

# Each to Their Own: Exploring the Optimal Embedding in RAG

Shiting Chen<sup>1\*</sup>, Zijian Zhao<sup>2</sup>, Jinsong Chen<sup>1\*</sup>

<sup>1</sup>Faculty of Education, University of Hong Kong, Hong Kong, China

<sup>2</sup>Department of Civil and Environmental Engineering,

The Hong Kong University of Science and Technology, Hong Kong, China

## Abstract

Recently, as Large Language Models (LLMs) have fundamentally impacted various fields, the methods for incorporating up-to-date information into LLMs or adding external knowledge to construct domain-specific models have garnered wide attention. Retrieval-Augmented Generation (RAG), serving as an inference-time scaling method, is notable for its low cost and minimal effort for parameter tuning. However, due to heterogeneous training data and model architecture, the variant embedding models used in RAG exhibit different benefits across various areas, often leading to different similarity calculation results and, consequently, varying response quality from LLMs. To address this problem, we propose and examine two approaches to enhance RAG by combining the benefits of multiple embedding models, named Mixture-Embedding RAG and Confident RAG. Mixture-Embedding RAG simply sorts and selects retrievals from multiple embedding models based on standardized similarity; however, it does not outperform vanilla RAG. In contrast, Confident RAG generates responses multiple times using different embedding models and then selects the responses with the highest confidence level, demonstrating average improvements of approximately 10% and 5% over vanilla LLMs and RAG, respectively. The consistent results across different LLMs and embedding models indicate that Confident RAG is an efficient plug-and-play approach for various domains. We will release our code upon publication.

## 1 Introduction

Large language models (LLMs) have recently accelerated the pace of transformation across multiple fields, including transportation (Lyu et al., 2025), arts (Zhao et al., 2025), and education (Gao et al., 2024), through various paradigms such as direct answer generation, training from scratch on different types of data, and fine-tuning on target domains. However, the hallucination problem (Henkel et al., 2024) associated with LLMs has confused people for a long time, stemming from multiple factors such as a lack of knowledge on the given prompt (Huang et al., 2025b) and a biased training process (Zhao, 2025).

Serving as a highly efficient solution, Retrieval-Augmented Generation (RAG) has been widely employed in constructing foundation models (Chen et al., 2024) and practical agents (Arslan et al., 2024). Compared to training methods like fine-tuning and prompt-tuning, its plug-and-play feature makes RAG an efficient, simple, and cost-effective approach. The main paradigm of RAG involves first calculating the similarities between a question and chunks in an external knowledge corpus, followed by incorporating the top  $K$  relevant chunks into the prompt to guide the LLMs (Lewis et al., 2020).

Despite the advantages of RAG, selecting the appropriate embedding models remains a crucial concern, as the quality of retrieved references directly influences the generation results of the LLM (Tu et al., 2025). Variations in training data and model architecture lead to different embedding models providing benefits across various domains. The differing similarity calculations across embedding models often leave researchers uncertain about how to choose the optimal one. Consequently, improving the accuracy of RAG from the perspective of embedding models continues to be an ongoing area of research.

---

\* Corresponding Author: Shiting Chen (u3011355@connect.hku.hk), Jinsong Chen (jinsong.chen@live.com)

To address this research gap, we propose two methods for improving RAG by combining the benefits of multiple embedding models. The first method is named Mixture-Embedding RAG, which sorts the retrieved materials from multiple embedding models based on normalized similarity and selects the top  $K$  materials as final references. The second method is named Confident RAG, where we first utilize vanilla RAG to generate answers multiple times, each time employing a different embedding model and recording the associated confidence metrics, and then select the answer with the highest confidence level as the final response. By validating our approach using multiple LLMs and embedding models, we illustrate the superior performance and generalization of Confident RAG, even though Mixture-Embedding RAG may lose to vanilla RAG. The main contributions of this paper can be summarized as follows:

- We first point out that in RAG, different embedding models operate within their own prior domains. To leverage the strengths of various embedding models, we propose and test two novel RAG methods: Mixture-Embedding RAG and Confident RAG. These methods effectively utilize the retrieved results from different embedding models to their fullest extent.
- While Mixture-Embedding RAG performs similarly to vanilla RAG, the Confident RAG method exhibits superior performance compared to both the vanilla LLM and vanilla RAG, with average improvements of 9.9% and 4.9%, respectively, when using the best confidence metric. Additionally, we discuss the optimal number of embedding models for the Confident RAG method based on the results.
- Our results reveal two outstanding confidence metrics: self-certainty and Distributional Perplexity (DP), both showing average improvements of approximately 10% compared to the vanilla LLM. Specifically, among the LLMs examined, self-certainty achieves a maximum increase of 10.4%, while DP demonstrates a maximum increase of 12.4% compared to the vanilla LLM. The reasons behind the better performance of these two metrics are discussed based on their formulas.

## 2 Preliminary and Related Work

### 2.1 Retrieval Augmented Generation (RAG)

When constructing domain-specific foundation models or private agents, modifying a trained LLM is often necessary. However, the computational and resource costs associated with full fine-tuning render it impractical for a large number of users. Even though techniques such as Low-Rank Adaptation (LoRA) and prefix-tuning provide lighter training paradigms, RAG offers its own advantages in terms of plug-and-play capabilities. Furthermore, RAG allows external knowledge to remain independent of LLMs, facilitating the replacement of databases with other domains or more up-to-date versions.

Currently, a common paradigm of vanilla RAG is illustrated in Fig. 1. Given a question  $q$ , an external corpus  $C$  divided into chunks as  $C = [c_1, c_2, \dots, c_m]$ , and an embedding model  $g()$ , we first generate the embeddings for the question and the corpus, respectively, as follows:

$$\begin{aligned} e_q &= g(q), \\ e_{c_j} &= g(c_j), \end{aligned} \quad (1)$$

where  $e_q, e_{c_j} \in \mathbb{R}^{d_g}$  are the embedding results for question  $q$  and chunk  $c_j$ , and  $d_g$  is the output dimension of the model  $g()$ .

Based on the understanding that two sentences with higher similarity yield more similar embeddings, we calculate the similarity between the question and each chunk using the cosine similarity metric, defined as:

$$s_{q, c_j} = \frac{e_q \cdot e_{c_j}}{\|e_q\| \|e_{c_j}\|}, \quad (2)$$

where  $s_{q, c_j}$  represents the cosine similarity between question  $q$  and chunk  $c_j$ ,  $\cdot$  denotes the dot product, and  $\|\cdot\|$  denotes the Euclidean norm. Finally, we select  $k$  chunks with the highest similarity and incorporate them into the prompt, expressed as:

$$h(t, q, c_j | j \in \mathcal{K}), \quad (3)$$

where  $h()$  is a montage function that combines the question and the retrieved chunks using a specified template  $t$ , and  $\mathcal{K}$  represents the index set of the selected  $k$  chunks. In this way, we construct a prompt that includes external knowledge from the corpus  $C$  that is most similar to question  $q$ .

Previous studies have validated a wide range of modifications to the RAG pipeline, demonstrating

effectiveness in open-domain question answering by improving various components such as adaptive retrieval strategies (Tang et al., 2024; Guan et al., 2024), dynamic document selection (Kim and Lee, 2024; Hei et al., 2024), joint retriever-generator optimization (Siriwardhana et al., 2021), reinforcement learning (Huang et al., 2025a), and hybrid retrieval-reranking frameworks (Sawarkar et al., 2024; Zhang et al., 2025).

Since embedding models can produce different similarity calculations, many studies have innovated in retriever embedding models to improve the performance of the RAG pipeline. For example, Invar-RAG (Liu et al., 2024) constructs an LLM-based retriever with LoRA-based representation learning and an invariance loss, directly addressing retriever embedding variance and locality, and demonstrates improved retrieval and response accuracy. Blended RAG (Sawarkar et al., 2024) combines dense semantic embeddings with sparse retrievers, empirically showing that this hybrid embedding space achieves higher QA accuracy, especially as corpus size scales. W-RAG (Nian et al., 2024) further validates that enhancing retriever embeddings using weak signals from the downstream QA task offers gains comparable to methods requiring full human annotations.

However, most studies have focused on the innovation of a single embedding model. Due to variations in model architecture and training datasets, different embedding models operate effectively in different domains. Our approach will leverage the strengths of various embedding models to provide the most relevant information for LLMs.

## 2.2 Confidence of LLMs

The confidence level of the LLM output can indicate the response quality to some extent. Several approaches have been employed to measure model confidence. Some of them depend on model prompting, external similarity or model-expressed judgments. These include linguistic confidence measures (Shrivastava et al., 2023; Xiong et al., 2023), where the model is prompted to state its confidence explicitly as part of the output, as well as LLM-as-judge or self-evaluation strategies that involve prompting the model to rate or compare its own answers in scalar or relative terms (Shrivastava et al., 2025; Ren et al., 2023; Hager et al., 2025). However, these methods are often poorly

calibrated, not directly related to the true predictive uncertainty of the model, and can be biased or inconsistent due to their reliance on language generation rather than statistical likelihood.

In contrast, probability-based metrics directly ground into models’ token probability output (Chen and Mueller, 2023; Jiang et al., 2021; Lin et al., 2024). These include average log-probability as the mean token log-probabilities, entropy-based measures that quantify prediction uncertainty via sequence or per-token entropy, self-certainty (Kang et al., 2025) that utilizes the intrinsic probability distribution of LLM outputs to assess the model confidence.

Currently, generative LLMs operate based on an auto-regressive process, formulated as:

$$y_i \sim f(x, y_{<i}) , \quad (4)$$

where  $f()$  is an LLM,  $x$  is the prompt,  $y_{<i}$  represents the first generated  $i-1$  tokens, and  $y_i$  denotes the  $i^{th}$  token. The output of  $f()$  is the probability distribution over each word in the vocabulary, from which  $y_i$  is sampled according to this probability vector. This process can be viewed as a Markov Decision Process (MDP), where the state consists of the prompt  $x$  along with the currently generated tokens  $y_{<i}$ , and the action space corresponds to the vocabulary, allowing each token to be chosen with a probability given by

$$p(v_j|x, y_{<i}; \pi_f) = f(x, y_{<i})[j] , \quad (5)$$

where  $v_j$  is the  $j^{th}$  element in the vocabulary  $v$ , and  $\pi_f$  represents the strategy. For simplicity, we will neglect the  $\pi_f$  term in the following discussion.

Many studies have observed that the probability of each token reflects the LLMs’ confidence. Intuitively, if the selected token has a higher chosen probability  $p(y_i|x, y_{<i})$ , or if the generated probability vector  $f(x, y_{<i})$  exhibits a less uniform distribution, the LLMs demonstrate higher confidence in their generated results. Our study will use the confidence measurement metrics summarized in (Kang et al., 2025):

- **Average Log-Probability (AvgLogP):** This metric computes the mean logarithmic likelihood of the generated tokens, reflecting the model’s confidence in the entire sequence:

$$\text{AvgLogP} = \frac{1}{n} \sum_{i=1}^n \log p(y_i|x, y_{<i}) , \quad (6)$$

where  $n$  represents the answer length. Higher log-probabilities suggest more reliable predictions.

- **Gini Impurity (Gini):** This metric measures the concentration of the predicted distribution:

$$\text{Gini} = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^{|v|} (p(v_j|x, y_{<i}))^2, \quad (7)$$

where  $|v|$  represents the vocabulary size. Higher values indicate more peaked distributions, reflecting greater certainty in the model's predictions.

- **Entropy:** This metric computes the uncertainty of the predicted distribution:

$$\text{Entropy} = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^{|v|} -p(j|x, y_{<i}) \log p(v_j|x, y_{<i}). \quad (8)$$

Lower entropy signifies higher confidence, as it indicates a more deterministic output distribution.

- **Distributional Perplexity (DP):** This metric generalizes perplexity to the entire predicted distribution:

$$\text{DP} = \frac{1}{n} \sum_{i=1}^n \exp \left( - \sum_{j=1}^{|v|} p(v_j|x, y_{<i}) \log p(v_j|x, y_{<i}) \right). \quad (9)$$

Lower values of distributional perplexity indicate that the model is more confident about its predictions across the entire distribution, as it reflects a clearer understanding of token relationships.

- **Self-Certainty:** This metric assesses the KL-divergence from a uniform distribution:

$$\text{Self-certainty} = - \frac{1}{n|v|} \sum_{i=1}^n \sum_{j=1}^{|v|} \log(|v| \cdot p(v_j|x, y_{<i})) . \quad (10)$$

Higher self-certainty values reflect greater confidence in the predictions, as they indicate a more concentrated distribution of outputs.

The Cumulative Distribution Function (CDF) relationship between these metrics and accuracy is provided in Fig. 2, where we utilize the negative of Entropy and DP.

### 3 Methodology

Based on the understanding that different embedding models excel in their respective domains, our study aims to develop a comprehensive RAG method that combines the strengths of various embedding models. In this work, given a question  $q$ , we utilize a single LLM  $f()$  for answer generation, while  $N$  embedding models  $\{g_1(), g_2(), \dots, g_N()\}$  assist in retrieving question-related information from a corpus, which has been divided into chunks  $C = [c_1, c_2, \dots, c_m]$ . Inspired by the stacking and bagging paradigms of ensemble learning, we propose two distinct methods: (1) Mixture-Embedding RAG, which first combines the retrieved results from multiple embedding models and incorporates them into the prompt for the LLM, and (2) Confident RAG, which employs vanilla RAG with different embedding models for initial answer generation and subsequently selects the answer with the highest confidence as the final response.

As shown in Fig. 1, the two proposed methods share the same preliminary process: first, they utilize different embedding models to identify the  $k$  chunks with the highest similarity, similar to the approach taken by the vanilla RAG described in Section 2.1. We define the retrieved chunks as  $R = \{r_{i,j} \mid i = 1, 2, \dots, N; j = 1, 2, \dots, k\}$ , where  $r_{i,j}$  is defined as  $c_{\mathcal{K}_j^i}$  and  $\mathcal{K}^i$  is the index set of the selected  $k$  chunks by embedding model  $g_i$ . Additionally, we define the similarity between the question  $q$  and chunk  $r_{i,j}$ , calculated by embedding model  $g_i$ , as  $w_{i,j}$ . The differences between the Mixture-Embedding RAG and Confident RAG are detailed as follows.

#### 3.1 Mixture-Embedding RAG

In this method, we aim to directly find the chunks with the highest similarity to the question among the candidates selected by all embedding models. Intuitively, we can select chunks using the similarity score  $w_{i,j}$ . However, this approach can lead to two problems. First, the candidates from different embedding models may have repetitions. Therefore, we select chunks without repetition to avoid any performance loss. Second, due to the different structures of embedding models, they may have varying ranges of similarity, so that simply comparing similarity across different models may introduce potential bias. To address this, we propose

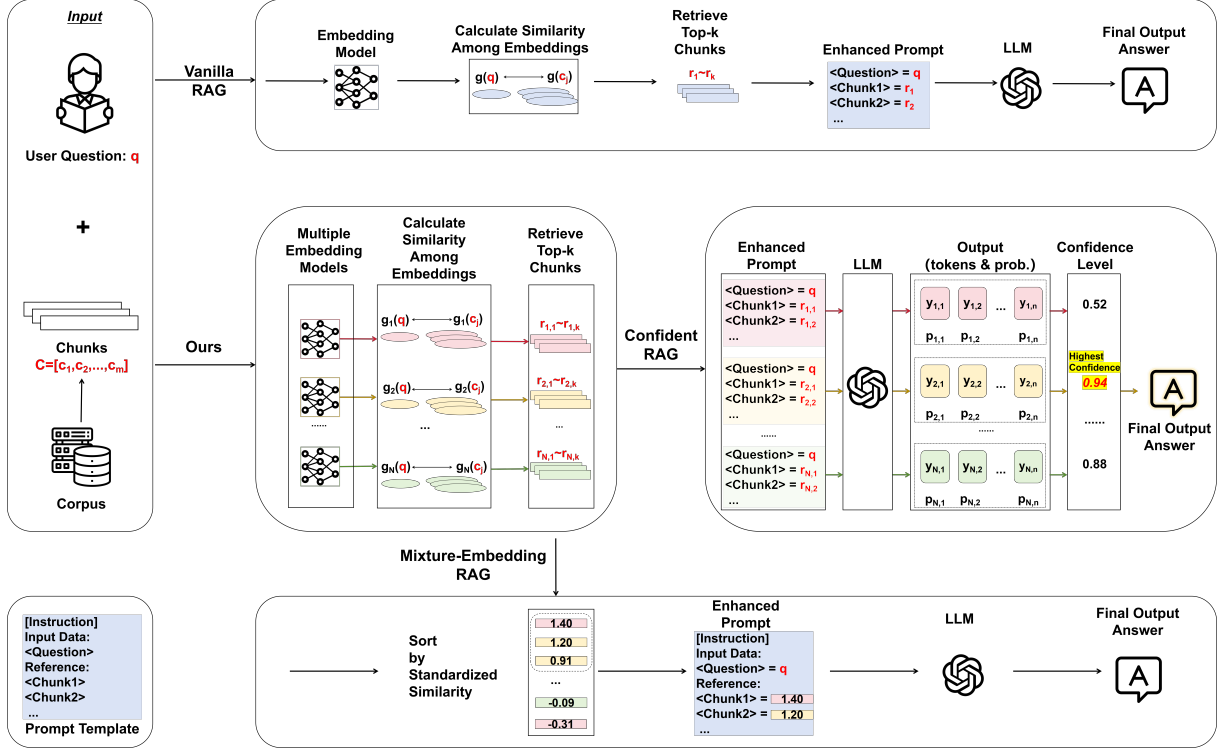


Figure 1: Workflow: This figure illustrates the workflows of vanilla RAG, as well as our proposed Mixture-Embedding RAG and Confident RAG.

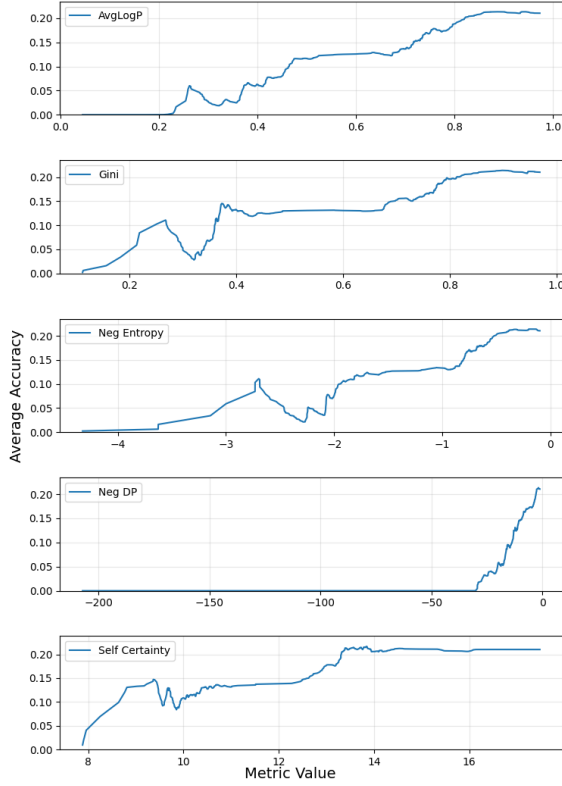


Figure 2: CDF of Accuracy for Different Metrics: The lines have been smoothed using a Gaussian filter.

to standardize the similarity score using Z-scores:

$$\hat{w}_{i,j} = \frac{w_{i,j} - \mu_i}{\sigma_i}, \quad (11)$$

where  $\mu_i$  and  $\sigma_i$  represent the mean and standard deviation of all similarity scores from embedding model  $g_i$ .

### 3.2 Confident RAG

Inspired by (Kang et al., 2025), we first let the LLM generate  $N$  separate answers, using the retrieved chunks from the  $N$  embedding models separately. Then, we evaluate the confidence of the  $N$  answers using the metrics mentioned in Section 2.2. Finally, we select the answer with the highest confidence as our final answer for each question.

## 4 Experiment

### 4.1 Experiment Setup

This section provides an overview of the experiment setup, including the model configuration, dataset, corpus, embedding models, and large language models.

- **Model configuration:** Utilizing six Nvidia A2 GPUs, it took approximately fourteen days to obtain the initial results from our



experiment. These included vanilla LLM, vanilla RAG with different embedding models, mixture-embedding RAG with various combinations of embedding models, and the combined results of confident RAG method.

- **Dataset:** Gsm8k (Cobbe et al., 2021) is a dataset of approximately 8,500 high-quality, linguistically diverse grade school math word problems designed to support question answering on basic mathematical tasks. These problems typically require 2 to 8 steps of reasoning, mainly involving elementary arithmetic operations, and are solvable by middle school students. The solutions are provided in natural language, emphasizing sequential reasoning rather than complex concepts or variables. The train dataset of Gsm8k was used as the retrieved corpus for RAG, while the first 500 items from the test dataset were used as user questions.
- **Corpus:** Two types of corpora were selected for retrieval: (1) Mathematics textbooks from OpenStax<sup>1</sup> covering domains such as calculus, algebra and trigonometry. The content was segmented into sub-section. (2) Math QA items from the Gsm8k train dataset. One textbook sub-section and three QA items will be retrieved according to the similarity based on their similarity to the user-provided math question each time.
- **Embedding models:** We chose four embedding models for encoding, including: all-MiniLM-L6-v2 (Face), ModernBERT-large (Warner et al., 2024), MathBERT (Peng et al., 2021) and stsb-roberta-large (Reimers and Gurevych, 2019) models.
- **LLMs:** To validate our proposed method and ensure repeatability, we selected three LLMs, following (Shao et al., 2025): Qwen2.5-Math-7B (Yang et al., 2024), Llama-3.1-8B (Grattafiori et al., 2024) and OLMo-2-1124-7B (OLMo et al., 2024).

## 4.2 Experiment Result

### 4.2.1 Vanilla LLM and RAG

Table 1 presents a comparison of three different LLMs with and without RAG using different embedding models. All models achieve significant improvements in accuracy after applying RAG method, with an average improvement of

5%. The accuracy of Qwen2.5-Math-7B increases from 75.2% to 80.5%, Llama-3.1-8B from 16.6% to 21.3%, and OLMo-2-1124-7B from 21.0% to 26.0%. This demonstrates that incorporating RAG methods can effectively enhance overall model performance on the task.

### 4.2.2 Results of Mixture-Embedding RAG method

Regarding the mixture-embedding RAG method (see Table 2), the average accuracies of the three LLMs did not perform well, compared to vanilla RAG. Llama-3.1-8B and OLMo-2-1124-7B demonstrated similar accuracy as vanilla RAG, with the differences ranging from -0.2% to 0.5%. Surprisingly, Qwen2.5-Math-7B showed a decrease of 5.5% compared to vanilla RAG.

### 4.2.3 Results of Confident RAG method

While using different embedding models for RAG in multiple rounds, the results of Confident RAG method are shown in Table 3. Four main findings were obtained:

(1) The accuracy of Confident RAG method surpasses that of both the vanilla RAG (with an average improvement ranging from 3.2% to 4.9%) and vanilla LLM (with an average improvement ranging from 8.1% to 9.9%).

(2) For each LLM, after applying the Confident RAG method, the accuracy improved by nearly 10% when using the best confidence metric compared to the vanilla LLM.

(3) Among all confidence metrics, self certainty and distributional perplexity demonstrated the best performance, with average improvements of 9.9% and 9.7%, respectively, over the vanilla LLM. These two metrics also performed well across different LLMs. For instance, there was an increase of 9.1% for Qwen2.5-Math-7B using self certainty as the confidence metric, an increase of 10.4% for Llama-3.1-8B with both metrics, and an increase of 12.3% for OLMo-2-1124-7B using distributional perplexity.

(4) Regarding the optimal number of embedding models (N) for multi-rounds, no evidence suggests that a larger n yields better accuracy. In our experiments, the accuracy when N=3 was always larger than when N=2. Meanwhile, the accuracy for N=3 and N=4 was similar, with a maximum difference of 1% across the three LLMs. Considering factors such as time cost, GPU capacity, and other constraints, this finding suggests that

<sup>1</sup><https://openstax.org/subjects/math>

LLM Model	Vanilla LLM	Vanilla RAG					
		Emb1	Emb2	Emb3	Emb4	Avg	Improvement
<b>Qwen2.5-Math-7B</b> (Yang et al., 2024)	75.2%	81.0%	83.0%	77.6%	80.2%	80.5%	5.3%↑
<b>Llama-3.1-8B</b> (Grattafiori et al., 2024)	16.6%	21.8%	19.8%	19.4%	24.0%	21.3%	4.7%↑
<b>OLMo-2-1124-7B</b> (OLMo et al., 2024)	21.0%	26.0%	28.4%	25.0%	24.6%	26.0%	5.0%↑
<b>Average</b>	37.6%	42.9%	43.7%	40.7%	42.9%	42.6%	5.0%↑

Table 1: Model Performance Comparison Across Different Embedding Methods: The four embedding models (1–4) correspond to all-MiniLM-L6-v2 (Face), ModernBERT-large (Warner et al., 2024), MathBERT (Peng et al., 2021), and stsb-roberta-large (Reimers and Gurevych, 2019), respectively, which will remain consistent in the following tables. The last two columns represent the average accuracy of RAG using the four embedding models and the improvement compared to the vanilla LLM. These details will remain the same in the subsequent tables.

LLM Model	Mix-Embedding RAG					
	2 Embs	3 Embs	4 Embs	Avg	v.s. Vanilla LLM	v.s. Vanilla RAG
<b>Qwen2.5-Math-7B</b> (Yang et al., 2024)	74.20%	76.00%	74.80%	75.0%	-0.2%↓	-5.5%↓
<b>Llama-3.1-8B</b> (Grattafiori et al., 2024)	20.90%	19.60%	22.80%	21.1%	4.5%↑	-0.2%↓
<b>OLMo-2-1124-7B</b> (OLMo et al., 2024)	26.20%	26.80%	26.40%	26.5%	5.5%↑	0.5%↑
<b>Average</b>	40.4%	40.8%	41.3%	40.9%	3.3%↑	-1.7%↓

Table 2: Performance of Mixture-Embedding RAG: This table illustrates the performance when using between 2 and 4 different embedding models randomly (denoted as 2 Embs to 4 Embs), comparing the results with those of the vanilla LLM and vanilla RAG.

researchers should seek a suitable trade-off based on their own condition when choosing the optimal N.

### 4.3 Analysis

#### 4.3.1 Mixture-embedding RAG

For general LLMs (e.g., Llama-3.1-8B and OLMo-2-1124-7B) without math fine-tuning, their internal math knowledge is limited, leading to lower accuracy in direct answer generation. In these cases, even noisy references retrieved by RAG are more reliable than the LLMs’ own outputs, as RAG at least provides partially correct information. However, while the mixture-embedding RAG method may optimize the retrieval ranking process and improve the quality of the references, the general LLMs’ capabilities prevent them from fully leveraging higher-quality references, resulting in performance similar to vanilla RAG. Additionally, if different embedding models return highly diverse references, directly combining the top-ranked documents may cause information overload or contextual confusion, negating the potential benefits of mixture-embedding method. Therefore, the performance of general LLMs matches that of vanilla RAG rather than surpassing it.

On the other hand, for the LLMs that have been fine-tuned based on math corpora, vanilla

RAG may result in smaller improvements, and lower-quality references can lead to poorer answer performance. For these types of LLMs, the mixture-embedding method may introduce additional noise in mathematical contexts, resulting in lower accuracy compared to vanilla RAG. The decline in accuracy may be caused by several factors: (1) Mathematical symbols and formulas may vary drastically across embedding models, making similarity calculations unstable. (2) Different models may encode math terms differently, causing the top-ranked reference to be suboptimal. (3) An embedding model might incorrectly rank an irrelevant math material highly, while better references from other models are ignored. If the LLM generates hallucinated answers based on incorrect references, its performance can degrade below that of the vanilla LLM. (4) Information overload or contextual confusion may also occur, similar to what happens with general LLMs.

#### 4.3.2 Confident RAG

As shown in Figure 2, there exists a positive correlation between confidence and accuracy. Therefore, the Confident RAG method improves overall accuracy by integrating multiple embedding models to generate answers and selecting the highest-confidence results using the most effective metric. This process effectively filters out low-confidence

LLM	Emb. Model	AvgLogP	Self-certainty	Gini	Entropy	DP
Qwen2.5-Math-7B (Yang et al., 2024)	1,2	82.0%	<b>85.0%</b>	81.8%	82.8%	83.6%
	1,3	79.4%	<b>83.4%</b>	79.2%	79.6%	81.0%
	1,4	81.4%	<b>84.6%</b>	81.6%	81.8%	82.2%
	2,3	81.8%	<b>84.4%</b>	81.4%	81.8%	82.4%
	2,4	81.4%	<b>83.6%</b>	81.0%	81.4%	82.2%
	3,4	79.2%	<b>82.2%</b>	79.8%	79.8%	80.6%
	Avg (n=2)	80.9%	<b>83.9%</b>	80.8%	81.2%	82.0%
	1,2,3	79.4%	<b>85.0%</b>	79.0%	79.8%	81.6%
	1,2,4	79.6%	<b>85.0%</b>	79.6%	80.6%	82.0%
	1,3,4	79.0%	<b>84.8%</b>	79.6%	80.0%	80.8%
	2,3,4	79.2%	<b>84.4%</b>	79.0%	79.8%	81.0%
	Avg (n=3)	79.3%	<b>84.8%</b>	79.3%	80.1%	81.4%
	1,2,3,4	78.2%	<b>84.8%</b>	78.4%	79.2%	80.6%
	Avg (n=2,3,4)	80.1%	<b>84.3%</b>	80.0%	80.6%	81.6%
Llama-3.1-8B (Grattafiori et al., 2024)	v.s. Vanilla RAG	0.4%↓	<b>3.8%</b> ↑	0.4%↓	0.1%↑	1.2%↑
	v.s. Vanilla LLM	4.9%↑	<b>9.1%</b> ↑	4.8%↑	5.4%↑	6.4%↑
	1,2	<b>27.2%</b>	<b>27.2%</b>	27.0%	27.0%	<b>27.2%</b>
	1,3	26.8%	26.8%	26.6%	27.0%	<b>27.2%</b>
	1,4	<b>26.8%</b>	26.4%	26.4%	26.4%	26.4%
	2,3	23.2%	<b>24.0%</b>	<b>24.0%</b>	<b>24.0%</b>	23.6%
	2,4	26.6%	<b>26.8%</b>	26.4%	26.0%	26.6%
	3,4	25.8%	26.4%	<b>26.6%</b>	<b>26.6%</b>	26.4%
	Avg (n=2)	26.1%	<b>26.3%</b>	26.2%	26.2%	26.2%
	1,2,3	28.6%	<u>29.2%</u>	<u>28.8%</u>	<u>29.4%</u>	<u>29.4%</u>
	1,2,4	<b>28.8%</b>	28.6%	28.0%	28.0%	28.2%
	1,3,4	27.4%	27.4%	27.2%	<b>27.6%</b>	<b>27.6%</b>
	2,3,4	25.8%	<b>26.4%</b>	26.0%	26.0%	<b>26.4%</b>
	Avg (n=3)	27.7%	<b>27.9%</b>	27.5%	27.8%	<b>27.9%</b>
OLMo-2-1124-7B (OLMo et al., 2024)	1,2,3,4	<b>27.8%</b>	27.6%	27.0%	27.4%	27.6%
	Avg (n=2,3,4)	26.8%	<b>27.0%</b>	26.7%	26.9%	<b>27.0%</b>
	v.s. Vanilla RAG	5.6%↑	<b>5.7%</b> ↑	5.5%↑	5.6%↑	<b>5.7%</b> ↑
	v.s. Vanilla LLM	10.2%↑	<b>10.4%</b> ↑	10.1%↑	10.3%↑	<b>10.4%</b> ↑
	1,2	31.0%	31.2%	30.4%	31.2%	32.2%
	1,3	29.8%	29.6%	29.8%	<b>30.0%</b>	<b>30.0%</b>
	1,4	29.4%	29.2%	28.6%	29.6%	<b>31.0%</b>
	2,3	31.6%	30.8%	30.2%	31.4%	<b>32.8%</b>
	2,4	32.6%	32.0%	31.0%	<b>33.0%</b>	33.8%
	3,4	29.6%	30.2%	29.6%	30.6%	<b>31.4%</b>
	Avg (n=2)	30.7%	30.5%	29.9%	31.0%	<b>31.9%</b>
	1,2,3	32.6%	32.0%	31.4%	32.2%	<b>33.2%</b>
	1,2,4	<u>32.8%</u>	<u>32.8%</u>	<u>31.6%</u>	33.2%	<b>34.8%</b>
	1,3,4	31.2%	31.2%	30.6%	32.0%	<b>34.6%</b>
Average	2,3,4	32.6%	32.2%	30.8%	<u>33.6%</u>	<b>36.6%</b>
	Avg (n=3)	32.3%	32.1%	31.1%	32.8%	<b>34.8%</b>
	1,2,3,4	<u>32.8%</u>	32.4%	31.0%	33.2%	<b>35.8%</b>
	Avg (n=2,3,4)	31.5%	31.2%	30.5%	31.8%	<b>33.3%</b>
	v.s. Vanilla RAG	5.5%↑	5.2%↑	4.5%↑	5.8%↑	<b>7.3%</b> ↑
	v.s. Vanilla LLM	10.5%↑	10.2%↑	9.5%↑	10.8%↑	<b>12.3%</b> ↑
	Avg (n=2)	45.9%	<b>46.9%</b>	45.6%	46.1%	46.7%
	Avg (n=3)	46.4%	<b>48.3%</b>	46.0%	46.9%	48.0%
	Avg (n=4)	46.3%	<b>48.3%</b>	45.5%	46.6%	48.0%
	Avg (n=2,3,4)	46.1%	<b>47.5%</b>	45.7%	46.4%	47.3%
	v.s. Vanilla RAG	3.5%↑	<b>4.9%</b> ↑	3.2%↑	3.9%↑	4.7%↑
	v.s. Vanilla LLM	8.5%↑	<b>9.9%</b> ↑	8.1%↑	8.8%↑	9.7%↑

Table 3: Accuracy Comparison Across Multi-RAG with Different Embedding Models: Avg( $n$ ) denotes the average accuracy across different combinations of  $n$  embedding models. Each line uses underline to indicate the best embedding combination within each LLM. Each row uses **bold** to signify the best metric for confidence evaluation. For performance comparison with Vanilla RAG, we use the average accuracy of all single embedding models as the baseline.



incorrect answers. Meanwhile, the combined effect of RAG and confidence filtering enhances the robustness, leading to significant improvements compared to vanilla LLMs. When employing the optimal confidence metric, all LLMs achieved an accuracy increase of nearly 10%, demonstrating the method’s universality. In the experiments, when the number of embedding models  $N > 3$ , the accuracy improvement became limited, likely due to redundant or noisy retrievals introduced by additional models. At  $N=3$ , the method achieved an optimal balance between diversity and computational efficiency. Increasing the number of models further yields only marginal benefits.

Self-Certainty and DP outperform other metrics since they directly measure the concentration and divergence of the probability distribution. Specifically, Self-Certainty measures how far the predicted distribution deviates from uniform. By scaling the probabilities by  $|v|$  and taking the negative logarithm, it heavily penalizes uniform-like distributions, favoring sharp peaks. This makes it highly discriminative for high-confidence answers. Additionally, DP is an exponential version of entropy. The exponentiation amplifies differences in entropy, making it more sensitive to the sharpness of the distribution. Low DP values indicate tightly clustered high-probability tokens, which strongly correlate with correct answers. In contrast, other metrics are less sensitive because they either average out uncertainties (AvgLogP) or lack normalization across vocabularies (Gini). While entropy can be useful, it is linear and less discriminative compared to DP’s exponential scaling. Therefore, Self-Certainty and DP are more sensitive to subtle variations in model confidence.

## 5 Conclusion

In this paper, we proposed to enhance the response performance of LLMs with RAG utilizing multiple embedding models. We introduced two related approaches: Mixture-Embedding RAG and Confident RAG. Mixture-Embedding RAG sorts and selects retrieved references from different embedding models based on their standardized similarity before guiding LLM. However, its results showed no improvement over vanilla RAG, possibly due to factors such as information overload, contextual confusion and the capabilities of LLMs. In contrast, Confident RAG selects the highest-confidence response from multiple RAG

outputs, each generated using a different embedding model. This approach demonstrated superior performance, with an average improvement of approximately 10% while utilizing the best confidence metric, compared to vanilla LLM. Self-certainty and distributional perplexity (DP) were identified as the most effective metrics for measuring model confidence. The number of embedding models ( $N$ ) can be determined by researchers based on their specific conditions, as no evidence suggesting that a larger  $N$  necessarily yields a better performance. The consistent and stable results from Confident RAG indicate that it can serve as a plug-and-play method for enhancing LLM performance across various domains.

## References

- Muhammad Arslan, Saba Munawar, and Christophe Cruz. 2024. Sustainable digitalization of business with multi-agent rag and llm. *Procedia Computer Science*, 246:4722–4731.
- Haolong Chen, Hanzhi Chen, Zijian Zhao, Kaifeng Han, Guangxu Zhu, Yichen Zhao, Ying Du, Wei Xu, and Qingjiang Shi. 2024. An overview of domain-specific foundation model: key technologies, applications and challenges. *SCIENCE CHINA Information Sciences*.
- Jiuhai Chen and Jonas Mueller. 2023. Quantifying uncertainty in answers from any language model and enhancing their trustworthiness. *arXiv preprint arXiv:2308.16175*.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Hugging Face. [sentence-transformers/all-minilm-l6-v2](https://huggingface.co/sentence-transformers/all-minilm-l6-v2).
- Lin Gao, Jing Lu, Zekai Shao, Ziyue Lin, Shengbin Yue, Chiokit Ieong, Yi Sun, Rory James Zainer, Zhongyu Wei, and Siming Chen. 2024. Fine-tuned large language model for visualization system: A study on self-regulated learning in education. *IEEE Transactions on Visualization and Computer Graphics*.

- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Wenbo Guan, Jiyu Lu, Qinyu Feng, Xiaoqian Li, and Jun Zhou. 2024. External knowledge is not always needed: An adaptive retrieval augmented generation method. In *2024 4th International Symposium on Artificial Intelligence and Intelligent Manufacturing (AIIM)*, pages 882–886. IEEE.
- Sophia Hager, David Mueller, Kevin Duh, and Nicholas Andrews. 2025. Uncertainty distillation: Teaching language models to express semantic confidence. *arXiv preprint arXiv:2503.14749*.
- Zijian Hei, Weiling Liu, Wenjie Ou, Juyi Qiao, Junming Jiao, Guowen Song, Ting Tian, and Yi Lin. 2024. Dr-rag: Applying dynamic document relevance to retrieval-augmented generation for question-answering. *arXiv preprint arXiv:2406.07348*.
- Owen Henkel, Zach Levonian, Chenglu Li, and Millie Postle. 2024. Retrieval-augmented generation to improve math question-answering: Trade-offs between groundedness and human preference. In *Proceedings of the 17th International Conference on Educational Data Mining*, pages 315–320.
- Jerry Huang, Siddarth Madala, Risham Sidhu, Cheng Niu, Hao Peng, Julia Hockenmaier, and Tong Zhang. 2025a. Rag-rl: Advancing retrieval-augmented generation via rl and curriculum learning. *arXiv preprint arXiv:2503.12759*.
- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. 2025b. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *ACM Transactions on Information Systems*, 43(2):1–55.
- Zhengbao Jiang, Jun Araki, Haibo Ding, and Graham Neubig. 2021. How can we know when language models know? on the calibration of language models for question answering. *Transactions of the Association for Computational Linguistics*, 9:962–977.
- Zhewei Kang, Xuandong Zhao, and Dawn Song. 2025. Scalable best-of-n selection for large language models via self-certainty. *arXiv preprint arXiv:2502.18581*.
- Kiseung Kim and Jay-Yoon Lee. 2024. Re-rag: Improving open-domain qa performance and interpretability with relevance estimator in retrieval-augmented generation. *arXiv preprint arXiv:2406.05794*.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in neural information processing systems*, 33:9459–9474.
- Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. 2024. Contextualized sequence likelihood: Enhanced confidence scores for natural language generation. *arXiv preprint arXiv:2406.01806*.
- Ziwei Liu, Liang Zhang, Qian Li, Jianghua Wu, and Guangxu Zhu. 2024. Invar-rag: Invariant llm-aligned retrieval for better generation. *arXiv preprint arXiv:2411.07021*.
- Tengfei Lyu, Siyuan Feng, Hao Liu, and Hai Yang. 2025. Llm-oddr: A large language model framework for joint order dispatching and driver repositioning. *arXiv preprint arXiv:2505.22695*.
- Jinming Nian, Zhiyuan Peng, Qifan Wang, and Yi Fang. 2024. W-rag: Weakly supervised dense retrieval in rag for open-domain question answering. *arXiv preprint arXiv:2408.08444*.
- Team OLMo, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, et al. 2024. 2 olmo 2 furious. *arXiv preprint arXiv:2501.00656*.
- Shuai Peng, Ke Yuan, Liangcai Gao, and Zhi Tang. 2021. Mathbert: A pre-trained model for mathematical formula understanding. *arXiv preprint arXiv:2105.00377*.

- Nils Reimers and Iryna Gurevych. 2019. [Sentence-bert: Sentence embeddings using siamese bert-networks](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics.
- Jie Ren, Yao Zhao, Tu Vu, Peter J Liu, and Balaji Lakshminarayanan. 2023. Self-evaluation improves selective generation in large language models. In *Proceedings on*, pages 49–64. PMLR.
- Kunal Sawarkar, Abhilasha Mangal, and Shivam Raj Solanki. 2024. Blended rag: Improving rag (retriever-augmented generation) accuracy with semantic search and hybrid query-based retrievers. In *2024 IEEE 7th International Conference on Multimedia Information Processing and Retrieval (MIPR)*, pages 155–161. IEEE.
- Rulin Shao, Shuyue Stella Li, Rui Xin, Scott Geng, Yiping Wang, Sewoong Oh, Simon Shaolei Du, Nathan Lambert, Sewon Min, Ranjay Krishna, Yulia Tsvetkov, Hannaneh Hajishirzi, Pang Wei Koh, and Luke Zettlemoyer. 2025. [Spurious rewards: Rethinking training signals in rlvr](#).
- Vaishnavi Shrivastava, Ananya Kumar, and Percy Liang. 2025. Language models prefer what they know: Relative confidence estimation via confidence preferences. *arXiv preprint arXiv:2502.01126*.
- Vaishnavi Shrivastava, Percy Liang, and Ananya Kumar. 2023. Llamas know what gpts don't show: Surrogate models for confidence estimation. *arXiv preprint arXiv:2311.08877*.
- Shamane Siriwardhana, Rivindu Weerasekera, Eliott Wen, and Suranga Nanayakkara. 2021. Fine-tune the entire rag architecture (including dpr retriever) for question-answering. *arXiv preprint arXiv:2106.11517*.
- Xiaqiang Tang, Qiang Gao, Jian Li, Nan Du, Qi Li, and Sihong Xie. 2024. Mba-rag: a bandit approach for adaptive retrieval-augmented generation through question complexity. *arXiv preprint arXiv:2412.01572*.
- Yiteng Tu, Weihang Su, Yujia Zhou, Yiqun Liu, and Qingyao Ai. 2025. Rbft: Robust fine-tuning for retrieval-augmented generation against retrieval defects. *arXiv preprint arXiv:2501.18365*.
- Benjamin Warner, Antoine Chaffin, Benjamin Clavié, Orion Weller, Oskar Hallström, Said Taghadouini, Alexis Gallagher, Raja Biswas, Faisal Ladhak, Tom Aarsen, Nathan Cooper, Griffin Adams, Jeremy Howard, and Iacopo Poli. 2024. [Smarter, better, faster, longer: A modern bidirectional encoder for fast, memory efficient, and long context finetuning and inference](#).
- Miao Xiong, Zhiyuan Hu, Xinyang Lu, Yifei Li, Jie Fu, Junxian He, and Bryan Hooi. 2023. Can llms express their uncertainty? an empirical evaluation of confidence elicitation in llms. *arXiv preprint arXiv:2306.13063*.
- An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong Tu, Jingren Zhou, Junyang Lin, et al. 2024. Qwen2. 5-math technical report: Toward mathematical expert model via self-improvement. *arXiv preprint arXiv:2409.12122*.
- Zhuocheng Zhang, Yang Feng, and Min Zhang. 2025. Levelrag: Enhancing retrieval-augmented generation with multi-hop logic planning over rewriting augmented searchers. *arXiv preprint arXiv:2502.18139*.
- Zijian Zhao. 2025. Let network decide what to learn: Symbolic music understanding model based on large-scale adversarial pre-training. In *Proceedings of the 2025 International Conference on Multimedia Retrieval*, pages 2128–2132.
- Zijian Zhao, Dian Jin, Zijing Zhou, and Xiaoyu Zhang. 2025. Automatic stage lighting control: Is it a rule-driven process or generative task? *arXiv preprint arXiv:2506.01482*.