

# Beyond Task-Level Context: Class-Conditional Context Vectors for Implicit In-Context Learning

Jianxin Zhang<sup>1\*</sup>, Yilu Hu<sup>1\*</sup>, Teng Liu<sup>1</sup>, Pei Guo<sup>1</sup>, Juntao Li<sup>1,2†</sup>

<sup>1</sup>School of Computer Science and Technology, Soochow University

<sup>2</sup>Key Laboratory of Data Intelligence and Advanced Computing, Soochow University  
{jxzhang2020, ylhuhy1}@stu.suda.edu.cn  
ljt@suda.edu.cn

## Abstract

Implicit In-Context Learning compresses demonstration examples into a single context vector and injects it into the model’s activation space, achieving few-shot-level performance at near zero-shot inference cost. However, existing approaches typically aggregate demonstrations from all classes into a shared, task-level context vector, capturing global task information but without explicitly preserving fine-grained, class-conditional semantic distinctions. In this work, we propose **Class-Conditional Context Vectors (C<sup>3</sup>V)**, a simple yet effective extension to implicit in-context learning that explicitly models class-specific contextual information by constructing separate context vectors for each class, without relying on explicit prompts. These class-conditional context vectors are additively injected into the model’s residual streams in a single forward pass, enabling contextual contributions to be modulated in a class-aware manner while keeping the backbone frozen. We evaluate C<sup>3</sup>V on multiple text classification benchmarks across several families of large language models. Experimental results demonstrate that C<sup>3</sup>V consistently outperforms task-level context vector baselines, and achieves higher average accuracy than conventional few-shot learning on most models.

## 1 Introduction

In-context learning (ICL) enables large language models (LLMs) to adapt to new tasks by conditioning on a small set of labeled demonstrations without updating model parameters (Honda and Oka, 2025; Saglam et al., 2025; Zhang et al., 2025b). However, explicitly concatenating demonstrations with each query incurs non-trivial inference overhead (Qin et al., 2025), motivating recent work on implicit in-context learning, which compresses demonstrations into latent context vectors

and injects them directly into the model’s activation space (Li et al., 2025; Todd et al., 2024; Liu et al., 2024; Hu et al., 2025; Wang et al., 2025a; Hendel et al., 2023; Rimsky et al., 2024). This formulation achieves performance comparable to few-shot learning while retaining near zero-shot inference efficiency.

Among existing implicit in-context learning approaches (Hu et al., 2025; Li et al., 2025; Todd et al., 2024; Liu et al., 2024; Wang et al., 2025a; Hendel et al., 2023; Rimsky et al., 2024), I2CL (Li et al., 2025) adopts a simple and effective formulation that blends context vectors with query activations via linear combination. In doing so, I2CL constructs a task-level context vector by aggregating demonstrations from all classes into a single shared representation. This design is effective for capturing global task information and achieves strong empirical performance across a range of benchmarks while retaining high inference efficiency (Li et al., 2025). However, by treating all demonstrations uniformly, it does not explicitly preserve class-specific structure. In classification tasks, demonstrations from different classes often exhibit distinct semantic characteristics, as shown in Section 4.3 and Appendix D, raising the question of whether a single task-level context vector is sufficient to represent such heterogeneous information. Given its simplicity and empirical success, we take I2CL as a natural foundation and investigate how its task-level context vectorization can be extended to better model class-specific structure.

Motivated by this consideration, we propose a simple class-conditional extension to implicit in-context learning. Specifically, we group demonstrations by class, construct a separate context vector for each label, and jointly inject all class-conditional context vectors into the model’s residual streams via learnable linear coefficients, which are trained using a small labeled support set, enabling class-aware conditioning without multiple

\*Equal contribution

†Corresponding author

inference passes or explicit prompting during inference.

We refer to this approach as Class-Conditional Context Vectors ( $C^3V$ ).  $C^3V$  provides a more fine-grained formulation than task-level context vectorization by retaining class-specific contextual information in the latent space. Experiments on multiple text classification benchmarks and 6 LLMs show that this class-conditional formulation improves classification accuracy over task-level context vector baselines in most cases. We further analyze the resulting representations to better understand how class-conditional context vectors influence latent activation structure. Our main contributions are summarized as follows:

- We highlight a potential limitation of task-level context vectorization in implicit in-context learning, which aggregates demonstrations from different classes into a single shared representation without explicitly modeling class-specific structure.
- We propose  $C^3V$ , a lightweight and general class-conditional extension that introduces class-conditional context vectors at the cost of a modest increase in inference time.
- We demonstrate the effectiveness of  $C^3V$  on multiple classification benchmarks and provide representation-level analyses to better understand its behavior.

## 2 Related Work

**In-Context Learning.** In-context learning enables LLMs to perform new tasks by conditioning on a small set of demonstrations provided in the input, without parameter updates (Honda and Oka, 2025). Prior work has studied various aspects of ICL, including prompt design, demonstration selection, and ordering (Liu et al., 2025; Sun et al., 2025; Kothapalli et al., 2025). Despite its effectiveness, explicit demonstration concatenation incurs non-trivial inference overhead, motivating more efficient alternatives (Qin et al., 2025).

**Implicit In-Context Learning and Context Vectors.** To reduce inference cost, recent approaches propose implicit in-context learning methods that compress demonstration information into latent context vectors and inject them directly into the model’s internal activations (Zhang et al., 2025a; Wang et al., 2025b; Nguyen et al., 2025; Yousefpour et al., 2025; Su et al., 2025; Konen

et al., 2024; Scalena et al., 2024; Zhao et al., 2025; Yang et al., 2025; Stolfo et al., 2025; Li et al., 2025; Todd et al., 2024; Liu et al., 2024; Hu et al., 2025; Wang et al., 2025a; Hendel et al., 2023; Rimsky et al., 2024). These methods achieve strong performance while avoiding explicit prompt-based conditioning. However, most existing context-vector-based approaches operate at the task level by aggregating demonstrations into a single shared context vector (Li et al., 2025; Todd et al., 2024; Liu et al., 2024; Hu et al., 2025; Wang et al., 2025a; Hendel et al., 2023; Rimsky et al., 2024). While effective for capturing global task structure, such aggregation may obscure fine-grained, class-conditional semantic information. In contrast to these approaches, our method introduces class-conditional structure directly into implicit context representations, without requiring explicit prompts.

**Class-Aware Conditioning and Representation Analysis.** Prior work has shown that incorporating label or class information into model inputs can play an important role in guiding model behavior (Yao et al., 2024; Wang et al., 2023), particularly in classification settings. These approaches suggest that explicitly exposing models to class-related cues can help align internal representations with task-relevant semantics.

Our work builds on this insight by introducing class-conditional structure directly into implicit context representations. We further analyze the representational effects of class-conditional context vectors through intra-class and inter-class distance measurements, complementing prior representation analyses of implicit in-context learning.

## 3 Method

We propose  $C^3V$ , a simple extension to implicit ICL that captures fine-grained class-level semantics. As illustrated in Figure 1, instead of constructing a single task-level context vector,  $C^3V$  builds separate context vectors for each class from labeled demonstrations and jointly injects them into the model via linear interventions. Algorithm 1 provides the pseudocode of the  $C^3V$  algorithm.

### 3.1 Preliminaries: I2CL

We briefly review I2CL, which forms the basis of our approach. Given a classification task with a small set of labeled demonstrations

$$\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N, \quad (1)$$

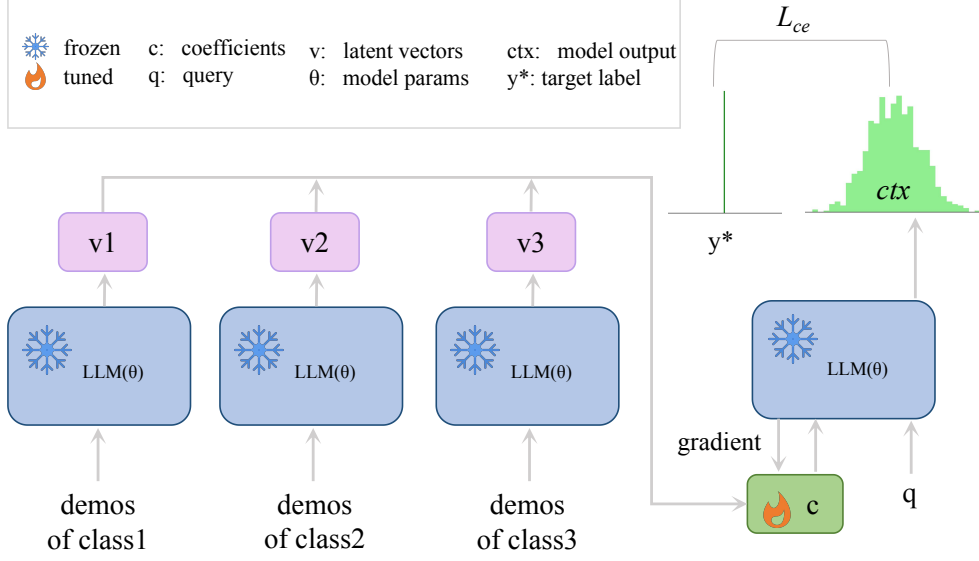


Figure 1: C<sup>3</sup>V overview. Demonstrations are grouped by class and encoded by a frozen model into class-conditional latent vectors. These vectors are jointly injected with the query using learnable coefficients in a single forward pass, which are optimized with cross-entropy while the backbone remains frozen.

For each demonstration  $(x_i, y_i)$ , I2CL extracts attention and MLP outputs at selected token positions (e.g., the end-of-sequence token  $e$ ) across all layers  $l = 1, \dots, L$ , denoted as  $\{a_{i,l}^e, m_{i,l}^e\}_{l=1}^L$ . These activations are averaged across all demonstrations to obtain a task-level context vector  $\{\bar{a}_l^e, \bar{m}_l^e\}_{l=1}^L$ , which is linearly injected into the residual stream during inference to bias the model toward task-relevant directions while preserving near zero-shot inference efficiency.

### 3.2 Class-Conditional Context Vector Construction

To capture class-specific semantics, we construct *class-conditional context vectors* by grouping demonstrations according to their labels. Let  $\mathcal{C}$  denote the set of class labels, and  $\mathcal{D}_c = \{(x_i, y_i) \in \mathcal{D} \mid y_i = c\}$  denote the demonstration subset for class  $c$ . For each  $(x_i, y_i) \in \mathcal{D}_c$ , we extract the same attention and MLP activations as in I2CL:

$$d_i = \{a_{i,l}^e, m_{i,l}^e\}_{l=1}^L. \quad (2)$$

Aggregating these activations within each class yields the class-conditional context vector

$$\bar{d}_c = \{\bar{a}_{c,l}^e, \bar{m}_{c,l}^e\}_{l=1}^L, \quad (3)$$

where

$$\bar{a}_{c,l}^e = \frac{1}{|\mathcal{D}_c|} \sum_{i \in \mathcal{D}_c} a_{i,l}^e,$$

$$\bar{m}_{c,l}^e = \frac{1}{|\mathcal{D}_c|} \sum_{i \in \mathcal{D}_c} m_{i,l}^e.$$

### 3.3 Class-Conditional Context Injection

Let  $\mathbf{r}_l^t$  denote the residual representation at layer  $l$  and token position  $t$ . We inject all class-conditional context vectors into the residual stream via a linear intervention:

$$\mathbf{r}_l^t = \mathbf{r}_{l-1}^t + \sum_{c \in \mathcal{C}} \lambda_{c,l}^a \bar{a}_{c,l}^e + \beta_l^a a_l^t + \sum_{c \in \mathcal{C}} \lambda_{c,l}^m \bar{m}_{c,l}^e + \beta_l^m m_l^t. \quad (4)$$

Here,  $\lambda_{c,l}^a$  and  $\lambda_{c,l}^m$  control the contribution of class-conditional attention and MLP components, while  $\beta_l^a$  and  $\beta_l^m$  preserve the original residual pathways of the backbone model. All class-conditional vectors are jointly injected within a single forward pass, resulting in a superposition of class-level semantic biases in the residual stream.

Rather than treating these coefficients as fixed hyperparameters, we learn

$$\mathbf{c} = \{\lambda_{c,l}^a, \lambda_{c,l}^m, \beta_l^a, \beta_l^m\}_{c \in \mathcal{C}, l=1}^L \quad (5)$$

via gradient-based optimization on the same few-shot demonstrations used to construct the class-conditional context vectors. In I2CL,  $\lambda$  is initialized to 0.1. Since our formulation injects class-conditional context vectors additively over classes, we scale the initialization by the number of classes and set it to  $0.1/|\mathcal{C}|$  to keep the overall injection magnitude comparable. Following I2CL, we initialize the residual scaling coefficients  $\beta$  to 1. All

coefficients are optimized by minimizing the negative log-likelihood of the ground-truth labels:

$$\mathcal{L} = -\frac{1}{|\mathcal{D}|} \sum_{(x,y) \in \mathcal{D}} \log P(y | x, \{\bar{d}_c\}_{c \in \mathcal{C}}, \mathbf{c}), \quad (6)$$

while keeping the backbone model parameters fixed.

### 3.4 Inference and Prediction

During inference, prediction is performed with a single forward pass in which all class-conditional context vectors are jointly injected into the model. Given a query input  $x$ , the model produces a unified hidden representation that integrates class-specific semantic information through the learned coefficients. Class scores are computed from the predicted probabilities associated with each candidate class label, and the predicted class is selected as

$$\hat{y} = \arg \max_{c \in \mathcal{C}} P(y_c | x, \{\bar{d}_c\}_{c \in \mathcal{C}}, \mathbf{c}), \quad (7)$$

where  $y_c$  denotes the label token corresponding to class  $c$ .

## 4 Experiments

### 4.1 Experimental Setup

**Implementation Details.** We outline the experimental configuration and setup used for all experiments unless otherwise specified. We construct the class-conditional context vectors by sampling five demonstration examples per class for each task. Demonstration examples are sampled from the training set and fixed across methods to ensure fair comparison. Input sequences are constructed using simple yet effective manually designed templates, which are summarized in Table 11 in the Appendix H. Each task is evaluated on a test set consisting of 500 data points, consistent with the baseline setup. For the training process, we employ the AdamW (Loshchilov and Hutter, 2019) optimizer with a learning rate of  $1 \times 10^{-2}$  and weight decay of  $1 \times 10^{-3}$ . The learning rate is adjusted using a cosine learning rate scheduler. The training is performed over 100 epochs with a batch size of 2. To ensure statistical reliability, each experiment is independently run five times under the same configuration with different random seeds, each corresponding to a different set of sampled demonstrations. All experiments were conducted on NVIDIA A5000 GPUs.

**Datasets.** We evaluate our method and baseline approaches on nine widely used text classification benchmarks. Following prior work (e.g., I2CL), we adopt the same benchmark suite to ensure fair comparison. Specifically, the evaluated datasets include SST-2 and SST-5 (Socher et al., 2013) for binary and fine-grained sentiment analysis, MR (Pang and Lee, 2005) for movie review sentiment classification, TREC (Voorhees and Tice, 2000) for question type classification, AGNews and DBpedia (Zhang et al., 2015) for topic and multi-class document classification, Subj (Pang and Lee, 2004) for subjectivity detection, EmoC (Chatterjee et al., 2019) for emotion classification, and HateSpeech18 (de Gibert et al., 2018) for hate speech detection.

**Models.** We select three famous model families: LLaMA, GPT and Qwen. Specifically, we choose LLaMA2-7B (Touvron et al., 2023), LLaMA3-8B (AI@Meta, 2024), LLaMA3.2-1B (AI@Meta, 2024), GPT-J-6B (Wang and Komatsuzaki, 2021), GPT2-XL (Radford et al., 2019) and Qwen3-8B (Team, 2025). The models were chosen based on their alignment with the I2CL and the availability of our computational resources. We only provide a detailed report of the results for LLaMA2-7B, with the results for other LLMs available in Appendix C.

**Baselines.** We compare our method against the following baselines: (1) Label Anchor (Wang et al., 2023), which uses anchor tokens in the input for task adaptation without fine-tuning model parameters. (2) Task Vector (Hendel et al., 2023), which injects a task-specific vector derived from demonstration examples into the model’s optimal replacement layer to guide task-specific behavior. (3) ICV (Liu et al., 2024), which leverages positive and negative demonstration pairs to generate the steering vector, mainly for open-ended generation tasks. (4) I2CL (Li et al., 2025), which condenses demonstration examples into a context vector and uses linear operations to fuse information. In addition to these, several tuning-based approaches are also included, such as (5) AutoCompressors (Chevalier et al., 2023), which finetunes the model to compress knowledge representations. (6) ICAE (Ge et al., 2024), which uses a LoRA-adapted encoder to encode context into memory slots. (7) CEPE (Yen et al., 2024), which utilizes demonstrations from the encoder context to improve performance. (8) Prompt-tuning (Lester et al., 2021), which adds learnable prompt tokens to adapt the model without modifying its parameters. (9) LoRA (Hu et al.,

Method	SST-2	SST-5	TREC	AGNews	Subj	HateS	DBPedia	EmoC	MR	Avg.	# cached P.
Zero-shot	83.00	27.00	50.00	70.20	51.40	54.20	72.00	41.80	73.60	58.13	0
Few-shot	94.44 <sub>1.44</sub>	41.72 <sub>3.68</sub>	77.32 <sub>4.41</sub>	85.68 <sub>2.00</sub>	52.56 <sub>3.09</sub>	70.24 <sub>5.80</sub>	96.64 <sub>0.48</sub>	75.48 <sub>1.63</sub>	93.24 <sub>0.50</sub>	76.37	2MDL
Label-anchor	83.32 <sub>5.95</sub>	27.68 <sub>4.21</sub>	77.48 <sub>3.49</sub>	83.72 <sub>1.04</sub>	53.00 <sub>2.95</sub>	64.52 <sub>8.09</sub>	81.40 <sub>3.67</sub>	59.12 <sub>10.60</sub>	84.40 <sub>5.89</sub>	68.29	2(M/K)DL
Task-vector	81.44 <sub>4.73</sub>	25.96 <sub>0.59</sub>	65.68 <sub>1.93</sub>	79.68 <sub>4.07</sub>	58.56 <sub>4.91</sub>	67.68 <sub>3.70</sub>	89.48 <sub>2.58</sub>	44.64 <sub>3.53</sub>	82.32 <sub>5.37</sub>	66.16	D
ICV	86.28 <sub>0.55</sub>	33.48 <sub>0.65</sub>	63.84 <sub>0.15</sub>	72.40 <sub>0.37</sub>	56.56 <sub>0.70</sub>	60.56 <sub>1.50</sub>	73.64 <sub>0.88</sub>	49.16 <sub>1.24</sub>	84.04 <sub>1.10</sub>	64.44	DL
I2CL	<b>87.68</b> <sub>2.47</sub>	39.12 <sub>2.69</sub>	78.56 <sub>5.32</sub>	85.48 <sub>1.16</sub>	<b>73.84</b> <sub>3.84</sub>	69.88 <sub>5.67</sub>	90.16 <sub>1.86</sub>	63.72 <sub>1.37</sub>	87.68 <sub>2.26</sub>	75.12	2DL
C <sup>3</sup> V	87.06 <sub>2.01</sub>	<b>44.64</b> <sub>2.65</sub>	<b>78.86</b> <sub>4.56</sub>	<b>85.52</b> <sub>2.49</sub>	70.40 <sub>7.06</sub>	<b>74.80</b> <sub>1.64</sub>	<b>93.16</b> <sub>0.97</sub>	<b>63.94</b> <sub>2.97</sub>	<b>88.48</b> <sub>3.12</sub>	<b>76.32</b>	2 C DL
AutoComp.	92.44 <sub>3.29</sub>	25.80 <sub>4.80</sub>	62.52 <sub>9.34</sub>	86.36 <sub>1.03</sub>	60.16 <sub>0.32</sub>	53.20 <sub>6.10</sub>	92.68 <sub>2.86</sub>	29.56 <sub>5.07</sub>	82.76 <sub>7.34</sub>	63.94	2(M/K)DL
ICAE	91.64 <sub>1.69</sub>	38.80 <sub>1.56</sub>	50.92 <sub>8.38</sub>	80.48 <sub>2.35</sub>	50.52 <sub>9.17</sub>	65.48 <sub>7.18</sub>	62.08 <sub>1.86</sub>	54.04 <sub>4.69</sub>	89.48 <sub>1.45</sub>	64.83	2(M/K)DL
CEPE	74.28 <sub>3.90</sub>	36.20 <sub>0.56</sub>	55.48 <sub>3.42</sub>	78.00 <sub>3.49</sub>	59.12 <sub>1.60</sub>	61.72 <sub>5.26</sub>	87.24 <sub>1.20</sub>	42.28 <sub>3.31</sub>	82.36 <sub>1.61</sub>	64.08	Md <sub>enc</sub>

Table 1: Comparison of methods on various classification benchmarks using LLaMA2-7B. The best results are highlighted in **bold**. In addition to evaluating the memory consumption, we also conduct an analysis of the cached parameters. Specifically,  $M$ ,  $D$ , and  $L$  denote the number of demonstration tokens, the model hidden dimensionality, and the number of architectural layers, respectively. Furthermore,  $1/K$  represents the compression ratio of the respective context-compression technique,  $|C|$  denotes the number of classes, and  $d_{enc}$  denotes the hidden dimension of the auxiliary encoder (when applicable). A detailed explanation of “# cached P.” can be found in Appendix A. Due to space constraints, the signs for the standard deviation were omitted, but it should be expressed as mean  $\pm$  standard deviation.

Method	# trainable params. (K)	SST-2	SST-5	TREC	AGNews	Subj	HateS	DBPedia	EmoC	MR	Avg.
Prompt-tuning	4.10	56.24 <sub>6.99</sub>	24.24 <sub>2.96</sub>	55.20 <sub>4.14</sub>	78.00 <sub>7.60</sub>	57.40 <sub>4.93</sub>	49.56 <sub>6.96</sub>	74.40 <sub>6.43</sub>	35.08 <sub>5.29</sub>	54.32 <sub>1.76</sub>	54.94
LoRA (r=8)	4194.30	84.80 <sub>6.59</sub>	39.87 <sub>4.33</sub>	75.97 <sub>10.77</sub>	83.80 <sub>2.32</sub>	70.47 <sub>10.68</sub>	75.32 <sub>2.88</sub>	91.40 <sub>3.54</sub>	53.67 <sub>16.27</sub>	83.07 <sub>0.25</sub>	73.15
LoRA (r=1)	524.29	82.44 <sub>0.77</sub>	38.88 <sub>3.47</sub>	73.08 <sub>6.56</sub>	83.32 <sub>2.76</sub>	60.68 <sub>0.68</sub>	63.32 <sub>2.83</sub>	84.52 <sub>4.14</sub>	57.88 <sub>2.66</sub>	87.96 <sub>1.68</sub>	70.23
IA3	262.14	89.40 <sub>2.08</sub>	46.93 <sub>0.81</sub>	75.41 <sub>4.94</sub>	84.43 <sub>1.45</sub>	56.67 <sub>3.07</sub>	62.54 <sub>5.58</sub>	93.91 <sub>0.49</sub>	59.75 <sub>3.67</sub>	88.00 <sub>1.88</sub>	73.00
C <sup>3</sup> V	0.064( C  + 1)	87.06 <sub>2.01</sub>	44.64 <sub>2.65</sub>	78.86 <sub>4.56</sub>	85.52 <sub>2.49</sub>	70.40 <sub>7.06</sub>	74.80 <sub>1.64</sub>	93.16 <sub>0.97</sub>	63.94 <sub>2.97</sub>	88.48 <sub>3.12</sub>	76.32

Table 2: Comparison of different PEFT methods on various datasets. Due to space constraints, the signs for the standard deviation were omitted, but it should be expressed as mean  $\pm$  standard deviation. Here,  $|C|$  denotes the number of class labels in the corresponding dataset.

2022) and (10) IA3 (Liu et al., 2022), which are PEFT methods from HuggingFace PEFT library<sup>1</sup>. All baseline configurations refer to the settings in I2CL.

**Evaluation Metrics.** We report classification accuracy as the primary evaluation metric. For multi-class datasets, we report macro-averaged accuracy. All results are averaged over five random seeds. All results are reported in percentage (%).

## 4.2 Main Results

**C<sup>3</sup>V Achieves Few-shot Level Accuracy with Moderate Parameter Caching.** Table 1 reports results on LLaMA2-7B. C<sup>3</sup>V achieves performance comparable to few-shot inference while requiring only moderate additional cached parameters. It attains 76.32% average accuracy, outperforming the zero-shot baseline by 18.19 points. **Notably, as shown in Appendix C, C<sup>3</sup>V even surpasses few-shot performance on average on the remaining 5 LLMs.**

**Comparison with Methods without Explicit Demonstrations.** Label Anchor, Task Vector, ICV, and I2CL perform inference without relying on explicit demonstration examples. As shown in Table 1, C<sup>3</sup>V achieves the best overall performance among these approaches, reaching the highest average accuracy. In particular, C<sup>3</sup>V outperforms the strongest baseline I2CL by 1.2 points on average, with gains on challenging datasets such as SST-5, HateS, and DBPedia. These results suggest that explicitly modeling class-conditional contextual information can offer a more effective alternative to task-level context representations, while avoiding any reliance on demonstration examples during inference.

**Comparison with Methods for Prompt Compression during Inference.** We compare C<sup>3</sup>V with representative prompt compression methods, including AutoCompressors (Chevalier et al., 2023), ICAE (Ge et al., 2024), and CEPE (Yen et al., 2024). As shown in Table 1, these approaches yield mixed improvements over zero-shot baselines and exhibit inconsistent performance across datasets.

<sup>1</sup><https://github.com/huggingface/peft>

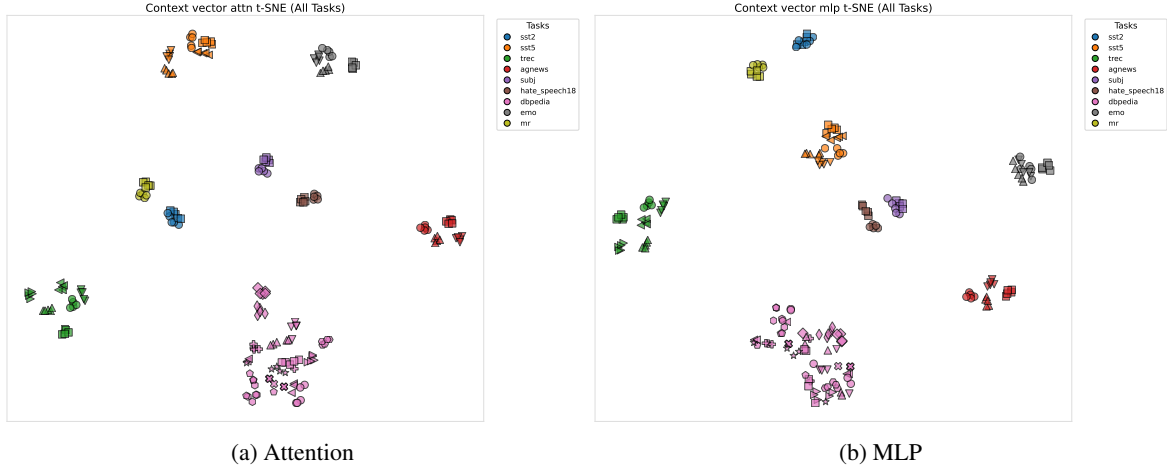


Figure 2: t-SNE visualization of class-conditional latent vectors on LLaMA-2-7B. Left: attention-based latent vectors; right: MLP-based latent vectors. Points are colored by dataset and shaped by class label. Latent vectors from five random seeds are included for each class.

In contrast,  $C^3V$  does not rely on a pre-trained universal compression model; instead, it constructs class-conditional context vectors and learns task-specific injection coefficients, enabling more targeted adaptation. We note that prompt compression methods primarily aim to reduce input length for long-context inference, making them complementary to  $C^3V$ , which focuses on improving task adaptation with moderate parameter caching.

### Comparison with Test-time PEFT Methods.

We compare  $C^3V$  with widely used test-time PEFT methods on LLaMA2-7B, including Prompt Tuning (Lester et al., 2021), LoRA (Hu et al., 2022), and IA3 (Liu et al., 2022). As shown in Table 2, Prompt Tuning underperforms the zero-shot baseline under data-scarce conditions, while LoRA and IA3 achieve stronger performance but require substantially more trainable parameters. In contrast,  $C^3V$  attains competitive adaptation performance with a number of learnable parameters that scales linearly with the number of classes.

## 4.3 Representation Analysis

To better understand the representational effects of class-conditional context vectors, we analyze the structure and separability of the latent representations using both visualization and distance-based metrics based on cosine similarity.

### 4.3.1 Visualization of Class-Conditional Latent Vectors

For each dataset, we construct five sets of class-conditional latent vectors, each obtained from a run

with a different random seed, corresponding to different sampled demonstrations. For visualization, we concatenate the layer-wise attention or MLP vectors of each class-conditional context into a single high-dimensional representation and project all such representations into two dimensions using t-SNE. As shown in Figure 2, for LLaMA2-7B, the visualization reveals a hierarchical organization: latent vectors from different datasets form well-separated global clusters, while vectors corresponding to different class labels further organize into distinct sub-clusters within each dataset. Despite this overall hierarchical structure, the visualization also exhibits cases with less clear class separation. For example, in Figure 2a, the two classes in the SST-2 dataset are relatively difficult to distinguish, while in Figure 2b, classes in the Emoc dataset (e.g., circles and inverted triangles) show substantial overlap. **Class-conditional latent vectors obtained from other models are presented in Appendix D.**

These patterns suggest that class-conditional context vectors capture generally stable and discriminative class-level semantics that are robust to randomness in demonstration sampling.

### 4.3.2 Inter- and Intra-Class Distance Analysis

We further quantify representation separability by computing inter- and intra-class distances within each dataset. Intra-class distance measures the variation among latent vectors of the same class across different random seeds, while inter-class distance measures the separation between latent vectors of different classes. All distances are computed us-

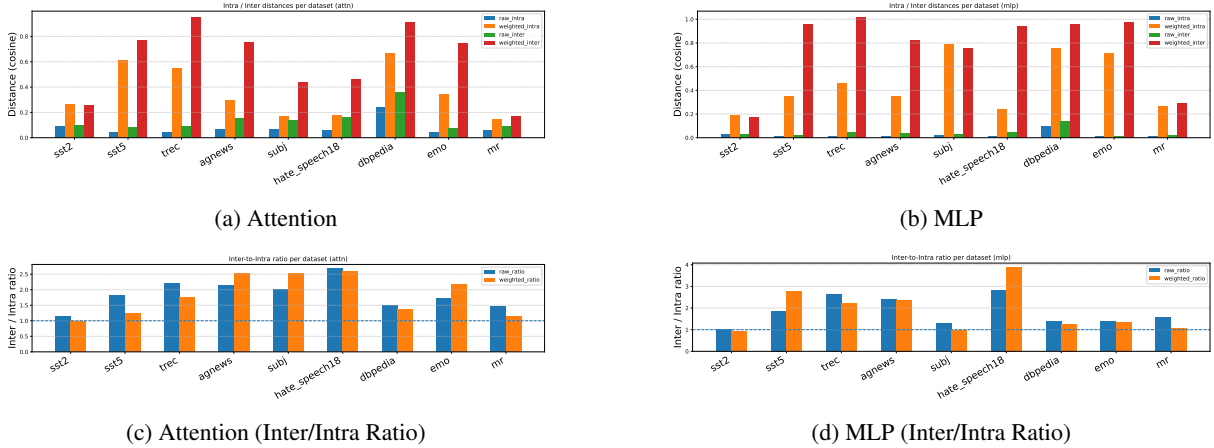


Figure 3: Inter- and intra-class distance analysis of latent representations across datasets on LLaMA-2-7B. (a,b) Inter- and intra-class distances for attention-based and MLP-based representations, where each panel includes distances computed both without coefficient scaling (bars labeled with the raw- prefix in the legend) and after multiplication by trained coefficients (bars labeled with the weight- prefix). (c,d) Corresponding inter/intra ratios for attention and MLP representations, comparing the settings without (raw-) and with (weight-) trained coefficient scaling.

ing cosine distance. We perform this analysis for both attention- and MLP-based representations, before and after modulation with the learned coefficients. The detailed formulations of the distance metric, as well as the definitions of intra-class and inter-class distances and the inter/intra ratio, are provided in Appendix B. As shown in Figure 3, we report the inter- and intra-class distance analysis for LLaMA2-7B; **corresponding results for other models are provided in Appendix E.**

Across most datasets, both before and after modulation with the learned coefficients, inter-class distances are larger than intra-class distances (Figures 3a–3b), indicating inherent class separability. After modulation, both inter-class and intra-class distances generally increase, indicating that, under cosine distance, pairwise distances between latent representations exhibit an overall consistent increase. To assess relative separability, we compute the inter/intra ratio (Figures 3c–3d), where values greater than 1 indicate favorable class discrimination. In most cases, the ratio remains above 1. Notably, for a small number of datasets (e.g., SST-2 and Subj under MLP-based modulated representations), the ratio falls below 1, which aligns with cases where C<sup>3</sup>V underperforms I2CL in the main results. These analyses suggest that class-conditional context vectors induce relatively stable and discriminative representations, while their effectiveness varies across different datasets.

Finally, we visualize the learned layer-wise coefficients to qualitatively analyze how context injection

is modulated across layers (see Appendix F).

#### 4.4 Ablation

We analyze the impact of each module in our framework using LLaMA2-7B (unless otherwise specified), with the reported results averaged across five random seeds and all nine datasets. **Due to space constraints, detailed ablation studies on target modules, target layers, injection positions, extraction positions, and context vector aggregation strategies are deferred to Appendix G.**

**Training and Inference Efficiency.** We compare the training and inference efficiency of C<sup>3</sup>V with relevant baselines on LLaMA3.2-1B, as summarized in Table 3. While both C<sup>3</sup>V and I2CL introduce learnable scaling coefficients, C<sup>3</sup>V incurs higher training time due to the introduction of additional parameters that scale with the number of classes. This effect is particularly pronounced on datasets with more classes, such as DBpedia, SST-5, TREC, and AGNews, where modeling class-conditional contextual information leads to a noticeable increase in training cost. On average, the training time increases from 104.60 seconds for I2CL to 146.93 seconds for C<sup>3</sup>V. At inference time, zero-shot and cached few-shot baselines achieve the lowest latency, while uncached few-shot inference is substantially more expensive. Compared to I2CL, C<sup>3</sup>V introduces moderate additional inference cost on most datasets, resulting in an average inference time of 22.11 seconds versus 13.60 sec-

Method	SST-2	SST-5	TREC	AGNews	Subj	HateS	DBPedia	EmoC	MR	Avg.
<i>Training Time</i>										
I2CL	55.27 <sub>1.74</sub>	146.17 <sub>1.01</sub>	177.07 <sub>1.24</sub>	134.32 <sub>4.15</sub>	66.03 <sub>0.99</sub>	64.47 <sub>0.85</sub>	99.93 <sub>5.46</sub>	134.59 <sub>0.53</sub>	63.59 <sub>2.44</sub>	104.60
C <sup>3</sup> V	65.49 <sub>1.91</sub>	201.37 <sub>4.79</sub>	265.27 <sub>2.92</sub>	175.95 <sub>26.25</sub>	66.26 <sub>3.10</sub>	63.13 <sub>3.01</sub>	273.45 <sub>40.80</sub>	146.19 <sub>4.03</sub>	65.28 <sub>1.19</sub>	146.93
<i>Inference Time</i>										
ZS	6.42	6.43	4.66	9.26	6.84	6.70	12.11	5.54	6.61	7.17
FS (Cache)	6.87 <sub>0.05</sub>	7.53 <sub>0.04</sub>	4.99 <sub>0.01</sub>	11.24 <sub>0.11</sub>	7.48 <sub>0.07</sub>	7.24 <sub>0.06</sub>	14.92 <sub>0.21</sub>	6.32 <sub>0.05</sub>	7.21 <sub>0.05</sub>	8.20
FS (No Cache)	26.09 <sub>3.46</sub>	85.46 <sub>2.70</sub>	64.67 <sub>2.53</sub>	141.98 <sub>6.18</sub>	45.97 <sub>6.48</sub>	34.96 <sub>5.13</sub>	131.87 <sub>14.67</sub>	58.09 <sub>4.56</sub>	36.93 <sub>2.30</sub>	69.56
I2CL	11.96 <sub>1.27</sub>	12.16 <sub>1.33</sub>	9.75 <sub>0.14</sub>	15.95 <sub>3.45</sub>	12.08 <sub>0.10</sub>	11.75 <sub>0.15</sub>	24.86 <sub>1.87</sub>	11.89 <sub>0.75</sub>	12.03 <sub>0.55</sub>	13.60
C <sup>3</sup> V	14.02 <sub>1.10</sub>	19.69 <sub>1.38</sub>	19.05 <sub>1.63</sub>	22.60 <sub>9.58</sub>	11.51 <sub>2.66</sub>	8.58 <sub>0.03</sub>	86.22 <sub>22.11</sub>	8.98 <sub>0.06</sub>	8.37 <sub>0.11</sub>	22.11

Table 3: Training and Inference Time (seconds) on LLaMA3.2-1B. Values are reported as mean with standard deviation shown in subscript. Due to space constraints, the  $\pm$  symbol is omitted.

Direction	Accuracy
Zero-shot	58.13
+ Noise Vector (Repl.)	69.73
+ Class-Conditional (No-Train, init = $0.1/ \mathcal{C} $ )	60.98
+ Class-Conditional (Train, init = 0.1)	71.94
Train, init = $0.1/ \mathcal{C} $ (Ours)	76.32

Table 4: Noise Injection, Training, and Coefficient Initialization

onds for I2CL. In particular, datasets with a larger number of classes, such as DBPedia (14 classes), exhibit a more pronounced increase in inference time. These results highlight a clear trade-off between efficiency and performance: C<sup>3</sup>V requires modestly higher training and inference time than task-level context vector methods, but achieves improved classification performance by explicitly modeling class-conditional structure, while remaining substantially more efficient than uncached few-shot inference.

### Scaling with the Number of Demonstrations.

We analyze how model performance scales with the number of demonstrations per class. As shown in Figure 4, C<sup>3</sup>V exhibits a clear and consistent improvement as more demonstrations are provided. Notably, across all shot settings shown in the figure, C<sup>3</sup>V consistently achieves higher accuracy than I2CL. In contrast, standard ICL shows limited gains beyond a small number of demonstrations and degrades as the context length grows. These results suggest that C<sup>3</sup>V scales more favorably with increasing demonstration size, avoiding the degradation observed in standard ICL at larger shot counts.

**Noise Injection, Training, and Coefficient Initialization.** We conduct an ablation study to disentangle the effects of noise-based injection, coeffi-

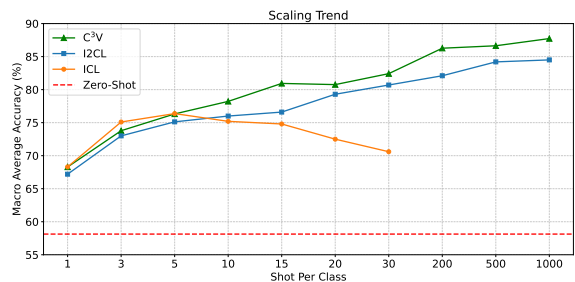


Figure 4: Scaling trend of C<sup>3</sup>V.

cient training, and initialization strategies, as summarized in Table 4. Replacing class-conditional context vectors with noise yields only limited improvement over the zero-shot baseline, indicating that the performance gains cannot be attributed solely to the class-conditional context vectors themselves and that the trainable coefficients play an important role. Injecting class-conditional context with fixed (non-trainable) coefficients yields only modest gains, suggesting that class-aware initialization ( $0.1/|\mathcal{C}|$ ) can introduce class-conditional signals into the model, but remains insufficient without learning. Training the coefficients with a uniform initialization (0.1) leads to substantial performance improvements; however, it remains inferior to the  $0.1/|\mathcal{C}|$  initialization, underscoring the importance of appropriate initialization of the injection parameters. Our full method (**Ours**: combining class-conditional context vectors with trainable coefficients and coefficient initialization of  $0.1/|\mathcal{C}|$ ) achieves the best performance, demonstrating that both coefficient training and initialization are crucial for fully leveraging class-conditional context vectors in C<sup>3</sup>V.

## 5 Conclusion

We introduce C<sup>3</sup>V, a lightweight extension to implicit ICL that explicitly models class-specific con-

textual information via class-conditional context injection. Experiments and representation analyses show that  $C^3V$  consistently improves classification performance over task-level context vectors by enhancing class discriminability, while preserving the efficiency of implicit ICL. Future work will explore extending  $C^3V$  to broader task settings and multi-modal scenarios.

## Limitations

While our method,  $C^3V$ , demonstrates promising results in improving ICL, there are several limitations that warrant further investigation.

1. **Task-Specific Applicability.** Currently,  $C^3V$  is evaluated primarily on text classification tasks. Although we show significant improvements in this domain, the effectiveness of class-conditional context vectors in more complex tasks, such as open-ended generation or multi-label classification, remains to be explored. Tasks involving long-range dependencies, multi-step reasoning, or generation may require further adaptations to the method or a more sophisticated use of context vectors.
2. **Parameter Growth with the Number of Classes.**  $C^3V$  employs class-conditional context vectors whose number increases with the size of the label set. As the label space grows, the number of class-related parameters to be learned or stored correspondingly increases, which may introduce additional memory and computational overhead. While this design facilitates modeling class-level semantic information, its scalability in settings with a large number of classes warrants further investigation. Future work may explore parameter-sharing mechanisms or more compact representations to improve scalability.

## Acknowledgments

We want to thank all the anonymous reviewers for their valuable comments. This work was supported by the National Science Foundation of China (NSFC No. 62576232), Key Laboratory of Data Intelligence and Advanced Computing, Soochow University.

## References

AI@Meta. 2024. [The llama 3 herd of models](#).

Ankush Chatterjee, Kedhar Nath Narahari, Meghana Joshi, and Puneet Agrawal. 2019. [Semeval-2019 task 3: Emocontext contextual emotion detection in text](#). In *Proceedings of the 13th International Workshop on Semantic Evaluation, SemEval@NAACL-HLT 2019, Minneapolis, MN, USA, June 6-7, 2019*, pages 39–48. Association for Computational Linguistics.

Alexis Chevalier, Alexander Wettig, Anirudh Ajith, and Danqi Chen. 2023. [Adapting language models to compress contexts](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 3829–3846. Association for Computational Linguistics.

Ona de Gibert, Naiara Pérez, Aitor García-Pablos, and Montse Cuadros. 2018. [Hate speech dataset from a white supremacy forum](#). In *Proceedings of the 2nd Workshop on Abusive Language Online, ALW@EMNLP 2018, Brussels, Belgium, October 31, 2018*, pages 11–20. Association for Computational Linguistics.

Tao Ge, Jing Hu, Lei Wang, Xun Wang, Si-Qing Chen, and Furu Wei. 2024. [In-context autoencoder for context compression in a large language model](#). In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.

Roe Hendel, Mor Geva, and Amir Globerson. 2023. [In-context learning creates task vectors](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, pages 9318–9333. Association for Computational Linguistics.

Ukyo Honda and Tatsushi Oka. 2025. [Exploring explanations improves the robustness of in-context learning](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2025, Vienna, Austria, July 27 - August 1, 2025*, pages 23693–23714. Association for Computational Linguistics.

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. [Lora: Low-rank adaptation of large language models](#). In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net.

Xinyan Hu, Kayo Yin, Michael I. Jordan, Jacob Steinhardt, and Lijie Chen. 2025. [Understanding in-context learning of addition via activation subspaces](#). *ArXiv*, abs/2505.05145v3.

Kai Konen, Sophie Jentzsch, Diaoulé Diallo, Peer Schütt, Oliver Bensch, Roxanne El Baff, Dominik Opitz, and Tobias Hecking. 2024. [Style vectors for steering generative large language models](#). In *Findings of the Association for Computational Linguistics: EACL 2024*, pages 782–802, St. Julian’s, Malta. Association for Computational Linguistics.

- Vignesh Kothapalli, Hamed Firooz, and Maziar Sanjabi. 2025. [Cot-icl lab: A synthetic framework for studying chain-of-thought learning from in-context demonstrations](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL 2025, Vienna, Austria, July 27 - August 1, 2025, pages 14620–14642. Association for Computational Linguistics.
- Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. [The power of scale for parameter-efficient prompt tuning](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, pages 3045–3059. Association for Computational Linguistics.
- Zhuowei Li, Zihao Xu, Ligong Han, Yunhe Gao, Song Wen, Di Liu, Hao Wang, and Dimitris N. Metaxas. 2025. [Implicit in-context learning](#). In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*. OpenReview.net.
- Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mohata, Tenghao Huang, Mohit Bansal, and Colin Raffel. 2022. [Few-shot parameter-efficient fine-tuning is better and cheaper than in-context learning](#). In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*.
- Hui Liu, Wenya Wang, Hao Sun, Chris Xing Tian, Chenqi Kong, Xin Dong, and Haoliang Li. 2025. [Unraveling the mechanics of learning-based demonstration selection for in-context learning](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL 2025, Vienna, Austria, July 27 - August 1, 2025, pages 2623–2641. Association for Computational Linguistics.
- Sheng Liu, Haotian Ye, Lei Xing, and James Y. Zou. 2024. [In-context vectors: Making in context learning more effective and controllable through latent space steering](#). In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net.
- Ilya Loshchilov and Frank Hutter. 2019. [Decoupled weight decay regularization](#). In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.
- Duy Nguyen, Archiki Prasad, Elias Stengel-Eskin, and Mohit Bansal. 2025. [Multi-attribute steering of language models via targeted intervention](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL 2025, Vienna, Austria, July 27 - August 1, 2025, pages 20619–20634. Association for Computational Linguistics.
- Bo Pang and Lillian Lee. 2004. [A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts](#). In *Proceedings of the 42nd Annual Meeting of the Association for Computational Linguistics, 21-26 July, 2004, Barcelona, Spain*, pages 271–278. ACL.
- Bo Pang and Lillian Lee. 2005. [Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales](#). In *ACL 2005, 43rd Annual Meeting of the Association for Computational Linguistics, Proceedings of the Conference, 25-30 June 2005, University of Michigan, USA*, pages 115–124. The Association for Computer Linguistics.
- Chengwei Qin, Wenhan Xia, Fangkai Jiao, Chen Chen, Yuchen Hu, Bosheng Ding, Ruirui Chen, and Shafiq Joty. 2025. [Beyond output matching: Bidirectional alignment for enhanced in-context learning](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL 2025, Vienna, Austria, July 27 - August 1, 2025, pages 32732–32758. Association for Computational Linguistics.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, and 1 others. 2019. [Language models are unsupervised multitask learners](#). *OpenAI blog*, 1(8):9.
- Nina Rimskey, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Matt Turner. 2024. [Steering llama 2 via contrastive activation addition](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL 2024, Bangkok, Thailand, August 11-16, 2024, pages 15504–15522. Association for Computational Linguistics.
- Baturay Saglam, Xinyang Hu, Zhuoran Yang, Dionysis Kalogerias, and Amin Karbasi. 2025. [Learning task representations from in-context learning](#). In *Findings of the Association for Computational Linguistics, ACL 2025, Vienna, Austria, July 27 - August 1, 2025*, pages 6634–6663. Association for Computational Linguistics.
- Daniel Scalena, Gabriele Sarti, and Malvina Nissim. 2024. [Multi-property steering of large language models with dynamic activation composition](#). In *Proceedings of the 7th BlackboxNLP Workshop: Analyzing and Interpreting Neural Networks for NLP*, pages 577–603, Miami, Florida, US. Association for Computational Linguistics.
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Y. Ng, and Christopher Potts. 2013. [Recursive deep models for semantic compositionality over a sentiment treebank](#). In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing, EMNLP 2013, 18-21 October 2013, Grand Hyatt Seattle, Seattle, Washington, USA, A meeting of SIG-DAT, a Special Interest Group of the ACL*, pages 1631–1642. ACL.

- Alessandro Stolfo, Vidhisha Balachandran, Safoora Yousefi, Eric Horvitz, and Besmira Nushi. 2025. [Improving instruction-following in language models through activation steering](#). In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*. OpenReview.net.
- Jingran Su, Jingfan Chen, Hongxin Li, Yuntao Chen, Li Qing, and Zhaoxiang Zhang. 2025. [Activation steering decoding: Mitigating hallucination in large vision-language models through bidirectional hidden state intervention](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12964–12974, Vienna, Austria. Association for Computational Linguistics.
- Chenghao Sun, Zhen Huang, Yonggang Zhang, Le Lu, Houqiang Li, Xinmei Tian, Xu Shen, and Jieping Ye. 2025. [Interpret and improve in-context learning via the lens of input-label mappings](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL 2025, Vienna, Austria, July 27 - August 1, 2025, pages 3873–3895. Association for Computational Linguistics.
- Qwen Team. 2025. [Qwen3 technical report](#). Preprint, arXiv:2505.09388.
- Eric Todd, Millicent L. Li, Arnab Sen Sharma, Aaron Mueller, Byron C. Wallace, and David Bau. 2024. [Function vectors in large language models](#). In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, and 49 others. 2023. [Llama 2: Open foundation and fine-tuned chat models](#). ArXiv, abs/2307.09288v2.
- Ellen M. Voorhees and Dawn M. Tice. 2000. [Building a question answering test collection](#). In *SIGIR 2000: Proceedings of the 23rd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, July 24-28, 2000, Athens, Greece*, pages 200–207. ACM.
- Ben Wang and Aran Komatsuzaki. 2021. Gpt-j-6b: A 6 billion parameter autoregressive language model.
- Futing Wang, Jianhao Yan, Yue Zhang, and Tao Lin. 2025a. [ELICIT: LLM augmentation via external in-context capability](#). In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*. OpenReview.net.
- 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](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 9840–9855. Association for Computational Linguistics.
- Mengru Wang, Ziwen Xu, Shengyu Mao, Shumin Deng, Zhaopeng Tu, Huajun Chen, and Ningyu Zhang. 2025b. [Beyond prompt engineering: Robust behavior control in llms via steering target atoms](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL 2025, Vienna, Austria, July 27 - August 1, 2025, pages 23381–23399. Association for Computational Linguistics.
- Tianyun Yang, Ziniu Li, Juan Cao, and Chang Xu. 2025. [Understanding and mitigating hallucination in large vision-language models via modular attribution and intervention](#). In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*. OpenReview.net.
- Hantao Yao, Rui Zhang, and Changsheng Xu. 2024. [TCP: textual-based class-aware prompt tuning for visual-language model](#). In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024*, pages 23438–23448. IEEE.
- Howard Yen, Tianyu Gao, and Danqi Chen. 2024. [Long-context language modeling with parallel context encoding](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL 2024, Bangkok, Thailand, August 11-16, 2024, pages 2588–2610. Association for Computational Linguistics.
- Ashkan Yousefpour, Taeheon Kim, Ryan Sungmo Kwon, Seungbeen Lee, Wonje Jeung, Seungju Han, Alvin Wan, Harrison Ngan, Youngjae Yu, and Jonghyun Choi. 2025. [Representation bending for large language model safety](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL 2025, Vienna, Austria, July 27 - August 1, 2025, pages 24073–24098. Association for Computational Linguistics.
- Jinghao Zhang, Yuting Liu, Wenjie Wang, Qiang Liu, Shu Wu, Liang Wang, and Tat-Seng Chua. 2025a. [Personalized text generation with contrastive activation steering](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL 2025, Vienna, Austria, July 27 - August 1, 2025, pages 7128–7141. Association for Computational Linguistics.
- Xiang Zhang, Junbo Jake Zhao, and Yann LeCun. 2015. [Character-level convolutional networks for text classification](#). In *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, pages 649–657.

Zheng Zhang, Shaocheng Lan, Lei Song, Jiang Bian, Yexin Li, and Kan Ren. 2025b. [Learning to select in-context demonstration preferred by large language model](#). In *Findings of the Association for Computational Linguistics, ACL 2025, Vienna, Austria, July 27 - August 1, 2025*, pages 11345–11360. Association for Computational Linguistics.

Yu Zhao, Alessio Devoto, Giwon Hong, Xiaotang Du, Aryo Pradipta Gema, Hongru Wang, Xuanli He, Kam-Fai Wong, and Pasquale Minervini. 2025. [Steering knowledge selection behaviours in llms via sae-based representation engineering](#). In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL 2025 - Volume 1: Long Papers, Albuquerque, New Mexico, USA, April 29 - May 4, 2025*, pages 5117–5136. Association for Computational Linguistics.

## A # Cached Parameters During Inference

- $M$ : number of demonstration (or context) tokens,
- $D$ : model hidden dimension,
- $L$ : number of transformer layers,
- $1/K$ : compression ratio.
- $d_{\text{enc}}$ : hidden dimension of the auxiliary encoder (if used).
- $|\mathcal{C}|$ : the number of classes (if used).

Specifically, the term “ $2MDL$ ” refers to the size of the standard key-value (KV) cache used in autoregressive inference. In the KV caching mechanism, for each of the  $M$  input tokens and across all  $L$  Transformer layers, the model stores:

- a  $D$ -dimensional Key ( $K$ ) vector, and
- a  $D$ -dimensional Value ( $V$ ) vector.

This results in a total of  $2 \times M \times D \times L = 2MDL$  cached parameters. The KV cache is maintained to avoid re-computing attention over the context at every generation step, but it directly contributes to memory footprint.

Few-shot method follows standard autoregressive inference. When the input prompt contains  $M$  demonstration tokens, the model caches their key-value representations across all  $L$  Transformer layers, resulting in a KV cache of size  $2MDL$ .

In contrast, the I2CL framework do *not* cache the full KV states of all  $M$  demonstration tokens, thereby avoiding the  $2MDL$  cost. Instead, it operates as a latent ICL approach:

- During a pre-processing phase, the demonstrations  $d_1, \dots, d_N$  are passed through the frozen LLM, and for each demonstration, we extract the multi-head attention output  $a_{i,l}^e$  and MLP output  $m_{i,l}^e$  at the final token position  $e$  of each layer  $l$ .
- These activations are then averaged across demonstrations to form a compact context vector.
- The resulting context representation is  $\{\bar{a}_l^e, \bar{m}_l^e\}_{l=1}^L$ , which contains two  $D$ -dimensional vectors per layer (one from attention, one from MLP), totaling  $2DL$  parameters.

During inference, no demonstration tokens are included in the input prompt. Instead, the pre-computed context vector  $\{\bar{a}_l^e, \bar{m}_l^e\}_{l=1}^L$  is directly injected into the residual stream of the query. This eliminates the need for a large KV cache proportional to  $M$ , making the method significantly more memory-efficient—especially as the number of demonstrations grows.

Building upon the latent ICL formulation of I2CL, our  $\mathbf{C}^3\mathbf{V}$  framework further refines the context representation by modeling class-conditional semantics explicitly. Instead of aggregating all demonstrations into a single task-level context vector,  $\mathbf{C}^3\mathbf{V}$  maintains a separate context vector for each class.

Concretely, for a classification task with  $|\mathcal{C}|$  classes, demonstrations belonging to the same class  $c \in \mathcal{C}$  are first processed by the frozen LLM to extract the attention and MLP activations at the final token position of each layer. These activations are averaged *within each class*, yielding a set of class-conditional context vectors

$$\{\bar{a}_{c,l}^e, \bar{m}_{c,l}^e\}_{l=1}^L, \quad c \in \mathcal{C}.$$

As a result,  $\mathbf{C}^3\mathbf{V}$  caches  $|\mathcal{C}|$  such context representations, each consisting of two  $D$ -dimensional vectors per layer. The total number of cached parameters therefore scales as

$$2|\mathcal{C}|DL,$$

which is linear in the number of classes but remains independent of the number of demonstration tokens  $M$ .

During inference, no demonstration tokens are appended to the query. Instead, all class-conditional context vectors are jointly injected into the residual stream through a linear intervention, enabling fine-grained class-aware modulation of

the model’s internal representations. Compared to I2CL, C<sup>3</sup>V incurs a modest increase in cached parameters proportional to  $|\mathcal{C}|$ , while preserving the same zero-shot inference cost and avoiding any KV-cache growth with respect to  $M$ .

For Label Anchor method, the demonstrations are compressed into  $\frac{M}{K}$  context tokens, including both label and formatting tokens. During inference, the model still performs token-level attention over these compressed tokens. As a result, the cached parameters correspond to the KV cache of the compressed context, yielding a memory cost of  $2 \times \frac{M}{K} \times D \times L$ .

For Task Vector method, the demonstrations are compressed into a single task vector  $\theta(S) \in \mathbb{R}^D$ . During inference, no demonstration tokens are included in the input; instead,  $\theta(S)$  is patched into the residual stream at an intermediate layer. Therefore, the cached parameters consist of a single activation vector of dimension  $D$ , independent of the number of demonstrations.

For ICV method, demonstrations are summarized into a set of layer-wise activation vectors  $\{h_{ICV}^l\}_{l=1}^L$  derived from the MLP outputs. During inference, no demonstration tokens are included; instead, one  $D$ -dimensional activation vector is added at each layer. Therefore, the cached parameters scale as  $DL$ , independent of the number of demonstrations.

ICAE uses a LoRA-fine-tuned encoder to compress the original long context  $M$  into a small number of “Memory Slots”. The compression ratio is  $1/K$ , so during inference, the effective sequence length entering the model is no longer  $M$ , but the reduced  $\frac{M}{K}$ . The KV cache parameter calculation is  $2 \times \frac{M}{K} \times D \times L$ , where the number 2 represents the Key and Value matrices.

In AutoCompressor, the context is compressed into  $\frac{M}{K}$  summary vectors, which are used as soft prompts and injected into the input of every layer of the model. During inference, the Key and Value representations of these summary tokens are cached at each layer, forming the Compressed KV Cache. Therefore, the size of the cached parameters is:  $2 \times \frac{M}{K} \times D \times L$ .

Ignoring the caching overhead from the query during inference, and considering only the cached parameters for the additional context, CEPE stores  $M \cdot d_{\text{enc}}$  parameters. Specifically, it caches the **last-layer hidden representations** from the small bidirectional encoder  $Model_{\text{enc}}$ , where  $d_{\text{enc}} = 1024$  (matching the hidden dimension of RoBERTa-large,

as used in the paper). These  $d_{\text{enc}}$ -dimensional vectors are later **up-projected to the decoder’s hidden dimension**  $d_{\text{dec}}$  (e.g., 4096 for LLaMA2-7b) **within the cross-attention layers** at each decoder block, enabling the decoder to attend to the encoded context without caching  $d_{\text{dec}}$ -dimensional keys and values for the additional tokens.

## B Definition of Intra-Class and Inter-Class Distances

In this appendix, we formally define the intra-class and inter-class distance metrics used in our representation analysis.

### B.1 Distance Metric

All distances are computed in the original high-dimensional representation space using *cosine distance*, defined as

$$d(\mathbf{x}, \mathbf{z}) = 1 - \frac{\mathbf{x}^\top \mathbf{z}}{\|\mathbf{x}\| \|\mathbf{z}\|}. \quad (8)$$

Cosine distance captures the angular relationship between representations and is widely used in representation learning due to its ability to emphasize semantic direction while remaining invariant to vector magnitude.

### B.2 Intra-Class Distance

Let  $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  denote the set of context representations, where each representation  $\mathbf{x}_i$  is associated with a class label  $y_i \in \{1, \dots, C\}$ . For each class  $c$ , we define the subset of representations belonging to that class as

$$\mathcal{X}_c = \{\mathbf{x}_i \mid y_i = c\}, \quad (9)$$

with cardinality  $n_c = |\mathcal{X}_c|$ . In our experiments,  $\mathcal{X}_c$  consists of the class-conditional context vectors for class  $c$  obtained from multiple independent runs with different random seeds.

The *intra-class distance* is computed as the average cosine distance over all unordered pairs of representations belonging to the same class, aggregated across all classes:

$$\text{Intra} = \frac{\sum_{c=1}^C \sum_{\substack{i < j \\ \mathbf{x}_i, \mathbf{x}_j \in \mathcal{X}_c}} d(\mathbf{x}_i, \mathbf{x}_j)}{\sum_{c=1}^C \binom{n_c}{2}}. \quad (10)$$

This metric reflects the degree of dispersion of representations within each class.

### B.3 Inter-Class Distance

The *inter-class distance* is defined as the average cosine distance between representations drawn from different classes:

$$\text{Inter} = \frac{\sum_{c_1 < c_2} \sum_{\mathbf{x}_i \in \mathcal{X}_{c_1}} \sum_{\mathbf{x}_j \in \mathcal{X}_{c_2}} d(\mathbf{x}_i, \mathbf{x}_j)}{\sum_{c_1 < c_2} n_{c_1} n_{c_2}}. \quad (11)$$

This quantity captures the overall semantic separation between classes in the representation space.

### B.4 Inter-to-Intra Distance Ratio

To summarize class separability with a single scalar metric, we further report the *inter-to-intra distance ratio*:

$$\text{Ratio} = \frac{\text{Inter}}{\text{Intra}}. \quad (12)$$

A higher ratio indicates that inter-class distances increase more substantially than intra-class distances, suggesting improved discriminability of the learned representations.

We emphasize that all distances are computed directly in the original representation space, without relying on dimensionality reduction techniques such as t-SNE, which may distort pairwise distances. Moreover, this analysis does not depend on training an auxiliary classifier, allowing us to assess representation geometry in a model-agnostic manner.

## C C<sup>3</sup>V on Other Large Language Models

In this section, we evaluate the performance of C<sup>3</sup>V on a diverse set of large language models beyond LLaMA2-7B, including models from the LLaMA, GPT, and Qwen families. All experiments follow the same experimental protocol as described in the main text, including dataset splits, demonstration construction, coefficient training, and evaluation metrics. For brevity, we focus on reporting the final performance results and key trends, while detailed analyses are deferred to the main results on LLaMA2-7B. Due to the limited context window of GPT2-XL, the AGNews and DBPedia datasets were not evaluated under the I2CL setting on GPT2-XL. For the same reason, our method, C<sup>3</sup>V, was also not evaluated on these datasets when using GPT2-XL. As shown in Table 6–10, overall, C<sup>3</sup>V generally outperforms the baselines across all evaluated models, with a single exception: on LLaMA3.2-1B, the Soft Prompt baseline achieves better performance than C<sup>3</sup>V. This behavior may

be related to the limited capacity of LLaMA3.2-1B, under which prompt-tuning-based adaptation methods such as Soft Prompt can sometimes be more effective. It is worth noting that, compared to C<sup>3</sup>V, Soft Prompt involves a larger number of trainable parameters. Notably, even on the strong Qwen3-8B model, our method achieves an average improvement of nearly 0.8 points over the strongest baseline across all datasets.

## D Attention- and MLP-Based Class-Conditional Latent Vector Visualizations for Other Models

Appendix D provides additional visualizations of class-conditional latent vectors obtained from other language models. Following the same visualization procedure as in Section 4.3.1, we concatenate layer-wise attention or MLP representations for each class-conditional context and project them into two dimensions using t-SNE. These results complement the main analysis and demonstrate that the observed representational patterns generalize across different models. The corresponding visualizations are shown in Figures 5–9.

## E Inter- and Intra-Class Distance Analysis for Other Models

This appendix presents the inter- and intra-class distance analysis for additional language models beyond LLaMA2-7B. The analysis follows the same experimental protocol and distance definitions as described in Appendix B, using cosine distance for all computations. These results complement the main analysis and further assess the generality of the observed representational patterns across different models. The corresponding visualizations are shown in Figures 10–14.

## F Layer-wise Coefficients of LLaMA2-7B

We present the layer-wise coefficients of LLaMA2-7B across all datasets for completeness and visualization purposes. As shown in Figure 15, C<sup>3</sup>V provides finer control over information injection by learning layer-wise coefficients that are not strictly positive. The coefficients fluctuate across layers, implying that contextual information may be reduced rather than consistently amplified at different depths.

Name	Accuracy(%)
Zero-shot	58.13
MHA	68.90
MLP	73.13
Hidden state	59.92
MHA+MLP	76.32

(a) Target Modules

Name	Accuracy(%)
Zero-shot	58.13
Early	59.37
Middle	70.40
Late	67.54
All(ours)	76.32

(b) Target Layers

Name	Accuracy(%)
Zero-shot	58.13
First	67.78
Last	71.09
All(ours)	76.32

(c) Injection Pos

Name	Accuracy(%)
Zero-shot	58.13
First	34.12
Random	69.29
Last (ours)	76.32

(d) Extraction Pos

Name	Accuracy(%)
Zero-shot	58.13
PCA	75.08
Mean(ours)	76.32

(e) Post Fuse Method

Table 5: Comparison of Different Settings.

## G Additional Ablation Studies

**Target Module.** We study different target modules for activation extraction and injection, including MHA, MLP, hidden states, and their combination. As shown in Table 5a, intervening at either MHA or MLP improves performance over the baseline, with the combined MHA+MLP setting achieving the best results. In contrast, using hidden states provides only marginal gains, likely due to their mixed and less targeted representations.

**Target Layer.** We examine where to inject context by grouping model layers into early, middle, and late stages, each containing one-third of the total layers. As shown in Table 5b, injecting context into middle or late layers yields larger gains than early-layer injection, while applying context across all layers achieves the best performance. This suggests that injecting contextual information across layers helps refine representations throughout the network.

**Injection Position.** By default,  $C^3V$  injects context across all token positions. We compare this strategy with alternatives that inject context at a single position, including the first token or the last token. As shown in Table 5c, injecting at the last token outperforms the first setting, while injecting across all positions achieves the best results.

**Extraction Position.** We study the effect of extracting class-conditional context vectors from different token positions, including the first token, a randomly selected token, and the last token in the sequence. As shown in Table 5d, extracting context vectors from the last token yields the best perfor-

mance, substantially outperforming other positions. In contrast, using the first token leads to a significant performance drop. These findings suggest that extracting context vectors from the last token yields the most effective performance, while using the first token leads to a significant drop, and random extraction performs worse than extraction from the last token.

**Context Vector Calculation.** We analyze how  $C^3V$  aggregates class-conditional context vectors across layers, comparing mean pooling and PCA-based fusion. Mean fusion averages layer-wise context vectors, whereas PCA fusion performs dimensionality reduction before aggregation. As shown in Table 5e, mean fusion outperforms PCA, indicating that simple averaging better preserves class-relevant information for effective context injection in  $C^3V$ .

## H Dataset Templates and Labels

The dataset templates and label spaces used in the classification tasks for our method and baselines are shown in Table 11.

## I $C^3V$ Algorithm Pseudocode

Algorithm 1 provides the pseudocode of the  $C^3V$  algorithm. The pseudocode summarizes the overall training procedure, including the construction of class-conditional context vectors, coefficient initialization, and coefficient optimization.

Method	SST-2	SST-5	TREC	AGNews	Subj	HateS	DBPedia	EmoC	MR	Avg.
Zero-shot	91.00	30.60	53.80	66.60	45.00	51.00	63.40	59.40	89.00	61.09
Few-shot	94.00 <sub>0.40</sub>	46.72 <sub>1.31</sub>	84.96 <sub>3.05</sub>	85.48 <sub>1.57</sub>	50.04 <sub>0.08</sub>	52.72 <sub>3.18</sub>	87.20 <sub>0.78</sub>	66.36 <sub>3.96</sub>	92.48 <sub>0.55</sub>	73.33
Soft Prompt	50.20 <sub>0.08</sub>	23.93 <sub>0.66</sub>	80.00 <sub>1.80</sub>	82.27 <sub>8.13</sub>	59.80 <sub>0.40</sub>	52.40 <sub>0.96</sub>	34.80 <sub>4.57</sub>	54.30 <sub>9.84</sub>	51.00 <sub>0.40</sub>	54.30
Task-vector	94.00 <sub>0.85</sub>	35.44 <sub>2.06</sub>	57.48 <sub>1.27</sub>	83.44 <sub>4.56</sub>	45.76 <sub>2.05</sub>	51.56 <sub>0.64</sub>	69.52 <sub>1.62</sub>	62.28 <sub>0.52</sub>	87.60 <sub>0.33</sub>	65.23
Label-anchor	89.88 <sub>1.65</sub>	34.12 <sub>2.05</sub>	79.64 <sub>2.05</sub>	70.28 <sub>2.82</sub>	46.96 <sub>14.34</sub>	52.40 <sub>3.44</sub>	54.36 <sub>19.58</sub>	68.12 <sub>3.95</sub>	87.56 <sub>0.74</sub>	64.81
I2CL	90.36 <sub>2.44</sub>	35.92 <sub>5.99</sub>	81.68 <sub>3.69</sub>	84.00 <sub>2.25</sub>	85.00 <sub>3.34</sub>	72.40 <sub>2.44</sub>	86.80 <sub>3.08</sub>	68.04 <sub>3.30</sub>	86.92 <sub>1.52</sub>	76.79
C <sup>3</sup> V	89.08 <sub>2.40</sub>	38.56 <sub>4.21</sub>	86.80 <sub>3.53</sub>	85.20 <sub>2.19</sub>	82.96 <sub>2.54</sub>	70.52 <sub>3.73</sub>	91.40 <sub>0.03</sub>	67.52 <sub>5.13</sub>	86.32 <sub>0.60</sub>	77.60

Table 6: Classification accuracy (in %) on nine datasets using Qwen3-8B. Results are reported as  $\text{acc}_{\text{std}}$ , where std denotes the standard deviation over multiple runs.

Method	SST-2	SST-5	TREC	AGNews	Subj	HateS	DBPedia	EmoC	MR	Avg.
Zero-shot	60.60	28.00	46.80	79.80	55.60	50.80	49.20	39.20	56.20	51.80
Few-shot	64.32 <sub>9.69</sub>	20.56 <sub>0.93</sub>	71.16 <sub>3.19</sub>	86.00 <sub>1.23</sub>	50.00 <sub>0.00</sub>	50.00 <sub>0.00</sub>	83.84 <sub>2.21</sub>	25.44 <sub>0.78</sub>	60.84 <sub>5.81</sub>	56.91
Soft Prompt	58.60 <sub>4.66</sub>	34.60 <sub>1.59</sub>	65.84 <sub>6.75</sub>	83.04 <sub>3.71</sub>	58.32 <sub>4.21</sub>	59.64 <sub>4.07</sub>	82.60 <sub>1.78</sub>	50.28 <sub>8.61</sub>	67.16 <sub>7.90</sub>	62.23
Task-vector	69.88 <sub>2.83</sub>	24.52 <sub>2.92</sub>	65.76 <sub>1.63</sub>	86.08 <sub>2.52</sub>	51.36 <sub>0.74</sub>	49.28 <sub>4.96</sub>	75.32 <sub>3.72</sub>	48.08 <sub>0.90</sub>	66.00 <sub>3.37</sub>	59.59
Label-anchor	58.92 <sub>1.57</sub>	28.80 <sub>4.76</sub>	51.76 <sub>1.42</sub>	76.36 <sub>3.04</sub>	50.72 <sub>0.45</sub>	50.00 <sub>0.13</sub>	62.00 <sub>5.04</sub>	25.84 <sub>0.50</sub>	58.64 <sub>1.27</sub>	51.45
I2CL	54.96 <sub>2.80</sub>	26.20 <sub>2.73</sub>	50.00 <sub>8.57</sub>	81.52 <sub>3.26</sub>	63.04 <sub>3.93</sub>	54.48 <sub>7.47</sub>	66.52 <sub>4.73</sub>	45.36 <sub>2.92</sub>	70.48 <sub>4.21</sub>	56.95
C <sup>3</sup> V	57.80 <sub>6.86</sub>	26.56 <sub>2.50</sub>	63.88 <sub>9.61</sub>	83.60 <sub>2.69</sub>	60.32 <sub>8.13</sub>	55.48 <sub>7.97</sub>	87.12 <sub>2.33</sub>	50.64 <sub>4.17</sub>	69.08 <sub>5.74</sub>	61.61

Table 7: Classification accuracy (in %) on nine datasets using LLaMA3.2-1B. Results are reported as  $\text{acc}_{\text{std}}$ , where std denotes the standard deviation over multiple runs.

Method	SST-2	SST-5	TREC	AGNews	Subj	HateS	DBPedia	EmoC	MR	Avg.
Zero-shot	55.60	31.40	66.40	73.60	51.00	50.40	56.20	41.00	55.60	53.47
Few-shot	93.00 <sub>0.62</sub>	39.40 <sub>3.21</sub>	78.20 <sub>5.94</sub>	84.80 <sub>1.30</sub>	62.47 <sub>12.44</sub>	64.93 <sub>2.65</sub>	82.33 <sub>3.64</sub>	52.20 <sub>3.92</sub>	92.73 <sub>0.50</sub>	72.23
Noise vector	68.48 <sub>6.60</sub>	38.60 <sub>3.03</sub>	66.88 <sub>9.34</sub>	82.16 <sub>2.78</sub>	69.96 <sub>5.73</sub>	56.04 <sub>5.09</sub>	80.44 <sub>4.29</sub>	43.84 <sub>7.59</sub>	77.76 <sub>8.82</sub>	64.91
Soft Prompt	58.96 <sub>12.03</sub>	29.56 <sub>7.61</sub>	74.24 <sub>10.38</sub>	85.64 <sub>2.77</sub>	69.00 <sub>9.78</sub>	62.68 <sub>5.03</sub>	86.40 <sub>0.00</sub>	54.32 <sub>11.88</sub>	60.32 <sub>11.85</sub>	64.57
Task-vector	69.28 <sub>8.96</sub>	35.20 <sub>2.83</sub>	64.80 <sub>1.50</sub>	84.00 <sub>3.74</sub>	50.40 <sub>0.00</sub>	55.60 <sub>3.41</sub>	73.28 <sub>1.27</sub>	41.64 <sub>0.99</sub>	66.28 <sub>4.70</sub>	60.05
Label-anchor	52.32 <sub>1.84</sub>	29.00 <sub>4.99</sub>	43.40 <sub>0.54</sub>	74.64 <sub>0.85</sub>	50.20 <sub>0.70</sub>	66.24 <sub>8.73</sub>	72.04 <sub>2.00</sub>	56.20 <sub>3.01</sub>	58.28 <sub>3.27</sub>	55.81
I2CL	91.40 <sub>2.26</sub>	32.60 <sub>1.25</sub>	80.27 <sub>2.58</sub>	82.56 <sub>1.57</sub>	64.67 <sub>3.52</sub>	75.80 <sub>1.02</sub>	85.00 <sub>0.33</sub>	52.92 <sub>5.28</sub>	84.27 <sub>2.07</sub>	72.17
C <sup>3</sup> V	89.16 <sub>2.33</sub>	43.60 <sub>3.86</sub>	78.24 <sub>3.95</sub>	85.04 <sub>2.57</sub>	64.96 <sub>5.84</sub>	62.72 <sub>4.80</sub>	91.56 <sub>1.37</sub>	64.80 <sub>4.64</sub>	88.52 <sub>3.73</sub>	74.29

Table 8: Classification accuracy (in %) on nine datasets using LLaMA3-8B. Results are reported as  $\text{acc}_{\text{std}}$ , where std denotes the standard deviation over multiple runs.

Method	SST-2	SST-5	TREC	AGNews	Subj	HateS	DBPedia	EmoC	MR	Avg.
Zero-shot	77.76	25.60	68.20	71.60	62.40	59.92	65.56	44.68	76.88	61.40
Few-shot	89.44 <sub>2.60</sub>	39.65 <sub>4.57</sub>	67.76 <sub>2.11</sub>	83.18 <sub>2.03</sub>	50.20 <sub>0.22</sub>	53.44 <sub>6.84</sub>	93.30 <sub>1.19</sub>	47.62 <sub>8.62</sub>	88.66 <sub>1.23</sub>	68.14
Soft Prompt	68.32 <sub>11.99</sub>	37.68 <sub>2.91</sub>	68.24 <sub>3.83</sub>	82.28 <sub>6.22</sub>	63.72 <sub>7.50</sub>	63.48 <sub>5.93</sub>	86.24 <sub>2.50</sub>	49.24 <sub>7.64</sub>	73.84 <sub>9.23</sub>	65.89
Task-vector	59.80 <sub>4.47</sub>	31.20 <sub>2.82</sub>	67.32 <sub>0.32</sub>	80.12 <sub>2.23</sub>	66.32 <sub>2.31</sub>	70.12 <sub>2.22</sub>	77.64 <sub>4.63</sub>	45.00 <sub>4.20</sub>	81.36 <sub>3.58</sub>	64.32
Label-anchor	86.12 <sub>3.95</sub>	27.44 <sub>4.21</sub>	51.68 <sub>4.87</sub>	77.52 <sub>1.53</sub>	52.64 <sub>2.99</sub>	54.32 <sub>7.27</sub>	90.16 <sub>1.13</sub>	42.16 <sub>10.04</sub>	87.56 <sub>1.38</sub>	63.29
I2CL	85.48 <sub>1.18</sub>	37.32 <sub>3.11</sub>	63.84 <sub>7.58</sub>	81.56 <sub>3.13</sub>	65.56 <sub>8.33</sub>	62.32 <sub>5.76</sub>	81.84 <sub>4.50</sub>	50.32 <sub>4.68</sub>	84.40 <sub>2.45</sub>	68.07
C <sup>3</sup> V	84.72 <sub>4.18</sub>	38.20 <sub>2.02</sub>	68.44 <sub>7.68</sub>	82.08 <sub>3.48</sub>	66.24 <sub>7.87</sub>	63.88 <sub>4.50</sub>	90.36 <sub>1.67</sub>	57.16 <sub>3.74</sub>	83.64 <sub>3.17</sub>	70.52

Table 9: Classification accuracy (in %) on nine datasets using GPT-J-6B. Results are reported as  $\text{acc}_{\text{std}}$ , where std denotes the standard deviation over multiple runs.

Method	SST-2	SST-5	TREC	Subj	HateS	EmoC	MR	Avg.
Zero-shot	74.60	31.40	35.40	64.40	71.80	38.60	70.40	55.23
Few-shot	73.44 <sub>8.97</sub>	35.92 <sub>2.41</sub>	60.64 <sub>5.00</sub>	64.28 <sub>10.89</sub>	51.76 <sub>3.06</sub>	38.64 <sub>7.46</sub>	75.76 <sub>9.49</sub>	57.21
Soft Prompt	59.24 <sub>5.12</sub>	23.48 <sub>2.16</sub>	46.04 <sub>8.46</sub>	52.76 <sub>3.74</sub>	62.72 <sub>7.50</sub>	36.08 <sub>3.28</sub>	55.84 <sub>6.06</sub>	48.02
Task-vector	81.08 <sub>4.87</sub>	28.52 <sub>1.37</sub>	41.40 <sub>5.35</sub>	71.80 <sub>1.86</sub>	62.48 <sub>2.83</sub>	37.60 <sub>2.48</sub>	78.40 <sub>2.36</sub>	57.33
Label-anchor	62.36 <sub>8.35</sub>	24.28 <sub>2.10</sub>	68.60 <sub>7.99</sub>	55.36 <sub>8.09</sub>	52.56 <sub>4.20</sub>	36.20 <sub>4.00</sub>	61.72 <sub>9.30</sub>	51.58
I2CL	80.16 <sub>3.98</sub>	33.84 <sub>2.60</sub>	51.48 <sub>5.26</sub>	65.96 <sub>4.83</sub>	68.32 <sub>4.76</sub>	47.92 <sub>1.84</sub>	83.20 <sub>3.29</sub>	61.55
C <sup>3</sup> V	83.60 <sub>1.55</sub>	37.24 <sub>4.82</sub>	58.40 <sub>8.09</sub>	69.32 <sub>7.62</sub>	65.24 <sub>5.83</sub>	55.95 <sub>4.79</sub>	82.16 <sub>4.89</sub>	64.56

Table 10: Classification accuracy (in %) on seven datasets using GPT-XL. Results are reported as  $\text{acc}_{\text{std}}$ , where std denotes the standard deviation over multiple runs.

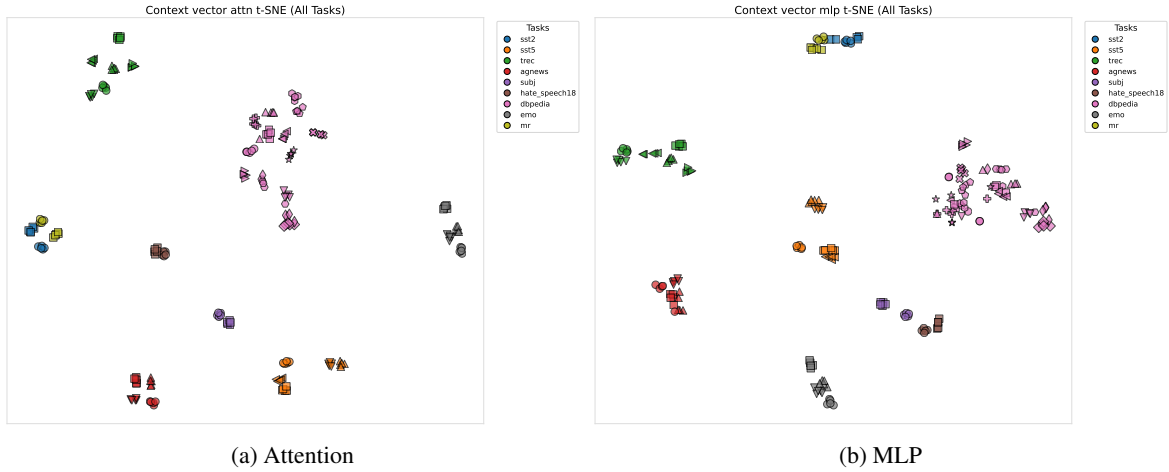


Figure 5: t-SNE visualization of class-conditional latent vectors on Qwen3-8B. Left: attention-based latent vectors; right: MLP-based latent vectors. Points are colored by dataset and shaped by class label. Latent vectors from five random seeds are included for each class.

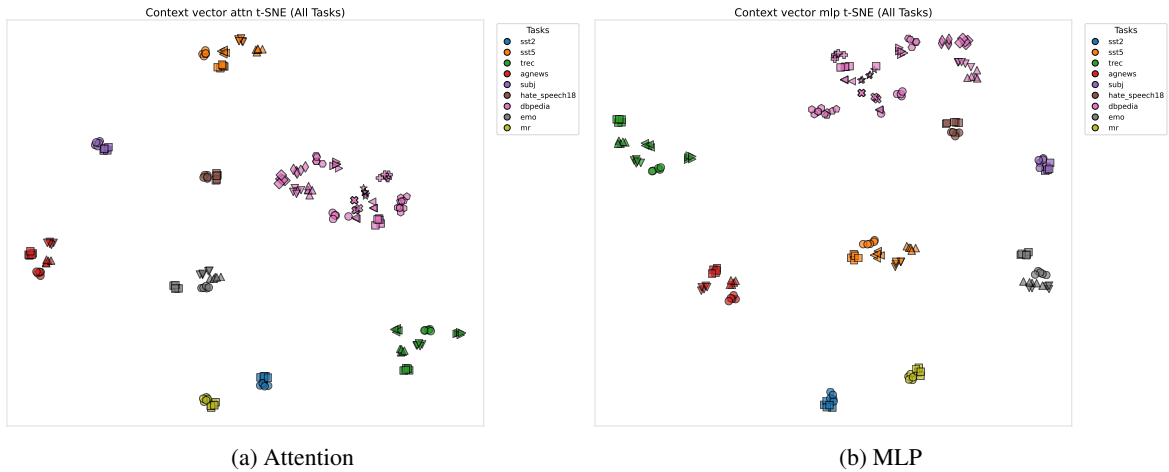


Figure 6: t-SNE visualization of class-conditional latent vectors on LLaMA3.2-1B. Left: attention-based latent vectors; right: MLP-based latent vectors. Points are colored by dataset and shaped by class label. Latent vectors from five random seeds are included for each class.

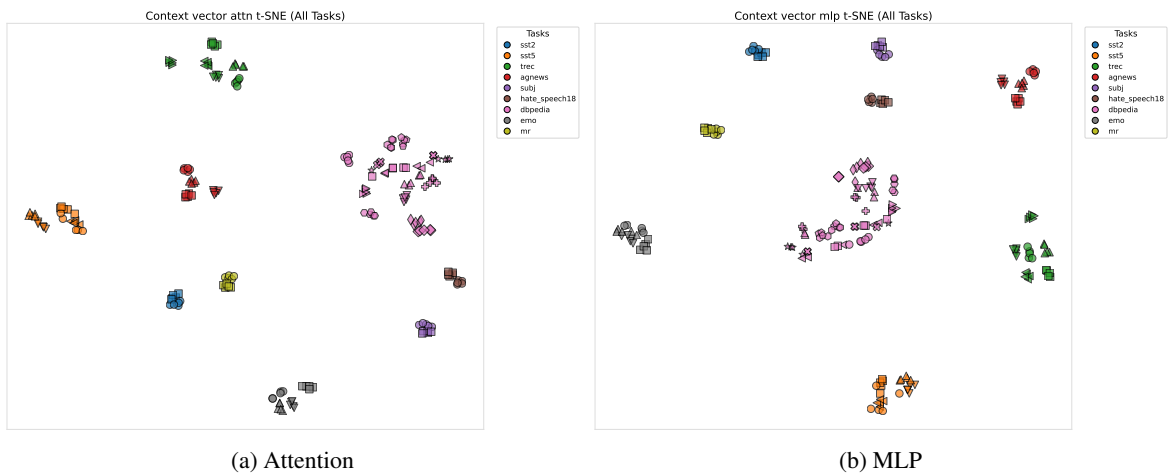


Figure 7: t-SNE visualization of class-conditional latent vectors on LLaMA3-8B. Left: attention-based latent vectors; right: MLP-based latent vectors. Points are colored by dataset and shaped by class label. Latent vectors from five random seeds are included for each class.

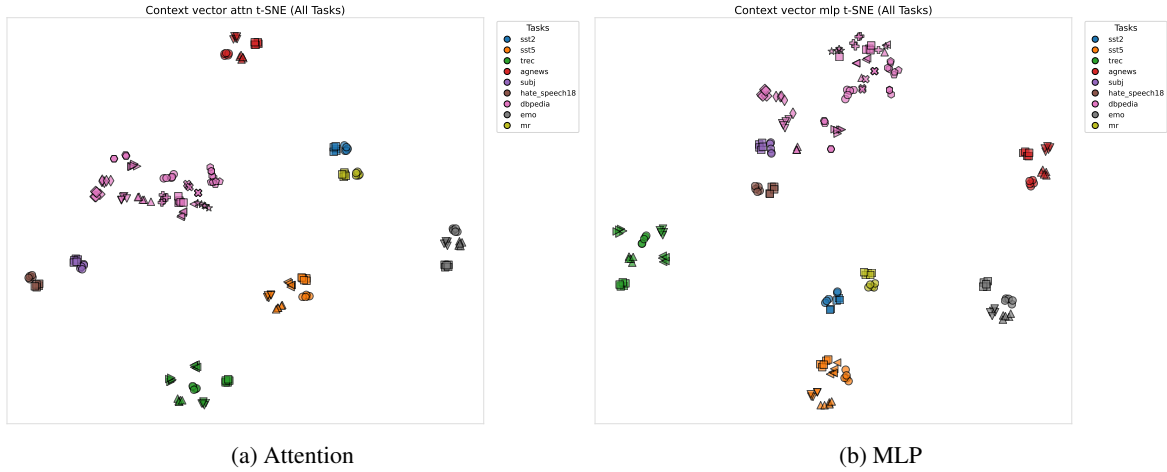


Figure 8: t-SNE visualization of class-conditional latent vectors on GPT-J-6B. Left: attention-based latent vectors; right: MLP-based latent vectors. Points are colored by dataset and shaped by class label. Latent vectors from five random seeds are included for each class.

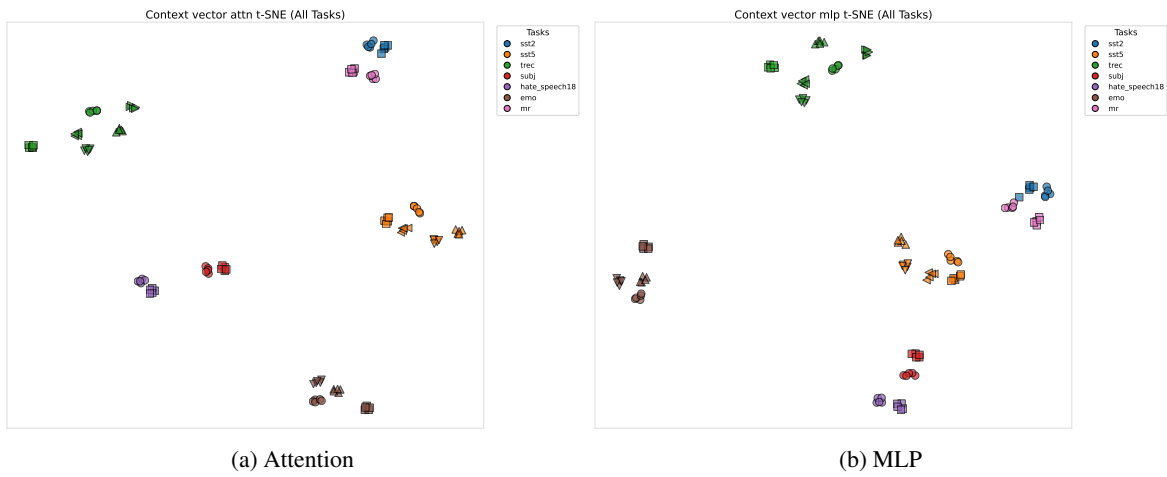


Figure 9: t-SNE visualization of class-conditional latent vectors on GPT2-XL. Left: attention-based latent vectors; right: MLP-based latent vectors. Points are colored by dataset and shaped by class label. Latent vectors from five random seeds are included for each class.

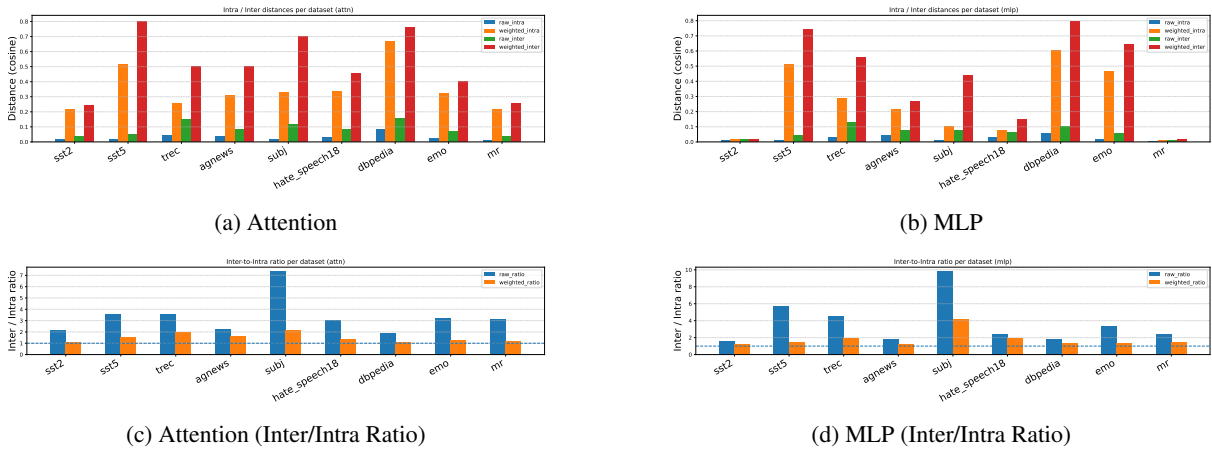
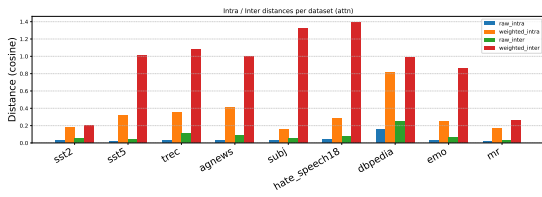
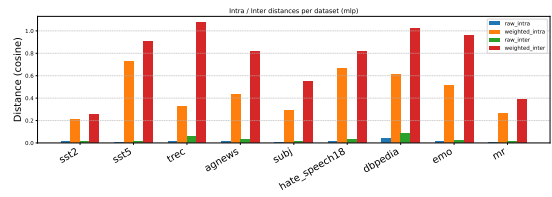


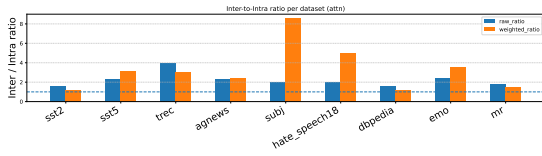
Figure 10: Inter- and intra-class distance analysis of latent representations across datasets on Qwen3-8b. (a,b) Inter- and intra-class distances for attention-based and MLP-based representations, where each panel includes distances computed both without coefficient scaling and after multiplication by trained coefficients. (c,d) Corresponding inter/intra ratios for attention and MLP representations, comparing the settings without and with trained coefficient scaling.



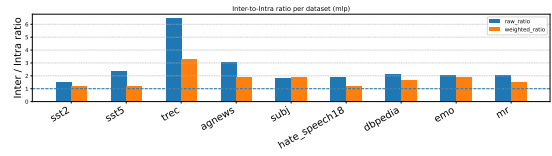
(a) Attention



(b) MLP

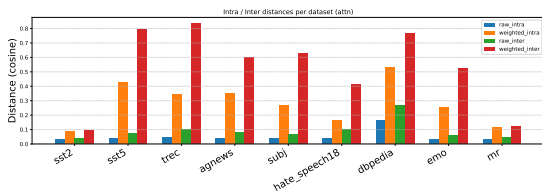


(c) Attention (Inter/Intra Ratio)

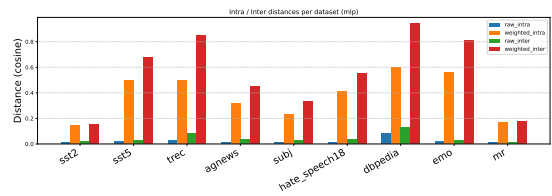


(d) MLP (Inter/Intra Ratio)

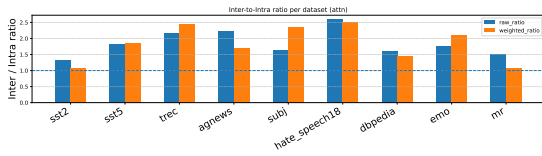
Figure 11: Inter- and intra-class distance analysis of latent representations across datasets on LLaMA3.2-1b. (a,b) Inter- and intra-class distances for attention-based and MLP-based representations, where each panel includes distances computed both without coefficient scaling and after multiplication by trained coefficients. (c,d) Corresponding inter/intra ratios for attention and MLP representations, comparing the settings without and with trained coefficient scaling.



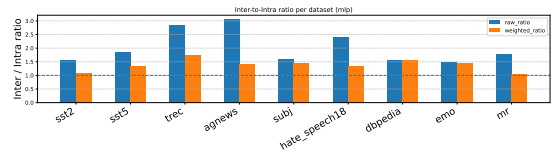
(a) Attention



(b) MLP

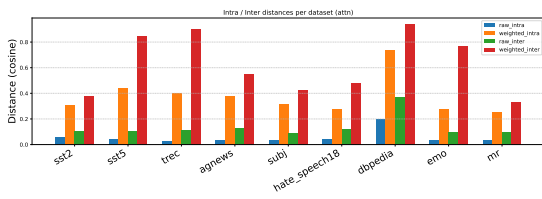


(c) Attention (Inter/Intra Ratio)

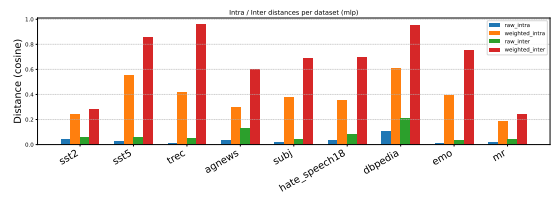


(d) MLP (Inter/Intra Ratio)

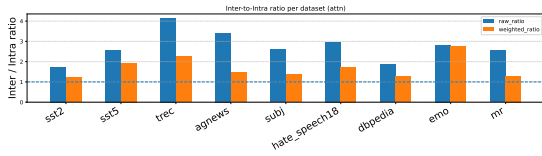
Figure 12: Inter- and intra-class distance analysis of latent representations across datasets on LLaMA3-8b. (a,b) Inter- and intra-class distances for attention-based and MLP-based representations, where each panel includes distances computed both without coefficient scaling and after multiplication by trained coefficients. (c,d) Corresponding inter/intra ratios for attention and MLP representations, comparing the settings without and with trained coefficient scaling.



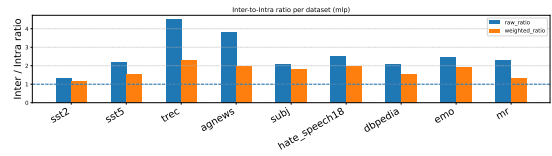
(a) Attention



(b) MLP

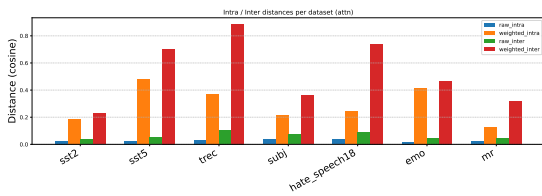


(c) Attention (Inter/Intra Ratio)

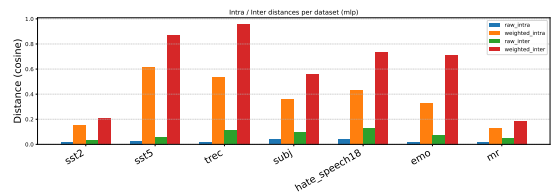


(d) MLP (Inter/Intra Ratio)

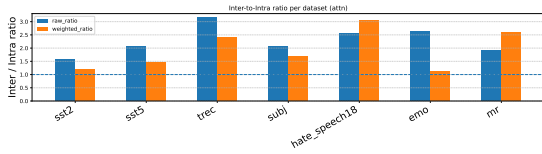
Figure 13: Inter- and intra-class distance analysis of latent representations across datasets on GPT-J-6B. (a,b) Inter- and intra-class distances for attention-based and MLP-based representations, where each panel includes distances computed both without coefficient scaling and after multiplication by trained coefficients. (c,d) Corresponding inter/intra ratios for attention and MLP representations, comparing the settings without and with trained coefficient scaling.



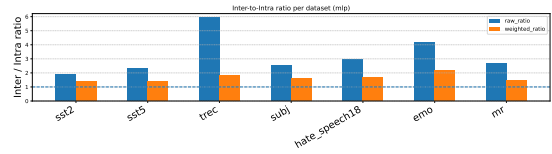
(a) Attention



(b) MLP



(c) Attention (Inter/Intra Ratio)



(d) MLP (Inter/Intra Ratio)

Figure 14: Inter- and intra-class distance analysis of latent representations across datasets on GPT2-XL. (a,b) Inter- and intra-class distances for attention-based and MLP-based representations, where each panel includes distances computed both without coefficient scaling and after multiplication by trained coefficients. (c,d) Corresponding inter/intra ratios for attention and MLP representations, comparing the settings without and with trained coefficient scaling.

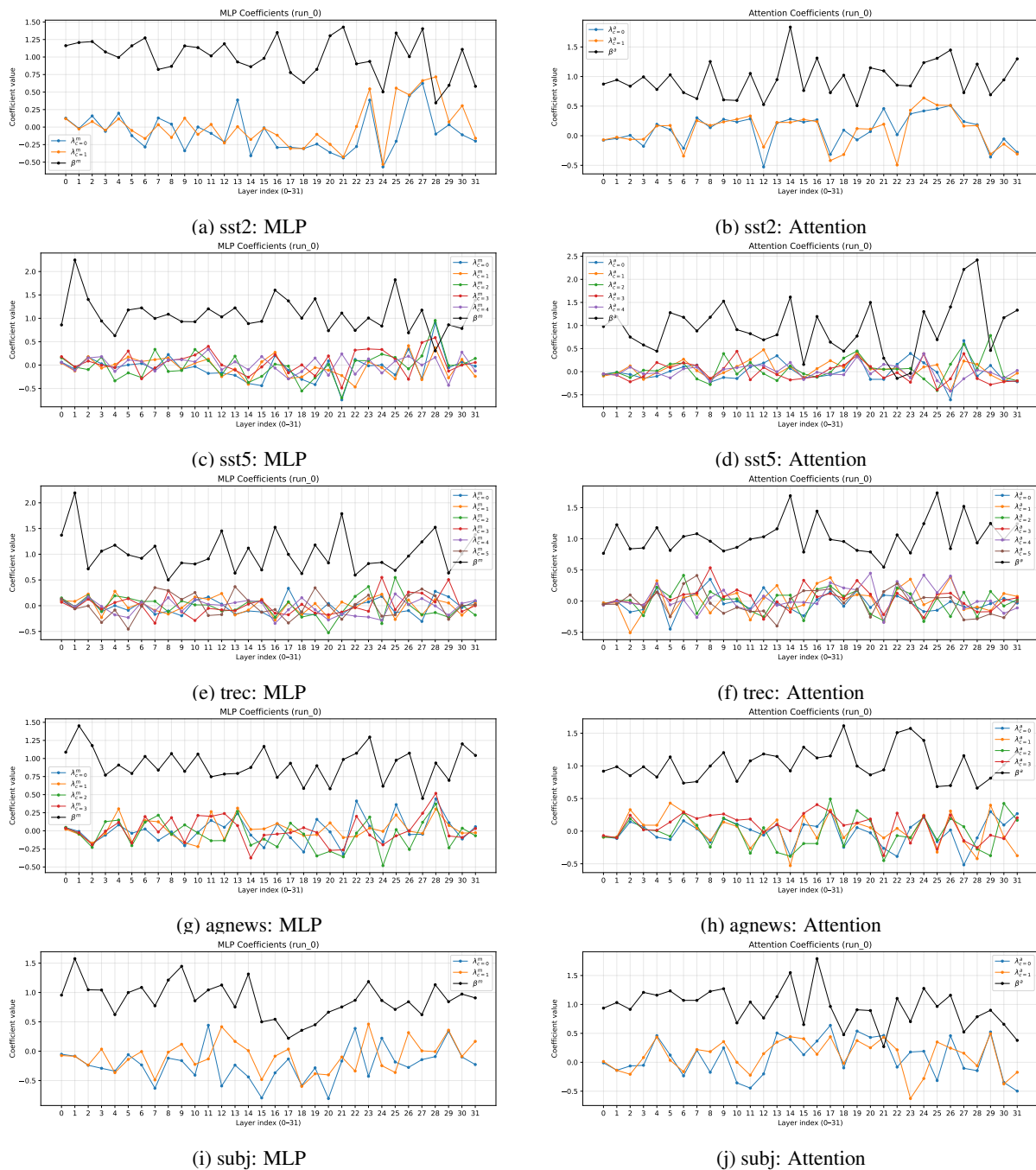
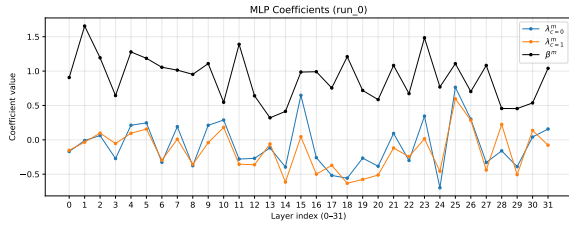
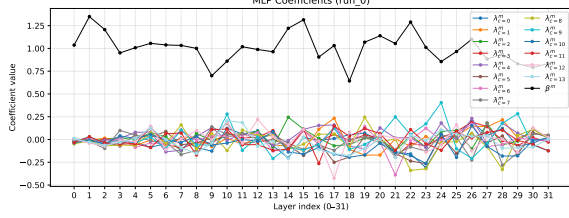


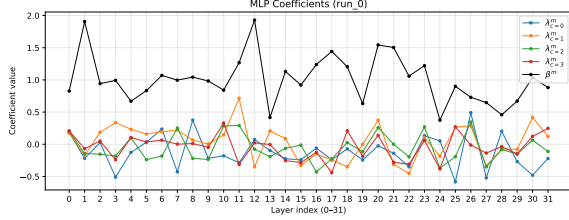
Figure 15: Layer-wise learned injection coefficients on LLaMA2-7B. Each row corresponds to a dataset; the left plot shows MLP coefficients, and the right plot shows attention coefficients. For each plot, the x-axis denotes the layer index, and the y-axis represents the learned coefficient values.



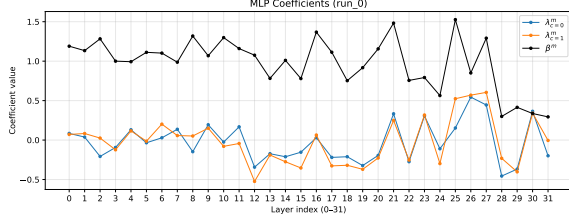
(k) hate\_speech18: MLP



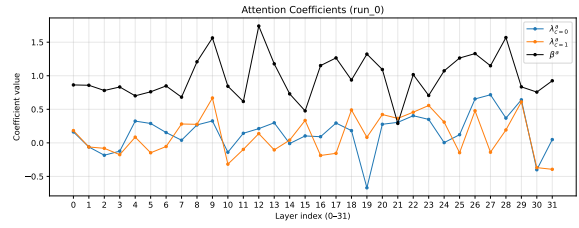
(m) dbpedia: MLP



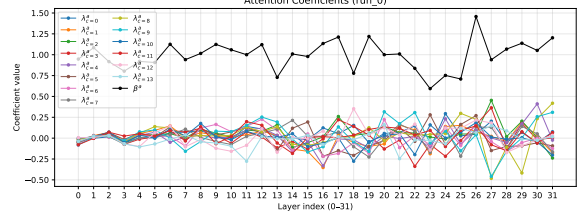
(o) emo: MLP



(q) mr: MLP



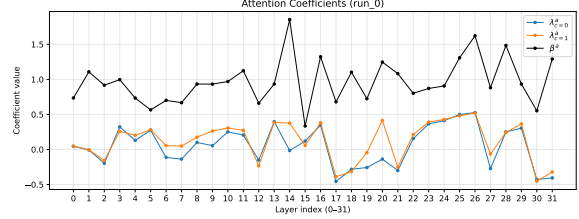
(l) hate\_speech18: Attention



(n) dbpedia: Attention



(p) emo: Attention



(r) mr: Attention

Figure 15: (continued)

<b>Dataset</b>	<b>Template</b>	<b>Label Space</b>
SST-2	Review: {Sentence} Sentiment: {Label}	negative / positive
SST-5	Sentence: {Sentence} Sentiment: {Label}	terrible / negative / neutral / positive / great
MR	Review: {Sentence} Sentiment: {Label}	negative / positive
Subj	Sentence: {Sentence} Label: {Label}	objective / subjective
DBPedia	Input: {Sentence} Label: {Label}	company / school / artist / athlete / politics / transportation / building / nature / village / animal / plant / album / film / book
AGNews	News: {Sentence} Type: {Label}	World / Sports / Business / Technology
TREC	Question: {Sentence} Answer Type: {Label}	Abbreviation / Entity / Person / Location / Number / Description
HateSpeech18	Text: {Sentence} Label: {Label}	neutral / hate
EmoC	Dialogue: {Sentence} Emotion: {Label}	others / happy / sad / angry

Table 11: Dataset templates and label spaces used in classification tasks.

---

**Algorithm 1: Class-Conditional Context Vectors ( $C^3V$ )**

---

**Input:** Frozen LM  $f_\theta$  with  $L$  layers; demonstrations  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$ ; class set  $\mathcal{C}$ ; token position for extraction  $e$ ; training steps  $TS$ ; learning rate  $\eta$

**Output:** Class-conditional context vectors  $\{\bar{d}_c\}_{c \in \mathcal{C}}$  and learned coefficients  $\mathbf{c}$

**Stage 1: Build class-conditional context vectors****foreach**  $c \in \mathcal{C}$  **do** $\mathcal{D}_c \leftarrow \{(x_i, y_i) \in \mathcal{D} \mid y_i = c\};$ **foreach**  $(x_i, y_i) \in \mathcal{D}_c$  **do**Run  $f_\theta$  on  $(x_i, y_i)$  and extract activations at position  $e$ ; $\{a_{i,l}^e, m_{i,l}^e\}_{l=1}^L;$ **for**  $l = 1$  **to**  $L$  **do** $\bar{a}_{c,l}^e \leftarrow \frac{1}{|\mathcal{D}_c|} \sum_{(x_i, y_i) \in \mathcal{D}_c} a_{i,l}^e;$  $\bar{m}_{c,l}^e \leftarrow \frac{1}{|\mathcal{D}_c|} \sum_{(x_i, y_i) \in \mathcal{D}_c} m_{i,l}^e;$  $\bar{d}_c \leftarrow \{\bar{a}_{c,l}^e, \bar{m}_{c,l}^e\}_{l=1}^L;$ **Stage 2: Initialize learnable coefficients**Initialize  $\lambda_{c,l}^a, \lambda_{c,l}^m \leftarrow 0.1/|\mathcal{C}|;$ Initialize  $\beta_l^a, \beta_l^m \leftarrow 1;$  $\mathbf{c} \leftarrow \{\lambda_{c,l}^a, \lambda_{c,l}^m, \beta_l^a, \beta_l^m\}_{c \in \mathcal{C}, l=1}^L;$ **Stage 3: Learn coefficients on demonstrations (backbone frozen)****for**  $ts = 1$  **to**  $TS$  **do**Sample mini-batch  $\mathcal{B} \subseteq \mathcal{D};$  $\mathcal{L} \leftarrow 0;$ **foreach**  $(x, y) \in \mathcal{B}$  **do**Run a *single* forward pass with class-conditional injection;**for**  $l = 1$  **to**  $L$  **do****foreach** token position  $t$  in the sequence **do** $\mathbf{r}_l^t \leftarrow \mathbf{r}_{l-1}^t + \sum_{c \in \mathcal{C}} \lambda_{c,l}^a \bar{a}_{c,l}^e + \beta_l^a a_l^t + \sum_{c \in \mathcal{C}} \lambda_{c,l}^m \bar{m}_{c,l}^e + \beta_l^m m_l^t;$ Compute  $P(y \mid x, \{\bar{d}_c\}, \mathbf{c})$  from the resulting logits; $\mathcal{L} \leftarrow \mathcal{L} - \log P(y \mid x, \{\bar{d}_c\}, \mathbf{c});$  $\mathcal{L} \leftarrow \mathcal{L}/|\mathcal{B}|;$ Update coefficients:  $\mathbf{c} \leftarrow \mathbf{c} - \eta \nabla_{\mathbf{c}} \mathcal{L};$ **Stage 4: Inference**Given query  $x$ , run one injected forward pass using  $\{\bar{d}_c\}$  and  $\mathbf{c}$ ;Compute class probabilities from label-token logits and output  $\hat{y} = \arg \max_{c \in \mathcal{C}} P(y_c \mid x, \{\bar{d}_c\}, \mathbf{c});$ 

---