

Revisiting RAG Ensemble: A Theoretical and Mechanistic Analysis of Multi-RAG System Collaboration

Yifei Chen

zhangboguodong@ruc.edu.cn
Gaoling School of Artificial Intelligence, Renmin
University of China
Haidian Qu, Beijing Shi, China

Yutao Zhu

Gaoling School of Artificial Intelligence, Renmin
University of China
Haidian Qu, Beijing Shi, China

Guanting Dong

Gaoling School of Artificial Intelligence, Renmin
University of China
Haidian Qu, Beijing Shi, China

Zhicheng Dou

Gaoling School of Artificial Intelligence, Renmin
University of China
Haidian Qu, Beijing Shi, China

Abstract

Retrieval-Augmented Generation (RAG) technology has been widely applied in recent years. However, despite the emergence of various RAG frameworks, a single RAG framework still cannot adapt well to a broad range of downstream tasks. Therefore, how to leverage the advantages of multiple RAG systems has become an area worth exploring. To address this issue, we have conducted a comprehensive and systematic investigation into ensemble methods based on RAG systems. Specifically, we have analyzed the RAG ensemble framework from both theoretical and mechanistic analysis perspectives. From the theoretical analysis, we provide the first explanation of the RAG ensemble framework from the perspective of information entropy. In terms of mechanism analysis, we have explored the RAG ensemble framework from both the pipeline and module levels. We carefully select four different pipelines (Branching, Iterative, Loop, and Agentic) and three different modules (Generator, Retriever, and Reranker) to solve seven different research questions. The experiments show that aggregating multiple RAG systems is both generalizable and robust, whether at the pipeline level or the module level. Our work lays the foundation for similar research on the multi-RAG system ensemble.

CCS Concepts

• **Information systems** → **Information integration.**

Keywords

Retrieval-Augmented Generation, Pipeline Ensemble, Module Ensemble, Model Preference

ACM Reference Format:

Yifei Chen, Guanting Dong, Yutao Zhu, and Zhicheng Dou. 2025. Revisiting RAG Ensemble: A Theoretical and Mechanistic Analysis of Multi-RAG

System Collaboration. In *Proceedings of CIKM'25*. ACM, New York, NY, USA, 13 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 Introduction

The emergence of Large Language Models (LLMs) has profoundly revolutionized many real-world tasks that rely on natural language [4, 45, 70]. However, when dealing with knowledge-intensive tasks, LLMs relying solely on their parametric knowledge often suffer from factual inconsistencies or hallucinations. To address these limitations, Retrieval-Augmented Generation (RAG) methods have been proposed, augmenting LLMs with dynamically retrieved external knowledge. This integration enhances response accuracy and reliability by grounding outputs in verifiable information sources.

As research in this field advances, more and more RAG methods have been proposed. Component Module RAG inserts various modules into the standard pipeline to better complete the retrieval task. For instance, the LongLLMLingua and RECOMP methods refine the retrieved knowledge with a refiner, and the SKR and Adaptive RAG methods distinguish the difficulty of questions by introducing a Judger [23, 27, 60, 65]. Pipeline Module RAG optimizes the whole process to improve the accuracy and efficiency. For example, RePlug method is suitable for tasks with varying difficulty levels, and methods such as Iter-RetGen and Self-RAG are suitable for solving multi-hop problems [2, 51, 52]. With the development of agent technology, the application of Agentic RAG technology is becoming increasingly widespread [35, 55, 62, 63]. For example, Search-o1 and WebThinker combine search and reasoning and perform well in deep search tasks [37, 38]. However, given the inherent complexity of tasks and the heterogeneity of RAG workflows, developing a universal RAG framework that generalizes effectively across diverse applications remains a significant challenge.

To further investigate this limitation, we analyze the framework classification proposed in FlashRAG [31] and select four representative RAG methods, each corresponding to a distinct pipeline type: Branching, Iterative, Loop, and Agentic. As illustrated in Figure 1, we evaluate these methods on four benchmark tasks and observe the following key findings:

(1) **Lack of generalizability in single RAG pipelines.** The upper part of Figure 1 presents the aggregated performance of each method across three datasets. Notably, Branching-based approaches (e.g., RePlug [52]) underperform in multi-hop reasoning tasks but

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
CIKM'25, Seoul, Korea

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 978-1-4503-XXXX-X/2018/06
<https://doi.org/XXXXXXX.XXXXXXX>

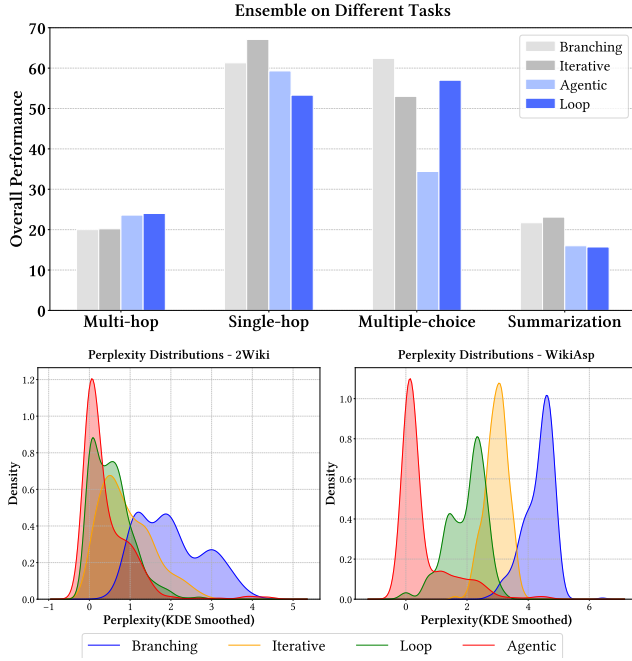


Figure 1: Overall Performance of Base RAG and Ensemble Methods. Top: the answers’ F1 scores of different methods. Bottom: the distribution of perplexity for the outputs.

excel in multiple-choice settings. A similar phenomenon can be observed in the Iterative method, reinforcing that each pipeline exhibits task-dependent performance biases.

(2) **Divergence in answer perplexity levels across pipelines.** By quantifying the perplexity of generated answers, we observe distinct density distributions for each method (lower part of Figure 1). The more concentrated the distribution, the more stable the generated results. For instance, Loop-based methods yield lower answer perplexity (greater stability) on 2Wiki, but higher perplexity on WikiASP, reflecting task-specific confidence disparities.

These results collectively demonstrate that single-RAG systems struggle with task generalization, whether measured by performance or output perplexity. This motivates our core research question: How can we aggregate multiple RAG systems to enhance generalization capability for complex, heterogeneous tasks?

To address this issue, one intuitive approach is to perform adaptive fine-tuning on the model to enhance its ability in RAG tasks. However, such methods may interfere with the model’s inherent capabilities and come with higher training costs. Another common strategy is to treat the model as a router, selecting the optimal single RAG system’s answer and discarding the remaining systems. However, we consider that the unselected answers may still contain valuable information for the task. Recent research has begun to explore component ensemble methods. Some studies suggest that meta-search engines, by aggregating results from multiple search engines, can provide more relevant information [41, 49]. Additionally, numerous studies focus on model-level ensemble strategies. We argue that, compared to routing methods, ensemble strategy can

better make full use of the useful information in each subsystem, improving the quality of the final results. However, existing methods mainly focus on multi-component ensemble on single level, while RAG tasks involve more complex input flows and system structures. Unfortunately, both in terms of theoretical modeling and mechanism explanation, there is still a significant lack of systematic research on ensemble across multiple RAG systems, which significantly limits its development and application.

To address this challenge, in this paper, we conduct a comprehensive and systematic study of ensemble methods based on RAG systems. Specifically, we perform an in-depth analysis of the RAG system ensemble method from theoretical analysis and mechanism analysis:

- (1) **From a theoretical analysis perspective:** We model the RAG system ensemble method on a non-Euclidean manifold. Through detailed derivation, we clarify the effectiveness of RAG system ensemble from the perspective of information entropy increase. As we know, this is the first work to model a system ensemble task from the perspective of information entropy.
- (2) **From a mechanism analysis perspective:** To achieve a comprehensive exploration of RAG ensemble, we conduct in-depth investigations of seven different research questions from both the pipeline and module levels. At the system level, we carefully select four different RAG pipelines (Branching, Iterative, Loop, and Agentic) for ensemble research. Additionally, we conduct ensemble experiments on closed-source RAG frameworks to further explore the characteristics of RAG ensemble. At the module level, we conduct experimental research on the retriever, reranker, and generator of the standard RAG framework. We carefully select three retrievers and five generation models for the experiments, and delve into the characteristics of applying generative rerankers to ensemble tasks. Moreover, our experiments cover a wide range of task sets, including single-hop tasks, multi-hop tasks, multiple-choice tasks, summarization tasks, and tasks in vertical domains, all of which have detailed ensemble analysis.

Our main findings include:

- RAG ensemble demonstrates clear advantages in both the framework type and the granularity of the ensemble. This reflects the good generalizability of the RAG ensemble method.
- In a significant portion of ensemble tasks, the RAG ensemble method exhibits scaling-up characteristics, meaning that increasing the external information has a notable positive impact on the final ensemble result. However, this characteristic also depends on the model’s strong resistance to information interference.
- The ensemble model shows a preference for certain groups of input information, and this preference becomes more pronounced as task difficulty increases.

2 Related Work

Retrieval-Augmented Generation. Retrieval-Augmented Generation (RAG) improves generation quality by integrating retrieved external knowledge [10, 17, 33, 73]. Mainstream RAG methods follow a “retrieve-then-read” paradigm, where the retrieval module supplies external knowledge as context for generation models to produce output [3, 48, 52]. Recently, numerous improved

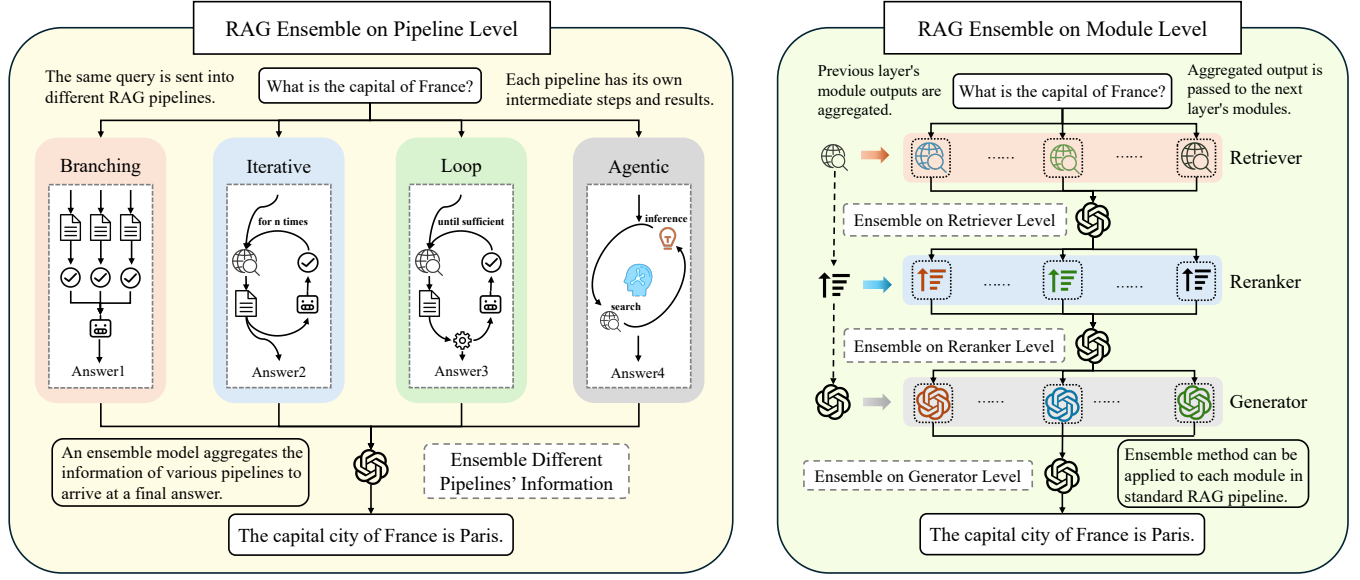


Figure 2: The overall architecture of RAG Ensemble. We conduct a detailed study from the perspective of the overall pipeline level and the module level. The left framework represents ensemble from pipeline level, while the right framework represents ensemble from module level.

RAG paradigms have emerged based on standard RAG technology [16, 31, 51]. Foundation studies have focused on capturing more relevant and high-quality retrieval documents through advanced query rewriting [15, 59] and fine-grained reranking techniques [65, 67, 68]. Other branch of efforts have attempted to introduce efficient fine-tuning strategies [11, 13, 42, 72, 74] to unlock the strong information-capturing capabilities of generators within RAG systems. Moreover, recent works aims to integrate self-correction and auto-evaluation methods into RAG domain [9].

Recently, with the rise in popularity of agentic search methods [30, 37], a series of approaches have attempted to further enhance large models' autonomy and deep information retrieval capabilities using reinforcement learning techniques. For example, Search-R1 and Research use reinforcement learning to empower the model with the ability to complete deep search tasks [6, 29]. Web-Thinker and Kimi-Researcher attempt to leverage reinforcement learning to enhance the model's report writing ability [38, 56]. Tool-Star and ARPO try to expand more tools to enhance the model's deep search capabilities [8, 12]. Furthermore, a series of studies have attempted to extend this progress to the multimodal domain [46, 61].

However, existing approaches mainly enhance individual RAG systems, lacking exploration of collaborative multiple RAG systems. Our paper aims to investigate the effectiveness of model ensemble in RAG scenarios and introduces the Ensemble-RAG paradigm, offering practical insights for real-world applications.

Model Ensemble in LLM. Ensemble of LLMs has significantly outperformed individual models by leveraging the strengths of different systems. Existing ensemble approaches can be mainly categorized into three types: (1) A series of studies have fine-tuned external routers to select the most suitable LLM for specific inputs, enabling model selection before inference [40, 53]; (2) Another

branch of efforts involves multiple models processing inputs incrementally and combining their outputs during decoding, showcasing strong collaborative potential [20, 36]; and (3) Some researchers focus on allowing each model to process inputs independently, then selecting the best response [5, 25]. To further improve the efficiency of LLMs ensemble, some works employ input compression [26, 34, 39] and speculative decoding [21, 69] to accelerate model inference. Unfortunately, these studies have not systematically examined the application of ensemble techniques in the RAG domain, and integrating information at the model level alone is insufficient to bridge system gaps. Our study makes the first attempt to integrate all external knowledge and outputs from different RAG systems to maximize performance.

3 Methodology and Theoretical Analysis

In this section, we introduce the framework of RAG ensemble, and then we theoretically analyze why combining multiple RAG systems can be effective.

3.1 RAG Ensemble

Given a set of RAG systems $\{S_1, S_2, \dots, S_n\}$, the ensemble framework aims to synthesize the inputs and outputs of these systems to generate a response Y . Specifically, each standard RAG system S_i includes a retriever R_i and a generator G_i .¹ Upon receiving a user input X , the retriever R_i retrieves relevant external knowledge, denoted as $D_i = R_i(X)$. Then, the generator G_i generates a response $Y_i = G_i(X, D_i)$. When multiple RAG systems exist, all inputs and

¹Many advanced RAG systems have designed several additional modules to improve the systems' performance. In theoretical analysis, we only consider the two fundamental components, without loss of generality, to make the formulation more clear. Henceforth, we will follow this principle to simplify the modules for better understanding.

outputs can be collected as follows:

$$S = \{S_1, S_2, \dots, S_n\}, \quad (1)$$

$$S_i = \{Y_i, D_i\}. \quad (2)$$

Finally, an ensemble model is employed to generate the final response as follows:

$$Y = f_\phi(X, S), \quad (3)$$

where $f_\phi(\cdot)$ denotes the ensemble model parameterized by ϕ .

The core idea of RAG ensemble lies in the ability of the model to aggregate the information from multiple RAG systems. This paper mainly focuses on the simplest method, which directly embeds the raw information from multiple RAG systems into a prompt and then inputs it to an LLM for ensemble. This process can be represented as follows:

$$Y = f_\phi(\text{Prompt}(X, S)). \quad (4)$$

The prompt we use for basic RAG ensemble is as follows:

RAG Ensemble Prompt

Here is a question and some external data from {num} systems' information:

System 1: {system 1's information}

System 2: {system 2's information}

System 3: {system 3's information}

.....

Question: {question}

Your task is to answer the question based on the given information. You should first output your reasoning process and then provide the final answer. The output format of reasoning process and final answer should be enclosed within `<think>` `</think>` and `<answer>` `</answer>` tags, respectively, and the final answer should be enclosed within boxed with latex format, i.e., "`<think>` reasoning process here `</think>``<answer>``a final answer here``</answer>`". Only output your reasoning process in `<think>``</think>` and your answer in `<answer>` boxed, and do not output any other words.

In the following, we will provide a detailed theoretical analysis of the rationale behind RAG ensemble.

3.2 Theoretical Analysis of RAG Ensemble

In this section, we provide a detailed analysis of RAG ensemble from the perspective of theoretical modeling. Given a probability distribution, information entropy describes the degree of uncertainty we have about an event. The higher the entropy, the more uncertain the system is, and the more information it contains. Mathematically, the entropy $H(X)$ of a continuous random variable X with probability distribution $P(x)$ can be shown as:

$$H(X) = - \int p(x) \log p(x) dx. \quad (5)$$

Inspired by the information bottleneck method [57, 71] in information theory, we consider that the key of the RAG ensemble

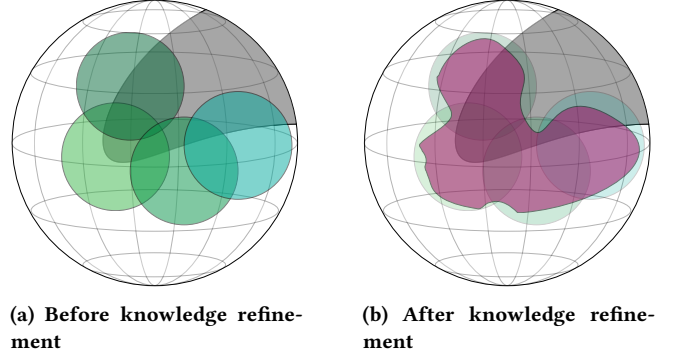


Figure 3: Comparison of multi-RAG before and after knowledge refinement. Aggregating and analyzing the information from multi-RAG system (Figure 3b) is the key operation of ensemble.

framework lies in its ensemble of information from multiple individual RAG systems, reducing the information entropy of the final answer. As shown in Figure 3, this process can be modeled on a non-Euclidean sphere.

As shown in Figure 3a, the entire white sphere represents the potential region for generating answers (i.e., the informational entropy of the answer), the gray area represents the useful information contained in the prompt, and each green circular area represents the external information introduced by an individual RAG system (including both useful and useless external knowledge). In Figure 3b, the purple area represents the useful information extracted by generator from all external knowledge. We consider that due to the introduction of useful knowledge from both the prompt and all individual RAG systems, the uncertainty in generating the answer is reduced. That is, RAG ensemble can aggregate useful information from multiple systems to generate a correct output. The useful knowledge of individual RAG systems and the useful knowledge after extraction from all systems can be obtained as follows:

$$e_i = g_\phi[q, (S_i)], \quad (6)$$

$$e^* = g_\phi[q, (S_1, S_2, \dots, S_n)], \quad (7)$$

where e_i represents individual RAG's useful knowledge, e^* represents all the useful knowledge extracted, and g_ϕ represents the process that the model ϕ extracts useful information. Then, the process by which ensemble model reduces the information entropy of the generated answer can be expressed by the following formula:

$$\begin{aligned} H(Y|X, K) &= H(Y|X, e^*) \\ &= H(Y) - \underbrace{I(X, e^*; Y)}_{\text{useful information}}, \\ I(X, e^*; Y) &= \underbrace{I(X, K; Y)}_{\text{all information}} - \underbrace{I(e_{\text{useless}})}_{\text{useless information}}, \end{aligned}$$

where K represents all the external knowledge provided, $H(Y|X, K)$ represents the conditional entropy of the generated response when both the user prompt and external knowledge are introduced, $H(Y)$

denotes the informational entropy of generating response. Moreover, $H(Y|X, K) = H(Y|X, e^*)$ means that only useful information can help reduce the target response's entropy. $I(X, e^*; Y)$, $I(X, K; Y)$ represents the mutual information between the input and useful reference knowledge, as well as the mutual information between the input and all external knowledge, respectively. $I(e_{\text{useless}})$ represents the useless information in the external knowledge that is discarded after model analysis. Therefore, when external knowledge is introduced as a condition, the information entropy of the generated answer decreases due to the useful knowledge extracted by the ensemble model (*i.e.*, the accuracy of generating answer improves).

For RAG ensemble system, if the external knowledge base does not contain conflicting information (*i.e.*, the external knowledge base cannot simultaneously contain knowledge such as "the sky is blue" and "the sky is green"), we can propose the following assumption:

Assumption

When performing ensemble tasks, an ideal model tends to refine the collected information in a direction that increases the amount of correct knowledge.

Suppose we have an ideal ensemble model ϕ^* that can perfectly extract useful information from the given external knowledge while ignoring all irrelevant information. In this case, the useful knowledge extracted from the input information of system i can be represented as follows:

$$e_i = g_{\phi^*}(q, d_i, a_i). \quad (8)$$

Meanwhile, the final useful information extracted from all sub-systems can be expressed as follows:

$$e^* = g_{\phi^*}(q, d_i, a_i, S_{\setminus i}), \quad (9)$$

here, $S_{\setminus i}$ represents the input information excluding system i .

When more system information is obtained, for the initial useful information e_i , the ensemble model can handle it in the following two ways:

(1) ϕ^* retains all the information e_i . Due to the introduction of information from other systems, the total information received by the ensemble model is at least not worse than the information received solely by the i -th system. In other words, e_i is included in e^* , which means:

$$\begin{aligned} I_1 &= I(q, e^*; a) \\ &= I(q, e_i; a) + I(q, e_{\setminus i}^*; a), \end{aligned} \quad (10)$$

$e_{\setminus i}^*$ represents the useful information from other systems excluding system i .

(2) If ϕ^* only extracts partial information from e_i , we define e_i^* represents the useful knowledge of system i after refinement, then the e_i^* can be represented as:

$$I(q, e_i^*; a) = I(q, e_i; a) - I(q, e_i^{\text{useless}}; a), \quad (11)$$

here, e_i^{useless} represents the part of the knowledge e_i that the model considers useless after final refinement. At this point, the information entropy of generating the answer can be expressed as:

$$\begin{aligned} I_2 &= I(q, e^*; a) \\ &= I(q, e_i^*; a) + I(q, e_{\setminus i}^*; a) \\ &= I(q, e_i; a) - I(q, e_i^{\text{useless}}; a) \\ &\quad + I(q, e_{\setminus i}^*; a). \end{aligned} \quad (12)$$

At the same time, based on the assumption before, if the model chooses to discard some of the information from system i , it must believe that other systems can provide more useful information, such that the total amount of effective information after refinement is greater than the information from system i alone, that is:

$$I(q, e_{\setminus i}^*; a) \geq I(q, e_i^{\text{useless}}; a). \quad (13)$$

Based on all the derivations above, we can conclude the following relationship:

$$\begin{aligned} H(a|q, e^*) &= H(a) - \min(I_1, I_2) \\ &\leq H(a) - I(q, e_i; a) \\ &= H(a|q, e_i), \end{aligned} \quad (14)$$

So the ensemble knowledge contains more useful information, and it helps reduce the information entropy in generating answers. Therefore, we consider that the process of ensemble can **introduce more helpful information than single system**. This is the core of RAG ensemble. In the following sections, we conduct a large number of experiments, clearly demonstrating the effectiveness of RAG ensemble.

4 Experiments

As shown in Figure 2, we conduct experiments to explore the RAG ensemble framework in depth. Based on the research in FlashRAG [31], we divide our experiments into ensemble at the pipeline level and at the module level. From the perspective of tasks, we carefully select datasets based on Wikipedia and MS MARCO [44]. We explore the following research questions regarding ensemble at the pipeline and module levels:

(1) Pipeline Level:

- **RQ1:** Does aggregating different pipelines effective? (§4.4)
- **RQ2:** Is the ensemble method still effective when aggregating closed - source model pipelines? (§4.4)
- **RQ3:** Is there a scaling - up phenomenon when aggregating? (§4.5)
- **RQ4:** Does the ensemble model show a preference when aggregating? (§4.6)

(2) Module Level:

- **RQ5:** When aggregating different generators' results, what is the performance of the ensemble framework? (§4.7)
- **RQ6:** Can aggregating different retrievers' results help improve the performance? (§4.8)
- **RQ7:** Is the ensemble method still effective when aggregating different rerankers? (§4.9)

4.1 Datasets

For the Wikipedia-based datasets, we carefully select four datasets: (1) TriviaQA [32], a reading comprehension dataset containing a lot of triples, (2) 2WikiMultiHopQA [19], a dataset containing multi-hop paths, (3) ARC [7], a multiple-choice dataset containing real scientific questions, and (4) WikiASP [18], a summarization generation dataset. They correspond to single-hop, multi-hop, multiple-choice, and open-domain summarization tasks, respectively. For each dataset, we randomly select 500 samples for evaluation. For the MS MARCO-based datasets, we follow the approach in RAG-Studio [43] and choose six vertical domains: Biomedical, Computing, Film, Finance, Law, and Music, with a maximum of 1000 test samples selected from each dataset. For the experiments on pipeline level, we use Wikipedia-based datasets, while for the module level, we use MS MARCO-based datasets. In the subsequent sections, we use “2Wiki” as the abbreviation for 2WikiMultiHopQA dataset.

4.2 Baselines

We select representative methods from each of four RAG frameworks: **Branching**, **Iterative**, **Loop** and **Agentic** methods. For the Agentic RAG technology, we conduct experiments with both prompting-based models and reinforcement learning models, respectively. For Branching framework, we choose Replug method [52]; for the Iterative framework, we choose Iter-RetGen method [51]; for the prompting-based Agentic framework, we choose Search-o1 method [37]; and for the RL-based Agentic framework, we choose R1-Searcher method [54]. Due to open-source models’ type limitations, for the Loop framework, we use Self-RAG with Llama3-8B-Instruct and FLARE with Qwen2.5-7B-Instruct as backbone models [2, 28].

In order to fully explore the characteristics of the ensemble method, in addition to RAG ensemble generation—method that directly generates an answer, we also set up a RAG ensemble selection method. Specifically, we use the same model as the one used for generation, but let it select the optimal answer from the given candidate answers. We call this method RAG ensemble selection.

4.3 Experimental Settings

In pipeline level experiments, we choose the Llama3-8B-Instruct[14] for the main experiment. In addition, we also select Llama3.1-8B-Instruct [14], Qwen2.5-7B-Instruct [47], Qwen2-7B-Instruct [66], and Mistral-7B-Instruct-v0.3 [24] for experiments related to base model ablation study. Due to the limitations of open-source models, for the Self-RAG method, we only use Llama3-8B-Instruct as base model. For retrieval step, we use e5-base [58] model as the retriever, and retrieve top-5 documents from Wikipedia dumps². We use exact match (EM) and F1 score for 2WikiMultiHopQA, TriviaQA, and ARC datasets as metrics, and use F1 score and ROUGE-L score for WikiASP dataset.

In the module level experiments, we use Llama3-8B-Instruct, Llama3.1-8B-Instruct, Qwen2-7B-Instruct, and Mistral-7B-Instruct-v0.3 for generator level and reranker level ensemble. For the experiment on retriever level ensemble, we choose e5-base model, contriever[22], and BM25[50] as our experiment’s retrievers. k1

²<https://archive.org/details/enwiki-20181220>

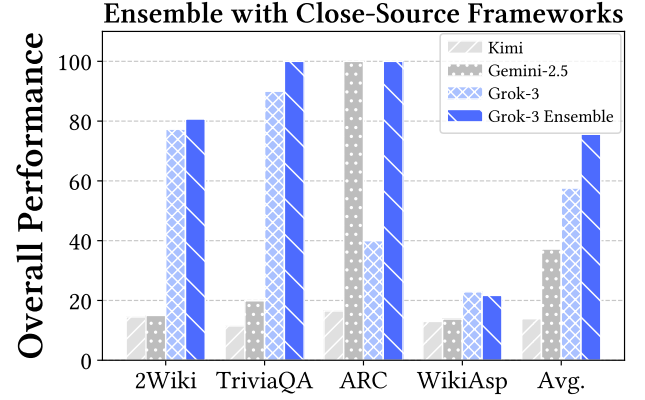


Figure 4: The performance of RAG Ensemble with different close-source frameworks.

parameter in BM25 algorithm is set to 1.5, and b parameter is set to 0.75.

4.4 Main Results on RAG Pipeline Ensemble

In this section, we perform a comparative performance analysis of RAG methods from both **open-source** and **closed-source** perspectives.

4.4.1 Analysis of Open-Source Methods. Table 1 presents the main results of our proposed RAG ensemble method. Overall, the ensemble method consistently demonstrates superior average performance across different backbone models in average performance. A more detailed analysis yields the following key observations:

- **Different methods excel in different tasks.** The effectiveness of different methods shows considerable variation depending on tasks. For instance, when using LLaMA as the backbone, RePlug achieves the best EM score on the ARC multiple-choice task, while Iter-RetGen outperforms RePlug and other variants on F1 score in the WikiASP summarization task. Similarly, R1-Searcher performs worse than RePlug on ARC and WikiASP but performs better on 2Wiki and TriviaQA. These results support our claim that no single RAG framework is universally optimal—different methods are more suited to specific task types.
- **The ensemble method’s stability and performance are superior to those of the baseline method on average.** Regardless of the backbone model, the performance of the ensemble method can reach the optimal or sub-optimal level on most metrics, and it achieves the best result in the average F1 score. This highlights the effectiveness of the ensemble method. In addition, whether integrating pure prompting-based RAG methods (such as Iter-RetGen and FLARE) or integrating RAG methods with training (such as Self-RAG and R1-Searcher), the ensemble method can also achieve good results. This highlights the stability of the ensemble method.
- **Generative ensemble outperforms selective ensemble in RAG.** Compared to selecting a single best candidate answer, fusing multiple answers through generative ensemble methods leads to superior performance. For instance, under the LLaMA-base

Table 1: Overall results on 4 QA datasets. The top two results in each backbone model’s baseline group are highlighted in bold and underlined. “Llama” means the backbone model is Llama3-8B-Instruct, and “Qwen” means the backbone model is Qwen2.5-7B-Instruct. The average score in the last column we report is the mean of F1 score across the four datasets. 2Wiki (2WikiMultiHopQA).

Method	Pipeline Type	2Wiki		TriviaQA		ARC		WikiASP		Avg.
		EM	F1	EM	F1	EM	F1	F1	Rouge-L	
Different RAG Baselines with Llama										
RePlug	Branching	13.6	20.0	53.4	61.3	84.4	62.4	21.7	11.6	41.4
Iter-RetGen	Iterative	10.2	20.2	58.0	67.1	76.2	53.0	23.1	12.3	40.9
Self-RAG	Loop	12.2	24.0	39.6	53.3	67.0	57.0	15.7	9.8	37.5
Search-o1	Agentic (Prompt)	15.2	23.6	50.6	59.3	26.8	34.4	16.0	5.8	33.3
R1-Searcher	Agentic (RL)	51.0	55.8	58.0	65.8	64.4	67.6	17.7	13.9	51.7
Multiple Systems Integrated Methods with Llama										
RAG Ensemble(Generation)		52.5	56.0	63.6	72.2	82.6	62.5	18.7	10.1	52.4
RAG Ensemble(Selection)		23.4	31.5	56.4	64.0	62.8	52.1	18.0	14.1	41.4
Different RAG Baselines with Qwen										
RePlug	Branching	23.8	27.5	47.4	54.5	91.0	68.2	18.9	11.3	42.3
Iter-RetGen	Iterative	30.2	36.2	59.2	66.8	90.4	68.0	23.1	12.0	48.5
FLARE	Loop	20.2	25.1	41.0	47.6	79.2	58.8	17.4	13.3	37.2
Search-o1	Agentic (Prompt)	32.0	39.0	55.0	63.2	93.0	70.0	18.4	10.2	47.7
R1-Searcher	Agentic (RL)	51.8	57.9	54.8	64.6	88.2	68.4	6.1	4.7	49.3
Multiple Systems Integrated Methods with Qwen										
RAG Ensemble(Generation)		54.8	64.5	55.2	65.8	91.8	70.4	19.8	14.3	55.1
RAG Ensemble(Selection)		37.2	43.5	54.5	63.2	90.4	69.0	17.1	12.6	48.2

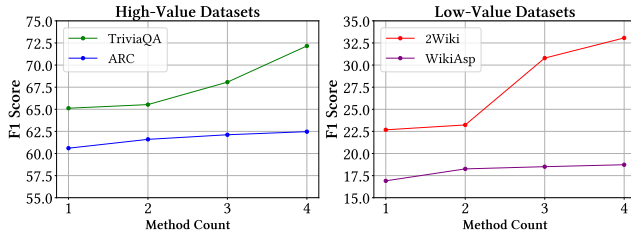


Figure 5: Ensemble of different scales at pipeline level.

setting, selection-based approaches underperform even some individual baselines, falling 27.6 points behind R1-Searcher on 2Wiki’s EM score and 21.6 points behind RePlug on ARC’s EM score. This performance gap may be attributed to the limited quality of selected candidates and inadequate comparative reasoning by the model. This result highlights the limitation of relying solely on a single RAG framework.

4.4.2 Analysis of Close-Source Methods. To comprehensively assess the applicability of the RAG Ensemble framework, we further conduct experiments on closed-source models. In this experiment, we select three closed-source models—Kimi, Gemini-2.5[1], and Grok-3—for RAG inference, using Kimi as the ensemble model.

Then, we randomly sample 20 subsets from each of four datasets for preliminary evaluation.

As shown in Figure 4, results demonstrate that the RAG Ensemble framework remains effective when applied to closed-source models. While individual models may achieve strong performance on specific tasks, they often underperform on others. This observation supports our claim in the introduction: the performance of a single RAG framework may vary significantly across tasks due to multiple factors. By aggregating outputs from different models, the RAG Ensemble effectively mitigates inter-system performance variance, thereby overcoming the limitations of any single model.

4.5 Ensemble System-Scale Scalability Analysis

Building on the theoretical analysis in section 3.2, the effectiveness of the RAG Ensemble framework primarily stems from the complementary information provided by different RAG systems, which helps mitigate the limitations of individual systems. To validate the hypothesis that incorporating more systems contributes to performance improvement, we investigate the scalability of the RAG system ensemble. In detail, we vary the number of aggregated systems from 1 to 4. For each setting, we enumerate all possible system combinations and report the average performance across combinations as the overall result for that scale.

As shown in Figure 5. Results indicate a clear upward trend in performance as the number of aggregated systems increases, revealing a strong positive correlation between ensemble size and

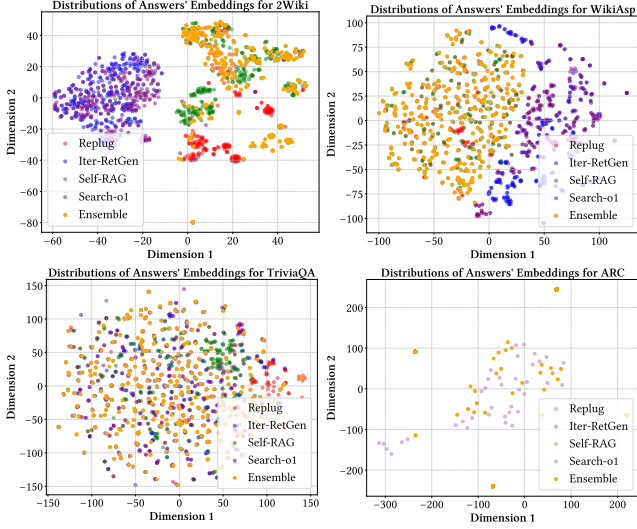


Figure 6: Distributions of ensemble answers and sub-answers.

reasoning effectiveness. This further supports our hypothesis: integrating information from more systems enhances final inference performance, demonstrating the framework’s strong scalability.

4.6 RAG Ensemble Preference Visualization

In preliminary experiments, we find that the ensemble method’s result is slightly lower than the best-performing baseline on certain datasets. We hypothesize that this may result from the ensemble model exhibiting preferences toward specific subsystems during answer generation. To validate this assumption, we further analyze the distributional differences between subsystem outputs and ensemble responses. Specifically, we use the BGE model [64] to encode answers generated by the ensemble and each subsystem, and apply t-SNE to visualize the embeddings. As shown in Figure 6, the visualization reveals two key insights:

- **The ensemble model exhibits clear preference.** On the 2Wiki and WikiASP datasets, the embedding distribution of ensemble answers is more closely aligned with those of Self-RAG and Replug, while differing significantly from Iter-RetGen and Search-o1. This indicates that, for these tasks, the ensemble model tends to rely more on information provided by Self-RAG and Replug.
- **The degree of ensemble preference correlates with subsystem performance.** On the TriviaQA and ARC datasets, where all subsystems achieve relatively high and comparable performance, the ensemble model shows no strong bias toward any particular method. In contrast, for more challenging tasks such as 2Wiki and WikiASP, where subsystem performance varies widely, the ensemble model tends to favor the stronger-performing methods. This suggests that ensemble preference is influenced, to some extent, by the performance distribution of the subsystems.

4.7 Generator-level Ensemble Analysis

Starting from this section, we will focus on exploring the experiments on RAG ensemble at the module level.

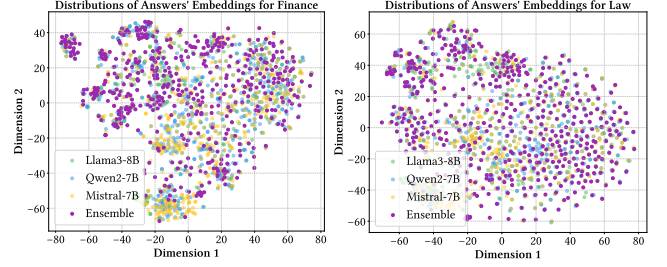


Figure 7: Distributions of ensemble answers and sub-answers on MS MARCO datasets.

Overall Results. To further evaluate the effectiveness of generator level’s ensemble in RAG systems, we design an experiment specifically targeting generator ensemble. For each question, the reference documents remain fixed, while only the answer-generating model is varied. All generated outputs are then fed into an ensemble model to produce a final answer. In this setup, Llama3-8B-Instruct, Qwen2-7B-Instruct, and Mistral-7B-Instruct-v0.3 serve as the generators to be aggregated, while Llama3.1-8B-Instruct and Qwen2.5-7B-Instruct are used as ensemble models. As shown in Table 2, our key findings are as follows:

- **Even under a fixed pipeline, aggregating outputs from different generators consistently yields strong performance gains.** In terms of the F1 score, ensemble using Llama3.1-8B-Instruct improves performance by 3.6% over the best-performing single-generator baseline (Llama3-8B-Instruct), while Qwen2.5-7B-Instruct achieves a 2.6% improvement. Across six vertical-domain tasks, the ensemble approach consistently reaches either the best or second-best results. This demonstrates ensemble framework’s robustness across domains.
- **Diversity among candidate answers plays a critical role in enhancing ensemble performance.** We compare the ensemble framework with a standard RAG setup where the same ensemble model directly generates answers using the same reference documents. Despite identical inputs, the ensemble framework, enriched with subsystem outputs, achieves superior results. For instance, using Llama3.1-8B-Instruct, the aggregated output yields a 6.7% improvement in average F1 score and a 4.1% gain in ROUGE-L compared to its standard RAG counterpart. This suggests that different generators offer complementary perspectives on the same evidence, which the ensemble model effectively synthesizes into more accurate answers. These results highlight the importance of answer diversity in improving final output quality.

Preference Visualization Analysis. Furthermore, we analyze the preference of the ensemble model in domain-specific tasks. We select the Finance and Law datasets for this analysis, with the results shown in the Figure 7. The results indicate that the ensemble model exhibits no clear preference in these two domains, as the embeddings of answers from individual systems and the aggregated outputs are relatively evenly distributed. This aligns with our earlier hypothesis that significant preference only emerges when the performance gap among subsystems is large. In the Finance

Table 2: Overall results for generator level’s ensemble on MS MARCO. The models we use here are all instruct versions. In the ensemble experiments, “Llama” means the ensemble model is Llama3-8B-Instruct, and “Qwen” means the ensemble model is Qwen2.5-7B-Instruct.

Backbone	Biomedical		Computing		Film		Finance		Law		Music		Avg.	
	F1	Rouge-L	F1	Rouge-L	F1	Rouge-L	F1	Rouge-L	F1	Rouge-L	F1	Rouge-L	F1	Rouge-L
<i>Different Base Models for Standard RAG</i>														
Llama3-8B	29.6	21.2	32.7	23.1	<u>45.9</u>	42.2	28.9	23.2	31.1	25.6	45.4	40.0	35.6	29.2
Qwen2-7B	29.5	21.2	30.1	24.0	38.8	33.7	29.2	22.7	31.7	21.8	40.6	34.5	33.3	26.3
Mistral-7B	26.5	19.6	31.2	<u>24.9</u>	38.4	34.6	26.5	20.6	30.1	24.5	36.7	31.5	31.6	26.0
Llama3.1-8B	27.5	21.2	28.9	23.1	44.8	41.9	28.3	23.8	29.1	24.1	43.3	37.7	33.7	28.6
Qwen2.5-7B	28.4	21.5	29.8	23.2	36.5	32.0	28.1	21.8	31.2	25.6	41.2	35.2	32.5	26.6
<i>Multiple Generators Ensemble</i>														
Ensemble with Llama	31.6	22.6	31.3	24.6	52.3	47.9	34.8	28.1	33.0	26.7	51.9	46.1	39.2	32.7
Ensemble with Qwen	<u>31.0</u>	<u>22.3</u>	<u>32.6</u>	25.6	52.3	<u>47.5</u>	<u>32.9</u>	<u>26.0</u>	<u>32.1</u>	<u>26.1</u>	<u>48.0</u>	<u>42.4</u>	<u>38.2</u>	<u>31.7</u>

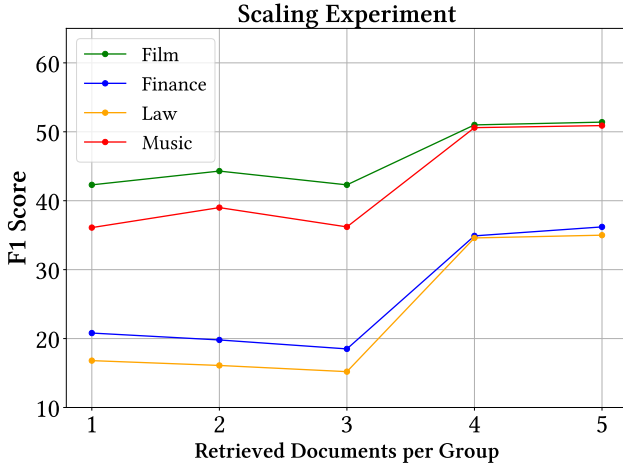


Figure 8: F1 Score for retriever level’s ensemble on MS MARCO.

and Law tasks, the performance of different generators is relatively similar, leading to low distinguishability among their outputs and, consequently, no strong preference in the ensemble model.

4.8 Retriever-level Ensemble Analysis

To fully explore the generalization of ensemble method in the retriever module, we conduct experiments using different retrievers on the MS MARCO datasets. We use three retrievers: BM25, Contriever, and E5, with each retriever recalling the top-5 documents. Additionally, we also perform scaling-up experiments at the document quantity level.

The experimental results are shown in Table 3. For all datasets, the ensemble method at the retriever level is still superior compared to the single-retriever RAG method. This highlights the strong generalization ability of the RAG ensemble at the retriever level. Subsequently, we continue with the scaling-up experiment on the number of retrieved documents, and the experimental results are shown in Figure 8.

Table 3: Different base models’ results and ensemble results under the same pipeline on TriviaQA dataset. All the base models are instruct versions.

Method	Biomedical	Computing	Film	Finance	Law	Music
BM25	17.3	16.0	38.0	18.1	15.2	34.9
Contriever	16.8	16.1	39.0	19.0	15.9	<u>37.5</u>
E5	<u>18.2</u>	<u>18.3</u>	<u>42.4</u>	<u>19.5</u>	<u>17.8</u>	37.3
Ensemble	33.1	34.4	51.4	36.2	35.0	50.9

In the early stages of documents increase, the improvement in ensemble performance is relatively small (Film, Music), and in some cases, the performance even decreases (Finance, Law). However, as the documents’ number continue to increase, a significant improvement in ensemble performance can be observed across all datasets. We consider that this is because the robustness of the model can have a significant impact on the effectiveness of ensemble task. In the early stages of documents increase, the ensemble model may not have been sufficiently robust due to the increase in redundant information. However, when the number of documents reaches a certain threshold, the valuable information might be more easily captured by the model due to its multi-perspective presentation, thus suppressing the noise and leading to a substantial performance improvement. This also highlights the importance of the stability of the ensemble model for this task.

4.9 Reranker-level Ensemble Analysis

In this section, we investigate the impact of reranker-level ensemble on RAG performance. The experimental setup is as follows: for each query, we retrieve ten relevant documents and use a base model to rank them, selecting the top five based on relevance.

Since the model may introduce hallucinations during ranking, we perform a post-check on the ranked list. Specifically, we remove duplicate document identifiers and filter out IDs outside the 1–10 range. For any missing identifiers after cleaning, we pad the sequence at the end to maintain a consistent length.

For the base model experiment, the same base model directly performs RAG over the top five ranked documents to generate an

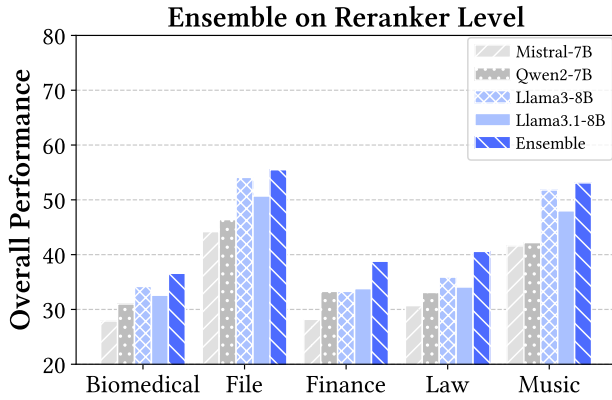


Figure 9: F1 Score for reranker level’s ensemble on MS MARCO.

answer. For the ensemble experiment, we provide the ensemble model with three sets of top-ranked documents, each produced by a different base model. We explicitly indicate that each set is ordered by descending relevance, guiding the model to generate a final answer based on all fifteen documents. Note that the ensemble model receives only the documents as input, without any intermediate answers. The output behavior of the ensemble model is summarized in the figure 9.

Experimental results show that ensemble at the reranker level remains effective. The aggregated performance consistently surpasses that of standard RAG with reranking. Additionally, the model demonstrates strong robustness in ensemble. Since different base models may rank the same document differently, the aggregated input may contain conflicting relevance signals. Nevertheless, the model is still able to produce accurate final answers, suggesting a degree of self-discrimination and resistance to noise during ensemble.

5 Conclusion

In this paper, we perform a thorough analysis of the method for aggregating information from multiple RAG systems to derive comprehensive answers. This is the first detailed analysis of the ensemble method about RAG ensemble framework. We establish a mathematical model that provides a solid formulation of the ensemble process. In addition, we conduct a lot of experiments at both the pipeline level and the module level, fully demonstrating the broad adaptability, effectiveness, and stability of the RAG ensemble framework. Moreover, we have drawn some important conclusions. For example, we find that the RAG system ensemble framework exhibits a scaling-up phenomenon, and that the ensemble model has different preferences for tasks of varying difficulty levels. These conclusions are supported by a series of our experiments. We hope this paper serves as a reference for research on the RAG system ensemble and encourages further work to optimize RAG system performance.

GenAI Usage Disclosure

In this paper, there is no use of GenAI tools whatsoever in any stage of the research.

References

- [1] Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M. Dai, Anja Hauth, Katie Millican, David Silver, Slav Petrov, Melvin Johnson, Ioannis Antonoglou, Julian Schrittwieser, Amelia Glaese, Jilin Chen, Emily Pitler, Timothy P. Lillicrap, Angeliki Lazaridou, Orhan Firat, James Molloy, Michael Isard, Paul Ronald Barham, Tom Hennigan, Benjamin Lee, Fabio Viola, Malcolm Reynolds, Yuanzhong Xu, Ryan Doherty, Eli Collins, Clemens Meyer, Eliza Rutherford, Erica Moreira, Kareem Ayoub, Megha Goel, George Tucker, Enrique Piqueras, Maxim Krikun, Iain Barr, Nikolay Savinov, Ivo Danihelka, Becca Roelofs, Anaïs White, Anders Andreassen, Tamara von Glehn, Lakshman Yagati, Mehran Kazemi, Lucas Gonzalez, Misha Khalman, Jakub Sygnowski, and et al. 2023. Gemini: A Family of Highly Capable Multimodal Models. *CoRR* abs/2312.11805 (2023). doi:10.48550/ARXIV.2312.11805 arXiv:2312.11805
- [2] Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2024. Self-RAG: Learning to Retrieve, Generate, and Critique through Self-Reflection. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7–11, 2024*. OpenReview.net. <https://openreview.net/forum?id=hSyW5go0v8>
- [3] Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George van den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, Diego de Las Casas, Aurelia Guy, Jacob Menick, Roman Ring, Tom Hennigan, Saffron Huang, Loren Maggiore, Chris Jones, Albin Cassirer, Andy Brock, Michela Paganini, Geoffrey Irving, Oriol Vinyals, Simon Osindero, Karen Simonyan, Jack W. Rae, Erich Elsen, and Laurent Sifre. 2022. Improving Language Models by Retrieving from Trillions of Tokens. In *International Conference on Machine Learning, ICML 2022, 17–23 July 2022, Baltimore, Maryland, USA (Proceedings of Machine Learning Research, Vol. 162)*, Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato (Eds.). PMLR, 2206–2240. <https://proceedings.mlr.press/v162/borgeaud22a.html>
- [4] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language Models are Few-Shot Learners. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6–12, 2020, virtual*, Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (Eds.). <https://proceedings.neurips.cc/paper/2020/hash/1457c0dbfcb4967418bfb8ac142f64a-Abstract.html>
- [5] Lingjiao Chen, Matei Zaharia, and James Zou. 2023. FrugalGPT: How to Use Large Language Models While Reducing Cost and Improving Performance. *CoRR* abs/2305.05176 (2023). doi:10.48550/ARXIV.2305.05176 arXiv:2305.05176
- [6] Mingyang Chen, Tianpeng Li, Haoze Sun, Yijie Zhou, Chenzheng Zhu, Haofen Wang, Jeff Z. Pan, Wen Zhang, Huajun Chen, Fan Yang, Zenan Zhou, and Weipeng Chen. 2025. ReSearch: Learning to Reason with Search for LLMs via Reinforcement Learning. arXiv:2503.19470 [cs.AI] <https://arxiv.org/abs/2503.19470>
- [7] Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have Solved Question Answering? Try ARC, the AI2 Reasoning Challenge. *CoRR* abs/1803.05457 (2018). arXiv:1803.05457 <http://arxiv.org/abs/1803.05457>
- [8] Guanting Dong, Yifei Chen, Xiaoxi Li, Jiajie Jin, Hongjin Qian, Yutao Zhu, Hangyu Mao, Guorui Zhou, Zhicheng Dou, and Ji-Rong Wen. 2025. Tool-Star: Empowering LLM-Brained Multi-Tool Reasoner via Reinforcement Learning. *CoRR* abs/2505.16410 (2025). doi:10.48550/ARXIV.2505.16410 arXiv:2505.16410
- [9] Guanting Dong, Jiajie Jin, Xiaoxi Li, Yutao Zhu, Zhicheng Dou, and Ji-Rong Wen. 2025. RAG-Critic: Leveraging Automated Critic-Guided Agentic Workflow for Retrieval Augmented Generation. 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*, Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (Eds.). Association for Computational Linguistics, 3551–3578. <https://aclanthology.org/2025.acl-long.179/>
- [10] Guanting Dong, Xiaoxi Li, Yuyao Zhang, and Mengjie Deng. 2025. Leveraging LLM-Assisted Query Understanding for Live Retrieval-Augmented Generation. *CoRR* abs/2506.21384 (2025). doi:10.48550/ARXIV.2506.21384 arXiv:2506.21384
- [11] Guanting Dong, Keming Lu, Chengpeng Li, Tingyu Xia, Bowen Yu, Chang Zhou, and Jingren Zhou. 2025. Self-play with Execution Feedback: Improving Instruction-following Capabilities of Large Language Models. In *The Thirtieth International Conference on Learning Representations, ICLR 2025, Singapore, April*

- 24-28, 2025. OpenReview.net. <https://openreview.net/forum?id=cRR0oDFEBC>
- [12] Guanting Dong, Hangyu Mao, Kai Ma, Licheng Bao, Yifei Chen, Zhongyuan Wang, Zhongxia Chen, Jiazhen Du, Huiyang Wang, Fuzheng Zhang, Guorui Zhou, Yutao Zhu, Ji-Rong Wen, and Zhicheng Dou. 2025. Agentic Reinforced Policy Optimization. arXiv:2507.19849 [cs.LG] <https://arxiv.org/abs/2507.19849>
 - [13] Guanting Dong, Xiaoshuai Song, Yutao Zhu, Runqi Qiao, Zhicheng Dou, and Ji-Rong Wen. 2024. Toward General Instruction-Following Alignment for Retrieval-Augmented Generation. *CoRR* abs/2410.09584 (2024). doi:10.48550/ARXIV.2410.09584 arXiv:2410.09584
 - [14] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurélien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Rozière, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Grégoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel M. Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelle van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Jo Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, and et al. 2024. The Llama 3 Herd of Models. *CoRR* abs/2407.21783 (2024). doi:10.48550/ARXIV.2407.21783 arXiv:2407.21783
 - [15] Luyu Gao, Xueguang Ma, Jimmy Lin, and Jamie Callan. 2023. Precise Zero-Shot Dense Retrieval without Relevance Labels. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL 2023, Toronto, Canada, July 9-14, 2023, Anna Rogers, Jordan L. Boyd-Graber, and Naoaki Okazaki (Eds.). Association for Computational Linguistics, 1762–1777. doi:10.18653/V1/2023.ACL-LONG.99
 - [16] Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Qianyu Guo, Meng Wang, and Haofen Wang. 2023. Retrieval-Augmented Generation for Large Language Models: A Survey. *CoRR* abs/2312.10997 (2023). doi:10.48550/ARXIV.2312.10997 arXiv:2312.10997
 - [17] Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Ming-Wei Chang. 2020. Retrieval Augmented Language Model Pre-Training. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event (Proceedings of Machine Learning Research, Vol. 119)*. PMLR, 3929–3938. <http://proceedings.mlr.press/v119/guu20a.html>
 - [18] Hiroaki Hayashi, Prashant Budania, Peng Wang, Chris Ackerson, Raj Neervannan, and Graham Neubig. 2021. WikiAsp: A Dataset for Multi-domain Aspect-based Summarization. *Trans. Assoc. Comput. Linguistics* 9 (2021), 211–225. doi:10.1162/TACL_A_00362
 - [19] Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. 2020. Constructing A Multi-hop QA Dataset for Comprehensive Evaluation of Reasoning Steps. In *Proceedings of the 28th International Conference on Computational Linguistics, COLING 2020, Barcelona, Spain (Online), December 8-13, 2020*, Doña Scott, Núria Bel, and Chengqing Zong (Eds.). International Committee on Computational Linguistics, 6609–6625. doi:10.18653/V1/2020.COLING-MAIN.580
 - [20] Hieu Hoang, Huda Khayrallah, and Marcin Junczys-Dowmunt. 2024. On-the-Fly Fusion of Large Language Models and Machine Translation. In *Findings of the Association for Computational Linguistics: NAACL 2024, Mexico City, Mexico, June 16-21, 2024*, Kevin Duh, Helena Gómez-Adorno, and Steven Bethard (Eds.). Association for Computational Linguistics, 520–532. doi:10.18653/V1/2024.FINDINGS-NAACL.35
 - [21] Kaixuan Huang, Xudong Guo, and Mengdi Wang. 2024. SpecDec++: Boosting Speculative Decoding via Adaptive Candidate Lengths. *CoRR* abs/2405.19715 (2024). doi:10.48550/ARXIV.2405.19715 arXiv:2405.19715
 - [22] Gautier Izcard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2022. Unsupervised Dense Information Retrieval with Contrastive Learning. *Trans. Mach. Learn. Res.* 2022 (2022). <https://openreview.net/forum?id=jKN1pXi7b0>
 - [23] Soyeong Jeong, Jinheon Baek, Sukmin Cho, Sung Ju Hwang, and Jong Park. 2024. Adaptive-RAG: Learning to Adapt Retrieval-Augmented Large Language Models through Question Complexity. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, NAACL 2024, Mexico City, Mexico, June 16-21, 2024, Kevin Duh, Helena Gómez-Adorno, and Steven Bethard (Eds.). Association for Computational Linguistics, 7036–7050. doi:10.18653/V1/2024.NAACL-LONG.389
 - [24] Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de Las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Léo Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7B. *CoRR* abs/2310.06825 (2023). doi:10.48550/ARXIV.2310.06825 arXiv:2310.06825
 - [25] Dongfu Jiang, Xiang Ren, and Bill Yuchen Lin. 2023. LLM-Blender: Ensembling Large Language Models with Pairwise Ranking and Generative Fusion. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL 2023, Toronto, Canada, July 9-14, 2023, Anna Rogers, Jordan L. Boyd-Graber, and Naoaki Okazaki (Eds.). Association for Computational Linguistics, 14165–14178. doi:10.18653/V1/2023.ACL-LONG.792
 - [26] Huiqiang Jiang, Qianhui Wu, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. 2023. LLMingua: Compressing Prompts for Accelerated Inference of Large Language Models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, Houda Bouamor, Juan Pino, and Kalika Bali (Eds.). Association for Computational Linguistics, 13358–13376. doi:10.18653/V1/2023.EMNLP-MAIN.825
 - [27] Huiqiang Jiang, Qianhui Wu, Xufang Luo, Dongsheng Li, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. 2024. LongLLMingua: Accelerating and Enhancing LLMs in Long Context Scenarios via Prompt Compression. 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, Lun-Wei Ku, Andre Martins, and Vivek Srikumar (Eds.). Association for Computational Linguistics, 1658–1677. doi:10.18653/V1/2024.ACL-LONG.91
 - [28] Zhengbao Jiang, Frank F. Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023. Active Retrieval Augmented Generation. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, Houda Bouamor, Juan Pino, and Kalika Bali (Eds.). Association for Computational Linguistics, 7969–7992. doi:10.18653/V1/2023.EMNLP-MAIN.495
 - [29] Bowen Jin, Hansi Zeng, Zhenrui Yue, Dong Wang, Hamed Zamani, and Jiawei Han. 2025. Search-R1: Training LLMs to Reason and Leverage Search Engines with Reinforcement Learning. *CoRR* abs/2503.09516 (2025). doi:10.48550/ARXIV.2503.09516 arXiv:2503.09516
 - [30] Jiajie Jin, Xiaoxi Li, Guanting Dong, Yuyao Zhang, Yutao Zhu, Zhao Yang, Hongjin Qian, and Zhicheng Dou. 2025. Decoupled Planning and Execution: A Hierarchical Reasoning Framework for Deep Search. *CoRR* abs/2507.02652 (2025). doi:10.48550/ARXIV.2507.02652 arXiv:2507.02652
 - [31] Jiajie Jin, Yutao Zhu, Xinyu Yang, Chenghao Zhang, and Zhicheng Dou. 2024. FlashRAG: A Modular Toolkit for Efficient Retrieval-Augmented Generation Research. *CoRR* abs/2405.13576 (2024). doi:10.48550/ARXIV.2405.13576 arXiv:2405.13576
 - [32] Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke Zettlemoyer. 2017. TriviaQA: A Large Scale Distantly Supervised Challenge Dataset for Reading Comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 1: Long Papers*, Regina Barzilay and Min-Yen Kan (Eds.). Association for Computational Linguistics, 1601–1611. doi:10.18653/V1/P17-1147
 - [33] Patrick S. H. Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Kuttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (Eds.). <https://proceedings.neurips.cc/paper/2020/hash/6b493230205f780e1bc26945df7481e5-Abstract.html>
 - [34] Jinyi Li, Yihui Lan, Lei Wang, and Hao Wang. 2024. PCToolkit: A Unified Plug-and-Play Prompt Compression Toolkit of Large Language Models. *CoRR* abs/2403.17411 (2024). doi:10.48550/ARXIV.2403.17411 arXiv:2403.17411
 - [35] Kuan Li, Zhongwang Zhang, Huifeng Yin, Liwen Zhang, Litu Ou, Jialong Wu, Wenbiao Yin, Baixuan Li, Zhengwei Tao, Xinyu Wang, Weizhou Shen, Junkai Zhang, Dingchu Zhang, Xixi Wu, Yong Jiang, Ming Yan, Pengjun Xie, Fei Huang, and Jingren Zhou. 2025. WebSailor: Navigating Super-human Reasoning for Web Agent. *CoRR* abs/2507.02592 (2025). doi:10.48550/ARXIV.2507.02592 arXiv:2507.02592
 - [36] Tianlin Li, Qian Liu, Tianyu Pang, Chao Du, Qing Guo, Yang Liu, and Min Lin. 2024. Purifying Large Language Models by Ensembling a Small Language Model. *CoRR* abs/2402.14845 (2024). doi:10.48550/ARXIV.2402.14845 arXiv:2402.14845
 - [37] Xiaoxi Li, Guanting Dong, Jiajie Jin, Yuyao Zhang, Yujia Zhou, Yutao Zhu, Peitian Zhang, and Zhicheng Dou. 2025. Search-o1: Agentic Search-Enhanced Large Reasoning Models. *CoRR* abs/2501.05366 (2025). doi:10.48550/ARXIV.2501.05366 arXiv:2501.05366
 - [38] Xiaoxi Li, Jiajie Jin, Guanting Dong, Hongjin Qian, Yutao Zhu, Yongkang Wu, Ji-Rong Wen, and Zhicheng Dou. 2025. WebThinker: Empowering Large Reasoning Models with Deep Research Capability. *CoRR* abs/2504.21776 (2025). doi:10.48550/ARXIV.2504.21776 arXiv:2504.21776

- [39] Junyi Liu, Liangzhi Li, Tong Xiang, Bowen Wang, and Yiming Qian. 2023. TCRA-LLM: Token Compression Retrieval Augmented Large Language Model for Inference Cost Reduction. In *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, Houda Bouamor, Juan Pino, and Kalika Bali (Eds.). Association for Computational Linguistics, 9796–9810. doi:10.18653/V1/2023.FINDINGS-EMNLP.655
- [40] Keming Lu, Hongyi Yuan, Runji Lin, Junyang Lin, Zheng Yuan, Chang Zhou, and Jingren Zhou. 2024. Routing to the Expert: Efficient Reward-guided Ensemble of Large Language Models. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, NAACL 2024, Mexico City, Mexico, June 16-21, 2024, Kevin Duh, Helena Gómez-Adorno, and Steven Bethard (Eds.). Association for Computational Linguistics, 1964–1974. doi:10.18653/V1/2024.NAACL-LONG.109
- [41] Yiyao Lu, Weiyei Meng, Liangcai Shu, Clement T. Yu, and King-Lup Liu. 2005. Evaluation of Result Merging Strategies for Metasearch Engines. In *Web Information Systems Engineering - WISE 2005, 6th International Conference on Web Information Systems Engineering*, New York, NY, USA, November 20-22, 2005, *Proceedings (Lecture Notes in Computer Science, Vol. 3806)*, Anne H. H. Ngu, Masaru Kitsuregawa, Erich J. Neuhold, Jen-Yao Chung, and Quan Z. Sheng (Eds.). Springer, 53–66. doi:10.1007/11581062_5
- [42] Kelong Mao, Zheng Liu, Hongjin Qian, Fengran Mo, Chenlong Deng, and Zhicheng Dou. 2024. RAG-Studio: Towards In-Domain Adaptation of Retrieval Augmented Generation Through Self-Alignment. In *Findings of the Association for Computational Linguistics: EMNLP 2024, Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (Eds.)*. Association for Computational Linguistics, Miami, Florida, USA, 725–735. doi:10.18653/v1/2024.findings-emnlp.41
- [43] Kelong Mao, Zheng Liu, Hongjin Qian, Fengran Mo, Chenlong Deng, and Zhicheng Dou. 2024. RAG-Studio: Towards In-Domain Adaptation of Retrieval Augmented Generation Through Self-Alignment. In *Findings of the Association for Computational Linguistics: EMNLP 2024, Miami, Florida, USA, November 12-16, 2024, Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (Eds.)*. Association for Computational Linguistics, 725–735. <https://aclanthology.org/2024.findings-emnlp.41>
- [44] Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, and Li Deng. 2016. MS MARCO: A Human Generated Machine Reading Comprehension Dataset. *CoRR abs/1611.09268* (2016). arXiv:1611.09268 <http://arxiv.org/abs/1611.09268>
- [45] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. 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*, Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (Eds.). http://papers.nips.cc/paper_files/paper/2022/hash/b1efde53be364a73914f58805a001731-Abstract-Conference.html
- [46] Runqi Qiao, Qiuna Tan, Guanting Dong, Minhui Wu, Chong Sun, Xiaoshuai Song, Jiapeng Wang, Zhuoma Gongque, Shanglin Lei, YiFan Zhang, Zhe Wei, Miaoquan Zhang, Runfeng Qiao, Xiao Zong, Yida Xu, Peiqing Yang, Zhimin Mao, Muxi Diao, Chen Li, and Honggang Zhang. 2025. We-Math: Does Your Large Multimodal Model Achieve Human-like Mathematical Reasoning?. 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, Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (Eds.). Association for Computational Linguistics, 20023–20070. <https://aclanthology.org/2025.acl-long.983/>
- [47] Qwen, ., An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. 2025. Qwen2.5 Technical Report. arXiv:2412.15115 [cs.CL] <https://arxiv.org/abs/2412.15115>
- [48] Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay, Amnon Shashua, Kevin Leyton-Brown, and Yoav Shoham. 2023. In-Context Retrieval-Augmented Language Models. *Trans. Assoc. Comput. Linguistics* 11 (2023), 1316–1331. doi:10.1162/TACL_A_00605
- [49] M. Elena Renda and Umberto Straccia. 2003. Web Metasearch: Rank vs. Score Based Rank Aggregation Methods. In *Proceedings of the 2003 ACM Symposium on Applied Computing (SAC), March 9-12, 2003, Melbourne, FL, USA*, Gary B. Lamont, Hisham Haddad, George A. Papadopoulos, and Brajendra Panda (Eds.). ACM, 841–846. doi:10.1145/952532.952698
- [50] Stephen E. Robertson, Steve Walker, and Micheline Hancock-Beaulieu. 1995. Large Test Collection Experiments on an Operational, Interactive System: Okapi at TREC. *Inf. Process. Manag.* 31, 3 (1995), 345–360. doi:10.1016/0306-4573(94)00051-4
- [51] Zhihong Shao, Yeyun Gong, Yelong Shen, Minlie Huang, Nan Duan, and Weizhu Chen. 2023. Enhancing Retrieval-Augmented Large Language Models with Iterative Retrieval-Generation Synergy. In *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, Houda Bouamor, Juan Pino, and Kalika Bali (Eds.). Association for Computational Linguistics, 9248–9274. doi:10.18653/V1/2023.FINDINGS-EMNLP.620
- [52] Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Richard James, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. 2024. REPLUG: Retrieval-Augmented Black-Box Language Models. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, NAACL 2024, Mexico City, Mexico, June 16-21, 2024, Kevin Duh, Helena Gómez-Adorno, and Steven Bethard (Eds.). Association for Computational Linguistics, 8371–8384. doi:10.18653/V1/2024.NAACL-LONG.463
- [53] Tal Shnitzer, Anthony Ou, Mirian Silva, Kate Soule, Yuekai Sun, Justin Solomon, Neil Thompson, and Mikhail Yurochkin. 2023. Large Language Model Routing with Benchmark Datasets. *CoRR abs/2309.15789* (2023). doi:10.48550/ARXIV.2309.15789 arXiv:2309.15789
- [54] Huatong Song, Jinhao Jiang, Yingqian Min, Jie Chen, Zhipeng Chen, Wayne Xin Zhao, Lei Fang, and Ji-Rong Wen. 2025. R1-Searcher: Incentivizing the Search Capability in LLMs via Reinforcement Learning. *CoRR abs/2503.05592* (2025). doi:10.48550/ARXIV.2503.05592 arXiv:2503.05592
- [55] Zhengwei Tao, Jialong Wu, Wenbiao Yin, Junkai Zhang, Baixuan Li, Haiyang Shen, Kuan Li, Liwen Zhang, Xinyu Wang, Yong Jiang, Pengjun Xie, Fei Huang, and Jingren Zhou. 2025. WebShaper: Agentically Data Synthesizing via Information-Seeking Formalization. *CoRR abs/2507.15061* (2025). doi:10.48550/ARXIV.2507.15061 arXiv:2507.15061
- [56] Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijie Chen, Yanru Chen, Yuankun Chen, Yutian Chen, Zhuofu Chen, Jialei Cui, Hao Ding, Mengnan Dong, Angang Du, Chenzhuang Du, Dikang Du, Yulun Du, Yu Fan, Yichen Feng, Kelin Fu, Bofei Gao, Hongcheng Gao, Peizhong Gao, Tong Gao, Xinran Gu, Longyu Guan, Haiqing Guo, Jianhang Guo, Hao Hu, Xiaoru Hao, Tianhong He, Weiran He, Wenyang He, Chao Hong, Yangyang Hu, Zhenxing Hu, Weixiao Huang, Zhiqi Huang, Zihao Huang, Tao Jiang, Zhejun Jiang, Xinyi Jin, Yongsheng Kang, Guokun Lai, Cheng Li, Fang Li, Haoyang Li, Ming Li, Wentao Li, Yanhao Li, Yiwei Li, Zhaowei Li, Zheming Li, Hongzhan Lin, Xiaohan Lin, Zongyu Lin, Chengyin Liu, Chenyu Liu, Hongzhang Liu, Jingyuan Liu, Junqi Liu, Liang Liu, Shaowei Liu, T. Y. Liu, Tianwei Liu, Weizhou Liu, Yangyang Liu, Yibo Liu, Yiping Liu, Yue Liu, Zhengying Liu, Enzhe Lu, Lijun Lu, Shengling Ma, Xinyu Ma, Yingwei Ma, Shaoguang Mao, Jie Mei, Xin Men, Yibo Miao, Siyuan Pan, Yebo Peng, Ruoyu Qin, Bowen Qu, Zeyu Shang, Lidong Shi, Shengyuan Shi, Feifan Song, Jianlin Su, Zhengyuan Su, Xinjie Sun, Flood Sung, Heyi Tang, Jiawen Tao, Qifeng Teng, Chensi Wang, Dinglu Wang, Feng Wang, Haiming Wang, Jianzhou Wang, Jiaxing Wang, Jinhong Wang, Shengjie Wang, Shuyi Wang, Yao Wang, Yejie Wang, Yiqin Wang, Yuxin Wang, Yuzhi Wang, Zhaoji Wang, Zhengtao Wang, Zhexu Wang, Chu Wei, Qianqian Wei, Wenhao Wu, Xingzhe Wu, Yuxin Wu, Chenjun Xiao, Xiaotong Xie, Weimin Xiong, Boyu Xu, Jing Xu, Jinjing Xu, L. H. Xu, Lin Xu, Suting Xu, Weixin Xu, Xinran Xu, Yangchuan Xu, Ziyao Xu, Junjie Yan, Yuzi Yan, Xiaofei Yang, Ying Yang, Zhen Yang, Zhilin Yang, Zonghan Yang, Haotian Yao, Xingcheng Yao, Wenjie Ye, Zhuorui Ye, Bohong Yin, Longhui Yu, Enming Yuan, Hongbang Yuan, Mengjie Yuan, Haobing Zhan, Dehao Zhang, Hao Zhang, Wanlu Zhang, Xiaobin Zhang, Yangkun Zhang, Yizhi Zhang, Yongting Zhang, Yu Zhang, Yutao Zhang, Yutong Zhang, Zheng Zhang, Haotian Zhao, Yikai Zhao, Huabin Zheng, Shaojie Zheng, Jianren Zhou, Xinyu Zhou, Zaida Zhou, Zhen Zhu, Weiye Zhuang, and Xinxing Zu. 2025. Kimi K2: Open Agentic Intelligence. arXiv:2507.20534 [cs.LG] <https://arxiv.org/abs/2507.20534>
- [57] Naftali Tishby, Fernando C. N. Pereira, and William Bialek. 2000. The information bottleneck method. *CoRR physics/0004057* (2000). <http://arxiv.org/abs/physics/0004057>
- [58] Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. 2022. Text Embeddings by Weakly-Supervised Contrastive Pre-training. *CoRR abs/2212.03533* (2022). doi:10.48550/ARXIV.2212.03533 arXiv:2212.03533
- [59] Liang Wang, Nan Yang, and Furu Wei. 2023. Query2doc: Query Expansion with Large Language Models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, Houda Bouamor, Juan Pino, and Kalika Bali (Eds.). Association for Computational Linguistics, 9414–9423. doi:10.18653/V1/2023.EMNLP-MAIN.585
- [60] Yile Wang, Peng Li, Maosong Sun, and Yang Liu. 2023. Self-Knowledge Guided Retrieval Augmentation for Large Language Models. In *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, Houda Bouamor, Juan Pino, and Kalika Bali (Eds.). Association for Computational Linguistics, 10303–10315. doi:10.18653/V1/2023.FINDINGS-EMNLP.691
- [61] Jiming Wu, Zihao Deng, Wei Li, Yiding Liu, Bo You, Bo Li, Zejun Ma, and Ziwei Liu. 2025. MMSearch-R1: Incentivizing LLMs to Search. arXiv:2506.20670 [cs.CV] <https://arxiv.org/abs/2506.20670>

- [62] Jialong Wu, Baixuan Li, Runnan Fang, Wenbiao Yin, Liwen Zhang, Zhengwei Tao, Dingchu Zhang, Zekun Xi, Yong Jiang, Pengjun Xie, Fei Huang, and Jingren Zhou. 2025. WebDancer: Towards Autonomous Information Seeking Agency. *CoRR* abs/2505.22648 (2025). doi:10.48550/ARXIV.2505.22648 arXiv:2505.22648
- [63] Jialong Wu, Wenbiao Yin, Yong Jiang, Zhenglin Wang, Zekun Xi, Runnan Fang, Linhai Zhang, Yulan He, Deyu Zhou, Pengjun Xie, and Fei Huang. 2025. WebWalker: Benchmarking LLMs in Web Traversal. 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*, Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (Eds.). Association for Computational Linguistics, 10290–10305. <https://aclanthology.org/2025.acl-long.508/>
- [64] Shitao Xiao, Zheng Liu, Peitian Zhang, Niklas Muennighoff, Defu Lian, and Jian-Yun Nie. 2024. C-Pack: Packed Resources For General Chinese Embeddings. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2024, Washington DC, USA, July 14-18, 2024*, Grace Hui Yang, Hongning Wang, Sam Han, Claudia Hauff, Guido Zuccon, and Yi Zhang (Eds.). ACM, 641–649. doi:10.1145/3626772.3657878
- [65] Fangyuan Xu, Weijia Shi, and Eunsol Choi. 2024. RECOMP: Improving Retrieval-Augmented LMs with Context Compression and Selective Augmentation. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net. <https://openreview.net/forum?id=mlJLVigNhp>
- [66] An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jianxin Yang, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Xuejing Liu, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuyang Liu, Zeyu Cui, Zhenru Zhang, Zhifang Guo, and Zhihao Fan. 2024. Qwen2 Technical Report. *CoRR* abs/2407.10671 (2024). doi:10.48550/ARXIV.2407.10671 arXiv:2407.10671
- [67] Ori Yoran, Tomer Wolfson, Ori Ram, and Jonathan Berant. 2024. Making Retrieval-Augmented Language Models Robust to Irrelevant Context. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net. <https://openreview.net/forum?id=ZS4m74kZpH>
- [68] Wenhao Yu, Hongming Zhang, Xiaoman Pan, Peixin Cao, Kaixin Ma, Jian Li, Hongwei Wang, and Dong Yu. 2024. Chain-of-Note: Enhancing Robustness in Retrieval-Augmented Language Models. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024*, Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (Eds.). Association for Computational Linguistics, 14672–14685. <https://aclanthology.org/2024.emnlp-main.813>
- [69] Hongyi Yuan, Keming Lu, Fei Huang, Zheng Yuan, and Chang Zhou. 2024. Speculative Contrastive Decoding. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics, ACL 2024 - Short Papers, Bangkok, Thailand, August 11-16, 2024*, Lun-Wei Ku, Andre Martins, and Vivek Srikumar (Eds.). Association for Computational Linguistics, 56–64. <https://aclanthology.org/2024.acl-short.5>
- [70] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. 2023. A Survey of Large Language Models. *CoRR* abs/2303.18223 (2023). doi:10.48550/ARXIV.2303.18223 arXiv:2303.18223
- [71] Kun Zhu, Xiaocheng Feng, Xiyuan Du, Yuxuan Gu, Weijiang Yu, Haotian Wang, Qianglong Chen, Zheng Chu, Jingchang Chen, and Bing Qin. 2024. An Information Bottleneck Perspective for Effective Noise Filtering on Retrieval-Augmented Generation. 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*, Lun-Wei Ku, Andre Martins, and Vivek Srikumar (Eds.). Association for Computational Linguistics, 1044–1069. doi:10.18653/V1/2024.ACL-LONG.59
- [72] Yutao Zhu, Zhaoheng Huang, Zhicheng Dou, and Ji-Rong Wen. 2024. One Token Can Help! Learning Scalable and Pluggable Virtual Tokens for Retrieval-Augmented Large Language Models. *CoRR* abs/2405.19670 (2024). doi:10.48550/ARXIV.2405.19670 arXiv:2405.19670
- [73] Yutao Zhu, Huaying Yuan, Shuting Wang, Jiongnan Liu, Wenhan Liu, Chenlong Deng, Zhicheng Dou, and Ji-Rong Wen. 2023. Large Language Models for Information Retrieval: A Survey. *CoRR* abs/2308.07107 (2023). doi:10.48550/ARXIV.2308.07107 arXiv:2308.07107
- [74] Yutao Zhu, Peitian Zhang, Chenghao Zhang, Yifei Chen, Binyu Xie, Zheng Liu, Ji-Rong Wen, and Zhicheng Dou. 2024. INTERS: Unlocking the Power of Large Language Models in Search with Instruction Tuning. 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*, Lun-Wei Ku, Andre
- Martins, and Vivek Srikumar (Eds.). Association for Computational Linguistics, 2782–2809. doi:10.18653/V1/2024.ACL-LONG.154