARER (Saravia et al., 2018) and BBCnews (Greene and Cunningham, 2006). Similar to conventional text classification methodologies, we treat the training sets as support sets and the test sets as query sets, while disregarding development sets if they exist. The detailed data statistics are provided in Appendix A and summarized in Table 3.

### 5.1.2 Baselines

We compare CLL-RetICL with the zero-shot approach as well as various RetICL methods.

**Zero-shot** predicts  $argmax_{y \in Y} P(y|q)$  without using any demonstrations (Radford et al., 2019; Brown et al., 2020). This method utilizes LLMs and linguistic label information to enhance text classification.

**Z-ICL** leverages physical neighbors to avoid selecting demonstrations that are overly similar to the query. Furthermore, it introduces the use of synonymous labels to mitigate the copying effect, highlighting the potential for effectively utilizing the linguistic meaning of labels (Lyu et al., 2022).

**KATE** employs a standard KNN approach to retrieve demonstrations, which remains the most widely used method in RetICL (Liu et al., 2021).

**NwayKshot** is a clustering-based retrieval method designed to tackle the challenge of homogeneity in demonstrations (Li et al., 2024).

**Cluster-TopN** builds on NwayKshot but applies  $k$ -means clustering to identify the cluster centers. It then selects the demonstration closest to the center from each sub-group (Zhdanov, 2019; Hassan and Alikhani, 2023).

**Votek** selects  $k$  representatives from  $N$  sub-groups through a voting mechanism to best represent the group (Su et al., 2022).

### 5.1.3 Experimental Details

**LLMs.** We conduct experiments using three LLMs: Gemini (Team et al., 2024), Llama (Dubey et al., 2024) and Mistral (Jiang et al., 2024). Specifically, we utilize fixed versions of these models, namely Gemini 1.5 Flash, Llama 3.2-90b-Vision,

LLM	Zero-shot		Z-ICL		KATE		Cluster-TopN		Votek		NwayKshot		CLL-RetICL	
	ACC	F1	ACC	F1	ACC	F1	ACC	F1	ACC	F1	ACC	F1	ACC	F1
SST2														
Gemini	93.29 <sub>.56</sub>	.933 <sub>.002</sub>	92.31 <sub>.29</sub>	.923 <sub>.005</sub>	94.17 <sub>.42</sub>	.941 <sub>.004</sub>	94.93 <sub>.47</sub>	.950 <sub>.003</sub>	94.16 <sub>.33</sub>	.942 <sub>.004</sub>	94.67 <sub>.47</sub>	.947 <sub>.002</sub>	<b>95.17<sub>.37</sub></b>	<b>.952<sub>.004</sub></b>
Llama	94.83 <sub>.62</sub>	.948 <sub>.004</sub>	<b>96.21<sub>.54</sub></b>	<b>.962<sub>.006</sub></b>	94.78 <sub>.56</sub>	.948 <sub>.004</sub>	93.61 <sub>.63</sub>	.936 <sub>.005</sub>	94.77 <sub>.47</sub>	.948 <sub>.004</sub>	90.82 <sub>.71</sub>	.908 <sub>.003</sub>	<u>95.06<sub>.43</sub></u>	<u>.951<sub>.004</sub></u>
Mistral	90.08 <sub>.32</sub>	.901 <sub>.002</sub>	90.72 <sub>.51</sub>	.906 <sub>.003</sub>	93.78 <sub>.21</sub>	.938 <sub>.003</sub>	94.88 <sub>.37</sub>	<u>.949<sub>.004</sub></u>	94.34 <sub>.46</sub>	.943 <sub>.002</sub>	94.34 <sub>.29</sub>	.943 <sub>.003</sub>	<b>95.60<sub>.21</sub></b>	<b>.956<sub>.003</sub></b>
Avg.	92.73	.927	93.08	.930	94.24	.942	<u>94.47</u>	<u>.945</u>	94.42	.944	93.27	.933	<b>95.28</b>	<b>.953</b>
CoLA														
Gemini	68.26 <sub>.56</sub>	.663 <sub>.008</sub>	60.21 <sub>.67</sub>	.583 <sub>.007</sub>	70.08 <sub>.84</sub>	.641 <sub>.008</sub>	80.32 <sub>.49</sub>	.765 <sub>.005</sub>	81.43 <sub>.63</sub>	.783 <sub>.007</sub>	<u>82.74<sub>.72</sub></u>	<u>.795<sub>.008</sub></u>	<b>83.60<sub>.91</sub></b>	<b>.801<sub>.006</sub></b>
Llama	61.74 <sub>.89</sub>	.585 <sub>.006</sub>	52.34 <sub>.72</sub>	.511 <sub>.007</sub>	68.36 <sub>.83</sub>	.650 <sub>.005</sub>	71.62 <sub>.76</sub>	.711 <sub>.007</sub>	61.42 <sub>.69</sub>	.607 <sub>.004</sub>	<u>74.52<sub>.71</sub></u>	<u>.686<sub>.008</sub></u>	<b>77.66<sub>.84</sub></b>	<b>.742<sub>.003</sub></b>
Mistral	74.30 <sub>.34</sub>	.697 <sub>.006</sub>	71.52 <sub>.56</sub>	.666 <sub>.003</sub>	78.71 <sub>.56</sub>	.752 <sub>.004</sub>	84.29 <sub>.56</sub>	.811 <sub>.005</sub>	84.48 <sub>.56</sub>	.821 <sub>.006</sub>	<u>85.23<sub>.56</sub></u>	<u>.816<sub>.007</sub></u>	<b>85.52<sub>.56</sub></b>	<b>.828<sub>.004</sub></b>
Avg.	68.10	.648	61.36	.587	72.38	.681	78.74	.762	75.78	.737	<u>80.83</u>	<u>.766</u>	<b>82.26</b>	<b>.790</b>
CARER														
Gemini	59.20 <sub>.51</sub>	.493 <sub>.004</sub>	65.85 <sub>.60</sub>	.607 <sub>.002</sub>	<u>70.85<sub>.52</sub></u>	<u>.621<sub>.002</sub></u>	61.65 <sub>.49</sub>	.533 <sub>.004</sub>	59.95 <sub>.67</sub>	.541 <sub>.004</sub>	66.25 <sub>.41</sub>	.596 <sub>.002</sub>	<b>72.65<sub>.67</sub></b>	<b>.669<sub>.005</sub></b>
Llama	56.75 <sub>.31</sub>	.488 <sub>.005</sub>	<u>65.70<sub>.60</sub></u>	<u>.594<sub>.003</sub></u>	61.95 <sub>.49</sub>	.537 <sub>.006</sub>	57.35 <sub>.32</sub>	.499 <sub>.002</sub>	59.50 <sub>.71</sub>	.526 <sub>.003</sub>	64.25 <sub>.54</sub>	.579 <sub>.004</sub>	<b>69.15<sub>.32</sub></b>	<b>.635<sub>.002</sub></b>
Mistral	56.50 <sub>.41</sub>	.506 <sub>.002</sub>	67.10 <sub>.48</sub>	.617 <sub>.004</sub>	68.89 <sub>.37</sub>	.601 <sub>.003</sub>	60.25 <sub>.29</sub>	.515 <sub>.003</sub>	58.75 <sub>.50</sub>	.498 <sub>.002</sub>	<u>72.10<sub>.43</sub></u>	<u>.670<sub>.003</sub></u>	<b>76.85<sub>.20</sub></b>	<b>.717<sub>.004</sub></b>
Avg.	57.48	.495	66.22	.606	67.23	.586	59.75	.516	59.40	.521	<u>67.53</u>	<u>.615</u>	<b>72.88</b>	<b>.674</b>
BBCNews														
Gemini	87.00 <sub>.31</sub>	.869 <sub>.013</sub>	87.70 <sub>.45</sub>	.872 <sub>.007</sub>	<b>90.99<sub>.21</sub></b>	<b>.909<sub>.005</sub></b>	85.30 <sub>.64</sub>	.850 <sub>.010</sub>	86.20 <sub>.35</sub>	.858 <sub>.011</sub>	88.60 <sub>.56</sub>	.884 <sub>.008</sub>	89.50 <sub>.37</sub>	.892 <sub>.006</sub>
Llama	94.89 <sub>.56</sub>	.948 <sub>.008</sub>	93.43 <sub>.41</sub>	.933 <sub>.004</sub>	94.70 <sub>.31</sub>	.946 <sub>.006</sub>	93.60 <sub>.41</sub>	.935 <sub>.004</sub>	96.00 <sub>.21</sub>	.960 <sub>.008</sub>	<u>96.10<sub>.52</sub></u>	<u>.960<sub>.007</sub></u>	<b>96.80<sub>.50</sub></b>	<b>.967<sub>.005</sub></b>
Mistral	<u>91.70<sub>.26</sub></u>	<u>.915<sub>.005</sub></u>	90.60 <sub>.31</sub>	.903 <sub>.002</sub>	<b>92.99<sub>.29</sub></b>	<b>.929<sub>.006</sub></b>	83.10 <sub>.46</sub>	.826 <sub>.017</sub>	83.00 <sub>.41</sub>	.825 <sub>.007</sub>	87.20 <sub>.29</sub>	.872 <sub>.010</sub>	88.10 <sub>.49</sub>	.879 <sub>.009</sub>
Avg.	91.20	.910	90.57	.902	<b>92.89</b>	<b>.928</b>	87.33	.870	88.40	.881	90.63	.905	<u>91.47</u>	<u>.912</u>

Table 1: Text classification results evaluated on four datasets using three LLMs. **Bold** indicates the best result and underline indicates the result worse than the best result.

Method	Gemini		Llama		Mistral		Avg.	
	ACC	F1	ACC	F1	ACC	F1	ACC	F1
SST2								
Baseline	94.67 <sub>.47</sub>	.947 <sub>.002</sub>	90.82 <sub>.71</sub>	.908 <sub>.003</sub>	94.34 <sub>.29</sub>	.943 <sub>.003</sub>	93.27	.932
PLA	<b>95.44<sub>.28</sub></b>	<b>.954<sub>.003</sub></b>	93.46 <sub>.31</sub>	<u>.934<sub>.004</sub></u>	94.34 <sub>.36</sub>	.943 <sub>.005</sub>	94.41	.943
NLP	<u>95.38<sub>.30</sub></u>	<b>.954<sub>.003</sub></b>	92.31 <sub>.16</sub>	.922 <sub>.004</sub>	<b>96.37<sub>.46</sub></b>	<b>.963<sub>.002</sub></b>	<u>94.68</u>	<u>.946</u>
CTL	<b>95.44<sub>.35</sub></b>	<b>.954<sub>.002</sub></b>	91.65 <sub>.62</sub>	.916 <sub>.004</sub>	95.11 <sub>.28</sub>	.951 <sub>.003</sub>	94.06	.940
Ours	95.17 <sub>.37</sub>	.952 <sub>.004</sub>	<b>95.06<sub>.43</sub></b>	<b>.951<sub>.004</sub></b>	95.60 <sub>.21</sub>	.956 <sub>.004</sub>	<b>95.28</b>	<b>.953</b>
CoLA								
Baseline	82.74 <sub>.72</sub>	.795 <sub>.008</sub>	64.52 <sub>.71</sub>	.586 <sub>.008</sub>	85.23 <sub>.56</sub>	.816 <sub>.007</sub>	77.49	.732
PLA	<u>83.31<sub>.54</sub></u>	<u>.798<sub>.006</sub></u>	73.53 <sub>.86</sub>	<u>.656<sub>.008</sub></u>	<u>85.31<sub>.75</sub></u>	<b>.832<sub>.008</sub></b>	<u>80.72</u>	<u>.762</u>
NLP	82.45 <sub>.43</sub>	.791 <sub>.005</sub>	64.05 <sub>.79</sub>	.579 <sub>.008</sub>	85.04 <sub>.64</sub>	.823 <sub>.005</sub>	77.18	.731
CTL	82.74 <sub>.86</sub>	.794 <sub>.007</sub>	62.79 <sub>.62</sub>	.579 <sub>.004</sub>	85.04 <sub>.95</sub>	.824 <sub>.011</sub>	76.86	.732
Ours	<b>83.60<sub>.91</sub></b>	<b>.801<sub>.006</sub></b>	<b>77.66<sub>.84</sub></b>	<b>.742<sub>.003</sub></b>	<b>85.52<sub>.58</sub></b>	<u>.828<sub>.004</sub></u>	<b>82.26</b>	<b>.790</b>
CARER								
Baseline	66.25 <sub>.41</sub>	.596 <sub>.002</sub>	64.25 <sub>.54</sub>	.579 <sub>.004</sub>	<u>72.10<sub>.43</sub></u>	<u>.670<sub>.003</sub></u>	<u>67.53</u>	<u>.615</u>
PLA	65.75 <sub>.64</sub>	.598 <sub>.005</sub>	61.65 <sub>.52</sub>	.556 <sub>.011</sub>	65.59 <sub>.61</sub>	.596 <sub>.008</sub>	64.32	.583
NLP	67.35 <sub>.39</sub>	<u>.619<sub>.004</sub></u>	64.40 <sub>.25</sub>	.583 <sub>.007</sub>	70.00 <sub>.38</sub>	.644 <sub>.005</sub>	67.25	<u>.615</u>
CTL	66.90 <sub>.45</sub>	.605 <sub>.007</sub>	65.40 <sub>.50</sub>	.586 <sub>.005</sub>	67.80 <sub>.44</sub>	.615 <sub>.007</sub>	66.70	.602
Ours	<b>72.65<sub>.67</sub></b>	<b>.669<sub>.005</sub></b>	<b>69.15<sub>.32</sub></b>	<b>.635<sub>.002</sub></b>	<b>76.85<sub>.20</sub></b>	<b>.717<sub>.004</sub></b>	<b>72.88</b>	<b>.673</b>
BBCNews								
Baseline	88.60 <sub>.56</sub>	.884 <sub>.008</sub>	96.10 <sub>.52</sub>	.960 <sub>.007</sub>	87.20 <sub>.29</sub>	.872 <sub>.010</sub>	90.63	.905
PLA	89.40 <sub>.35</sub>	.891 <sub>.005</sub>	96.70 <sub>.60</sub>	<u>.966<sub>.003</sub></u>	<b>89.50<sub>.29</sub></b>	<b>.895<sub>.002</sub></b>	<u>91.86</u>	<u>.917</u>
NLP	89.00 <sub>.37</sub>	.889 <sub>.002</sub>	96.40 <sub>.56</sub>	.964 <sub>.004</sub>	88.40 <sub>.42</sub>	.883 <sub>.006</sub>	91.20	.875
CTL	<b>90.30<sub>.54</sub></b>	<b>.901<sub>.003</sub></b>	96.50 <sub>.71</sub>	.964 <sub>.003</sub>	<u>89.40<sub>.63</sub></u>	<u>.893<sub>.006</sub></u>	<b>92.06</b>	<b>.919</b>
Ours	89.50 <sub>.37</sub>	.892 <sub>.006</sub>	<b>96.80<sub>.50</sub></b>	<b>.967<sub>.005</sub></b>	88.10 <sub>.45</sub>	.879 <sub>.009</sub>	91.47	.912

Table 2: A comparative analysis of various linguistic label retrieval methods across four datasets.

and Mistral Large. These recently developed models demonstrate strong performance and exceptional generalization across a variety of tasks.

**Similarity function.** We define a similarity function, *sim*, as the cosine similarity between two sentence embeddings. These embeddings are generated using the all-MiniLM-L6-v2 model from the SBERT (Reimers and Gurevych, 2019).

**Implementation details.** For all LLMs, we use two random seeds and report the average results. We set the default number of demonstrations  $k$  per class to 3 for all experiments. We adopt the typical prompt design methodology proposed by (Luo et al., 2024). To ensure accurate and consistent results in text classification tasks, we employ fixed

hyperparameters for LLMs, thereby minimizing variability and limiting creative outputs. Further details are provided in Appendix B.

## 5.2 Experimental Results

### 5.2.1 Main results

Table 1 presents the results obtained using various retrieval strategies across three LLMs. The zero-shot approach, which does not rely on retrieving relevant demonstrations from the support set, leverages only the semantic understanding of labels. This strategy enables LLMs to achieve a baseline level of accuracy without additional context. Although Z-ICL mitigates the Copying Effect by leveraging physical neighbors and synonym labels, it only marginally outperforms the zero-shot baseline. However, it lags behind other methods, likely due to the inherent complexity and challenges associated with selecting appropriate synonym labels. KATE achieves better performance than zero-shot and Z-ICL by utilizing the most similar demonstrations to the query. However, it is susceptible to errors caused by misleading similarities. As a result, KATE still struggles to perform well on the CoLA and CARER datasets. To mitigate the effects of misleading similarities, NwayKshot generally outperforms KATE in most scenarios. However, as noted earlier, NwayKshot still struggles to identify an optimal combination of demonstrations. VoteK attempts to further select more effective and relevant demonstrations from the support set. However, this method still fails to utilize label information effectively. On the other hand, Cluster-TopN leverages label information from a distributional

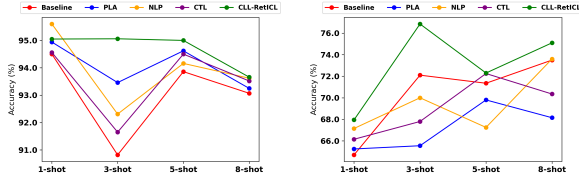


Figure 4: A comparison of the performance of various shot configurations is presented across a baseline and four linguistic label retrieval strategies. Evaluations for the SST2 task (using Llama) are on the left, while results for the CARER task (using Mistral) appear on the right.

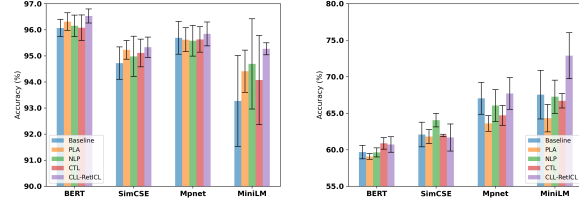


Figure 5: A comparison of the performance of various sentence embedding models is presented, with evaluations conducted on SST2 on the left and CARER on the right.

perspective but does not account for the linguistic meaning of the labels. While both VoteK and Cluster-TopN show improvements in accuracy for certain tasks, they fall short in addressing a fundamental issue: the importance of linguistic label meaning in text classification tasks. This oversight leads to inconsistent performance and highlights their inherent weaknesses. Finally, our proposed method, CLL-RetICL, significantly outperforms all baseline approaches. On average, CLL-RetICL improves RetICL’s performance by an absolute margin of 2–15% over the zero-shot strategy and by 0.57–13.48% over existing RetICL-based methods. These results demonstrate consistent performance gains across all datasets and LLMs by effectively leveraging the relationships between linguistic labels and their corresponding sentences.

**Comparison to Variants of Label-Related RetICL.** We use the NwayKshot method as our baseline, a retrieval-based approach that does not utilize linguistic label information. To enhance performance, we evaluate four proposed strategies that incorporate linguistic label related retrieval methods, with the results summarized in Table 2. All four strategies outperform the baseline across all datasets and LLMs, demonstrating the benefits of leveraging label information. Among these, CLL-RetICL consistently delivers the best performance, achieving an average absolute improvement of 0.8–5.3% over the NwayKshot method. While

PLA, NLP, and CTL also surpass the baseline, they show minor performance drops on certain tasks. In contrast, CLL-RetICL not only outperforms these methods in most tasks but also achieves consistent gains in classification accuracy.

### 5.3 Ablation Study

We conduct detailed ablation studies to analyze the significance of each component in CLL-RetICL. In our ablation study, the NwayKshot approach serves as the baseline, as shown in the following tables and figures.

**Effect of the number of shots.** The number of shots significantly impacts the performance of LLMs. We explore experiments comparing four different shot configurations for each label class: 1-shot, 3-shot, 5-shot, and 8-shot. Figure 4 presents partial results, while the complete results are provided in Appendix D.1. The results in Figure 4 demonstrate that CLL-RetICL consistently outperforms the baseline methods across different values of  $k$ . While some alternative strategies occasionally achieve better performance than CLL-RetICL, they lack robustness and often fall short of both CLL-RetICL and the baselines. This indicates that CLL-RetICL delivers more stable performance across a range of scenarios. Based on the experimental results, we selected  $k = 3$  as the hyperparameter for the number of shots, as CLL-RetICL demonstrated higher improvement with a 3-shot configuration.

**Effect of sentence embedding model.** Pre-trained sentence embeddings play a crucial role in ICL. The objective is to evaluate the effectiveness of the proposed methods by comparing them against four off-the-shelf sentence embedding models. Figure 5 illustrates the average performance of three LLMs across two datasets. CLL-RetICL consistently outperforms the baseline and the other three strategies across all sentence embedding models, with the exception of SimCSE (Gao et al., 2021) in the CARER dataset. We attribute the relatively lower performance of our method with SimCSE to the fact that SimCSE has already employed contrastive learning to fine-tune the pre-trained sentence embedding model. This suggests that our approach is generally more effective for pre-trained sentence embeddings that do not utilize contrastive learning strategies. Compared to other sentence embedding models, MiniLM demonstrates the greatest improvement over the baseline;

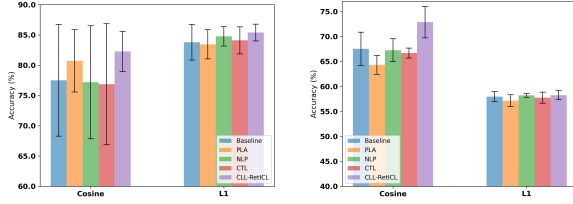


Figure 6: A comparison of the performance of various similarity functions is presented, with evaluations conducted on CoLA on the left and CARER on the right.

therefore, we have chosen it as our default. Full results are presented in Appendix D.2.

**Effect of similarity function.** To evaluate the effect of the similarity function in our CLL-RetICL model, we compare its performance using another similarity function, L1, as described in (Winata et al., 2023). The results are presented in Figure 6 with detailed results provided in Appendix D.3.

CLL-RetICL performs effectively with both cosine and L1 similarity functions. However, experiments show that cosine similarity outperforms the L1 function, suggesting that it better leverages CLL-RetICL’s potential. Consequently, we use cosine similarity as the default.

**Effect of w/o similarity between demonstration and query.** Because our proposed additional component can serve as a scoring criterion for selecting demonstrations, the question arises whether the similarity score between demonstrations and the query should be included in CLL-RetICL.

We evaluate the problem and present the results in Figure 7. Our findings indicate that the performance without the component addressing the similarity between queries and sentences is consistently lower than that of linguistically labeled RetICL. In fact, it performs even worse than the baseline. These results highlight that the similarity component between queries and sentences is an essential part of the retrieval process. Detailed results are presented in Appendix D.4.

**Effect of trade-off hyperparameters.** We use a trade-off approach to balance the impact between sentences and their label set. Based on the results of the previous experiment, sentence-query similarity remains a crucial factor in selecting relevant demonstrations. This raises an important question: how should we trade off between the original method, which retrieves the closest demonstrations to the query, and our approach? To address this question, we evaluate the effects of various hyperparameter

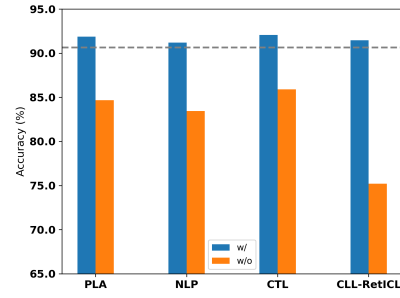


Figure 7: A comparison of the retrieval process with and without incorporating the similarity score between the query and the sentence is illustrated on BBCNews dataset. The baseline is represented by a dashed line.

settings. Specifically, we focus on hyperparameters lower than 1.0, as previous research has consistently shown that closer demonstrations generally outperform those that are further away. We maintain the principle that proximity to the query remains a core factor in our approach. Based on these observations in Appendix D.5, we found that the trade-off hyperparameter has some influence on the final results. However, their impact on PLA, NLP, and CTL methods is relatively small. Interestingly, we observed that a trade-off hyperparameter value of 1.0 yields the best performance for our CLL-RetICL method. Consequently, we adopt 1.0 as the default hyperparameter.

## 6 Conclusion

This paper introduces a new paradigm Contrastive Linguistic Label Retrieval-based In-Context Learning. Unlike existing approaches that universally sample demonstrations without considering the linguistic label information, we propose a general framework for identifying more effective and relevant demonstrations. This framework enhances the capabilities of LLMs to produce more accurate text classification results. Additionally, we design a universal prompt that is adaptable to all text classification tasks. Empirical evaluation on four datasets demonstrates that CLL-RetICL significantly outperforms conventional practices in RetICL by incorporating the similarity between linguistic labels and sentences. This highlights the promising performance of CLL-RetICL and opens up several intriguing research opportunities for further methodological exploration.

## 7 Limitations

**Requiring Semantic Labels.** Our approach focuses exclusively on the semantic label text clas-



sification task. Certain text classification scenarios, however, may involve ambiguous label classes, such as class0, class1, . . . . Ambiguities in labeling may introduce additional challenges, as discussed in Appendix E.1. Addressing these issues remains an open area for future research.

**Better Descriptive Labels** Recently, the use of class-label synonyms has become a popular and compelling topic of research (Pawar et al., 2024). In our work, we also present results using class-label synonyms on the SST-2 dataset, as shown in Table 17. Our findings indicate that CLL-RetICL consistently performs well across different label synonym settings. However, the overall performance with label synonym settings is lower compared to using the original labels. These results suggest that more accurate, semantic, and suitable labels could further enhance the effectiveness of our method.

Moreover, some classification tasks include explanations for the meaning of each label. Using more descriptive sentences and designing multi-label descriptors can help reduce the risk of bias and support effective mitigation strategies. In this work, we did not utilize those explanations. Incorporating these explanations into the classification process is left as a direction for future work.

**Enhance prompt clarity.** In previous work, researchers observed that well-crafted prompts can lead to better results. However, in this study, we did not compare the effects of different prompt formats. Determining how to construct optimal prompts to leverage the potential of our CLL-RetICL framework fully remains an open question and is left for future exploration.

## References

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Youngjin Chae and Thomas Davidson. 2023. Large language models for text classification: From zero-shot learning to fine-tuning. *Open Science Foundation*.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Aleksandra Edwards and Jose Camacho-Collados. 2024. Language models for text classification: Is in-context learning enough? *arXiv preprint arXiv:2403.17661*.

Wenqi Fan, Yajuan Ding, Liangbo Ning, Shijie Wang, Hengyun Li, Dawei Yin, Tat-Seng Chua, and Qing Li. 2024. A survey on rag meeting llms: Towards retrieval-augmented large language models. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 6491–6501.

John Fields, Kevin Chovanec, and Praveen Madiraju. 2024. A survey of text classification with transformers: How wide? how large? how long? how accurate? how expensive? how safe? *IEEE Access*.

Lingyu Gao, Aditi Chaudhary, Krishna Srinivasan, Kazuma Hashimoto, Karthik Raman, and Michael Bendersky. 2023. Ambiguity-aware in-context learning with large language models. *arXiv preprint arXiv:2309.07900*.

Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. Simcse: Simple contrastive learning of sentence embeddings. *arXiv preprint arXiv:2104.08821*.

Derek Greene and Pádraig Cunningham. 2006. Practical solutions to the problem of diagonal dominance in kernel document clustering. In *Proceedings of the 23rd international conference on Machine learning*, pages 377–384.

Qi Guo, Leiyu Wang, Yidong Wang, Wei Ye, and Shikun Zhang. 2024. What makes a good order of examples in in-context learning. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 14892–14904.

Sabit Hassan and Malihe Alikhani. 2023. D-calm: A dynamic clustering-based active learning approach for mitigating bias. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 5540–5553.

Eduard Hovy, Laurie Gerber, Ulf Hermjakob, Chin-Yew Lin, and Deepak Ravichandran. 2001. [Toward semantics-based answer pinpointing](#). In *Proceedings of the First International Conference on Human Language Technology Research*.

Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. 2024. Mixtral of experts. *arXiv preprint arXiv:2401.04088*.

Hongpeng Jin, Wenqi Wei, Xuyu Wang, Wenbin Zhang, and Yanzhao Wu. 2023. Rethinking learning rate tuning in the era of large language models. In *2023 IEEE 5th International Conference on Cognitive Machine Intelligence (CogMI)*, pages 112–121. IEEE.

Guozheng Li, Peng Wang, Jiajun Liu, Yikai Guo, Ke Ji, Ziyu Shang, and Zijie Xu. 2024. Meta in-context learning makes large language models better zero

- and few-shot relation extractors. *arXiv preprint arXiv:2404.17807*. 709 710
- Jing Li, Billy Chiu, Shanshan Feng, and Hao Wang. 2020. Few-shot named entity recognition via meta-learning. *IEEE Transactions on Knowledge and Data Engineering*, 34(9):4245–4256. 711 712 713 714
- Junlong Li, Jinyuan Wang, Zhuosheng Zhang, and Hai Zhao. 2022. Self-prompting large language models for zero-shot open-domain qa. *arXiv preprint arXiv:2212.08635*. 715 716 717 718
- Xiaonan Li, Kai Lv, Hang Yan, Tianyang Lin, Wei Zhu, Yuan Ni, Guotong Xie, Xiaoling Wang, and Xipeng Qiu. 2023. Unified demonstration retriever for in-context learning. *arXiv preprint arXiv:2305.04320*. 719 720 721 722
- Xiaonan Li and Xipeng Qiu. 2023. Finding support examples for in-context learning. *arXiv preprint arXiv:2302.13539*. 723 724 725
- Xin Li and Dan Roth. 2002. [Learning question classifiers](#). In *COLING 2002: The 19th International Conference on Computational Linguistics*. 726 727 728
- Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2021. What makes good in-context examples for gpt-3? *arXiv preprint arXiv:2101.06804*. 729 730 731 732
- Man Luo, Xin Xu, Yue Liu, Panupong Pasupat, and Mehran Kazemi. 2024. In-context learning with retrieved demonstrations for language models: A survey. *arXiv preprint arXiv:2401.11624*. 733 734 735 736
- Xinxi Lyu, Sewon Min, Iz Beltagy, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2022. Z-icl: Zero-shot in-context learning with pseudo-demonstrations. *arXiv preprint arXiv:2212.09865*. 737 738 739 740
- Costas Mavromatis, Balasubramaniam Srinivasan, Zhengyuan Shen, Jiani Zhang, Huzefa Rangwala, Christos Faloutsos, and George Karypis. 2023. Which examples to annotate for in-context learning? towards effective and efficient selection. *arXiv preprint arXiv:2310.20046*. 741 742 743 744 745 746
- Aristides Milios, Siva Reddy, and Dzmitry Bahdanau. 2023. In-context learning for text classification with many labels. In *Proceedings of the 1st GenBench Workshop on (Benchmarking) Generalisation in NLP*, pages 173–184. 747 748 749 750 751
- Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, et al. 2022. In-context learning and induction heads. *arXiv preprint arXiv:2209.11895*. 752 753 754 755 756
- Sachin Pawar, Nitin Ramrakhiyani, Anubhav Sinha, Manoj Apte, and Girish Palshikar. 2024. Why generate when you can discriminate? a novel technique for text classification using language models. In *Findings of the association for computational linguistics: EACL 2024*, pages 1099–1114. 757 758 759 760 761 762
- Branislav Pecher, Ivan Srba, Maria Bielikova, and Joaquin Vanschoren. 2024. Automatic combination of sample selection strategies for few-shot learning. *arXiv preprint arXiv:2402.03038*. 763 764 765 766
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9. 767 768 769 770
- Nils Reimers and Iryna Gurevych. 2019. [Sentence-bert: Sentence embeddings using siamese bert-networks](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics. 771 772 773 774 775
- Hamidreza Rouzegar and Masoud Makrehchi. 2024. Enhancing text classification through llm-driven active learning and human annotation. *arXiv preprint arXiv:2406.12114*. 776 777 778 779
- Elvis Saravia, Hsien-Chi Toby Liu, Yen-Hao Huang, Junlin Wu, and Yi-Shin Chen. 2018. Carer: Contextualized affect representations for emotion recognition. In *Proceedings of the 2018 conference on empirical methods in natural language processing*, pages 3687–3697. 780 781 782 783 784 785
- Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, Ed H Chi, Nathanael Schärli, and Denny Zhou. 2023. Large language models can be easily distracted by irrelevant context. In *International Conference on Machine Learning*, pages 31210–31227. PMLR. 786 787 788 789 790 791
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pages 1631–1642. 792 793 794 795 796 797 798
- Hongjin Su, Jungo Kasai, Chen Henry Wu, Weijia Shi, Tianlu Wang, Jiayi Xin, Rui Zhang, Mari Ostendorf, Luke Zettlemoyer, Noah A Smith, et al. 2022. Selective annotation makes language models better few-shot learners. *arXiv preprint arXiv:2209.01975*. 799 800 801 802 803
- Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*. 804 805 806 807 808 809
- Lewis Tunstall, Nils Reimers, Unso Eun Seo Jo, Luke Bates, Daniel Korat, Moshe Wasserblat, and Oren Pereg. 2022. Efficient few-shot learning without prompts. *arXiv preprint arXiv:2209.11055*. 810 811 812 813
- Lean Wang, Lei Li, Damai Dai, Deli Chen, Hao Zhou, Fandong Meng, Jie Zhou, and Xu Sun. 2023. Label words are anchors: An information flow perspective for understanding in-context learning. *arXiv preprint arXiv:2305.14160*. 814 815 816 817 818

- Zhiqiang Wang, Yiran Pang, and Yanbin Lin. 2024. Smart expert system: Large language models as text classifiers. *arXiv preprint arXiv:2405.10523*.
- Alex Warstadt, Amanpreet Singh, and Samuel R Bowman. 2018. Neural network acceptability judgments. *arXiv preprint arXiv:1805.12471*.
- Genta Indra Winata, Liang-Kang Huang, Soumya Vadamannati, and Yash Chandarana. 2023. Multilingual few-shot learning via language model retrieval. *arXiv preprint arXiv:2306.10964*.
- Zhiyong Wu, Yaoxiang Wang, Jiacheng Ye, and Lingpeng Kong. 2022. Self-adaptive in-context learning: An information compression perspective for in-context example selection and ordering. *arXiv preprint arXiv:2212.10375*.
- Hao Yu, Zachary Yang, Kellin Pelrine, Jean Francois Godbout, and Reihaneh Rabbany. 2023. Open, closed, or small language models for text classification? *arXiv preprint arXiv:2308.10092*.
- Yazhou Zhang, Mengyao Wang, Chenyu Ren, Qiuchi Li, Prayag Tiwari, Benyou Wang, and Jing Qin. 2024. Pushing the limit of llm capacity for text classification. *arXiv preprint arXiv:2402.07470*.
- Yiming Zhang, Shi Feng, and Chenhao Tan. 2022a. Active example selection for in-context learning. *arXiv preprint arXiv:2211.04486*.
- Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. 2022b. Automatic chain of thought prompting in large language models. *arXiv preprint arXiv:2210.03493*.
- Fedor Zhdanov. 2019. Diverse mini-batch active learning. *arXiv preprint arXiv:1901.05954*.

## A Data Statistics

We take the four text classification tasks including SST2, CoLA, CARER, and BBCNews. See the descriptions and statistics in Table 3.

We use the original SST-2 dataset that only comprises the complete sentences that are not labeled neutral, and its original split is 6920/872/1821 (Socher et al., 2013).

CoLA comprises 10,657 sentences sourced from 23 linguistics publications. Each sentence has been expertly annotated for acceptability (i.e., grammaticality) by the original authors (Warstadt et al., 2018). CoLA is divided into two subsets: a training set and a development set. In our work, we treat the development set as the test set.

CARER is a dataset of English Twitter messages with six basic emotions: anger, fear, joy, love, sadness, and surprise (Saravia et al., 2018). The original CARER dataset has been split into trainsets, validation, and test sets.

The BBC News Topic Classification dataset consists of 2,225 articles published on the BBC News website between 2004 and 2005. Each article is labeled under one of 5 categories: business, entertainment, politics, sport, or tech (Greene and Cunningham, 2006). The original BBCNews dataset has been split into trainsets and test sets.

As stated in the main text, we exclude the validation set.

## B Implementation Details

All implementations are done in PyTorch.

**Prompt template.** We adopt the prompt used in the CARER task as a template. Following the approach of Luo et al. (2024), we design our prompt for universal text classification tasks, as shown in Table 5. (*demo\_1*), (*demo\_2*), (*demo\_3*) are selected demonstration from support set. (*query*) is the current query sentence.

**Budget.** We conducted experiments on LLMs across four public datasets, utilizing APIs to compute the results. To ensure consistency and avoid generating creative outputs, we fixed the LLMs’ hyperparameters, as detailed in Table 4. The total cost of running these experiments through the APIs amounted to approximately 1,000 US dollars.

## C Linguistic Label Retrieval Hypothesis

To explore the question: Are the clustering centers reliable, we analyze the distribution of clustering



Figure 8: An example illustrating the distribution of queries, linguistic labels, and clustering centers in a pre-trained sentence embedding model using t-SNE. 10%\_N and 10%\_P represent a pair, indicating that 10% of the support set is used to estimate the clustering center. In this notation, "N" refers to the Negative-labeled clustering center, while "P" denotes the Positive-labeled clustering center.

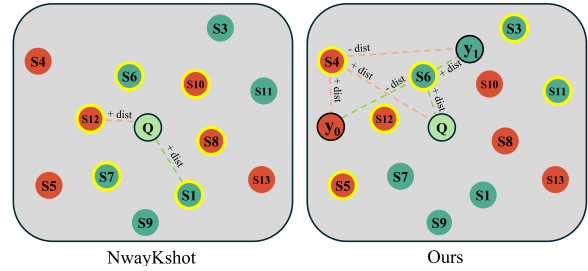


Figure 9: A comparison of the retrieved demonstrations between NwayKshot and our method in the sentence embedding vector space. A yellow circle indicates a selected sentence.

centers, as shown in Figure 8.

We present five examples in Table 6 to illustrate our findings. The experiment was conducted on the SST2 dataset, where we treated the training set as the support set and the test set as the query set. For each query example in the test set, we provide its index and the indices of the selected demonstrations from the training set. Additionally, we report the similarity scores calculated using cosine distance within the vector space of a pre-trained sentence embedding model (Reimers and Gurevych, 2019).

## D Additional Results

### D.1 Effect of the number of shots.

We present the detailed results in Table 7.



Dataset	Trainset	Testset	Label
SST2	6,920	1,821	"Negative", "Positive"
CoLA	8,551	1,043	"Unacceptable", "Acceptable"
CARER	16,000	2,000	"Sadness", "Joy", "Love", "Anger", "Fear", "Surprise"
BBCNews	1,225	1,000	"Business", "Entertainment", "Politics", "Sport", "Tech"

Table 3: Statistics of datasets as well as their labels.

Configure	Gemini	Llama	Mistral
"temperature"	0.2	0.01	0.01
"top_p"	0.9	0.9	0.5
"top_k"	1	1	1
"max_output_tokens"	2	2	2

Table 4: generation\_config of hyperparameters in various LLMs.

## D.2 Effect of sentence embedding model.

We present the detailed results in Table 8.

## D.3 Effect of similarity function.

We present the detailed results in Table 9.

## D.4 Effect of w/o similarity between demonstration and query.

We present the detailed results in Table 12.

## D.5 Effect of trade-off hyperparameters.

We present the detailed results in Table 13, Table 14, Table 15, Table 10, Table 11.

## E Limitations

### E.1 Requiring Semantic Labels.

Handling complex classification tasks with ambiguous labels presents additional challenges for our method, as CLL-RetICL relies heavily on semantic label representations. To illustrate this issue, we use the TREC dataset (Li and Roth, 2002; Hovy et al., 2001), which provides both abbreviated and full-form class labels. In our analysis, we adopt the coarse-label scheme and specifically compare the abbreviated class labels with their corresponding full descriptive labels. The abbreviated labels include [ABBR, ENTY, DESC, HUM, LOC, NUM], while their full counterparts are [Abbreviation, Entity, Description and abstract concept, Human being, Location, Numeric value]. The results, shown in Table 16, demonstrate that using abbreviated labels weakens the performance of our method compared to using the full descriptive labels.

### E.2 Better Descriptive Labels

System message	"You are given a task where there are multiple classes, and for each class, a few labeled examples are provided. Based on these examples, you need to classify a new unseen instance. Choose ONLY one tag and output the tag. Do Not output others."
CARER	
Prompt	Class0: sadness 1. Example 1: (demo_1) -> "sadness" 2. Example 2: (demo_2) -> "sadness" 3. Example 3: (demo_3) -> "sadness" Class1: joy 1. Example 1: (demo_1) -> "joy" 2. Example 2: (demo_2) -> "joy" 3. Example 3: (demo_3) -> "joy" Class2: love 1. Example 1: (demo_1) -> "love" 2. Example 2: (demo_2) -> "love" 3. Example 3: (demo_3) -> "love" Class3: anger 1. Example 1: (demo_1) -> "anger" 2. Example 2: (demo_2) -> "anger" 3. Example 3: (demo_3) -> "anger" Class4: fear 1. Example 1: (demo_1) -> "fear" 2. Example 2: (demo_2) -> "fear" 3. Example 3: (demo_3) -> "fear" Class5: surprise 1. Example 1: (demo_1) -> "surprise" 2. Example 2: (demo_2) -> "surprise" 3. Example 3: (demo_3) -> "surprise" Query: (query) Prediction:

Table 5: Designed universal prompt for all text classification tasks.

query	label	Not Correct						Correct					
		index	sentence	label_N	label_P	center_N	center_P	index	sentence	label_N	label_P	center_N	center_P
36	N	3465	0.501	0.111	<b>0.187</b>	0.576	<b>0.626</b>	1344	0.379	<b>0.173</b>	0.087	<b>0.432</b>	0.417
		5169	0.447	0.095	<b>0.162</b>	0.434	<b>0.454</b>	5012	0.399	<b>0.124</b>	0.063	<b>0.656</b>	0.613
		4441	0.436	0.096	<b>0.131</b>	0.477	<b>0.526</b>	6432	0.399	<b>0.188</b>	0.131	<b>0.461</b>	0.450
	P	5529	0.603	<b>0.232</b>	0.187	0.393	<b>0.479</b>	4310	0.580	0.170	<b>0.247</b>	0.363	<b>0.444</b>
		4310	0.580	<b>0.247</b>	0.170	0.363	<b>0.444</b>	5529	0.603	0.187	<b>0.231</b>	0.393	<b>0.479</b>
		5879	0.507	<b>0.175</b>	0.106	0.561	<b>0.646</b>	6723	0.427	0.129	<b>0.291</b>	0.442	<b>0.591</b>
49	N	4084	0.694	<b>0.022</b>	−0.033	<b>0.525</b>	0.458	4084	0.694	<b>0.022</b>	−0.033	<b>0.525</b>	0.458
		3465	0.564	0.111	<b>0.187</b>	0.576	<b>0.626</b>	6331	0.562	<b>0.122</b>	0.015	<b>0.647</b>	0.569
		6331	0.562	<b>0.122</b>	0.015	<b>0.647</b>	0.569	4290	0.543	<b>0.093</b>	0.044	<b>0.583</b>	0.569
	P	1625	0.672	0.072	<b>0.134</b>	0.531	<b>0.567</b>	1625	0.672	0.072	<b>0.134</b>	0.531	<b>0.567</b>
		4273	0.576	<b>0.040</b>	0.038	0.533	<b>0.579</b>	1268	0.506	−0.024	<b>0.136</b>	0.428	<b>0.482</b>
		1936	0.565	<b>0.669</b>	0.103	0.510	<b>0.550</b>	543	0.516	0.126	<b>0.272</b>	0.379	<b>0.448</b>
1690	N	1613	0.569	<b>-0.029</b>	−0.053	<b>0.341</b>	0.321	1613	0.569	<b>-0.029</b>	−0.053	<b>0.341</b>	0.321
		5550	0.520	<b>0.025</b>	0.015	<b>0.283</b>	0.247	4127	0.497	<b>0.037</b>	−0.033	<b>0.271</b>	0.250
		4127	0.497	<b>0.037</b>	−0.033	<b>0.271</b>	0.250	801	0.466	<b>0.127</b>	0.047	<b>0.464</b>	0.396
	P	3043	0.600	0.019	<b>0.056</b>	0.194	<b>0.217</b>	3043	0.600	0.019	<b>0.056</b>	0.194	<b>0.217</b>
		444	0.502	<b>0.112</b>	0.093	0.334	<b>0.338</b>	4941	0.401	−0.029	<b>0.091</b>	0.464	<b>0.562</b>
		1856	0.480	−0.008	<b>-0.004</b>	0.327	<b>0.363</b>	1856	0.480	−0.008	<b>-0.004</b>	0.327	<b>0.363</b>
1694	N	2433	0.545	−0.054	<b>0.038</b>	0.434	<b>0.441</b>	3367	0.485	<b>0.092</b>	0.061	<b>0.586</b>	0.540
		3367	0.485	<b>0.092</b>	0.061	<b>0.586</b>	0.540	4925	0.478	<b>-0.002</b>	−0.041	<b>0.507</b>	0.456
		4925	0.478	<b>-0.002</b>	−0.041	<b>0.507</b>	0.456	3643	0.428	<b>0.057</b>	−0.026	<b>0.557</b>	0.543
	P	613	0.455	−0.031	<b>-0.011</b>	0.356	<b>0.405</b>	6713	0.433	0.110	<b>0.258</b>	0.549	<b>0.615</b>
		324	0.447	<b>0.181</b>	0.105	0.572	<b>0.623</b>	5337	0.424	0.056	<b>0.188</b>	0.337	<b>0.478</b>
		5135	0.446	0.114	<b>0.198</b>	0.470	<b>0.572</b>	5135	0.446	0.114	<b>0.198</b>	0.470	<b>0.572</b>
1809	N	5557	0.512	<b>0.168</b>	0.138	<b>0.321</b>	0.292	5557	0.512	<b>0.168</b>	0.138	<b>0.321</b>	0.292
		2756	0.405	0.045	<b>0.061</b>	<b>0.328</b>	0.336	4071	0.388	<b>0.114</b>	0.099	<b>0.402</b>	0.376
		2690	0.397	0.088	<b>0.910</b>	<b>0.566</b>	0.507	1430	0.382	<b>0.038</b>	0.023	0.331	<b>0.332</b>
	P	2193	0.480	0.090	<b>0.130</b>	<b>0.570</b>	0.567	2193	0.480	0.090	<b>0.130</b>	<b>0.570</b>	0.567
		6385	0.465	<b>0.094</b>	0.046	0.470	<b>0.477</b>	296	0.347	0.018	<b>0.135</b>	0.295	<b>0.417</b>
		679	0.391	<b>0.149</b>	0.095	0.427	<b>0.444</b>	897	0.359	0.162	<b>0.254</b>	0.394	<b>0.398</b>

Table 6: Five examples comparing incorrect demonstration combinations with their correct counterparts, as evaluated on SST2 task. In this notation, "N" refers to the "Negative" label, while "P" denotes the "Positive" label.

Method	1-shot			3-shot			5-shot			8-shot		
	Gemini	Llama	Mistral	Gemini	Llama	Mistral	Gemini	Llama	Mistral	Gemini	Llama	Mistral
SST2												
Baseline	94.67	94.50	94.61	94.67	90.82	94.34	95.16	93.86	95.00	<u>95.60</u>	93.07	95.00
PLA	94.78	94.94	94.34	<b>95.44</b>	<u>93.46</u>	94.34	95.00	<u>94.62</u>	95.22	95.33	93.35	95.00
NLP	<b>95.38</b>	<b>95.60</b>	<b>96.59</b>	<u>95.38</u>	92.31	<b>96.37</b>	<b>95.94</b>	94.16	<u>95.44</u>	95.00	<u>93.62</u>	<u>95.10</u>
CTL	94.28	94.56	94.28	<b>95.44</b>	91.65	95.11	<u>95.28</u>	94.50	94.89	95.00	93.52	<b>95.16</b>
CLL-RetlCL	<u>94.89</u>	<u>95.05</u>	<u>95.33</u>	95.17	<b>95.06</b>	<u>95.60</u>	<u>95.28</u>	<b>95.00</b>	<b>95.71</b>	<b>96.21</b>	<b>93.66</b>	<b>95.16</b>
CARER												
Baseline	<u>66.55</u>	63.45	64.70	66.25	64.25	<u>72.10</u>	68.23	<u>70.95</u>	71.35	69.50	69.35	73.50
PLA	63.80	60.30	65.25	65.75	61.65	65.55	67.70	62.10	69.80	67.10	65.40	68.15
NLP	66.30	<u>64.60</u>	<u>67.15</u>	<u>67.35</u>	64.40	70.00	<u>70.31</u>	65.35	67.25	<u>69.69</u>	<b>72.30</b>	<u>73.60</u>
CTL	64.30	<u>61.20</u>	<u>66.15</u>	66.90	<u>65.40</u>	67.80	68.61	68.05	<u>72.25</u>	68.30	68.65	70.35
CLL-RetlCL	<b>66.75</b>	<b>65.50</b>	<b>67.95</b>	<b>72.65</b>	<b>69.15</b>	<b>76.85</b>	<b>69.05</b>	<b>74.45</b>	<b>72.30</b>	<b>70.75</b>	<u>71.35</u>	<b>75.10</b>

Table 7: Full results of various shots effect in our proposed methods.

Approach	Bert				Simcse				Mpnet				MiniLM			
	Gemini	Llama	Mistral	Avg.	Gemini	Llama	Mistral	Avg.	Gemini	Llama	Mistral	Avg.	Gemini	Llama	Mistral	Avg.
SST2																
Baseline	95.93	95.76	96.54	96.07 <sub>0.33</sub>	95.11	93.85	95.21	94.72 <sub>0.62</sub>	95.66	94.93	<b>96.48</b>	95.69 <sub>0.63</sub>	94.67	90.82	94.34	93.27 <sub>1.74</sub>
PLA	95.93	96.26	<u>96.76</u>	<u>96.31</u> <sub>0.34</sub>	94.78	<b>95.27</b>	<u>95.66</u>	<u>95.23</u> <sub>0.36</sub>	<b>95.71</b>	95.02	96.15	95.62 <sub>0.46</sub>	<b>95.44</b>	<u>93.46</u>	94.34	94.41 <sub>0.81</sub>
NLP	<u>96.26</u>	95.60	96.59	96.15 <sub>0.41</sub>	<u>95.38</u>	93.90	<u>95.66</u>	94.98 <sub>0.77</sub>	95.55	94.87	<u>96.32</u>	95.58 <sub>0.59</sub>	<u>95.38</u>	92.31	<b>96.37</b>	<u>94.68</u> <sub>1.73</sub>
CTL	95.88	95.60	<u>96.76</u>	96.08 <sub>0.49</sub>	94.94	94.56	<b>95.82</b>	95.11 <sub>0.53</sub>	95.40	<u>95.18</u>	<u>96.32</u>	95.63 <sub>0.49</sub>	<b>95.44</b>	91.65	95.11	94.06 <sub>1.71</sub>
CLL-RetICL	<b>96.32</b>	<b>96.37</b>	<b>96.92</b>	<b>96.53</b> <sub>0.27</sub>	<b>95.77</b>	<u>94.83</u>	95.39	<b>95.33</b> <sub>0.39</sub>	95.44	<b>95.60</b>	<b>96.48</b>	<b>95.84</b> <sub>0.46</sub>	95.17	<b>95.06</b>	95.60	<b>95.28</b> <sub>0.23</sub>
CARER																
Baseline	58.65	59.45	60.90	59.67 <sub>0.93</sub>	<b>64.45</b>	60.60	61.20	62.08 <sub>1.69</sub>	<u>63.95</u>	<u>68.30</u>	<u>68.85</u>	67.03 <sub>2.19</sub>	66.25	64.25	<u>72.10</u>	67.53 <sub>3.33</sub>
PLA	58.75	58.85	59.65	59.08 <sub>0.40</sub>	62.65	62.30	60.45	61.80 <sub>0.96</sub>	62.50	63.25	65.05	63.60 <sub>1.07</sub>	65.75	61.65	65.55	64.32 <sub>1.88</sub>
NLP	59.10	59.30	60.50	59.63 <sub>0.62</sub>	<u>63.10</u>	<b>63.75</b>	<b>65.30</b>	<b>64.05</b> <sub>0.92</sub>	63.00	67.15	68.00	66.05 <sub>2.18</sub>	<u>67.35</u>	64.40	70.00	67.25 <sub>2.29</sub>
CTL	<u>59.80</u>	<b>61.12</b>	<u>61.70</u>	<b>60.87</b> <sub>0.79</sub>	61.90	61.75	62.20	61.95 <sub>0.18</sub>	62.90	65.05	66.15	64.70 <sub>1.35</sub>	66.90	<u>65.40</u>	67.80	66.70 <sub>0.99</sub>
CLL-RetICL	<b>59.90</b>	<u>60.05</u>	<b>62.20</b>	<u>60.72</u> <sub>1.05</sub>	59.05	<u>62.80</u>	<u>63.15</u>	61.67 <sub>1.86</sub>	<b>64.65</b>	<b>69.15</b>	<b>69.30</b>	<b>67.70</b> <sub>2.16</sub>	<b>72.65</b>	<b>69.15</b>	<b>76.85</b>	<b>72.88</b> <sub>3.14</sub>

Table 8: A Comparison of various pre-trained sentence embedding models.

Approach	Cosine				L1			
	Gemini	Llama	Mistral	Avg.	Gemini	Llama	Mistral	Avg.
CoLA								
Baseline	82.74	64.52	85.23	77.50 <sub>9.23</sub>	<u>85.43</u>	79.65	86.28	83.79 <sub>2.95</sub>
PLA	<u>83.31</u>	<u>73.53</u>	<u>85.31</u>	80.72 <sub>5.14</sub>	84.66	80.09	85.62	83.46 <sub>2.41</sub>
NLP	82.45	64.05	85.04	77.18 <sub>9.34</sub>	85.13	<u>82.64</u>	<u>86.57</u>	<u>84.78</u> <sub>1.62</sub>
CTL	82.74	62.79	85.04	76.86 <sub>9.99</sub>	84.56	81.17	<u>86.57</u>	84.10 <sub>2.23</sub>
CLL-RetICL	<b>83.60</b>	<b>77.66</b>	<b>85.52</b>	<b>82.26</b> <sub>3.34</sub>	<b>85.81</b>	<b>83.51</b>	<b>86.86</b>	<b>85.39</b> <sub>1.39</sub>
CARER								
Baseline	66.25	64.25	<u>72.10</u>	67.53 <sub>3.33</sub>	<u>57.30</u>	57.20	<u>59.40</u>	57.97 <sub>1.01</sub>
PLA	65.75	61.65	65.55	64.32 <sub>1.88</sub>	55.95	56.75	58.75	57.15 <sub>1.18</sub>
NLP	<u>67.35</u>	64.40	70.00	67.25 <sub>2.29</sub>	<b>57.95</b>	<b>58.75</b>	57.95	<u>58.22</u> <sub>0.38</sub>
CTL	66.90	<u>65.40</u>	67.80	66.70 <sub>0.99</sub>	56.75	57.20	59.30	57.75 <sub>1.11</sub>
CLL-RetICL	<b>72.65</b>	<b>69.15</b>	<b>76.85</b>	<b>72.88</b> <sub>3.14</sub>	<b>57.95</b>	<u>57.30</u>	<b>59.50</b>	<b>58.25</b> <sub>0.92</sub>

Table 9: A Comparison of Similar Function Methods

CTL	ACC								
LLM	(0.3, 0.3)	(0.3, 0.5)	(0.3, 1.0)	(0.5, 0.3)	(0.5, 0.5)	(0.5, 1.0)	(1.0, 0.3)	(1.0, 0.5)	(1.0, 1.0)
CoLA									
Gemini	82.92	83.30	<b>84.35</b>	83.39	83.40	84.05	83.84	83.76	83.60
Llama	62.24	63.44	62.73	63.26	63.44	63.53	65.45	<b>64.04</b>	62.79
Mistral	<b>85.91</b>	85.33	85.90	85.71	85.62	85.33	85.71	85.33	85.04
CARER									
Gemini	65.55	65.80	65.45	<b>67.00</b>	66.50	65.35	63.65	63.65	66.90
Llama	64.35	62.95	63.30	<b>68.25</b>	67.75	63.15	63.10	63.05	65.40
Mistral	68.30	67.75	67.10	<b>71.10</b>	70.90	66.60	65.10	64.95	67.80

Table 10: A comparison of classification accuracy (%) to assess the impact of various trade-off hyperparameters in the CTL strategy.

CTL	F1								
LLM	(0.3, 0.3)	(0.3, 0.5)	(0.3, 1.0)	(0.5, 0.3)	(0.5, 0.5)	(0.5, 1.0)	(1.0, 0.3)	(1.0, 0.5)	(1.0, 1.0)
CoLA									
Gemini	0.791	0.795	<b>0.809</b>	0.798	0.797	0.806	0.802	0.801	0.801
Llama	0.552	0.584	0.570	0.581	0.581	0.589	<b>0.603</b>	0.588	0.579
Mistral	<b>0.825</b>	0.817	0.824	0.823	0.821	0.819	0.820	0.816	0.824
CARER									
Gemini	0.594	0.596	0.595	<b>0.612</b>	0.607	0.590	0.575	0.573	0.605
Llama	0.576	0.559	0.567	<b>0.612</b>	0.602	0.570	0.570	0.566	0.586
Mistral	0.631	0.605	0.607	<b>0.648</b>	0.645	0.595	0.585	0.582	0.615

Table 11: A comparison of F1 score (%) to assess the impact of various trade-off hyperparameters in the CTL strategy.

Method	Gemini		Llama		Mistral		Avg.	
	ACC	F1	ACC	F1	ACC	F1	ACC	F1
CARER								
Baseline	66.25	0.596	64.25	0.579	72.10	0.670	67.53	0.615
PLA	<b>65.75</b>	<b>0.598</b>	<b>61.65</b>	<b>0.556</b>	<b>65.55</b>	<b>0.596</b>	<b>64.32</b>	<b>0.583</b>
w/o	57.80	0.514	56.75	0.497	56.10	0.488	56.88	0.499
NLP	<b>67.35</b>	<b>0.619</b>	<b>64.40</b>	<b>0.583</b>	<b>70.00</b>	<b>0.644</b>	<b>67.25</b>	<b>0.615</b>
w/o	57.05	0.502	58.90	0.512	59.70	0.529	58.55	0.514
CTL	<b>66.90</b>	<b>0.605</b>	<b>65.40</b>	<b>0.586</b>	<b>67.80</b>	<b>0.615</b>	<b>66.70</b>	<b>0.602</b>
w/o	58.40	0.521	56.30	0.494	57.75	0.505	57.48	0.507
Ours	<b>72.65</b>	<b>0.669</b>	<b>69.15</b>	<b>0.635</b>	<b>76.85</b>	<b>0.717</b>	<b>72.88</b>	<b>0.673</b>
w/o	52.25	0.468	55.50	0.486	51.70	0.463	53.15	0.472
BBCNews								
Baseline	88.60	0.884	96.10	0.960	87.20	0.872	90.63	0.905
PLA	<b>89.40</b>	<b>0.891</b>	<b>96.70</b>	<b>0.966</b>	<b>89.50</b>	<b>0.895</b>	<b>91.86</b>	<b>0.917</b>
w/o	79.50	0.777	94.20	0.940	80.30	0.796	84.67	0.837
NLP	<b>89.00</b>	<b>0.889</b>	<b>96.40</b>	<b>0.964</b>	<b>88.40</b>	<b>0.883</b>	<b>91.20</b>	<b>0.875</b>
w/o	84.60	0.843	80.20	0.801	85.50	0.854	83.43	0.832
CTL	<b>90.30</b>	<b>0.901</b>	<b>96.50</b>	<b>0.964</b>	<b>89.40</b>	<b>0.893</b>	<b>92.06</b>	<b>0.919</b>
w/o	83.40	0.822	94.20	0.942	80.10	0.792	85.90	0.852
Ours	<b>89.50</b>	<b>0.892</b>	<b>96.80</b>	<b>0.967</b>	<b>88.10</b>	<b>0.879</b>	<b>91.47</b>	<b>0.912</b>
w/o	77.00	0.750	70.10	0.698	78.50	0.770	75.20	0.739

Table 12: A comparison of the retrieval process with and without incorporating the similarity score between the query and sentence.

CLL-RetICL	0.3		0.5		0.7		1.0	
LLM	ACC	F1	ACC	F1	ACC	F1	ACC	F1
CoLA								
Gemini	82.92	0.791	83.21	0.794	83.01	0.793	<b>83.60</b>	<b>0.801</b>
Llama	77.60	0.746	<b>77.68</b>	<b>0.757</b>	76.53	0.737	77.66	0.742
Mistral	85.43	0.818	85.33	0.817	<b>85.71</b>	0.822	85.52	<b>0.828</b>
CARER								
Gemini	69.10	0.636	69.85	0.640	70.10	0.640	<b>72.65</b>	<b>0.669</b>
Llama	68.50	0.625	68.65	0.625	<b>69.95</b>	<b>0.635</b>	69.15	<b>0.635</b>
Mistral	72.20	0.665	72.20	0.671	71.65	0.656	<b>76.85</b>	<b>0.717</b>

Table 13: A comparison of classification accuracy (%) and F1 score to assess the impact of various trade-off hyperparameters in CLL-RetICL strategy.

PLA	0.3		0.5		0.7		1.0	
LLM	ACC	F1	ACC	F1	ACC	F1	ACC	F1
CoLA								
Gemini	<b>83.57</b>	<b>0.799</b>	82.42	0.784	83.17	0.794	83.31	0.798
Llama	73.12	0.661	73.03	0.649	<b>74.29</b>	<b>0.681</b>	73.53	0.656
Mistral	<b>85.71</b>	0.821	85.42	0.817	85.23	0.813	85.31	<b>0.832</b>
CARER								
Gemini	66.60	0.602	<b>66.70</b>	<b>0.603</b>	65.85	0.598	65.75	0.598
Llama	<b>65.15</b>	<b>0.585</b>	64.90	<b>0.585</b>	62.50	0.562	61.65	0.556
Mistral	<b>65.80</b>	0.595	65.15	0.579	62.50	0.558	65.55	<b>0.596</b>

Table 14: A comparison of classification accuracy (%) and F1 score to assess the impact of various trade-off hyperparameters in PLA strategy.

NLP	0.3		0.5		0.7		1.0	
LLM	ACC	F1	ACC	F1	ACC	F1	ACC	F1
CoLA								
Gemini	83.51	0.798	83.31	<b>0.802</b>	<b>83.69</b>	0.801	82.45	0.791
Llama	<b>64.11</b>	0.553	62.73	0.539	63.10	0.542	64.05	<b>0.579</b>
Mistral	85.52	0.820	<b>85.53</b>	0.820	85.33	0.817	85.04	<b>0.823</b>
CARER								
Gemini	65.75	0.594	65.25	0.586	66.05	0.595	<b>67.35</b>	<b>0.619</b>
Llama	<b>64.60</b>	<b>0.586</b>	64.10	0.581	63.05	0.567	64.40	0.583
Mistral	<b>70.40</b>	0.636	68.65	0.625	69.70	0.632	70.00	<b>0.644</b>

Table 15: A comparison of classification accuracy (%) and F1 score to assess the impact of various trade-off hyperparameters in NLP strategy.

TREC	Abbr.		Full	
	ACC	F1	ACC	F1
Llama				
Nwaykshot	<b>57.40</b>	<b>0.591</b>	56.60	0.577
CLL-RetICL	56.40	0.603	<b>57.60</b>	<b>0.603</b>

Table 16: A comparison of classification accuracy (%) and F1 score to evaluate the impact of abbreviated labels versus full labels on the TREC dataset.

SST2	Positive/Negative		Great/Terrible		Good/Bad	
	ACC	F1	ACC	F1	ACC	F1
Llama						
Nwaykshot	90.82	0.908	92.48	0.927	91.98	0.923
CLL-RetICL	<b>95.06</b>	<b>0.951</b>	<b>93.90</b>	<b>0.939</b>	<b>94.23</b>	<b>0.944</b>

Table 17: A comparison of classification accuracy (%) and F1 score to evaluate the impact of synonym labels on the SST2 dataset. The synonym pairs used in this study are drawn from previously published work (Pawar et al., 2024).