

A-RAG: Scaling Agentic Retrieval-Augmented Generation via Hierarchical Retrieval Interfaces

Mingxuan Du¹, Benfeng Xu^{2†}, Chiwei Zhu¹, Shaohan Wang¹, Pengyu Wang¹
 Xiaorui Wang², Zhendong Mao^{1‡}

¹University of Science and Technology of China, Hefei, China

²Metastone Technology, Beijing, China
 dumingxuan@mail.ustc.edu.cn

Frontier language models have demonstrated strong reasoning and long-horizon tool-use capabilities. However, existing RAG systems fail to leverage these capabilities. They still rely on two paradigms: (1) designing an algorithm that retrieves passages in a single shot and concatenates them into the model’s input, or (2) predefining a workflow and prompting the model to execute it step-by-step. Neither paradigm allows the model to participate in retrieval decisions, preventing efficient scaling with model improvements. In this paper, we introduce **A-RAG**, an **Agentic RAG** framework that exposes hierarchical retrieval interfaces directly to the model. A-RAG provides three retrieval tools: `keyword_search`, `semantic_search`, and `chunk_read`, enabling the agent to adaptively search and retrieve information across multiple granularities. Experiments on multiple open-domain QA benchmarks show that A-RAG consistently outperforms existing approaches with comparable or lower retrieved tokens, demonstrating that A-RAG effectively leverages model capabilities and dynamically adapts to different RAG tasks. We further systematically study how A-RAG scales with model size and test-time compute. We will release our code and evaluation suite to facilitate future research. Code and evaluation suite are available at <https://github.com/Ayanami0730/arag>.

1 Introduction

The development of LLMs has entered a new phase, where the primary scaling direction is shifting from single-turn text understanding and generation toward complex reasoning and multi-step, tool-augmented interaction (OpenAI, 2025c; Anthropic, 2025b; Google, 2025; DeepSeek-AI et al., 2025; Shao et al., 2025; Yang et al., 2025a; Team et al., 2025; MiniMax AI, 2025). This transformation has significantly enhanced the capabilities and practicality of LLM-based agents, demonstrating remarkable progress in domains such as coding and deep research. By integrating frontier models, coding agents (Anysphere, Inc., 2023; Anthropic, 2025a) have substantially improved the productivity of software engineers, while deep research agents (OpenAI, 2025b; Google LLC, 2024; Tongyi Deep-Research Team, 2025) have greatly accelerated researchers’ ability to conduct surveys and gather information. This marks a paradigm shift. However, methods in the RAG domain have rarely addressed this transition.

Existing RAG methods primarily rely on two paradigms: (1) designing an algorithm (with or without graph structures) that retrieves multiple passages in a single shot and concatenates them into the model’s input (Yan et al., 2024; Sarthi et al., 2024; Gutiérrez et al., 2025a; Edge et al., 2025; Guo et al., 2025; Qian et al., 2025; Huang et al., 2025); (2) predefining a workflow and prompting the model to execute it step-by-step through multiple iterations (Jiang et al., 2023;

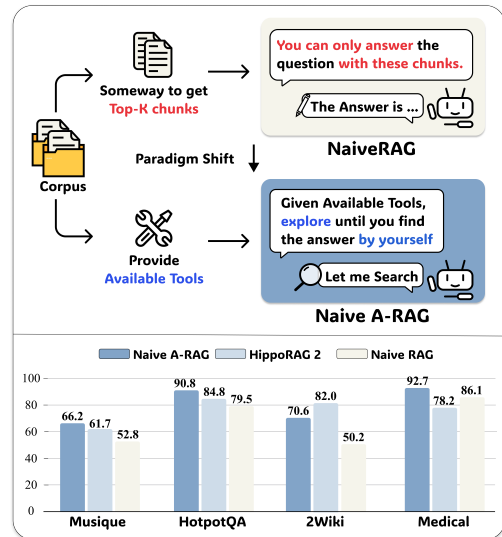


Figure 1: Two paradigms comparison and performance results.

[†]Project lead.

[‡]Corresponding author.

2024). Neither approach is truly agentic, as the model is not allowed to adapt the workflow based on the specific task, choose different interaction strategies, or decide when sufficient evidence has been gathered to provide an answer.

As illustrated in Figure 1, the key distinction between Naive RAG and Naive Agentic RAG lies in the agent’s autonomy, and our preliminary experiments show that even the simplest Naive Agentic RAG, equipped with only a single embedding-based tool to retrieve from the corpus, consistently outperforms Naive RAG and previous baselines. This result demonstrates the potential of the agentic RAG paradigm.

To address these limitations, we propose A-RAG, an Agentic RAG framework featuring hierarchical retrieval interfaces. Our key insight is that information within a corpus is inherently organized at multiple granularities, ranging from fine-grained keyword-level signals to coarser sentence-level and chunk-level representations. Accordingly, we design a suite of retrieval tools that enable the agent to access information across these granularities. We observe that when equipped with this hierarchical toolset, the agent spontaneously generalizes to diverse workflows tailored to various tasks, yielding consistent performance gains.

Comprehensive experiments across multiple benchmarks demonstrate that A-RAG substantially surpasses prior methods. Furthermore, we conduct systematic studies on Test-Time Scaling behavior, showing that A-RAG’s performance improves steadily with increased computational resources, indicating that our framework scales efficiently alongside advances in model capabilities. In summary, our contributions include:

- We identify the paradigm shift from static LLM pipelines to dynamic agent-based systems, and highlight the necessity of transforming RAG into an agentic framework.
- We introduce A-RAG, an agentic RAG framework with hierarchical retrieval interfaces. By conducting comprehensive experiments, we validate that multi-granularity tools are essential for unlocking stronger model performance.
- We present further scaling analyses across multiple dimensions, demonstrating that our framework scales efficiently alongside advances in model capabilities and test-time computation.

2 Related Work

We compare three RAG paradigms in Figure 2: Graph RAG, Workflow RAG, and Agentic RAG (A-RAG). We identify three principles that define true agentic autonomy, and demonstrate that A-RAG is the only paradigm satisfying all three. A detailed comparison across existing methods is provided in Appendix A.

2.1 Basic RAG

Early research demonstrated that retrieval can help models incorporate external knowledge to answer questions more accurately (Lewis et al., 2021). Subsequent work has continuously improved upon this foundation through query rewriting (Chan et al., 2024), adaptive routing strategies (Jeong et al., 2024), retrieval quality evaluation (Yan et al., 2024), and reranking mechanisms.

2.2 Graph RAG

In 2024, Microsoft introduced GraphRAG (Edge et al., 2025), which constructs entity-relation graphs from corpora to help models develop holistic understanding of large-scale knowledge bases. This approach has rapidly evolved into a mainstream RAG paradigm, with researchers advancing the frontier through innovations in knowledge graph structure design, semantic unit definition, and retrieval strategies (Guo et al., 2025; Shen et al., 2025; Yang et al., 2025b; Song et al., 2025b). Among these, RAPTOR (Sarhi et al., 2024) constructs hierarchical tree structures through recursive summarization for multi-level retrieval. LightRAG (Guo et al., 2025) combines knowledge graphs with vector retrieval for both local and global search. HippoRAG (Gutiérrez et al., 2025a,b) mimics hippocampal memory indexing using Personalized PageRank for efficient multi-hop reasoning. While these methods incorporate richer structure, they still rely on predefined retrieval algorithms rather than model-driven decisions. If the initially retrieved context is insufficient, the model cannot leverage its reasoning capabilities to iteratively gather more comprehensive and accurate information.

2.3 Workflow RAG

With the emergence of LLM-based agents, many works have explored agentic approaches to RAG. However, most rely on predefined agent-workflows that prompt models to execute fixed procedures step by step. So we refer to these methods as Workflow RAG. Some further employ SFT and RL to help models follow these workflows more robustly. Among the training-free methods, FLARE (Jiang et al., 2023) triggers retrieval when generation confidence drops, IRCoT (Trivedi et al., 2023) interleaves chain-of-thought reasoning with retrieval steps, and RA-ISF (Liu et al., 2024)

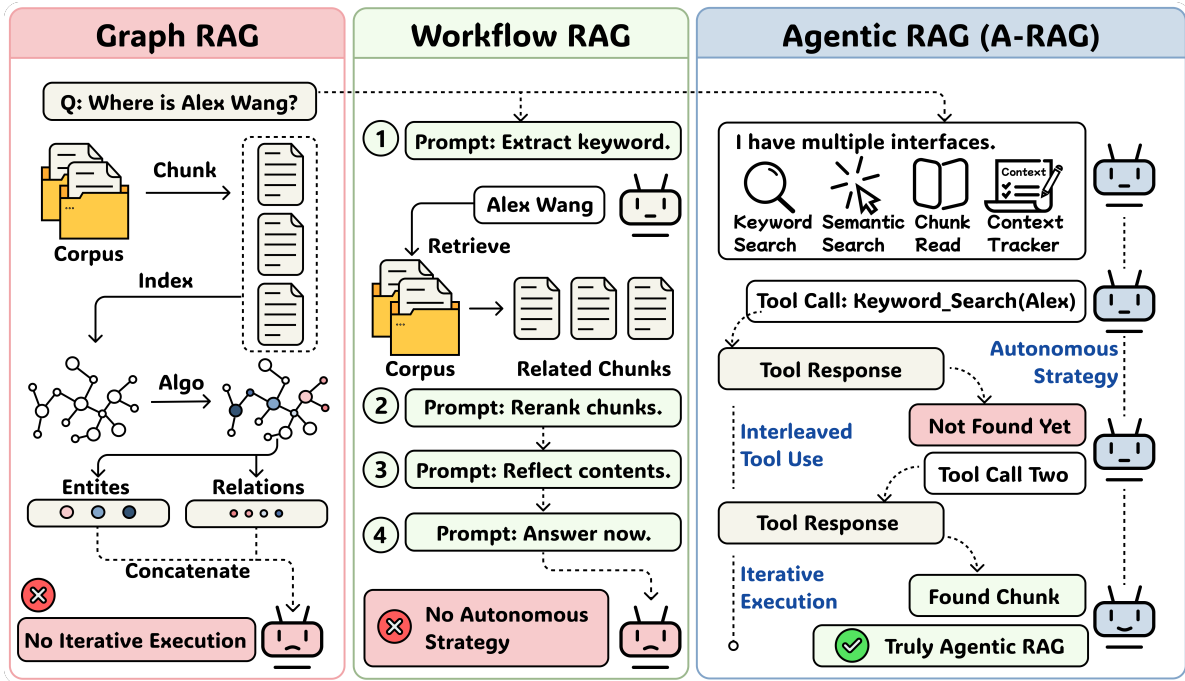


Figure 2: Comparison of three paradigms. We identify three principles of agentic autonomy: Autonomous Strategy, Iterative Execution, and Interleaved Tool Use. Only A-RAG satisfies all three, making it a truly agentic framework.

decomposes complex queries through iterative self-feedback. Multi-agent approaches further extend this paradigm: MA-RAG (Nguyen et al., 2025) coordinates specialized agents via collaborative chain-of-thought, RAGentA (Besrouer et al., 2025; Chang et al., 2024) combines hybrid retrieval with citation tracking for question answering. Training-based methods have demonstrated that even smaller models can learn effective retrieval strategies (Asai et al., 2023; Chan et al., 2024; Chen et al., 2025; Xiong et al., 2025; Jin et al., 2025; Song et al., 2025a; Luo et al., 2025). Despite their sophistication, these workflows remain fixed at design time: the model cannot adapt its strategy based on task characteristics. In contrast, we demonstrate that with agent-friendly hierarchical retrieval interfaces, models can autonomously adopt diverse interaction strategies without any predefined workflows, exhibiting stronger and more robust performance across varying task complexities.

3 Methodology

In this section, we present A-RAG, an agentic-RAG framework that exposes hierarchical retrieval interfaces to models. As illustrated in Figure 3, our approach consists of three key components: (i) a hierarchical index, (ii) a suite of retrieval tools, and (iii) a simple agent loop design to clearly demonstrate the effectiveness of A-RAG.

3.1 Hierarchical Index Construction

To enable efficient multi-granularity retrieval, we construct a hierarchical index that organizes corpus information at different levels of abstraction. Our indexing procedure is lightweight and consists of only two stages: chunking and embedding.

Chunking. Following the setup of LinearRAG (Zhuang et al., 2025), we partition the corpus into chunks of approximately 1,000 tokens each, ensuring that chunk boundaries align with sentence boundaries to preserve semantic coherence. Each chunk serves as a self-contained semantic unit that the agent can selectively access through dedicated retrieval interfaces, rather than being indiscriminately concatenated into the context as in conventional RAG approaches.

Embedding. For each chunk c_i , we decompose it into sentences $\{s_{i,1}, s_{i,2}, \dots, s_{i,n_i}\}$ using rule-based sentence segmentation. We then compute dense vector representations using a pre-trained sentence encoder f_{emb} : $\mathbf{v}_{i,j} = f_{\text{emb}}(s_{i,j})$. This sentence-level embedding enables fine-grained semantic matching while maintaining a mapping from sentences back to their parent chunks, allowing the agent to first identify relevant sentences and then read the complete chunk contexts.

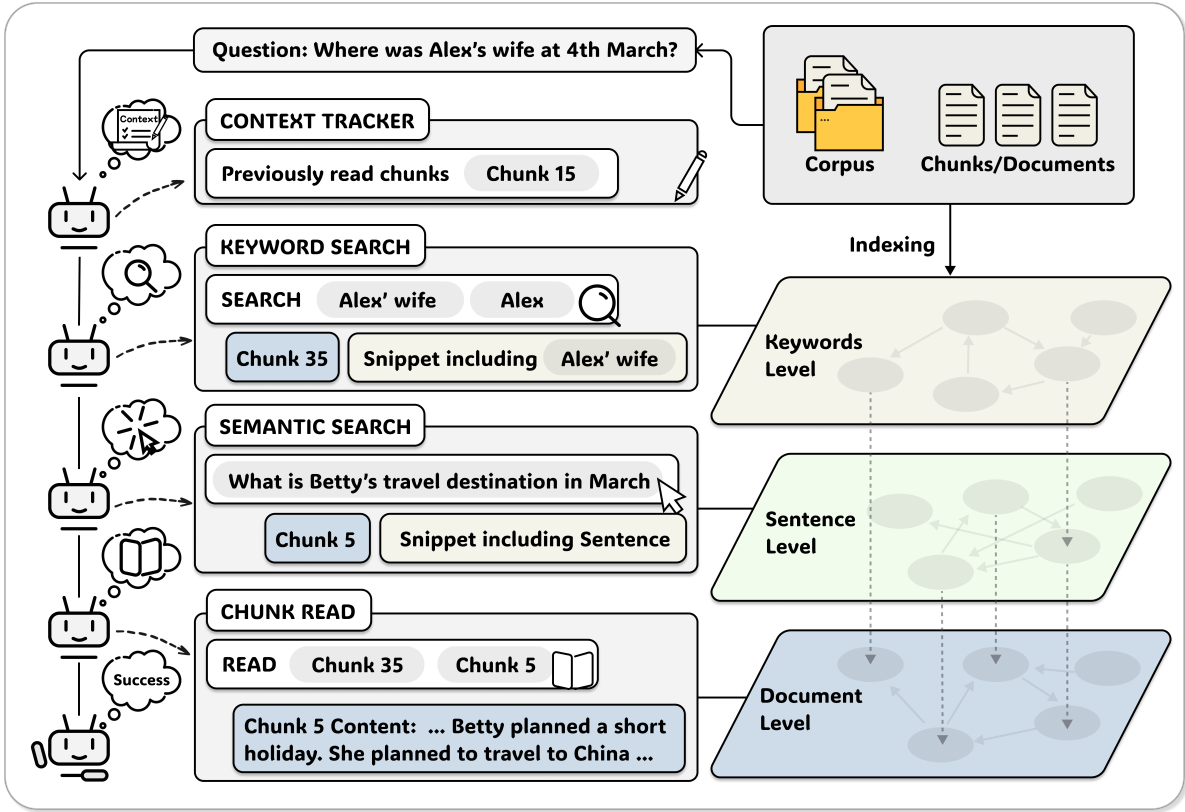


Figure 3: Overview of A-RAG framework. The agent iteratively uses hierarchical retrieval tools (keyword search, semantic search, chunk read) to gather information from the corpus and autonomously decides when to provide the final answer.

Keyword-Level. For the keyword-level information, we avoid pre-indexing. Instead of constructing inverted indices or knowledge graphs during the offline phase, we perform exact text matching directly at query time. This design choice significantly reduces both indexing time and computational cost compared to graph-based approaches. Through this lightweight indexing procedure, we obtain a three-level information representation: an implicit keyword-level for precise entity matching via runtime text search, sentence-level embeddings for semantic search, and chunk-level storage for full content access, which collectively support the hierarchical retrieval interfaces.

3.2 Hierarchical Retrieval Interfaces

We design three retrieval tools that operate at different granularities, enabling the agent to adaptively choose the most suitable search strategy based on the characteristics of each question.

Keyword Search. This tool performs exact lexical matching to locate chunks containing specific terms. The agent provides a keyword list $\mathcal{K} = \{k_1, k_2, \dots, k_m\}$ and a parameter k specifying the number of results to return. The relevance score of chunk c_i is computed as:

$$\text{Score}_{\text{kw}}(c_i, \mathcal{K}) = \sum_{k \in \mathcal{K}} \text{count}(k, T_i) \cdot |k| \quad (1)$$

where $\text{count}(k, T_i)$ denotes the frequency of keyword k in chunk text T_i , and $|k|$ is the character length of the keyword (longer keywords are weighted higher as they are typically more specific). For each matched chunk, we construct an abbreviated snippet by extracting sentences that contain at least one keyword:

$$\text{Snippet}(c_i, \mathcal{K}) = \{s \in \text{Sent}(c_i) \mid \exists k \in \mathcal{K}, k \subseteq s\} \quad (2)$$

where $\text{Sent}(c_i)$ denotes the set of sentences in chunk c_i . The tool returns the top- k chunk IDs along with their snippets, allowing the agent to autonomously decide the next action.

Semantic Search. This tool finds semantically similar passages using dense retrieval. Given a natural language query q , we encode it into a query embedding $\mathbf{v}_q = f_{\text{emb}}(q)$ and compute cosine similarity with all sentence embeddings:

$$\text{Score}_{\text{sem}}(s_{i,j}, q) = \frac{\mathbf{v}_{i,j}^T \mathbf{v}_q}{\|\mathbf{v}_{i,j}\| \|\mathbf{v}_q\|} \quad (3)$$

We retrieve the top-ranked sentences and aggregate them by their parent chunks. Each chunk’s relevance score is determined by its highest-scoring sentence. The tool returns the top- k chunk IDs along with the matched sentences within each chunk as snippets, allowing the agent to autonomously decide the next action.

Chunk Read. Based on the snippets returned by keyword search and semantic search, the agent can determine which chunks require full reading and use this tool to access their complete content. The agent can also read adjacent chunks to gather additional context when needed.

This hierarchical design is inherently agent-friendly, allowing the agent to access corpus information at different granularities based on its own judgment. Rather than loading large amounts of context indiscriminately, the agent can incrementally retrieve information on-demand, minimizing context overhead while maintaining the flexibility to gather comprehensive evidence when needed.

3.3 Agent Loop

Since our method primarily focuses on interface design and investigating test-time scaling behavior in A-RAG, we deliberately adopt the simplest agent loop backbone to minimize confounding factors from complex orchestration mechanisms.

Agent Loop. We adopt the ReAct-like framework (Yao et al., 2023), where the model iteratively performs reasoning and tool calling in an interleaved manner. At each iteration, the agent selects *one* tool to call, observes the result, and decides the next action. We intentionally avoid parallel tool calling and other sophisticated designs to facilitate clean observation of how different interface configurations influence agent behavior. When the maximum iteration budget is reached without producing an answer, we prompt the agent to synthesize a response based on the information gathered so far.

Context Tracker. To prevent redundant information retrieval and unnecessary token consumption, we maintain a context tracker that records which chunks have been read during the retrieval process. Specifically, we track a set $\mathcal{C}^{\text{read}} = \{c_{i_1}, c_{i_2}, \dots, c_{i_k}\}$, where each c_{i_j} denotes the ID of a previously accessed chunk. When the agent attempts to read a chunk $c_i \in \mathcal{C}^{\text{read}}$, instead of returning the full text again, the chunk read tool returns a notification message “This chunk has been read before”, consuming zero additional tokens. This mechanism not only reduces computational cost but also encourages the agent to explore diverse parts of the corpus rather than repeatedly examining the same passages.

This straightforward design allows us to cleanly isolate and analyze the impact of hierarchical interfaces on agent behavior and retrieval performance.

4 Experiments

In this section, we conduct comprehensive experiments to evaluate the effectiveness of A-RAG across multiple benchmarks and analyze its test-time scaling behavior.

4.1 Experimental Setting

Datasets. We evaluate A-RAG on four widely-used multi-hop QA datasets: HotpotQA (Yang et al., 2018), 2WikiMultiHopQA (Ho et al., 2020), MuSiQue (Trivedi et al., 2022), and GraphRAG-Bench (Xiang et al., 2025). Following the experimental setup of LinearRAG (Zhuang et al., 2025), we use the same corpus and questions to ensure fair comparison across different methods.

Baselines. We organize all compared methods into two groups: (i) **Vanilla Baselines:** including direct zero-shot LLM inference and Naive RAG method; (ii) **Graph-RAG and Workflow RAG:** we benchmark against representative graph-enhanced approaches including GraphRAG (Edge et al., 2025), HippoRAG2 (Gutiérrez et al., 2025a), and LinearRAG (Zhuang et al., 2025), as well as workflow-based methods including FaithfulRAG (Zhang et al., 2025a), MA-RAG (Nguyen et al., 2025), and RAGentA (Besrouer et al., 2025; Chang et al., 2024). We compare these baselines against our A-RAG (Naive), equipped with only a single embedding search tool, and A-RAG (Full).

Table 1: Results (%) of baselines and A-RAG on benchmark datasets in terms of LLM-Evaluation Accuracy(LLM-Acc) and Contain-Match Accuracy(Cont-Acc). The best result for each backbone LLM is highlighted in **bold**, while the second result is indicated with an underline.

Method	MuSiQue		HotpotQA		2Wiki		Med.	Novel
	LLM	Cont	LLM	Cont	LLM	Cont	LLM	LLM
<i>GPT-4o-mini</i>								
<i>Vanilla Baselines</i>								
Direct Answer	18.3	13.9	45.4	40.7	30.3	49.7	68.6	45.3
Naive RAG	38.6	36.1	74.5	<u>72.9</u>	42.6	59.0	75.3	68.5
<i>Graph-RAG and Workflow RAG</i>								
GraphRAG	26.4	20.8	33.2	33.3	18.4	47.2	51.3	28.8
HippoRAG2	40.6	38.4	80.7	69.7	64.7	68.5	72.0	<u>70.1</u>
LinearRAG	34.8	26.3	72.0	60.5	<u>62.9</u>	<u>62.3</u>	53.1	45.4
FaithfulRAG	28.8	22.6	60.5	52.5	38.8	38.1	42.5	33.3
MA-RAG	34.1	27.4	60.6	54.4	51.0	53.4	62.3	44.5
RAGentA	32.2	29.9	63.0	62.4	27.7	50.3	67.7	61.3
<i>A-RAG (Ours)</i>								
A-RAG (Naive)	43.8	<u>38.5</u>	76.6	70.7	52.3	62.4	<u>79.0</u>	70.0
A-RAG (Full)	46.1	39.6	<u>77.1</u>	74.0	60.2	63.7	79.4	72.7
<i>GPT-5-mini</i>								
<i>Vanilla Baselines</i>								
Direct Answer	35.8	26.5	63.6	53.5	51.3	54.0	90.5	45.1
Naive RAG	52.8	48.7	81.2	79.5	50.2	66.5	86.1	70.6
<i>Graph-RAG and Workflow RAG</i>								
GraphRAG	48.3	39.1	82.5	74.9	66.5	70.7	87.3	<u>77.1</u>
HippoRAG2	61.7	52.5	84.8	75.0	82.0	79.7	78.2	54.3
LinearRAG	62.4	51.8	86.2	77.6	<u>87.2</u>	<u>84.8</u>	79.2	54.7
FaithfulRAG	52.9	52.8	76.9	75.3	51.8	56.6	75.4	60.7
MA-RAG	40.0	31.6	67.1	57.9	54.7	54.3	68.3	45.1
RAGentA	38.3	37.4	61.2	65.0	24.0	53.5	73.7	60.2
<i>A-RAG (Ours)</i>								
A-RAG (Naive)	<u>66.2</u>	<u>59.7</u>	<u>90.8</u>	<u>85.3</u>	70.6	76.9	<u>92.7</u>	80.4
A-RAG (Full)	74.1	65.3	94.5	88.0	89.7	88.9	93.1	85.3

Evaluation Metrics. Following LinearRAG, we employ two metrics for end-to-end QA assessment: (1) LLM-Evaluation Accuracy (LLM-Acc, corresponding to GPT-Acc in LinearRAG), an LLM-based metric that determines semantic equivalence between predictions and ground-truth answers, and (2) Contain-Match Accuracy (Contain-Acc), which verifies whether the ground-truth answer appears within the generated response. For HotpotQA, 2WikiMultiHopQA, and MuSiQue with short-form answers, we report both metrics. For GraphRAG-Bench with long-form descriptive answers, we report LLM-Acc only, as lengthy ground-truth answers rarely appear verbatim in generated responses, making Contain-Acc uninformative.

Implementation. We evaluate all methods using both GPT-4o-mini and GPT-5-mini (OpenAI, 2025a) as backbone LLMs. For dense retrieval, all methods except LinearRAG utilize Qwen3-Embedding-0.6B (Zhang et al., 2025b) with $k = 5$ for top- k results; LinearRAG uses its original embedding model due to incompatibility between its NER module and Qwen3-Embedding. We intentionally include both earlier and frontier reasoning models to provide a comprehensive view of how RAG methods perform across different capability levels. For LLM-based evaluation, we use GPT-5-mini as the judge, which demonstrates improved accuracy and stability based on our human verification. Detailed configuration and hyperparameters are provided in Appendix B.

Table 2: Ablation study results (%) on benchmark datasets.

Method	MuSiQue		HotpotQA		2Wiki		Med.	Novel
	LLM	Cont	LLM	Cont	LLM	Cont	LLM	LLM
A-RAG (Full)	74.1	65.3	94.5	88.0	89.7	88.9	93.1	85.3
w/o KW Search	72.6	65.3	93.0	87.4	88.9	88.1	93.2	85.0
w/o Semantic	69.4	63.3	93.9	88.4	89.1	88.0	92.1	85.2
w/o Chunk Read	73.6	67.0	93.6	88.8	89.0	87.9	93.3	85.1

4.2 Main Results

Table 1 presents the main experimental results across all benchmarks. We highlight three key observations from our experiments:

Vanilla retrieval method remain robust baseline. Under our unified evaluation setting with GPT-5-mini as the judge and Qwen3-Embedding for dense retrieval, vanilla baselines demonstrate robust performance across both GPT-4o-mini and GPT-5-mini backbones. Existing Graph-RAG and Workflow RAG methods fail to consistently outperform these simple baselines across all datasets.

Naive A-RAG establishes a new strong baseline for agentic RAG. As a simplified variant equipped with only a single embedding-based retrieval tool, A-RAG (Naive) surpasses existing Graph-RAG and Workflow RAG methods on multiple datasets, demonstrating the inherent advantages of the agentic paradigm. This advantage becomes more pronounced when switching to GPT-5-mini as the backbone. This result suggests that granting models greater autonomy in retrieval decisions yields better performance than relying on fixed retrieval algorithms, even without sophisticated multi-granularity tools.

A-RAG outperforms existing RAG methods through hierarchical retrieval interfaces. A-RAG is designed for reasoning models with tool-use capabilities, aligning with the current development trend of the LLM field. With GPT-4o-mini as the backbone, A-RAG (Full) achieves the best performance on 3 out of 5 datasets. When switching to GPT-5-mini with stronger reasoning and tool-calling capabilities, A-RAG (Full) achieves superior results across all benchmarks. The consistent improvements of A-RAG over both baseline methods and Naive A-RAG demonstrate that the A-RAG framework is agent-friendly. It allows models to leverage their reasoning capabilities to dynamically adjust strategies and orchestrate different interfaces based on task requirements, thereby achieving better performance.

4.3 Ablation Study

To investigate the contribution of each retrieval tool, we conduct ablation experiments by systematically removing individual components from A-RAG (Full). We evaluate three ablation variants: (i) w/o Keyword Search and w/o Semantic Search, which directly remove the corresponding retrieval tool from the agent’s toolkit; (ii) w/o Chunk Read, which replaces the snippet-based responses of keyword and semantic search with complete chunk texts and removes the chunk read tool entirely.

As shown in Table 2, the full hierarchical configuration achieves optimal overall performance. A-RAG (Full) consistently achieves the best results on most benchmarks. Removing either semantic search or keyword search leads to performance degradation, highlighting the importance of multi-granularity information for multi-hop retrieval tasks. The inferior performance of w/o Chunk Read compared to A-RAG (Full) demonstrates that our progressive information acquisition design allows the agent to make autonomous judgments and precisely read the most relevant content. This design not only enhances agent autonomy but also enables the model to selectively read only the most relevant chunks in full, avoiding the noise introduced by irrelevant content.

5 Analysis and Discussion

To understand the advantages and characteristics of A-RAG as a new paradigm, we conduct further experiments and analyses in this section.

5.1 Test-Time Scaling Analysis

Since A-RAG grants LLMs greater autonomy in retrieval decisions, increasing computational resources at test time can further scale the framework’s performance. As shown in Figure 4, we conduct experiments on the first 300 tasks

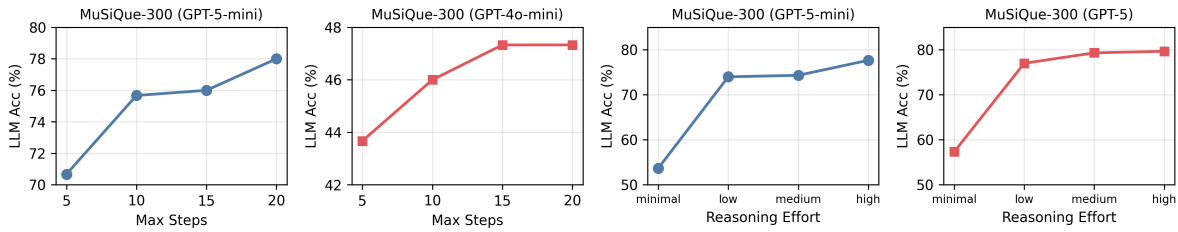


Figure 4: Test-time scaling analysis on MuSiQue-300. Left two: LLM-Acc vs. max steps with GPT-5-mini and GPT-4o-mini. Right two: LLM-Acc vs. reasoning effort with GPT-5-mini and GPT-5.

of MuSiQue and find that both increasing max-step and reasoning effort effectively scale model performance. When scaling from 5 to 20 steps, GPT-5-mini improves by approximately 8% while GPT-4o-mini improves by only about 4%, indicating that stronger reasoning models are better equipped for longer-horizon exploration. When scaling reasoning effort from minimal to high, both GPT-5-mini and GPT-5 achieve substantial improvements of approximately 25%. These results demonstrate that A-RAG effectively leverages test-time compute, positioning it as a promising paradigm for future development.

5.2 Context Efficiency Analysis

Context efficiency is crucial for integrating RAG into complex agentic systems. We analyze the tokens retrieved from the corpus to measure how efficiently each method utilizes context (Table 3).

Table 3: Retrieved tokens across methods (GPT-5-mini backbone). Lower values indicate higher efficiency.

Method	MuSi.	Hotpot.	2Wiki	Med.	Novel
Naive RAG	5,387	5,358	5,506	5,418	4,997
HippoRAG2	5,411	5,380	5,538	5,447	5,019
GraphRAG	9,234	8,744	4,201	9,391	9,318
LinearRAG	5,418	5,353	5,518	5,427	4,998
FaithfulRAG	5,342	5,310	5,419	5,410	4,994
MA-RAG	9,566	8,007	8,857	6,858	6,101
A-RAG (Naive)	56,360	27,455	45,406	23,657	22,391
A-RAG (Full)	5,663	2,737	2,930	7,678	6,087

A-RAG achieves superior accuracy with higher context efficiency. Contrary to the intuition that more retrieved content leads to better performance, A-RAG (Full) retrieves comparable or fewer tokens than traditional RAG methods while achieving superior accuracy.

Hierarchical interfaces are key to context efficiency. Comparing A-RAG (Naive) and A-RAG (Full) reveals a striking pattern: A-RAG (Naive) retrieves more tokens than A-RAG (Full) but achieves lower performance. This validates our hierarchical interface design, the progressive information disclosure grants the model greater autonomy while avoiding irrelevant content.

5.3 Failure Mode Analysis

We manually reviewed the first 100 incorrect cases of A-RAG on MuSiQue and categorized them into 2-level error types (Figure 5). The majority of failures stem from reasoning chain errors. Among these, entity confusion is the most common, with substantial portions also attributed to wrong retrieval strategies and question misunderstanding. Detailed category definitions and analysis on other datasets are provided in Appendix D.

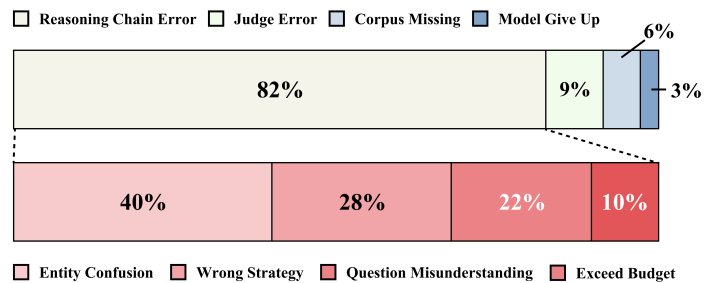


Figure 5: Failure mode distribution of A-RAG. Top: primary categories. Bottom: breakdown of reasoning chain errors.

6 Conclusion

In this work, we recognize agentic RAG as a fundamental paradigm shift in RAG. We introduce A-RAG, an agentic RAG framework featuring hierarchical retrieval interfaces that enable LLMs to autonomously access corpus information at keyword, sentence, and chunk levels. Extensive experiments demonstrate that A-RAG consistently outperforms existing Graph-RAG and Workflow RAG methods across diverse benchmarks, while our analysis validates its efficient test-time scaling behavior. Our findings suggest that future research should focus on designing agent-friendly interfaces rather than complex retrieval algorithms, and explore new interaction paradigms between language models and external knowledge sources.

Limitations

Our work primarily aims to highlight the paradigm shift from traditional RAG to agentic RAG and demonstrate hierarchical interfaces as a promising scaling direction. However, we do not exhaustively enumerate all possible tool designs or systematically compare different tool subsets and their impacts on agent behavior. A comprehensive ablation across diverse tool configurations could provide deeper insights into optimal interface design, which we leave for future work.

Due to computational resource constraints, we have not validated the framework on larger and more powerful models such as GPT-5, and Gemini-3. Given that A-RAG is specifically designed for reasoning models with strong tool-use capabilities, we anticipate that performance gains would be more pronounced with these frontier models, but empirical verification remains to be conducted.

Additionally, while we demonstrate strong results on multi-hop QA benchmarks, the generalization of A-RAG to other knowledge-intensive tasks such as fact verification, dialogue systems, and long-form generation warrants further investigation.

Ethical Considerations

All datasets used in this work are publicly available benchmarks that have been previously curated and processed by prior research with appropriate ethical considerations. Our work focuses on fundamental research for improving retrieval-augmented generation in large language models, and does not involve the collection of new data or human subjects. As a methodological contribution to RAG systems, our approach does not introduce additional ethical risks beyond those inherent to the underlying language models.

References

- Anthropic. 2025a. Claude code: Agentic coding assistant. <https://code.claude.com/docs/en/overview>. Accessed: 2025-12-04.
- Anthropic. 2025b. Introducing claude sonnet 4.5.
- Ansphere, Inc. 2023. Cursor: Ai code editor. <https://www.cursor.com/>. Accessed: 2025-12-04.
- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2023. Self-rag: Learning to retrieve, generate, and critique through self-reflection. *Preprint*, arXiv:2310.11511.
- Ines Besrou, Jingbo He, Tobias Schreieder, and Michael Färber. 2025. Ragenta: Multi-agent retrieval-augmented generation for attributed question answering. *Preprint*, arXiv:2506.16988.
- Chi-Min Chan, Chunpu Xu, Ruibin Yuan, Hongyin Luo, Wei Xue, Yike Guo, and Jie Fu. 2024. Rq-rag: Learning to refine queries for retrieval augmented generation. *Preprint*, arXiv:2404.00610.
- Chia-Yuan Chang, Zhimeng Jiang, Vineeth Rakesh, Menghai Pan, Chin-Chia Michael Yeh, Guanchu Wang, Mingzhi Hu, Zhichao Xu, Yan Zheng, Mahashweta Das, and Na Zou. 2024. Main-rag: Multi-agent filtering retrieval-augmented generation. *Preprint*, arXiv:2501.00332.
- Yiqun Chen, Lingyong Yan, Weiwei Sun, Xinyu Ma, Yi Zhang, Shuaiqiang Wang, Dawei Yin, Yiming Yang, and Jiaxin Mao. 2025. Improving retrieval-augmented generation through multi-agent reinforcement learning. *Preprint*, arXiv:2501.15228.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, and 1 others. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *Preprint*, arXiv:2501.12948.
- Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt, Dasha Metropolitan-sky, Robert Osazuwa Ness, and Jonathan Larson. 2025. From local to global: A graph rag approach to query-focused summarization. *Preprint*, arXiv:2404.16130.
- Google. 2025. A new era of intelligence with gemini 3.
- Google LLC. 2024. Gemini deep research: Your personal research assistant. <https://gemini.google.com/overview/deep-research/>. Accessed: 2025-12-04.
- Zirui Guo, Lianghao Xia, Yanhua Yu, Tu Ao, and Chao Huang. 2025. Lightrag: Simple and fast retrieval-augmented generation. *Preprint*, arXiv:2410.05779.
- Bernal Jiménez Gutiérrez, Yiheng Shu, Yu Gu, Michihiro Yasunaga, and Yu Su. 2025a. Hipporag: Neurobiologically inspired long-term memory for large language models. *Preprint*, arXiv:2405.14831.
- Bernal Jiménez Gutiérrez, Yiheng Shu, Weijian Qi, Sizhe Zhou, and Yu Su. 2025b. From rag to memory: Non-parametric continual learning for large language models. *Preprint*, arXiv:2502.14802.
- Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. 2020. Constructing a multi-hop qa dataset for comprehensive evaluation of reasoning steps. *Preprint*, arXiv:2011.01060.
- Haoyu Huang, Yongfeng Huang, Junjie Yang, Zhenyu Pan, Yongqiang Chen, Kaili Ma, Hongzhi Chen, and James Cheng. 2025. Retrieval-augmented generation with hierarchical knowledge. *Preprint*, arXiv:2503.10150.
- Soyeong Jeong, Jinheon Baek, Sukmin Cho, Sung Ju Hwang, and Jong C. Park. 2024. Adaptive-rag: Learning to adapt retrieval-augmented large language models through question complexity. *Preprint*, arXiv:2403.14403.
- 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. *Preprint*, arXiv:2305.06983.
- Bowen Jin, Hansi Zeng, Zhenrui Yue, Jinsung Yoon, Serkan Arik, Dong Wang, Hamed Zamani, and Jiawei Han. 2025. Search-r1: Training llms to reason and leverage search engines with reinforcement learning. *Preprint*, arXiv:2503.09516.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2021. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Preprint*, arXiv:2005.11401.

- Yanming Liu, Xinyue Peng, Xuhong Zhang, Weihao Liu, Jianwei Yin, Jiannan Cao, and Tianyu Du. 2024. Ra-isf: Learning to answer and understand from retrieval augmentation via iterative self-feedback. *Preprint*, arXiv:2403.06840.
- Haoran Luo, Haihong E, Guanting Chen, Qika Lin, Yikai Guo, Fangzhi Xu, Zemin Kuang, Meina Song, Xiaobao Wu, Yifan Zhu, and Luu Anh Tuan. 2025. Graph-r1: Towards agentic graphrag framework via end-to-end reinforcement learning. *Preprint*, arXiv:2507.21892.
- MiniMax AI. 2025. Minimax-m2: A model built for max coding & agentic workflows.
- Thang Nguyen, Peter Chin, and Yu-Wing Tai. 2025. Ma-rag: Multi-agent retrieval-augmented generation via collaborative chain-of-thought reasoning. *Preprint*, arXiv:2505.20096.
- OpenAI. 2025a. Gpt-5 system card. System card, OpenAI. Accessed: 2025-12-12.
- OpenAI. 2025b. Introducing deep research. <https://openai.com/index/introducing-deep-research/>. Accessed: 2025-12-04.
- OpenAI. 2025c. Introducing gpt-5.
- Hongjin Qian, Zheng Liu, Peitian Zhang, Kelong Mao, Defu Lian, Zhicheng Dou, and Tiejun Huang. 2025. Memorag: Boosting long context processing with global memory-enhanced retrieval augmentation. *Preprint*, arXiv:2409.05591.
- Parth Sarthi, Salman Abdullah, Aditi Tuli, Shubh Khanna, Anna Goldie, and Christopher D. Manning. 2024. Raptor: Recursive abstractive processing for tree-organized retrieval. *Preprint*, arXiv:2401.18059.
- Zhihong Shao, Yuxiang Luo, Chengda Lu, Z. Z. Ren, Jiewen Hu, Tian Ye, Zhibin Gou, Shirong Ma, and Xiaokang Zhang. 2025. Deepseekmath-v2: Towards self-verifiable mathematical reasoning. *Preprint*, arXiv:2511.22570.
- Zhili Shen, Chenxin Diao, Pavlos Vougiouklis, Pascual Merita, Shriram Piramanayagam, Enting Chen, Damien Graux, Andre Melo, Ruofei Lai, Zeren Jiang, Zhongyang Li, YE QI, Yang Ren, Dandan Tu, and Jeff Z. Pan. 2025. Gear: Graph-enhanced agent for retrieval-augmented generation. *Preprint*, arXiv:2412.18431.
- Huatong Song, Jinhao Jiang, Yingqian Min, Jie Chen, Zhipeng Chen, Wayne Xin Zhao, Lei Fang, and Ji-Rong Wen. 2025a. R1-searcher: Incentivizing the search capability in llms via reinforcement learning. *Preprint*, arXiv:2503.05592.
- Jinyeop Song, Song Wang, Julian Shun, and Yada Zhu. 2025b. Efficient and transferable agentic knowledge graph rag via reinforcement learning. *Preprint*, arXiv:2509.26383.
- Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, and 1 others. 2025. Kimi k2: Open agentic intelligence. *Preprint*, arXiv:2507.20534.
- Tongyi DeepResearch Team. 2025. Tongyi deepresearch technical report. *Preprint*, arXiv:2510.24701.
- Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022. Musique: Multihop questions via single-hop question composition. *Preprint*, arXiv:2108.00573.
- Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2023. Interleaving retrieval with chain-of-thought reasoning for knowledge-intensive multi-step questions. *Preprint*, arXiv:2212.10509.
- Zhishang Xiang, Chuanjie Wu, Qinggang Zhang, Shengyuan Chen, Zijin Hong, Xiao Huang, and Jinsong Su. 2025. When to use graphs in rag: A comprehensive analysis for graph retrieval-augmented generation. *Preprint*, arXiv:2506.05690.
- Guangzhi Xiong, Qiao Jin, Xiao Wang, Yin Fang, Haolin Liu, Yifan Yang, Fangyuan Chen, Zhixing Song, Dengyu Wang, Minjia Zhang, Zhiyong Lu, and Aidong Zhang. 2025. Rag-gym: Systematic optimization of language agents for retrieval-augmented generation. *Preprint*, arXiv:2502.13957.
- Shi-Qi Yan, Jia-Chen Gu, Yun Zhu, and Zhen-Hua Ling. 2024. Corrective retrieval augmented generation. *Preprint*, arXiv:2401.15884.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, and 1 others. 2025a. Qwen3 technical report. *Preprint*, arXiv:2505.09388.
- Cehao Yang, Xiaojun Wu, Xueyuan Lin, Chengjin Xu, Xuhui Jiang, Yuanliang Sun, Jia Li, Hui Xiong, and Jian Guo. 2025b. Graphsearch: An agentic deep searching workflow for graph retrieval-augmented generation. *Preprint*, arXiv:2509.22009.

- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W. Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. *Preprint*, arXiv:1809.09600.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2023. React: Synergizing reasoning and acting in language models. *Preprint*, arXiv:2210.03629.
- Qinggang Zhang, Zhishang Xiang, Yilin Xiao, Le Wang, Junhui Li, Xinrun Wang, and Jinsong Su. 2025a. Faithfulrag: Fact-level conflict modeling for context-faithful retrieval-augmented generation. *Preprint*, arXiv:2506.08938.
- Yanzhao Zhang, Mingxin Li, Dingkun Long, Xin Zhang, Huan Lin, Baosong Yang, Pengjun Xie, An Yang, Dayiheng Liu, Junyang Lin, Fei Huang, and Jingren Zhou. 2025b. Qwen3 embedding: Advancing text embedding and reranking through foundation models. *Preprint*, arXiv:2506.05176.
- Luyao Zhuang, Shengyuan Chen, Yilin Xiao, Huachi Zhou, Yujing Zhang, Hao Chen, Qinggang Zhang, and Xiao Huang. 2025. Linearrag: Linear graph retrieval augmented generation on large-scale corpora. *arXiv preprint arXiv:2510.10114*.

A Comparison of RAG Method Autonomy

We identify three key principles to determine whether a RAG method is truly agentic: (1) **Autonomous Strategy**: Whether the method allows the LLM to dynamically choose and organize high-level strategies (e.g., whether/when/how to retrieve, decompose, verify, re-plan) without being constrained to a single pre-specified workflow or being primarily decided by external rules/classifiers/evaluators. (2) **Iterative Execution**: Whether the method supports multi-round execution and can adapt the number of rounds based on intermediate results, rather than being strictly one-shot. (3) **Interleaved Tool Use**: Whether the method follows a ReAct-like action→observation→reasoning loop, where each tool call is conditioned on observations from previous tool outputs instead of a fixed toolchain that is always executed in the same order.

Table 4 compares existing RAG methods across these three dimensions. As shown, while existing methods may partially satisfy one or two principles, A-RAG is the only method that fully satisfies all three, making it a truly agentic RAG framework.

Table 4: Comparison of agentic characteristics across RAG methods. ✓ indicates the method clearly satisfies the criterion; ✗ indicates it does not; △ indicates a boundary case.

Method	Auto.	Iter.	Interl.
Naive RAG	✗	✗	✗
Self-RAG	✗	✓	✓
CRAG	✗	✗	✗
Adaptive-RAG	✗	△	△
FLARE	✗	△	△
IRCoT	✗	✓	△
RQ-RAG	✗	✗	✗
RA-ISF	✗	△	△
RAPTOR	✗	✗	✗
GraphRAG	✗	✗	✗
LightRAG	✗	✗	✗
MemoRAG	✗	✗	✗
HippoRAG2	✗	✗	✗
LinearRAG	✗	✗	✗
FaithfulRAG	✗	✗	✗
MA-RAG	△	✓	△
RAGentA	△	✓	△
A-RAG (Ours)	✓	✓	✓

B Baseline Reproduction Details

All baseline results are reproduced locally under a unified evaluation setting. All methods use top- $k=5$ for retrieval and max_tokens ≥ 16384 to prevent reasoning truncation. We briefly describe each baseline method below:

- **GraphRAG** (Edge et al., 2025): Constructs knowledge graphs from documents with hierarchical community structure, enabling both local entity-based and global community-based retrieval for query-focused summarization.
- **HippoRAG2** (Gutiérrez et al., 2025a): Mimics human hippocampal memory indexing using knowledge graphs and Personalized PageRank, enabling single-step multi-hop reasoning with improved efficiency.
- **LinearRAG** (Zhuang et al., 2025): Simplifies graph construction by replacing relation extraction with entity extraction, creating hierarchical graphs with two-stage retrieval.
- **FaithfulRAG** (Zhang et al., 2025a): Resolves knowledge conflicts between retrieved content and model’s parametric knowledge through self-fact mining, conflict identification, and reasoning integration.
- **MA-RAG** (Nguyen et al., 2025): Multi-agent framework with specialized agents (Planner, Step Definer, Extractor, QA) collaborating through chain-of-thought reasoning.
- **RAGentA** (Besrou et al., 2025): Multi-agent system with hybrid sparse-dense retrieval, iterative document filtering, and citation-attributed answer generation.

Table 5 summarizes the key reproduction configurations for each method.

Table 5: Baseline reproduction configurations.

Method	Configuration
GraphRAG	Local Search; chunk_size=1200; top_k_entities=10; Qwen3-Embedding-0.6B.
HippoRAG2	graph_type=facts_and_sim_passage_node; retrieval_top_k=200; qa_top_k=5.
LinearRAG	Official configs; all-mpnet-base-v2.
FaithfulRAG	3-step Fact Mining; chunk_topk=5; Qwen3-Embedding-0.6B.
MA-RAG	Plan Agent + Plan Executor; numpy cosine similarity.
RAGentA	4-Agent; alpha=0.65; Qwen3-Embedding + FAISS; BM25.

C Agent Loop Algorithm

Algorithm 1 A-RAG Agent Loop

Input: question q , tools \mathcal{T} , LLM \mathcal{M} , max iterations L

- 1: $\mathcal{M}_{\text{msg}} \leftarrow [\{q\}]$, $\mathcal{C}^{\text{read}} \leftarrow \emptyset$
- 2: **for** $\ell = 1$ to L **do**
- 3: response $\leftarrow \mathcal{M}(\mathcal{M}_{\text{msg}}, \mathcal{T})$
- 4: **if** response contains tool call (t, args) **then**
- 5: $\mathcal{M}_{\text{msg}}.\text{append}(\text{response})$
- 6: result $\leftarrow t.\text{execute}(\mathcal{C}, \text{args})$
- 7: $\mathcal{M}_{\text{msg}}.\text{append}(\text{result})$
- 8: **if** $t = \text{chunk_read}$ **then**
- 9: $\mathcal{C}^{\text{read}} \leftarrow \mathcal{C}^{\text{read}} \cup \text{args.chunk_ids}$
- 10: **end if**
- 11: **else**
- 12: **return** response
- 13: **end if**
- 14: **end for**
- 15: $\mathcal{M}_{\text{msg}}.\text{append}([\text{“Answer the question”}])$
- 16: **return** $\mathcal{M}(\mathcal{M}_{\text{msg}})$

Algorithm 1 presents the pseudocode for the A-RAG agent loop. The agent maintains a message history \mathcal{M}_{msg} and a set of read chunks $\mathcal{C}^{\text{read}}$ to track context. At each iteration, the LLM \mathcal{M} receives the message history and available tools \mathcal{T} , then decides whether to call a tool or return a final answer. If a tool is called, the result is appended to the message history. The loop continues until the agent produces an answer or reaches the maximum iteration limit L , at which point it is prompted to synthesize a response.

D Failure Mode Details

To understand how failure modes shift with the paradigm change from Naive RAG to Agentic RAG, we manually analyzed the first 100 incorrect cases from two settings: (1) GPT-4o-mini with Naive RAG on HotpotQA and MuSiQue, and (2) GPT-5-mini with A-RAG on MuSiQue and 2WikiMultiHopQA. This analysis aims to identify optimization opportunities for future research.

D.1 Naive RAG Failure Categories

For Naive RAG with GPT-4o-mini, we define the following failure modes:

- **Model Understanding:** The gold answer exists in retrieved documents, but the model fails to correctly understand or extract it.
- **Multi-hop Retrieval:** Gold exists in the corpus but single-pass retrieval fails to find it.
- **Judge Error:** The model provides a correct answer, but is misjudged as incorrect.

- **Top-K Insufficient:** Gold is not in the corpus, or k=5 cannot cover the complete answer chain.

Table 6: Naive RAG (GPT-4o-mini) failure mode distribution.

Failure Mode	HotpotQA	MuSiQue
Model Understanding	36%	35%
Multi-hop Retrieval	27%	36%
Judge Error	23%	14%
Top-K Insufficient	13%	15%
Gold Answer Error	1%	0%

D.2 A-RAG Failure Categories

For A-RAG with GPT-5-mini, we define a two-level taxonomy:

Primary Categories:

- **Reasoning Chain Error:** The model performs multiple retrieval rounds but makes errors in the reasoning chain, leading to incorrect final answers.
- **Judge Error:** The model provides a correct answer but is misjudged.
- **Model Gave Up:** The model exhausts retrieval rounds and claims “information not found”.
- **Corpus Missing:** The gold answer does not exist in the corpus.

Secondary Categories (within Reasoning Chain Error):

- **Entity Confusion:** The model reads chunks containing gold but is distracted by other information.
- **Wrong Strategy:** Incorrect search query construction.
- **Question Misunderstanding:** Complex question structure causes fundamental misunderstanding.
- **Exceed Budget:** Exhausts all retrieval rounds without finding the answer.

Table 7: A-RAG (GPT-5-mini) primary failure mode distribution.

Failure Mode	MuSiQue	2Wiki
Reasoning Chain Error	82%	45%
Model Gave Up	3%	33%
Judge Error	9%	19%
Corpus Missing	6%	3%

Table 8: A-RAG (GPT-5-mini) secondary failure mode distribution within reasoning chain errors.

Failure Mode	MuSiQue	2Wiki
Entity Confusion	40%	71%
Wrong Strategy	28%	29%
Question Misunderstanding	22%	0%
Exceed Budget	10%	0%

D.3 Analysis

Paradigm shift changes the bottleneck. For Naive RAG, approximately 50% of failures stem from retrieval limitations (multi-hop retrieval + top-k insufficient), indicating the core problem is “cannot find documents”. In contrast, A-RAG’s dominant failure mode (82% on MuSiQue) is reasoning chain errors, shifting the bottleneck to “found documents but reasoned incorrectly”.

Entity confusion is the primary challenge. Across both datasets, entity confusion is the largest secondary failure mode (40% on MuSiQue, 71% on 2Wiki), suggesting that improving the model’s ability to disambiguate and extract correct entities from retrieved context is a key optimization direction.

Dataset characteristics affect failure patterns. MuSiQue shows 22% question misunderstanding errors due to its complex multi-hop question structures, while 2Wiki shows 33% model gave up cases, indicating different optimization priorities for different task types

Naive RAG (Direct Answer)
You are a helpful assistant. Answer the question based on the provided context. [Context] {retrieved_chunks} [Question] {question} Please provide a direct answer based on the context above.
System Prompt: Naive A-RAG
You are a question-answering assistant with access to a document corpus through available tools. Your goal is to answer questions accurately by finding and analyzing relevant information from the provided documents. [Available Tools] - naive_embedding_search: Return top-k chunks by embedding similarity - chunk_read: Read the full content of a specific chunk [Strategy] Work iteratively: retrieve -> read -> answer. [When Answering] - Ground your response in the retrieved documents - Cite the specific chunks that support your answer - Provide clear, direct answers supported by evidence - Avoid speculation beyond what the documents support
System Prompt: A-RAG (Full)
You are a question-answering assistant with access to a document corpus through available tools. Your goal is to answer questions accurately by finding and analyzing relevant information from the provided documents. [Available Tools] - keyword_search: Find chunks by exact keyword matching - semantic_search: Find chunks by semantic similarity - chunk_read: Read the full content of a specific chunk [Strategy] Work iteratively: search -> read -> evaluate -> search -> read -> ... -> answer. For multi-hop questions, decompose the problem and tackle each sub-question step by step. [When Answering] - Ground your response in the retrieved documents - Cite the specific chunks that support your answer - Provide clear, direct answers supported by evidence - Avoid speculation beyond what the documents support

Figure 6: **System prompts for different RAG configurations.** Top: Naive RAG uses direct answer without tool calling. Middle: Naive A-RAG with single embedding tool. Bottom: A-RAG (Full) with hierarchical retrieval tools.

E Prompt Templates and Tool Descriptions

We deliberately use minimal system prompts to demonstrate the simplicity and effectiveness of the agentic RAG paradigm. As shown in Figure 6, all configurations share the same basic instruction structure, differing only in available tools and strategy descriptions. The complete tool descriptions provided to the agent are shown in Figure 7 and Figure 8.

<p>Tool: naive_embedding_search</p> <p>Return top-k chunks by embedding similarity (no keyword/BM25). Parameters: - query (string): User question to search relevant chunks - top_k (integer): Number of chunks to return (default: 5)</p>
<p>Tool: keyword_search</p> <p>Search for document chunks using keyword-based exact text matching (case-insensitive). Returns chunk IDs and abbreviated sentence snippets where the keywords appear. IMPORTANT: This tool matches keywords literally in the text. Use SHORT, SPECIFIC terms (1-3 words maximum). Each keyword is matched independently. Examples of GOOD keywords: - Entity names: "Albert Einstein", "Tesla", "Python", "Argentina" - Technical terms: "photosynthesis", "quantum mechanics" - Key concepts: "climate change", "GDP growth" Examples of BAD keywords (DO NOT use): - Long phrases: "the person who invented the telephone" -> use "Alexander Bell" instead - Questions: "when did World War 2 start" -> use "World War 2", "1939" instead - Descriptions: "the country between France and Spain" -> use "Andorra" instead - Full sentences: "how does the stock market work" -> use "stock market", "trading" instead RETURNS: Abbreviated snippets marked with "..." showing where keywords appear. These snippets help you identify relevant chunks, but you MUST use chunk_read to get the full text for answering questions. Parameters: - keywords (array[string]): List of keywords to search. Each keyword should be 1-3 words maximum (e.g., ['Einstein', 'relativity theory', '1905']). - top_k (integer): Number of top-ranked chunks to return (default: 5, max: 20)</p>

Figure 7: **Tool descriptions (Part 1)**. Top: naive_embedding_search for Naive A-RAG. Bottom: keyword_search for exact text matching.

Tool: semantic_search

Semantic search using embedding similarity. Matches your query against sentences in each chunk via vector similarity.

WHEN TO USE: - When keyword search fails to find relevant information - When exact wording in documents is unknown - For conceptual/meaning-based matching

RETURNS: Abbreviated snippets with matched sentences. Use chunk_read to get full text for answering.

Parameters: - query (string): Natural language query describing what information you're looking for - top_k (integer): Number of most relevant results to return (default: 5, max: 20)

Tool: chunk_read

Read the complete content of document chunks by their IDs.

This tool returns the full text of the specified chunks, allowing you to examine the complete context and details that are not visible in search snippets.

IMPORTANT: Search results (keyword_search and semantic_search) only show abbreviated snippets marked with "..." - they are NOT sufficient for answering questions. You MUST use chunk_read to get the full content before formulating your answer.

STRATEGY: - Always read promising chunks identified by your searches - Make sure to read the most relevant chunks to gather complete information - If information seems incomplete or truncated, read adjacent chunks (+/- 1) - Reading full text is essential for accurate answers

Note: Previously read chunks will be marked as already seen to avoid redundant information.

Parameters: - chunk_ids (array[string]): List of chunk IDs to retrieve (e.g., ['0', '24', '172'])

Figure 8: **Tool descriptions (Part 2)**. Top: semantic_search for meaning-based retrieval. Bottom: chunk_read for accessing full document content.