

Beyond Single Representations: Embedding Fusion for Stable Text Classification

Anonymous ACL submission

Abstract

Embedding fusion has emerged as a powerful technique for enhancing performance across various NLP tasks. While prior research suggests that different layers of language models encode distinct representations and that pooling strategies influence performance, there is a lack of systematic analysis regarding the practical significance of these differences or the impact of combining embeddings from multiple models. This study provides a rigorous evaluation of layer-wise fusion strategies to determine their actual contribution to classification performance. Our findings reveal that the effectiveness of individual layers is more dependent on dataset characteristics than on the model architecture itself. Furthermore, we demonstrate that fusing embeddings from multiple models yields more robust and consistent representations across tasks, with the influence of any single model diminishing as the number of integrated models increases. Notably, experiments on low-resource datasets show that embedding fusion provides particularly significant gains when training data is scarce, highlighting its robustness and adaptability in data-constrained environments.

1 Introduction

With recent advancements in large language models (LLMs), the representational capacity of decoder-based architectures (Brown et al., 2020; Touvron et al., 2023a,b) has drawn increasing attention for downstream NLP tasks (Zhang et al., 2022; Sun et al., 2023). Despite their strong zero- and few-shot performance, these models are primarily trained for next-token prediction, leading to layer-wise differences in how semantic and contextual information is encoded. While final-layer embeddings are commonly used (Brown et al., 2020), prior studies in encoder-based models (Devlin et al., 2019) suggest that intermediate layers may yield richer representations for classification

(Zhang et al., 2024).

The availability of generative LLMs and specialized embedding models (Lee et al., 2025; Wang et al., 2024) makes embedding fusion a viable strategy for leveraging complementary knowledge. However, systematic guidelines for strategic layer selection and multi-model synergy remain limited. Specifically, the benefit of intermediate layers over the final layer—and how fusion scales across generative and specialized models—remains unexplored. To address these gaps, our contributions are as follows:

Layer-Aware Representation Analysis We quantitatively examine how representations vary across layers and how different pooling strategies, when considered alongside layer selection, influence classification performance. Furthermore, we assess whether consistent patterns arise in both dedicated embedding models and generative LLMs, establishing a comprehensive foundation for embedding fusion.

Embedding Fusion We show that relying solely on individual layer performance does not guaranty consistent stability, while fusing embeddings from multiple models produces more robust and reliable outcomes.

Adaptability in Low-Resource Settings Through experiments on low-resource datasets, we demonstrate that embedding fusion can operate efficiently even with limited data availability. These findings suggest that fusion strategies not only enhance representational robustness but also offer practical advantages in data-scarce environments.

2 Related Work

2.1 Text Classification and Pretrained Language Models

Text classification has evolved from traditional statistical methods to deep learning-based paradigms. Early research primarily utilized representations

like TF-IDF and n -grams with models such as SVM (Pang et al., 2002) and logistic regression (Zhang et al., 2003). The paradigm shifted significantly with pre-trained language models. Representative models like BERT (Devlin et al., 2019) and RoBERTa (Liu et al., 2019) capture rich contextual information through bidirectional architectures and large-scale pre-training. Fine-tuning these models on downstream tasks consistently yields superior performance compared to traditional approaches (Youngmin et al., 2024).

2.2 Extended Applications of LLM

The evolution of Large Language Models (LLMs) based on decoder architectures has enabled zero-shot (Kojima et al., 2023) and few-shot (Zhang et al., 2022) learning for downstream tasks such as text classification, spurring research on prompt-based approaches such as chain-of-thought (CoT) (Wei et al., 2023) and CARP (Sun et al., 2023). Numerous studies have reported on the performance of LLMs in classification tasks, including models such as GPT-3 (Brown et al., 2020) and the LLaMA series (Touvron et al., 2023a,b), as well as empirical studies analyzing their behavior (Sarkar et al., 2023; Gretz et al., 2023). However, these methods exhibit considerable performance variability depending on prompt design (Cao et al., 2024; He et al., 2024).

Meanwhile, there is growing interest in using LLMs not only as generative models but also as providers of high-quality embeddings (Tao et al., 2025). Recent research suggests that relatively lightweight LLMs (e.g. up to 7B parameters) can produce strong embedding quality with efficient computational resources (Wang et al., 2024; Lee et al., 2025). Furthermore, several studies have emphasized that different embedding layers within a model can produce optimal representations for tasks, with a particular focus on the importance of layer selection (Zhang et al., 2024). These findings highlight the role of intermediate representations in understanding the encoding behavior of transformer models in various applications (Skean et al., 2024).

2.3 Embedding Fusion

As diverse pretrained models, including large language models (LLMs), continue to emerge, there has been increasing interest in combining embeddings extracted from multiple models (Shinnou et al., 2018; Blandfort et al., 2019). Previous studies

have reported performance improvements on tasks such as text classification and sentiment analysis through embedding fusion.

For example, LLMEmbed (Liu et al., 2024) demonstrates that combining embeddings from LLaMA2 (Touvron et al., 2023b) with those from BERT and RoBERTa can effectively leverage the distinctive representational characteristics of each model. Moreover, a variety of fusion strategies have been proposed not only in NLP, but also in other domains for instance, QUARC (Kumar et al., 2020) applies quaternion-based operations.

However, significant performance differences arise depending on the fusion method employed, and not all combinations lead to consistent improvements (Ko et al., 2024). Furthermore, since each embedding layer possesses a different representational capacity, understanding the layer-wise characteristics is critical for effective fusion (Kaushik et al., 2024).

3 Methodology

In this study, we investigate embedding-based text classification from three complementary perspectives. First, we revisit whether previously reported layer-wise differences in the representations of large language models also manifest in text classification tasks and further examine whether similar patterns can be observed in embedding models.

Second, we assess the stability of performance across layers and show that this limitation can be alleviated through multi-model embedding fusion.

Finally, we investigate the effectiveness of embedding fusion in low-resource settings, analyzing whether the combined representations provide a significant performance boost when labeled data is limited. We aim to determine if fusion strategies can offer a practical, data-efficient solution for high-quality text classification. Our experiments employ both pretrained generative LLMs (up to 14B parameters) and specialized embedding models.

3.1 Layer and Pooling Strategies

Layer-wise Representation Previous research has shown that different layers of decoder-based LLMs capture different aspects of semantic and contextual information. In this section, we verify that these results apply equally to LLMs with different model architectures, and also to embedding models explicitly optimized for representation

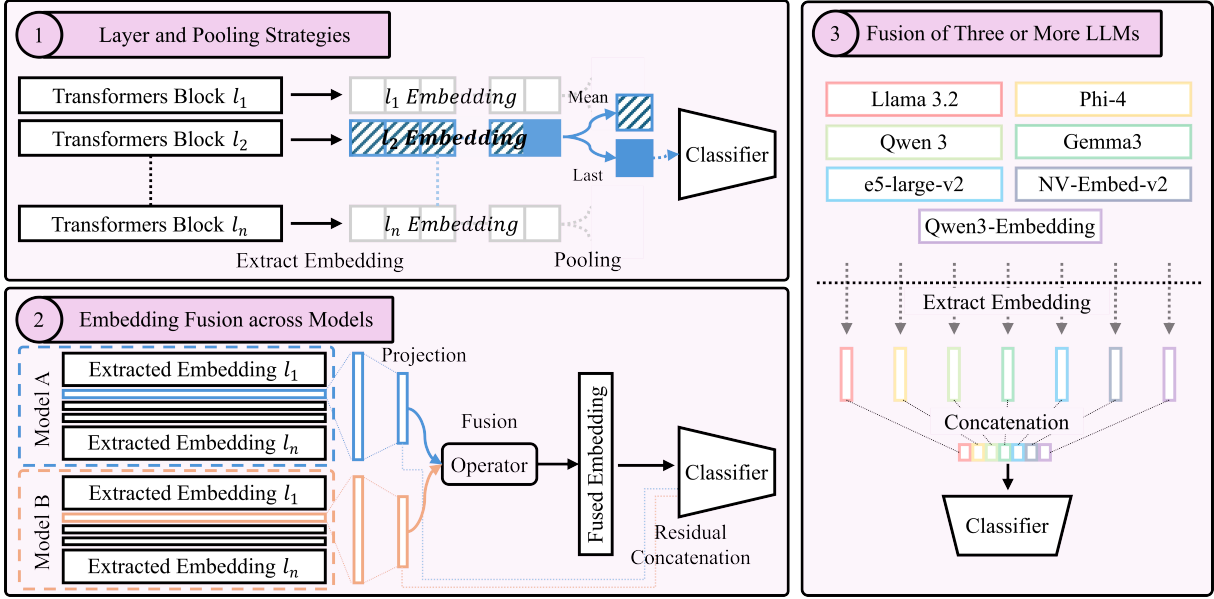


Figure 1: Overview of the experimental framework for evaluating embedding fusion. Here, l_n denotes the Transformer block layer, where n represents the final layer of each model. (1) For each model, embeddings extracted from every layer are aggregated using both mean and last pooling strategies and evaluated through a classifier to measure layer/pooling sensitivity. (2) Embeddings from different models are projected onto a shared dimensional space and fused through a predefined operator (e.g., concatenation, summation, or convolution) to produce unified representations. (3) All selected embeddings are concatenated, and the fused representations are classified to obtain the final performance for multi-model fusion analysis.

learning. To assess the layer-by-layer representational power for each text data set, we systematically evaluate embeddings extracted from all layers. This quantitative analysis, conducted on both generative LLMs and dedicated embedding models, validates the findings of previous studies (Skean et al., 2024).

Pooling Next, we examine how pooling strategies interact with layer selection. While embedding models explicitly define pooling strategies (Tang and Yang, 2024), generative models typically do not account for pooling when extracting embeddings. Thus, choosing an appropriate pooling strategy becomes essential. In particular, we compare last pooling, which directly uses the representation at the final token position, with mean pooling, which averages hidden states across the entire sequence. For each model and dataset, we systematically analyze whether the impact of pooling remains consistent across datasets within the same model and whether performance differences emerge.

3.2 Embedding Fusion across Models

Dimension Projection Embeddings extracted from multiple models may have different dimensions, so we apply a linear projection to unify them before fusion. Specifically, to project an embedding

$E \in \mathbb{R}^{d_1}$ to a target dimension d_2 , we use learnable parameters: a weight matrix $W \in \mathbb{R}^{d_2 \times d_1}$ and a bias vector $b \in \mathbb{R}^{d_2}$. These projected embeddings are then integrated using various fusion techniques.

The choice of d_2 involves a trade-off: larger dimensions preserve more information but increase computational cost, while smaller ones reduce overhead at the risk of information loss. In this study, we project onto the smaller dimension to improve efficiency.

Fusion Methods A key focus of this study is to enhance the representational capabilities of various models by combining embeddings extracted from different LLMs. While earlier sections analyze performance across different layers, here we fix the embeddings to those from the last layer in order to ensure consistency across models and to isolate the effect of different fusion strategies. Formally, we define the fusion process as a mapping

$$F : \{E_1, E_2, \dots, E_n\} \rightarrow E', \quad (1)$$

where $E_i \in \mathbb{R}^d$ are embeddings extracted from different models and $E' \in \mathbb{R}^{d'}$ is the fused embedding. The specific instantiations of F explored in this study are as follows:

- **Concatenation:** Directly concatenating the embedding vectors E along the matrix dimension to

form a single embedding:

$$E' = [E_1 \parallel E_2 \parallel \dots \parallel E_n] \quad (2)$$

- **Sum:** After aligning the dimensions of the embeddings, we form a new embedding by element-wise addition:

$$E' = \sum_{i=1}^n f(E_i) \quad (3)$$

- **Multiplication:** After aligning dimensions, embeddings are reshaped into 2D arrays of size $\sqrt{d} \times \sqrt{d}$ (where d must be a perfect square), and then combined via matrix multiplication:

$$\begin{aligned} E_1 \in \mathbb{R}^d &\rightarrow \text{reshape} \rightarrow E'_1 \in \mathbb{R}^{\sqrt{d} \times \sqrt{d}} \\ E_2 \in \mathbb{R}^d &\rightarrow \text{reshape} \rightarrow E'_2 \in \mathbb{R}^{\sqrt{d} \times \sqrt{d}} \\ E' &= E'_1 \cdot E'_2 \\ E'' &= \text{flatten}(E') \in \mathbb{R}^d \end{aligned} \quad (4)$$

- **Hadamard (Element-wise Product):** Embeddings are combined by multiplying corresponding elements at the same position across models.
- **Quaternion Fusion (Kumar et al., 2020):** Embeddings are treated as quaternion-valued vectors and fused using quaternion operations, preserving multidimensional inter-model relationships.
- **Residually Enhanced Fusion (Gardias et al., 2020):** The newly generated embedding is combined with the original one via residual connections to incorporate additional information while preserving the original expressiveness.

Layer-Aware Fusion Beyond the choice of fusion techniques, we further investigate whether selecting embeddings from dataset-specific optimal layers can enhance fusion performance. For each dataset, we first identify the most effective layers using the methodology described in Section 3.1. These layers are then used in fusion experiments involving single or multiple models. We compare the classification performance of these optimal-layer-based fusions against baselines that use default or last-layer embeddings. This analysis evaluates whether tailoring layer selection to each dataset improves accuracy and stability without requiring fine-tuning of the underlying models.

Table 1: Summary of experimental setups. The top section lists dataset statistics and the bottom section details the generation and embedding models used, including dimension (Dim), layers, and parameter size (Params).

Section 1: Dataset Statistics			
Dataset	Classes	Train	Test
SST-2	2	67,349	872
MR	2	40,000	10,000
R8	8	5,485	2,189
R52	52	6,532	2,568
Rotten Tomatoes	2	8,530	1,066
LegalBench	3	200	50
Section 2: Model Specifications			
Model Type	Dim	Layers	Params
<i>Generation</i>			
llama3.2	4096	28	3B
qwen3	4096	36	8B
gemma3	2304	48	12B
phi4	3584	40	14.7B
<i>Embedding</i>			
e5-large-v2	1024	24	0.335B
nv-embed-v2	4096	32	7.1B
qwen3-embedding	4096	24	8B

3.3 Fusion of Three or More LLMs

While the previous sections focused on combining embeddings from two models, this section investigates the performance impact when combining embeddings from three or more LLMs.

The motivation behind combining multiple models is to leverage their complementary strengths, potentially maximizing performance improvement. The key objectives of this experiment are twofold: first, to determine whether combining embeddings from three or more models leads to performance improvement; and second, to assess whether performance variance decreases as more models are combined, thus improving stability.

The fusion method used for combining embeddings in this experiment is primarily concatenation, and all possible combinations are tested. To eliminate external sources of variation, embeddings are fixed to the last layer, and mean pooling is applied consistently across all models.

4 Experiment

We evaluated our approach using six benchmark datasets: SST-2 (Socher et al., 2013), MR (Maas et al., 2011), and Rotten Tomatoes (Pang and Lee, 2005) for sentiment classification, R8 and R52 for topic classification, and a subset of LegalBench (Guha et al., 2023). Following the official guidelines, LegalBench was split into training and

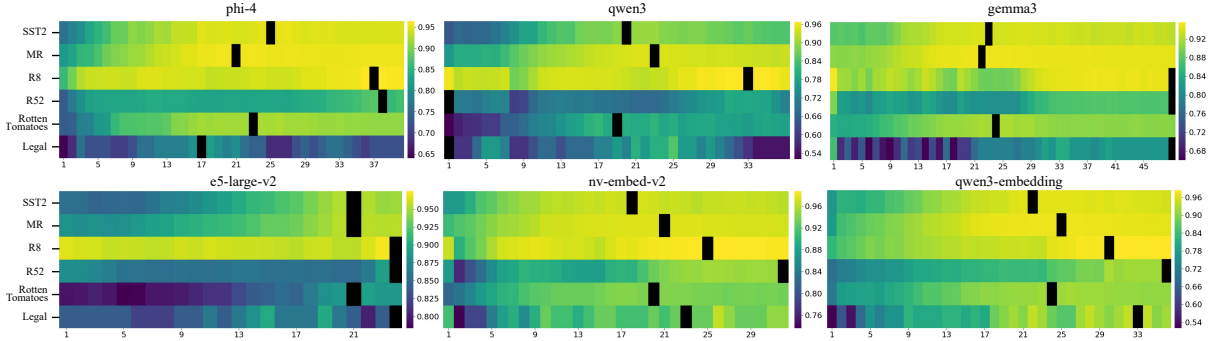


Figure 2: Comparison of Single-Layer Embedding Classification Performance in LLMs

Model	Layer	SST2		MR		R8		R52		Rotten Tomatoes		LegalBench			
		Acc	F_1	Acc	F_1	Acc	F_1	Acc	F_1	Acc	F_1	Acc	F_1		
Generation	llama3.2	last ₂₈	0.911	0.911	0.957	0.957	0.975	0.925	0.902	0.366	0.895	0.895	0.836	0.575	
		best _l	0.950 ₁₄	0.950 ₁₄	0.965 ₁₄	0.965 ₁₄	0.979 ₂₄	0.941 ₂₄	0.924 ₂₆	0.485 ₂₆	0.910 ₁₃	0.910 ₁₃	0.920 ₀	0.833 ₁₅	
	phi4	last ₄₀	0.944	0.944	0.960	0.960	0.978	0.938	0.908	0.406	0.909	0.909	0.860	0.592	
		best _l	0.954 ₂₅	0.954 ₂₅	0.966 ₂₁	0.966 ₂₁	0.980 ₃₉	0.943 ₁₃	0.933 ₃₀	0.583 ₃₀	0.923 ₂₁	0.923 ₂₁	0.960 ₁₉	0.912 ₁₉	
	qwen3	last ₃₆	0.943	0.943	0.958	0.958	0.967	0.867	0.858	0.196	0.911	0.911	0.780	0.535	
		best _l	0.951 ₂₄	0.951 ₂₄	0.965 ₂₂	0.965 ₂₂	0.975 ₃₂	0.916 ₃₀	0.896 ₀	0.389 ₀	0.919 ₂₅	0.919 ₂₅	0.920 ₀	0.634 ₂₄	
	gemma3	last ₄₈	0.953	0.953	0.959	0.959	0.976	0.927	0.920	0.480	0.923	0.923	0.888	0.764	
		best _l	0.964 ₂₃	0.964 ₂₃	0.967 ₂₁	0.967 ₂₁	0.978 ₄₈	0.931 ₄₈	0.922 ₄₈	0.492 ₄₈	0.926 ₂₆	0.926 ₂₆	0.900 ₄₈	0.778 ₄₈	
	Embedding	e5-large-v2	last ₂₄	0.945	0.945	0.953	0.953	0.973	0.922	0.895	0.357	0.891	0.891	0.916	0.629
			best _l	0.953 ₂₁	0.953 ₂₁	0.956 ₂₁	0.956 ₂₁	0.974 ₂₄	0.931 ₂₄	0.896 ₂₄	0.365 ₂₄	0.894 ₂₁	0.894 ₂₁	0.920 ₂₄	0.632 ₂₄
nv-embed-v2		last ₃₁	0.955	0.955	0.967	0.967	0.983	0.954	0.951	0.667	0.922	0.922	0.940	0.896	
		best _l	0.969 ₁₇	0.969 ₁₇	0.972 ₂₀	0.972 ₂₀	0.984 ₂₅	0.958 ₂₂	0.953 ₃₁	0.681 ₃₁	0.935 ₁₉	0.925 ₁₉	0.960 ₂₂	0.912 ₂₂	
qwen3-embedding		last ₃₆	0.949	0.949	0.964	0.964	0.982	0.946	0.950	0.664	0.906	0.906	0.920	0.884	
		best _l	0.964 ₂₂	0.964 ₂₂	0.972 ₂₅	0.972 ₂₅	0.983 ₃₀	0.952 ₃₁	0.952 ₃₆	0.699 ₃₆	0.922 ₂₄	0.922 ₂₄	0.960 ₃₃	0.911 ₃₃	

Table 2: Classification performance using embeddings extracted from specific layers. Each dataset reports both Accuracy and F1-score. Results for the last layer are averaged over five runs, while results for the best layer correspond to the best performance observed across the five runs.

test sets at an 8:2 ratio, while the original partitions were used for the other datasets. Detailed statistics of the datasets and the models employed are summarized in Table 1. The experiments involved a diverse range of generation and embedding models, including Llama3.2 (Grattafiori et al., 2024), Phi-4 (Abdin et al., 2024), Qwen3 (Yang et al., 2025), and Gemma3 (Team et al., 2025) as generative models, alongside e5-large-v2 (Wang et al., 2024), NV-Embed-v2 (Lee et al., 2025), and Qwen3-Embedding (Zhang et al., 2025). Specific configurations for each model, such as embedding dimensions (Dim), layer counts, and parameter scales ($Param$), are also detailed in the table.

Experiments were conducted on a system equipped with two NVIDIA RTX 4090 GPUs (24 GB memory each). A multilayer perceptron (MLP)-based classifier was applied for text classification with fused embedding vectors. Training was performed with batch size = 1024, learning rate = 1e-4, optimizer = Adam, and 120 epochs. Performance was evaluated using accuracy and F1-score.

4.1 Results by Layer Strategies

Layer-wise Representation Previous research has reported that embeddings extracted from intermediate layers often exhibit stronger representational power than those from the final layer. Building on the observation that generative language models become increasingly specialized for text generation in their final layers, we conducted a detailed quantitative evaluation to examine whether this pattern consistently holds across different layers and models. As shown in Figure 2, in some cases, generative and embedding models achieved peak performance in intermediate layers rather than in the final layer. However, this trend cannot be generalized. For example, in the case of Gemma3, despite being a generative model, the final layer performed best on three datasets. Interestingly, for Qwen3, the 0th layer achieved the highest accuracy on two datasets. Moreover, even in embedding models explicitly optimized for representation learning, the final layer did not always guarantee superior performance. Table 2 presents the numerical re-

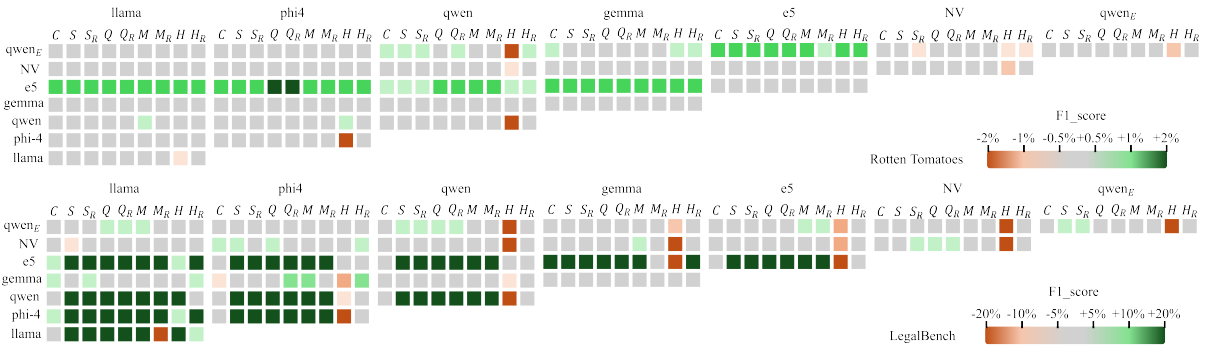


Figure 3: F1-score variations across layer combinations when combining two models on the Rotten Tomatoes and LegalBench datasets. The model abbreviations are as follows: llama(LLaMA-3.2), phi4(Phi-4), qwen(Qwen3), gemma(Gemma-12B-it), e5(e5-large-v2), NV(NV-Embed-v2), and qwen_E(Qwen3-Embedding). Fusion methods are denoted as C (Concatenation), S(Sum), S_R(Sum Residual), Q (Quaternion), Q_R(Quaternion Residual), M(Multiplication), M_R(Multiplication Residual), H(Hadamard), and H_R(Hadamard Residual).

	SST2	MR	RS	R52	Rotten Tomatoes	Legal Bench						
llama 3.2	0.942	0.942	0.952	0.952	0.962	0.881	0.873	0.291	0.895	0.895	0.812	0.554
	0.911	0.911	0.957	0.957	0.975	0.925	0.902	0.366	0.877	0.877	0.836	0.575
phi-4	0.944	0.944	0.955	0.955	0.960	0.883	0.860	0.247	0.909	0.909	0.720	0.489
	0.930	0.930	0.960	0.960	0.978	0.938	0.908	0.406	0.885	0.885	0.860	0.592
qwen3	0.943	0.943	0.954	0.954	0.938	0.703	0.820	0.158	0.911	0.911	0.708	0.484
	0.917	0.917	0.958	0.958	0.967	0.867	0.858	0.196	0.885	0.885	0.780	0.535
gemma3	0.953	0.953	0.958	0.958	0.958	0.874	0.871	0.305	0.923	0.923	0.780	0.547
	0.928	0.928	0.959	0.959	0.976	0.927	0.920	0.480	0.892	0.892	0.880	0.764
e5-large-v2	0.937	0.937	0.949	0.949	0.971	0.936	0.889	0.336	0.891	0.891	0.808	0.555
	0.946	0.946	0.953	0.953	0.973	0.922	0.895	0.357	0.882	0.882	0.916	0.629
nv-embed-v2	0.949	0.949	0.962	0.962	0.977	0.939	0.931	0.546	0.929	0.929	0.880	0.604
	0.955	0.955	0.967	0.967	0.984	0.954	0.952	0.679	0.923	0.923	0.940	0.896
qwen3 embedding	0.949	0.949	0.964	0.964	0.982	0.946	0.950	0.664	0.906	0.906	0.920	0.882
	0.928	0.928	0.961	0.961	0.977	0.932	0.929	0.524	0.888	0.888	0.920	0.884

Figure 4: Performance differences across pooling strategies for each model and dataset.

346 results, comparing accuracy and F1-scores between
 347 the final and the best-performing layers. Notably,
 348 selecting the optimal layer alone improved accuracy
 349 by up to 4% and F1-score by as much as 20%.

350 **Pooling** Experimental results for pooling
 351 strategies across models and datasets are pre-
 352 sented in Figure 4. Each block represents:
 353 $\frac{\text{last}/acc}{\text{mean}/acc} \mid \frac{\text{last}/F_1}{\text{mean}/F_1}$ Here, **last** refers to
 354 pooling from the final **token position** along the
 355 sequence dimension, not to the last layer of the
 356 model. During the experiments, tokenizer padding
 357 positions were fixed according to each model’s
 358 baseline configuration, and only the pooling
 359 strategy was varied. To ensure a fair comparison,
 360 all experiments used the final-layer representations.
 361 Green blocks indicate the pooling method that
 362 achieved higher performance, while purple blocks

363 represent cases where the performance gap
 364 between pooling strategies exceeded 6%. Orange
 365 blocks highlight unusual cases in which the choice
 366 of pooling method led to different trends between
 367 F1-score and accuracy.

368 On average, more than a 1% difference in per-
 369 formance was observed between pooling strategies,
 370 and in the LegalBench dataset, the gap exceeded
 371 20%. The most notable observation arises from
 372 the distributional differences in generative mod-
 373 els. Despite having distinct architectures, the four
 374 generative models did not exhibit model-specific
 375 consistency. Instead, the effect of pooling tended
 376 to generalize across datasets rather than across
 377 models. In contrast, embedding models showed
 378 consistent behavior across model types.

4.2 Fusion Strategy Experiments 379

380 This section evaluates various embedding fusion
 381 strategies using the models introduced earlier. To
 382 analyze the pure effect of fusion itself, embeddings
 383 were consistently extracted from the final layer of
 384 each model, as in the pooling experiments. The per-
 385 formance gain or loss of each fused embedding was
 386 assessed by comparing it with the higher F1-score
 387 obtained from the two individual models using their
 388 single-layer embeddings.

389 Figure 3 presents the results of model combi-
 390 nations on two datasets, Rotten Tomatoes, where
 391 accuracy has not yet saturated, and LegalBench,
 392 which suffers from limited data resources and
 393 shows instability in both accuracy and F1-score.
 394 As expected, simple concatenation between iden-
 395 tical models did not lead to any improvement in
 396 either dataset. In the Rotten Tomatoes dataset, em-
 397 bedding fusion with the e5 model achieved more

Table 3: Performance comparison between single-model baselines and two-model fusion methods across six text classification datasets. l indicates the optimal layer index for each dataset.

Dataset	Models	Layer	Fusion	ACC	F_1	$\Delta(\%)$
SST2	NV-Embed-v2	–	–	0.955	0.955	
	NV-Embed-v2, e5-large-v2	17, 21	Hadamard(R)	0.9748	0.9748	+1.98%
MR	Qwen3-Embedding	–	–	0.964	0.964	
	NV-Embed-v2, e5-large-v2	20, 21	Hadamard	0.9738	0.9738	+0.98%
R8	NV-Embed-v2	–	–	0.983	0.954	
	NV-Embed-v2, Qwen3-Embedding	25, 30	concatenation	0.9863	0.9607	+0.33%, +0.67%
R52	NV-Embed-v2	–	–	0.951	0.667	
	NV-Embed-v2, Qwen3-Embedding	31, 36	Hadamard(R)	0.9599	0.7702	+0.82%, +10.32%
Rotten Tomatoes	NV-Embed-v2	–	–	0.922	0.922	
	NV-Embed-v2, gemma3-12b-it	19, 26	Sum(R)	0.9371	0.9371	+1.51%
LegalBench	Qwen3-Embedding	–	–	0.940	0.896	
	NV-Embed-v2, e5-large-v2	22, 24	Multiplication(R)	0.9800	0.9270	+4.0%, +3.1%

Table 4: Classification accuracy and F1-score with multi-model fusion (≥ 3 models). $N(C)$ denotes the number of possible combinations of n models selected from the seven models, without repetition or order. All models use embeddings extracted from the last layer with mean pooling.

Models	SST2		MR		R8		R52		Rotten Tomatoes		LegalBench	
	Acc	F_1	Acc	F_1	Acc	F_1	Acc	F_1	Acc	F_1	Acc	F_1
Single	0.911-0.955	0.911-0.955	0.953-0.967	0.953-0.967	0.967-0.983	0.916-0.954	0.858-0.951	0.196-0.667	0.891-0.922	0.891-0.922	0.780-0.940	0.535-0.896
3(35)	0.930-0.961	0.930-0.961	0.961-0.971	0.961-0.971	0.978-0.984	0.932-0.956	0.921-0.956	0.461-0.721	0.891-0.925	0.891-0.924	0.860-0.980	0.592-0.927
4(35)	0.934-0.963	0.934-0.963	0.961-0.971	0.961-0.971	0.978-0.985	0.932-0.958	0.934-0.957	0.545-0.720	0.894-0.925	0.894-0.925	0.880-0.98	0.620-0.927
5(21)	0.941-0.962	0.941-0.962	0.962-0.971	0.962-0.971	0.980-0.984	0.942-0.955	0.939-0.956	0.589-0.722	0.899-0.925	0.899-0.925	0.920-0.980	0.884-0.927
6(7)	0.956-0.963	0.956-0.963	0.965-0.971	0.965-0.971	0.982-0.984	0.945-0.953	0.946-0.956	0.639-0.724	0.906-0.925	0.906-0.925	0.940-0.980	0.898-0.927
7(1)	0.958	0.958	0.971	0.971	0.982-0.983	0.949-0.951	0.956	0.697	0.924-0.925	0.924-0.925	0.940-0.980	0.898-0.927

than a 0.5% performance gain in all cases except for combinations involving NV embeddings. Other combinations also showed slight improvements, but in most cases, the performance increase was not substantial.

These findings suggest that, in embedding fusion, the improvement in classification performance arises not merely from the inherent strength of a single model but from the integration of models with distinct representational spaces, which produces richer and more stable feature representations. Meanwhile, in the LegalBench dataset, despite its unstable distribution, most fusion cases exhibited over 20% improvement in F1-score, highlighting that the fusion effect becomes more pronounced under limited-data conditions.

4.3 Optimal-Layer Fusion Analysis

According to the previous experiments, we confirmed that selecting an appropriate embedding layer alone can improve classification performance, and that combining embeddings from different models can further enhance it. However, since the optimal embedding layer varies depending on dataset characteristics and model architecture, it is

practically infeasible to universally determine and fuse the optimal layers in real-world applications. In this section, we therefore assume that the optimal layer of each model is known in advance and analyze the theoretical upper bound of performance improvement that embedding fusion can achieve under such ideal conditions.

Table 3 presents the results of combining the two embeddings that achieved the highest performance for each dataset, compared with the best-performing single embedding. Although the observed improvements are marginal, these results suggest that layer-aware embedding fusion can potentially overcome the performance plateau inherent to single-embedding models.

4.4 Fusion of Three or More Models

We extend our analysis to the fusion of three or more models to examine how classification performance scales with the number of integrated representations. Following the previous setting, all embeddings are consistently taken from the final layer to focus on the effect of fusion itself. For each configuration, we report the minimum and maximum accuracy and F1 scores. Table 4 summarizes the results for all multi-model fusion settings, where the

447 number in parentheses indicates the total number
448 of possible model combinations. The "Single" row
449 represents the range of performance observed from
450 individual models. As the number of fused models
451 increases, both accuracy and F1 score generally
452 improve across datasets, suggesting that integrat-
453 ing heterogeneous embeddings leads to more stable
454 and reliable representations.

455 In this context, stability refers not to the improve-
456 ment in maximum performance but to the notable
457 increase in the minimum performance. This effect
458 is most pronounced in the R52 and LegalBench
459 datasets, where the minimum F1 score increases
460 substantially from 0.196 to 0.697 and from 0.535 to
461 0.898, respectively. Similar patterns are observed
462 across the remaining datasets, consistently showing
463 that the lower bound of performance shifts upward
464 with larger fusion sets.

465 Overall, these findings indicate that multi-model
466 fusion enhances representational robustness by mit-
467 igating dependence on design choices such as layer
468 selection or pooling strategy, thereby producing
469 consistently improved results even without exten-
470 sive tuning.

471 5 Conclusion

472 Previous studies show that language models exhibit
473 distinct layer-wise representations, yet it remains
474 unclear if this consistently manifests in text clas-
475 sification across both generative and embedding-
476 oriented LLMs. Furthermore, the existence of a
477 universally optimal pooling strategy remains uncer-
478 tain.

479 Our systematic analysis of these two factors re-
480 veals that while certain layers capture richer fea-
481 tures, this does not always improve classification;
482 in some cases, initial layers even yield peak per-
483 formance. In generative models, pooling effectiveness
484 is driven more by dataset characteristics than by
485 model architecture.

486 Building on these findings, we demonstrate that
487 multi-model embedding fusion enhances stability
488 and generalization by alleviating dependence on
489 specific design choices, such as layer or pooling se-
490 lection. This efficiency is especially pronounced in
491 low-resource environments, where fusion mitigates
492 external factors and provides more reliable repre-
493 sentations. Ultimately, our results offer new empir-
494 ical evidence for the performance gains achievable
495 through layer-aware fusion, underscoring its impor-
496 tance for future model design.

6 Limitation

497 While this study provides a systematic analysis
498 of layer selection, pooling strategies, and embed-
499 ding fusion, several limitations should be acknowl-
500 edged. First, our investigation was limited to text
501 classification tasks. Future work could expand this
502 scope to other downstream objectives such as re-
503 trieval, clustering, or text generation, offering a
504 more comprehensive view of how embedding fu-
505 sion behaves across different task types. Second,
506 the proposed evaluation framework focused on
507 static, task-agnostic fusion rather than incorporat-
508 ing adaptive or learnable weighting mechanisms
509 that could dynamically optimize contributions from
510 multiple models. Finally, although our findings
511 demonstrate the effectiveness of embedding fusion
512 in low resource settings, a more extensive evalua-
513 tion on domain shifted and noisy data is needed to
514 further validate its robustness and generalization.
515

7 Future Work

516 Future research may explore adaptive fusion frame-
517 works that dynamically modulate the contribution
518 of each model or layer based on task characteristics
519 and data properties. Another promising direction is
520 to extend layer-aware embedding fusion to multi-
521 modal contexts by integrating textual, visual, and
522 structured representations, thereby enhancing repre-
523 sentational diversity and transferability. In addition,
524 investigating how the efficiency of embedding fu-
525 sion scales with model size and architecture diver-
526 sity could offer valuable insight into the theoretical
527 and practical limits of representation integration.
528

References

- 529
530 Marah Abdin, Jyoti Aneja, Harkirat Behl, Sébastien
531 Bubeck, Ronen Eldan, Suriya Gunasekar, Michael
532 Harrison, Russell J. Hewett, Mojan Javaheripi, Piero
533 Kauffmann, James R. Lee, Yin Tat Lee, Yuanzhi Li,
534 Weishung Liu, Caio C. T. Mendes, Anh Nguyen,
535 Eric Price, Gustavo de Rosa, Olli Saarikivi, and
536 8 others. 2024. [Phi-4 technical report](#). *Preprint*,
537 arXiv:2412.08905.
- 538 Philipp Blandfort, Tushar Karayil, Federico Raue, Jörn
539 Hees, and Andreas Dengel. 2019. [Fusion strategies
540 for learning user embeddings with neural networks](#).
541 *Preprint*, arXiv:1901.02322.
- 542 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie
543 Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind
544 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
545 Askell, and 1 others. 2020. [Language models are
546 few-shot learners](#). *Preprint*, arXiv:2005.14165.

547	Bowen Cao, Deng Cai, Zhisong Zhang, Yuexian Zou, and Wai Lam. 2024. On the worst prompt performance of large language models . <i>Preprint</i> , arXiv:2406.10248.	601
548		602
549		603
550		604
551	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding . <i>Preprint</i> , arXiv:1810.04805.	605
552		606
553		607
554		608
555	Przemek Gardias, Eric Arthur, and Huaming Sun. 2020. Enhanced residual networks for context-based image outpainting . <i>Preprint</i> , arXiv:2005.06723.	610
556		611
557		612
558	Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, and 542 others. 2024. The llama 3 herd of models . <i>Preprint</i> , arXiv:2407.21783.	613
559		614
560		615
561		616
562		617
563		618
564		619
565		620
566	Shai Gretz, Alon Halfon, Ilya Shnayderman, Orith Toledo-Ronen, Artem Spector, Lena Dankin, Yanis Katsis, Ofir Arviv, Yoav Katz, Noam Slonim, and Liat Ein-Dor. 2023. Zero-shot topical text classification with LLMs - an experimental study . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 9647–9676, Singapore. Association for Computational Linguistics.	621
567		622
568		623
569		624
570		625
571		626
572		627
573		628
574	Neel Guha, Julian Nyarko, Daniel E. Ho, Christopher Ré, Adam Chilton, Aditya Narayana, Alex Chohlas-Wood, Austin Peters, Brandon Waldon, Daniel N. Rockmore, Diego Zambrano, Dmitry Talisman, Enam Hoque, Faiz Surani, Frank Fagan, Galit Sarfaty, Gregory M. Dickinson, Haggai Porat, Jason Hegland, and 21 others. 2023. Legalbench: A collaboratively built benchmark for measuring legal reasoning in large language models . <i>Preprint</i> , arXiv:2308.11462.	629
575		630
576		631
577		632
578		633
579		634
580		635
581		636
582		637
583		638
584	Jia He, Mukund Rungta, David Koleczek, Arshdeep Sekhon, Franklin X. Wang, and Sadid Hasan. 2024. Does prompt formatting have any impact on llm performance? <i>Preprint</i> , arXiv:2411.10541.	639
585		640
586		641
587		642
588	Arjun Ramesh Kaushik, Sunil Rufus R. P, and Nalini Ratha. 2024. Enhancing authorship attribution through embedding fusion: A novel approach with masked and encoder-decoder language models . <i>Preprint</i> , arXiv:2411.00411.	643
589		644
590		645
591		646
592		647
593	Young Su Ko, Jonathan Parkinson, and Wei Wang. 2024. Benchmarking text-integrated protein language model embeddings and embedding fusion on diverse downstream tasks . <i>bioRxiv</i> .	648
594		649
595		650
596		651
597	Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2023. Large language models are zero-shot reasoners . <i>Preprint</i> , arXiv:2205.11916.	652
598		653
599		654
600		655
		656
	Deepak Kumar, Nalin Kumar, and Subhankar Mishra. 2020. Quarc: Quaternion multi-modal fusion architecture for hate speech classification . <i>Preprint</i> , arXiv:2012.08312.	
	Chankyu Lee, Rajarshi Roy, Mengyao Xu, Jonathan Raiman, Mohammad Shoeybi, Bryan Catanzaro, and Wei Ping. 2025. Nv-embed: Improved techniques for training llms as generalist embedding models . <i>Preprint</i> , arXiv:2405.17428.	
	Chun Liu, Hongguang Zhang, Kainan Zhao, Xinghai Ju, and Lin Yang. 2024. Llmembed: Rethinking lightweight llm’s genuine function in text classification . <i>Preprint</i> , arXiv:2406.03725.	
	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach . <i>Preprint</i> , arXiv:1907.11692.	
	Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. 2011. Learning word vectors for sentiment analysis . In <i>Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies</i> , pages 142–150, Portland, Oregon, USA. Association for Computational Linguistics.	
	Bo Pang and Lillian Lee. 2005. Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales . In <i>Proceedings of the ACL</i> .	
	Bo Pang, Lillian Lee, and Shivakumar Vaithyanathan. 2002. Thumbs up? sentiment classification using machine learning techniques . <i>Preprint</i> , arXiv:cs/0205070.	
	Souvika Sarkar, Dongji Feng, and Shubhra Kanti Kar-maker Santu. 2023. Zero-shot multi-label topic inference with sentence encoders and LLMs . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 16218–16233, Singapore. Association for Computational Linguistics.	
	Hiroyuki Shinnou, Xinyu Zhao, and Kanako Komiya. 2018. Domain adaptation using a combination of multiple embeddings for sentiment analysis . In <i>Proceedings of the 32nd Pacific Asia Conference on Language, Information and Computation</i> , Hong Kong. Association for Computational Linguistics.	
	Oscar SKEAN, Md Rifat Arefin, Yann LeCun, and Ravid Shwartz-Ziv. 2024. Does representation matter? exploring intermediate layers in large language models . <i>Preprint</i> , arXiv:2412.09563.	
	Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank . In <i>Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing</i> , pages	

