

Fishing for Answers: Exploring One-shot vs. Iterative Retrieval Strategies for Retrieval Augmented Generation

Huifeng Lin, Gang Su, Rui Zhao
SenseTime Research
zhaorui@sensetime.com

You Wu
HKUST, Hong Kong
wuyouscut@gmail.com

Jintao Liang
BUPT, Beijing
ljt2016@bupt.edu.cn

Ziyue Li
Technical University of Munich, Germany
ziyue.li@tum.de

Abstract

Retrieval-Augmented Generation (RAG) based on Large Language Models (LLMs) is a powerful solution to understand and query the industry’s closed-source documents. However, basic RAG often struggles with complex QA tasks in legal and regulatory domains, particularly when dealing with numerous government documents. The top- k strategy frequently misses golden chunks, leading to incomplete or inaccurate answers. To address these retrieval bottlenecks, we explore two strategies to improve evidence coverage and answer quality. The first is a One-SHOT retrieval method that adaptively selects chunks based on a token budget, allowing as much relevant content as possible to be included within the model’s context window. Additionally, we design modules to further filter and refine the chunks. The second is an iterative retrieval strategy built on a Reasoning Agentic RAG framework, where a reasoning LLM dynamically issues search queries, evaluates retrieved results, and progressively refines the context over multiple turns. We identify query drift and retrieval laziness issues and further design two modules to tackle them. Through extensive experiments on a dataset of government documents, we aim to offer practical insights and guidance for real-world applications in legal and regulatory domains. The code and dataset will be released upon acceptance.

1 Introduction

Retrieval-Augmented Generation (RAG) (Chen et al., 2024; Lewis et al., 2020; Gao et al., 2023) based on Large Language Models (LLMs) has been intensively deployed to conduct question-answering (QA) towards closed-source or internal documents. Although Large Language Models (LLMs) (Singh, 2023; Zhao et al., 2023; Zhu et al., 2024) have demonstrated remarkable capabilities in natural language understanding and generation, LLMs rely on static training data, making them prone to hallucinations and limiting their

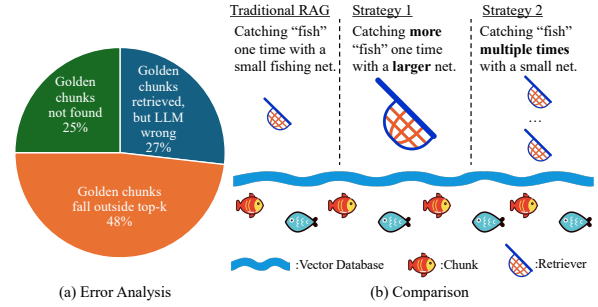


Figure 1: (a) Error analysis: 48% of traditional RAG failures is due to not finding the golden chunk in the top- k ; (b) Comparison of three retrieval methods: Traditional RAG is casting a small fishing net once; We propose either to use a “larger” net or a same net but “multi-times”.

ability to provide accurate, up-to-date information in knowledge-intensive tasks (Rawte et al., 2023; Zhang et al., 2023; Huang et al., 2025). By integrating relevant information from external knowledge bases or search engines, RAG enhances factual accuracy and broadens the model’s domain coverage (Zhao et al., 2024; Li et al., 2024).

Traditional RAG methods typically follow a one-round retrieval paradigm, where the user’s query is first embedded into a vector representation and matched against a vector-based index of document chunks stored in a knowledge base or internal repository. The retriever then selects the top-ranked chunks based on vector similarity, and these retrieved chunks, together with the original query, are provided to a generative LLM, which synthesizes the final answer using both the retrieved evidence and the query context.

Despite the effectiveness of basic RAG methods, they often struggle when applied to real-world scenarios involving complex and heterogeneous government documents. In our analysis of QA tasks on such documents (Fig. 1(a)), we find that nearly half of the incorrect answers result from the failure to retrieve the golden chunks using the traditional top- k selection strategy. This retrieval failure pre-

vents the generative model from accessing crucial evidence, ultimately degrading the quality of the final answer.

Since the one-round retrieval may miss the golden chunks, then the intuitive practical solution is: (1) either we increase the number of chunks retrieved, or (2) keep the same number of chunks and just retrieve multiple times. Take the analogy of fishing (Fig. 1(b)), the traditional RAG’s top- k retrieval can be seen as casting a small fishing net once (e.g., top-5 usually), hoping to capture all relevant evidence in a single pass. However, this “single-shot trawling” approach may miss the golden chunks altogether, especially in complex or diverse knowledge domains. The first solution is basically to dase a larger net, but it may catch too much irrelevant information (“noise” or “junk fish”); The second is to adopt a “multi-cast” fishing strategy, where smaller, more targeted nets are cast repeatedly, progressively refining the retrieval toward high-value chunks.

Motivated by the two retrieval ideas introduced above: casting a larger net and casting multiple times, we implement two corresponding strategies. The first is a One-SHOT retrieval strategy that removes the fixed top- k constraint and instead selects as many relevant chunks as fit within a predefined token budget, ranking them by relevance-per-token to maximize evidence density. To further enhance precision, we incorporate rule-based modules, such as *chunk filter* to filter out irrelevant chunks based on chunk meta information (e.g., year, location in our case). The second is an iterative retrieval strategy based on an Agentic RAG framework, where a reasoning-capable LLM manages retrieval over multiple turns by issuing intermediate queries, assessing retrieved results, and progressively refining the context. These strategies offer practical solutions to improve evidence coverage and robustness in complex QA scenarios.

The complex QA usually lies in two abilities of RAG system, which we refer to as the reasoning ability: (1) the ability to decompose the complex question (generation reasoning), and (2) the ability to retrieve multiple cues (retrieval reasoning). To systematically evaluate the performance of traditional RAG systems and our two proposed strategies on government documents, we design a benchmark that categorizes questions into four levels based on whether generation reasoning and retrieval reasoning are required.

This paper offers valuable insights (underlined)

through extensive experiments on real-world government documents, which we believe are highly beneficial for applications in legal and regulatory domains. The resulting RAG system has been continuously deployed in practice and serving various clients for over two years.

2 Related Works

2.1 Basic RAG

Retrieval-Augmented Generation (RAG) was introduced to overcome the static knowledge limitations of LLMs by integrating external retrieval mechanisms during inference (Chen et al., 2024; Gao et al., 2023). Naive RAG methods represent the earliest implementations, typically using sparse retrieval techniques like BM25 (Robertson et al., 2009) to fetch documents based on keyword overlap (Ma et al., 2023). While efficient for simple factoid queries, these approaches offered limited semantic understanding, thus often retrieving noisy or redundant content and failing to reason across multiple sources.

The emergence of Advanced RAG and Modular RAG was aimed at addressing key limitations of the Naive RAG, particularly in terms of retrieval precision, information integration, and system flexibility (Gao et al., 2023). Advanced RAG improves retrieval quality through techniques such as dense semantic matching, re-ranking, and multi-hop querying, while also introducing refined indexing strategies like fine-grained chunking and metadata-aware retrieval. Modular RAG rethinks the Naive RAG by breaking down the end-to-end process of indexing, retrieval, and generation into discrete, configurable modules. This design allows for greater architectural flexibility and enables system developers to incorporate diverse techniques into specific stages, such as enhancing retrieval with fine-tuned search modules (Lin et al., 2023). In response to specific task demands, various restructured and iterative module designs have also emerged. As a result, modular RAG has increasingly become a dominant paradigm in the field, supporting both serialized pipeline execution and end-to-end learning across modular components.

Despite these advances, many RAG systems remain constrained by static control logic, making them ill-suited for complex QA tasks where key evidence may be scattered or initially missed. Recent work on Agentic RAG introduces reasoning and tool use into retrieval, enabling more adap-

tive behavior (Ravuru et al., 2024; Li et al., 2025a; Wu et al., 2025). Motivated by this, we explore two complementary strategies: a One-SHOT retrieval method with token-aware evidence selection, and an agent-driven iterative framework. These approaches aim to improve retrieval robustness and adaptivity in real-world QA scenarios.

2.2 Agentic RAG

The year 2025 is marked as the year of agentic AI, with applications emerging such as agentic LLMs and so on (Ruan et al., 2023; Kong et al., 2024; Zhang et al.). Recent progress in RAG has moved beyond static, rule-based pipelines toward more dynamic, decision-driven systems broadly known as *Agentic RAG*. These systems embed retrieval decisions into the model’s reasoning flow, enabling LLMs to actively determine when and how to interact with external tools during generation.

A growing body of work has demonstrated the effectiveness of prompting large models to autonomously invoke tools, reformulate queries, or break down complex questions. For instance, methods like ReAct (Yao et al., 2023), Self-Ask (Press et al., 2023), and Search-o1 (Li et al., 2025a) allow models to interleave generation with retrieval by identifying knowledge gaps and issuing targeted queries. Built-in function calling APIs (Eleti et al., 2023) further support structured tool use in models like GPT and Gemini.

Beyond prompting, other approaches train models to improve their retrieval capabilities through learning signals. DeepRetrieval (Jiang et al., 2025), Search-R1 (Jin et al., 2025a), and ReZero (Dao and Le, 2025) integrate retrieval into the learning loop, optimizing search behaviors with reinforcement signals. More advanced systems like Deep-Researcher (Zheng et al., 2025) extend this idea to open-domain web environments, encouraging models to reason, search, and synthesize in a tightly coupled loop.

3 Methodology

This section will give the details on how we achieve the two strategies.

3.1 One-SHOT Strategy: “Casting A Bigger Net!”

Our One-SHOT (all capital to infer the bigger net) retrieval strategy implements a token-constrained RAG framework designed to maximize recall

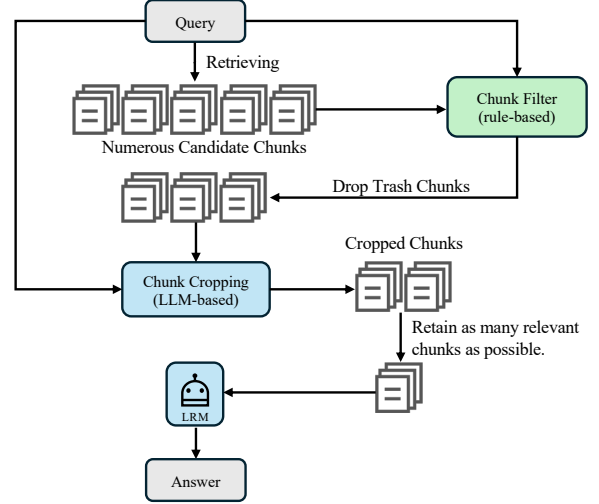


Figure 2: One-SHOT Strategy. This category includes methods such as Token-Constrained Top-K (retain as many relevant chunks as possible), along with techniques like chunk filter and chunk cropping.

within a fixed context window. Rather than selecting a fixed number of top- k (e.g., top-5) chunks, the system dynamically selects as many as possible relevant chunks that can fit within a given token budget, enhancing evidence coverage in a single retrieval, denoted as Token-Constrained Top- K_{max} . A few components are further designed:

3.1.1 Token-Constrained Top- K_{max}

To mitigate the issue of golden chunks being excluded under a fixed Top- k setting, we propose a Token-Constrained Top- K_{max} strategy that aims to **maximize recall within a predefined token budget** (e.g., 32,000 tokens). Each candidate chunk is associated with a relevance score r_i (e.g., embedding similarity) and a token count t_i , and the goal is to select a subset of chunks that together fit within the context window while maximizing the total relevance. Let $x_i \in \{0, 1\}$ be a binary variable indicating whether chunk i is selected, and T_{max} be the token limit. The selection problem can be formalized as:

$$\begin{aligned} & \max \sum_{i=1}^n r_i \cdot x_i \\ & \text{s.t.} \quad \sum_{i=1}^n t_i \cdot x_i \leq T_{max}, \quad x_i \in \{0, 1\} \end{aligned} \tag{1}$$

This formulation effectively prioritizes chunks with high relevance-per-token ratios, enabling dense and informative context construction within the model’s input limits.

3.1.2 Key Optimization Components

As illustrated in Fig. 1(b), such a greedy retrieval inevitably introduce some unrelated chunks. To further refine the evidence retrieved by the Token-Constrained Top- K_{max} strategy, we design two complementary optimization components: *chunks filter* and *chunks summary*.

The **chunk filter** component (rule-based) selectively removes irrelevant chunks and augments missing but potentially useful ones by analyzing chunk meta information such as temporal expressions, locations, and named entities. Specifically, the drop operation filters out chunks that lack alignment with query entities, while the add operation supplements the retrieved set with additional chunks containing matching or complementary elements. This rule-based filtering improves both precision and recall of the retrieval results.

The **chunk cropping** component (LLM-based) further crops the remaining related chunks. It leverages LLM reasoning to assess and refine the initial evidence holistically. Given a specific related chunk that survives the filtering process, the LLM evaluates the contents inside this chunk carefully with respect to the original query and crop out the information deemed irrelevant or redundant. Chunks are basically refined to be shorter after this process, resulting in a more condensed and informative evidence set for answer generation.

In practice, both chunk filter and chunk cropping are conducted first before the token length is satisfied, such that the recall is being possibly maximized.

3.2 Iterative Retrieval Strategy: “Casting A Small Net Multi-Times!”

Our iterative retrieval strategy implements a reasoning agentic RAG (Liang et al., 2025) that leverages the reasoning and function-calling capabilities of large reasoning language models (LRMs) to perform adaptive information retrieval.

3.2.1 Architecture Overview

As shown in Fig. 3, the iterative retrieval system follows a multi-turn paradigm where the language reasoning model acts as an intelligent agent capable of making informed decisions about when and how to retrieve additional information. Unlike traditional RAG systems that rely on explicit pipeline components, our approach embeds the core functionalities implicitly within the model’s reasoning

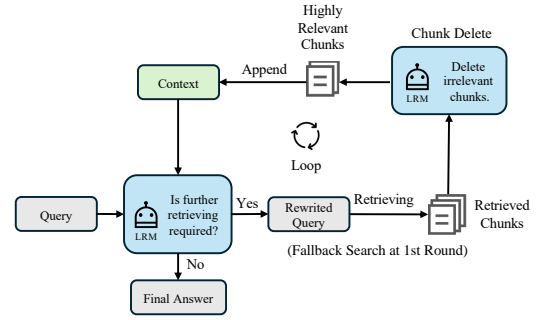


Figure 3: Iterative retrieval strategy, with fallback module to avoid query drift and chunk delete to avoid retrieval laziness.

process, managing context analysis, tool invocation, and information synthesis naturally through the model’s inherent reasoning abilities.

Similar to existing agentic RAG models, e.g., Search-O1 (Li et al., 2025b) and Search-R1 (Jin et al., 2025b), the system employs a structured reasoning process encapsulated within `<think>` and `</think>` tags, where the model evaluates information sufficiency, identifies knowledge gaps, and formulates search queries. When additional information is required, the model generates structured function calls using the `<tool_call>` format. The system provides two primary tools: `chunk_search` for information retrieval (still use the small “fishing net”) and `chunk_delete` for context refinement and information overload prevention.

3.2.2 Multi-Turn Retrieval Process

The iterative retrieval process operates within a configurable maximum turn limit (typically set to 5 turns) to balance thoroughness with efficiency. In each turn, the system performs the following steps:

- Reasoning Phase:** The model analyzes the current context, including the original query and all previously retrieved information, to determine if additional retrieval is necessary.
- Query Formulation:** If retrieval is deemed necessary, the model generates a focused search query designed to address specific information gaps identified during reasoning.
- Retrieval Execution:** The system executes the `chunk_search` function against the vector database to identify relevant document chunks. The function accepts the rewritten query and returns the top- k most relevant document chunks ($k = 5$) with metadata and text content.
- Context Integration:** Retrieved chunks are

integrated into the system context, providing the model with additional information for subsequent reasoning.

5. Iteration Decision: The model determines whether to continue searching or provide a final answer based on the accumulated information.

3.2.3 Key Optimization Components

Our complete iterative system incorporates several critical optimization mechanisms that address the inherent challenges of agentic retrieval:

Fallback Search Mechanism: To mitigate the risk of query drift and ensure comprehensive coverage, the system implements a strategic fallback mechanism during the initial retrieval phase. When the model performs its first chunk_search operation, the system automatically executes an additional retrieval using the original user query in parallel to the model’s reformulated query. This dual-retrieval approach is specifically applied only to the first turn to provide a strong foundation for subsequent iterations—ensuring the model starts with a good baseline of relevant information from the original query while still benefiting from its own query reformulation.

Adaptive Context Management: The system incorporates the chunk_delete function, which allows the model to remove irrelevant retrieved information from the working context by specifying the chunk IDs of unwanted chunks. This mechanism helps maintain focus on pertinent information and prevents the accumulation of noise that could degrade answer quality.

These optimization components are essential for addressing two key challenges:

- **Query drift:** where autonomous query reformulation by LRM drifts away from the original query and leads to less relevant results, as shown in Appx. A.
- **Retrieval laziness:** where in the iterative agentic retrieving, the model prematurely terminates search due to context overload: specifically, the heavier the context retrieved from previous round, the less likely the agent will initiate the next retrieval action due to “cognitive burden”, as shown in Appx. B. Such a cognitive burden is also observed in function calling (Yang et al., 2025b).

The effectiveness of these components is demonstrated in our experiments in Section 4.3.

4 Experiment

4.1 Experiment Settings

Datasets. Our dataset is composed of 1,000 carefully designed questions, all extracted from a corpus of over 40,000 real-world government documents. Specifically, the questions are categorized into four levels (ranging from Level 1 to Level 4) based on two key criteria: whether retrieval requires reasoning and whether answering requires reasoning. We ensure that all answers are retrievable from the original document corpus. Figure 4 illustrates the difference between each level.

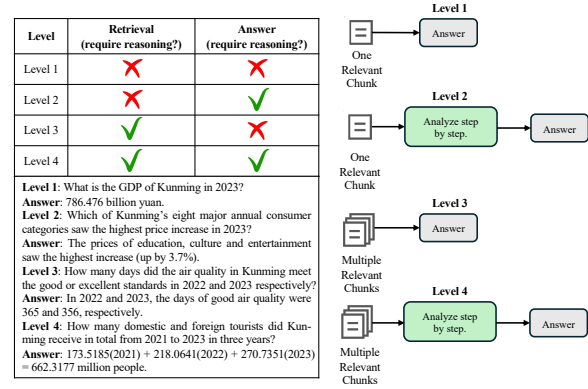


Figure 4: Categories of questions in our dataset. The table in the upper-left shows differences between questions of various levels. The lower-left provides examples of different-level questions and their reference answers. The diagrams on the right illustrate the solving processes for questions of different levels.

Basic RAG Setting. Our baseline system employs traditional single-turn retrieval with a fixed Top-5 strategy. This basic RAG approach retrieves the top-5 most relevant document chunks based on vector similarity and provides them to the language model for answer generation without any iterative refinement or optimization mechanisms. All experiments are conducted using Qwen3-32B as the language model.

Evaluation Setting. We employ an LLM-as-a-judge evaluation framework to assess answer quality across multiple dimensions including factual accuracy, completeness, and relevance. We adopt SenseChat-5 to perform answer evaluation. The evaluation scores are normalized to a 0-100 scale for consistent comparison across different approaches.

Method	L1	L2	L3	L4	Avg
Basic RAG (Top-5)	89.5	88.5	76.0	70.0	81.0
Token-Constrained Top- K_{max}	90.0	91.0	85.5	83.5	87.5
+ Chunk Filter	92.0	93.0	90.0	89.0	91.0
+ Chunk Cropping	91.0	92.5	86.5	84.0	88.5
+ Chunk Filter & Cropping	92.0	93.0	87.5	83.5	89.0
Improvement over Baseline	+2.5	+4.5	+14.0	+19.0	+10.0

Table 1: Ablation study of One-SHOT retrieval components. Each row builds incrementally on the previous configuration.

4.2 Results on One-SHOT Strategy

4.2.1 Ablation studies

We first conduct ablation studies on the One-SHOT retrieval strategy to evaluate the impact of each optimization component, as shown in Table 1. Starting from the baseline *Basic RAG (Top-5)*, our *Token-Constrained Top- K_{max}* strategy brings a notable improvement by relaxing the fixed top- k constraint and retrieving more relevant content within a token budget. This approach performs particularly well on L3 and L4 tasks, which demand deeper reasoning and broader evidence coverage. With the addition of *Chunk Filter*, we achieve the highest overall performance (91.0% average), amounting to a +10.0% gain over the baseline. This confirms the effectiveness of combining token-based recall expansion with relevance-aware filtering.

Adding the *Chunk Filter* before reaching the token budgets further boosts performance, achieving the best overall result (91.0% average). The largest gains appear on L3 (+14.0%) and L4 (+19.0%), where filtering irrelevant chunks mitigates noise from large retrieval sets and allows the model to focus on high-value evidence. The *Chunk Cropping* further improves performance by removing noisy chunks through LLM-based semantic judgments. However, its effect is more moderate compared to the rule-based *Chunk Filter*, possibly due to the trade-off between precision and recall introduced by aggressive semantic cropping. Combining *Chunk Filter* and *Chunk Cropping* yields mixed results, suggesting that their benefits are not strictly additive and require careful integration.

Overall, these results demonstrate that optimizing evidence selection through both token budgeting and relevance refinement is crucial for enhancing RAG performance in complex QA scenarios.

4.2.2 Chunk Retention Analysis

To better understand how the One-SHOT pipeline refines evidence through multi-stage processing, we analyze the *chunk retention ratio* at each stage: after the initial retrieval (Token-Constrained Top-

K_{max}), after applying the rule-based *Chunk Filter*, and finally after applying LLM-based *Chunk Cropping*. The results in Table 2 show that chunk retention varies across QA difficulty levels, with L2 consistently yielding the lowest retention. Both *Chunk Filter* and *Chunk Cropping* reduce chunk volume, with the latter being slightly more aggressive across all levels. Interestingly, L1 and L4 retain more chunks after refinement, suggesting that simple factoid queries (L1) and domain-specific regulatory queries (L4) allow clearer chunk-to-query alignment. When combined, the two methods further reduce context size while maintaining high relevance, confirming that structured and semantic filtering are complementary in eliminating irrelevant or low-utility content, and that refinement sensitivity varies by question type.

Methods	L1	L2	L3	L4
+ Chunk Filter	85	76	80	82
+ Chunk Cropping	83	73	75	79
+ Chunk Filter & Cropping	80	71	72	76

Table 2: Chunk retention ratio (%) after applying chunk refinement modules under different QA difficulty levels. Values indicate the percentage of retrieved chunks retained after each processing step.

Models	Metrics	Basic RAG (Top-5)	One-SHOT Strategy
SenseChat-5	Scores	77.5	83.5
	Mins	30	24
Qwen2.5-32B	Scores	80.5	91.5
	Mins	23	21
Qwen3-32B no think	Scores	77.0	91.0
	Mins	26	24
Qwen3-32B think	Scores	75.0	89.5
	Mins	51	83

Table 3: Evaluation scores and time costs on different LLMs for One-SHOT strategy.

4.2.3 Applications on Different LLMs

As One-SHOT strategy basically increase the recall, so theoretically it won’t behave too much differently on L1 to L4, instead, we would like to see how the different choice of LLMs will affect the overall result, thus, Table 3 only presents the overall evaluation scores and time costs of applying the One-SHOT strategy across different LLMs. Compared to the baseline *Basic RAG (Top-5)*, all models show substantial performance improvements when using “a big net”, i.e., the token-constrained One-SHOT strategy. For instance, *Qwen2.5-32B* achieves the highest score of 91.5 (+11.0), while *SenseChat-5* improves by 6.0 points to reach 83.5. In terms of efficiency, the One-SHOT strategy generally reduces

time costs for most models, as it avoids multi-turn reasoning and minimizes redundant computations. For example, the time for *SenseChat-5* drops from 30 to 24 minutes. However, for *Qwen3-32B* in the *think* configuration, we observe a significant increase in latency (from 51 to 83 minutes). This overhead is attributed to the additional reasoning steps performed within the `<think>` blocks during relevance evaluation. These results demonstrate the effectiveness and generality of the One-SHOT strategy across LLMs with different sizes and reasoning capabilities. They also highlight the trade-off between reasoning complexity and latency, suggesting that careful system design is needed to balance performance gains and computational efficiency.

4.3 Results on Iterative Retrieval Strategy

4.3.1 Ablation studies

We evaluated our iterative retrieval approach through a series of ablation studies, progressively building from a basic agentic system to our complete framework described in Section 3.2. Our experiments reveal several key insights about the challenges and opportunities in multi-turn retrieval systems. As shown in Table 4, the *Basic Agentic (Multi-turn)* can autonomously decide when and what to search across multiple turns, but it shows mixed results: while complex retrieval tasks (L3 and L4) benefit from the iterative approach, simpler questions (L1 and L2) experience slight performance degradation. This counterintuitive result reveals a critical challenge in agentic retrieval: **query drift**, where the model’s autonomous query reformulation may deviate from the user’s original intent, particularly for straightforward questions where the initial query is already well-formed. For example, when given the original query "I left my previous company 3 years ago, can I still recover my housing fund contributions?", the model might reformulate it as "housing fund withdrawal time limit regulations" and then further drift to "housing fund withdrawal conditions for resigned employees". This reformulation fundamentally changes the intent from recovering/reclaiming contributions to asking about withdrawal procedures, leading to irrelevant retrieval results.

+ *Fallback Search* benefits easy QAs more, i.e., L1 and L2 without retrieval reasoning need, vulnerable to query drift. Fallback search performs parallel retrieval using both the model’s reformulated query and the original user query during the

Method	L1	L2	L3	L4	Avg
Basic RAG (Top-5)	89.5	88.5	76.0	70.0	81.0
Basic Agentic (Multi-turn)	86.5	86.0	81.0	80.5	83.5
+ Fallback Search	90.0	89.5	84.0	82.5	87.5
+ Chunk Delete	91.5	91.0	88.5	89.0	90.0
Improvement over Baseline	+2.0	+2.5	+12.5	+19.0	+9.0

Table 4: Ablation study of iterative retrieval components. Each row builds incrementally on the previous configuration.

first retrieval turn, similar to a residual link design, significantly increases the model performance in L1 (+3.5%) and L2 (+3.5%), avoiding query drift. This boost is more evident than L3 (+3.0%) and L4 (+2.0%). This enhancement resulted in consistent improvements across all question levels, effectively mitigating the performance degradation observed in simple questions while maintaining benefits for complex queries.

+ *Chunk Delete* instead benefits hard QA more, i.e., L3 and L4 with retrieval reasoning need, vulnerable to retrieval laziness. As mentioned, we observed **retrieval laziness**: the agent is lazier to conduct the next retrieval when the previous round context is heavier. + *Chunk Delete* further boosts L3 (+4.5%) and L4 (+6.5%) performance, more than L1 (+1.5%) and L2 (+1.5%).

Method	L1	L2	L3	L4	Avg
Basic RAG (Top-5)	89.5 _{1.00}	88.5 _{1.00}	76.0 _{1.00}	70.0 _{1.00}	81.0 _{1.00}
Agentic (Qwen3)	91.5 _{1.22}	91.0 _{1.34}	88.5 _{1.58}	89.0 _{2.22}	90.0 _{1.59}
Agentic (DeepSeek)	83.0 _{1.86}	93.0 _{1.85}	76.0 _{2.02}	74.0 _{2.17}	81.5 _{1.97}
Agentic (DeepSeek*)	92.2 _{1.86}	96.9 _{1.85}	92.7 _{2.02}	90.2 _{2.17}	93.1 _{1.97}

Table 5: Comparison with DeepSeek-R1, with accuracy as the main result and the average retrieval times executed as the subscript. The original DeepSeek results are mostly lower than Qwen3 due to dozens of questions being rejected, so we report the relative percentage of correct answers based on the questions DeepSeek can generate results for, noted as DeepSeek* in grey color.

4.3.2 Applications on Different LLMs

As we can observe in Table 5, no matter based on *Qwen3* (Yang et al., 2025a) or *DeepSeek-R1* (Guo et al., 2025), our agentic iterative strategy can generally achieve better results than basic RAG, and this may be highly due to the more frequent retrieval: as we observe, *Agentic (Qwen3)* increases retrieval by +59%, and *Agentic (DeepSeek)* by +97%. Moreover, in the complex questions L3 and L4, we observe twice higher frequency of retrieval.

However, due to the high-sensitivity filtering and security protocol of *DeepSeek-R1*, a lot of questions are being rejected, resulting worse performance than *Qwen3* (as shown in *Agentic (DeepSeek)*); We exclude those questions and re-

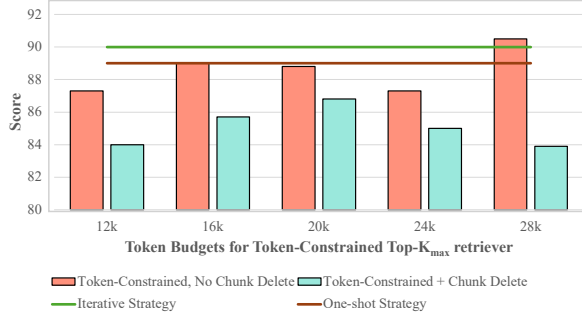


Figure 5: Performance comparison of combined strategies with different configurations.

calculate the relative accuracy rate as *Agentic (DeepSeek*)*, we can see that among the questions *DeepSeek* is allowed to answer, the performance is very promising. Yet, this still supports our decision of deployment based on *Qwen3*.

4.4 Combine both? “Cast A Big Net Multi-Times”

Overviewing the two strategies in Table 1 and Table 4, we can see that “different paths lead to the same goal”, both strategies, either “cast a bigger net” or “cast a small net multi-times”, reach the similar accuracy of 90% - 91%. For practitioners working with government documents and sufficient token budgets, “cast a bigger net” is recommended due to fewer retrieval rounds and runtime; for those with limited tokens and who are comfortable with more runtime, “cast a small net multi-times” is recommended.

Combining Both: To explore the potential synergies between One-SHOT and iterative approaches, we combined elements from both strategies. Specifically, we replaced the fixed Top-5 retriever in the Iterative Retrieval Strategy’s fallback search mechanism with the Token-Constrained Top-K retriever from our One-SHOT strategy, testing various context budgets to maximize golden chunk recall in the first retrieval round.

Surprisingly, all combined configurations underperformed the original iterative approach. Analysis revealed that the longer working context after the first retrieval round (due to Token-Constrained Top- K_{max}) caused the LLM-based chunk delete tool to perform poorly, frequently removing genuinely useful chunks that were critical for answering the questions (see Appendix C for a detailed case study).

When we removed the chunk delete from the combined strategy, performance improved compared to the combined strategy with chunk delete, but still fell short of the original iterative approach.

It revealed that without the chunk delete, the longer working context after the first retrieval round led to severe **retrieval laziness**. Even when information was incomplete, the model tended to avoid further retrieval operations, effectively degenerating into a One-SHOT strategy. This explains why the combined approach without chunk delete achieved scores nearly same to the pure One-SHOT.

Analysis and Implications: Our experiments demonstrate that the two strategies are not complementary in their current forms. The Token-Constrained Top- K_{max} approach, while effective in isolation, creates longer initial contexts that interfere with the iterative retrieval process in two ways: (1) it degrades the performance of the chunk delete tool, leading to the removal of useful information, and (2) it exacerbates retrieval laziness when the chunk delete tool is removed, preventing the model from conducting necessary follow-up searches. From cognitive science’s perspective, “People make better choices with 6 options but experience decision paralysis with 30 options” (Iyengar and Lepper, 2000): this is what is happening when combining the strategies: chunk delete works fine when given 5 chunk options in the iterative strategy, but fails when being combined with Top- K_{max} , which usually retrieves 30+ chunks.

These findings suggest that the one-shot and iterative strategies operate under different assumptions about optimal context management and retrieval behavior. The one-shot strategy benefits from comprehensive initial retrieval within a controlled context size, while the iterative strategy relies on focused, incremental information gathering with active context management. Combining these approaches requires more sophisticated mechanisms to balance context size, relevance filtering, and retrieval motivation—an area for future research.

5 Conclusion

In this paper, we explored two strategies, namely “casting a bigger net” and “casting a small net multiple times”, supported by our carefully designed retrieval and reasoning modules. Both strategies demonstrated significant improvements, achieving over +10% performance gains compared with basic RAG baselines when applied to complex government documents. Our findings highlight the importance of adopting adaptive retrieval and iterative reasoning to better handle lengthy, heterogeneous government texts.

References

- Jiawei Chen, Hongyu Lin, Xianpei Han, and Le Sun. 2024. Benchmarking large language models in retrieval-augmented generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17754–17762.
- Alan Dao and Thinh Le. 2025. [Rezero: Enhancing llm search ability by trying one-more-time](#). *Preprint*, arXiv:2504.11001.
- Atty Eleti, Jeff Harris, and Logan Kilpatrick. 2023. [Function calling and other api updates](#).
- Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yixin Dai, Jiawei Sun, Haofen Wang, and Haofen Wang. 2023. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997*, 2:1.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shitong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and 1 others. 2025. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *ACM Transactions on Information Systems*, 43(2):1–55.
- Sheena S Iyengar and Mark R Lepper. 2000. When choice is demotivating: Can one desire too much of a good thing? *Journal of personality and social psychology*, 79(6):995.
- Pengcheng Jiang, Jiacheng Lin, Lang Cao, Runchu Tian, Seongku Kang, Zifeng Wang, Jimeng Sun, and Jiawei Han. 2025. [Deepretrieval: Hacking real search engines and retrievers with large language models via reinforcement learning](#). *Preprint*, arXiv:2503.00223.
- Bowen Jin, Hansi Zeng, Zhenrui Yue, Jinsung Yoon, Sercan Arik, Dong Wang, Hamed Zamani, and Jiawei Han. 2025a. [Search-r1: Training llms to reason and leverage search engines with reinforcement learning](#). *Preprint*, arXiv:2503.09516.
- Bowen Jin, Hansi Zeng, Zhenrui Yue, Jinsung Yoon, Sercan Arik, Dong Wang, Hamed Zamani, and Jiawei Han. 2025b. Search-r1: Training llms to reason and leverage search engines with reinforcement learning. *arXiv preprint arXiv:2503.09516*.
- Yilun Kong, Jingqing Ruan, Yihong Chen, Bin Zhang, Tianpeng Bao, Shi Shiwei, Du Qing, Xiaoru Hu, Hangyu Mao, Ziyue Li, and 1 others. 2024. Tptu-v2: Boosting task planning and tool usage of large language model-based agents in real-world industry systems. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pages 371–385.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, and 1 others. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in neural information processing systems*, 33:9459–9474.
- Jiarui Li, Ye Yuan, and Zehua Zhang. 2024. Enhancing llm factual accuracy with rag to counter hallucinations: A case study on domain-specific queries in private knowledge-bases. *arXiv preprint arXiv:2403.10446*.
- Xiaoxi Li, Guanting Dong, Jiajie Jin, Yuyao Zhang, Yujia Zhou, Yutao Zhu, Peitian Zhang, and Zhicheng Dou. 2025a. [Search-o1: Agentic search-enhanced large reasoning models](#). *Preprint*, arXiv:2501.05366.
- Xiaoxi Li, Guanting Dong, Jiajie Jin, Yuyao Zhang, Yujia Zhou, Yutao Zhu, Peitian Zhang, and Zhicheng Dou. 2025b. Search-o1: Agentic search-enhanced large reasoning models. *arXiv preprint arXiv:2501.05366*.
- Jintao Liang, Gang Su, Huifeng Lin, You Wu, Rui Zhao, and Ziyue Li. 2025. Reasoning rag via system 1 or system 2: A survey on reasoning agentic retrieval-augmented generation for industry challenges. *arXiv preprint arXiv:2506.10408*.
- Xi Victoria Lin, Xilun Chen, Mingda Chen, Weijia Shi, Maria Lomeli, Richard James, Pedro Rodriguez, Jacob Kahn, Gergely Szilvasy, Mike Lewis, and 1 others. 2023. Ra-dit: Retrieval-augmented dual instruction tuning. In *The Twelfth International Conference on Learning Representations*.
- Xinbei Ma, Yeyun Gong, Pengcheng He, Hai Zhao, and Nan Duan. 2023. Query rewriting in retrieval-augmented large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5303–5315.
- Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A. Smith, and Mike Lewis. 2023. [Measuring and narrowing the compositionality gap in language models](#). *Preprint*, arXiv:2210.03350.
- Chidaksh Ravuru, Sagar Srinivas Sakhinana, and Venkataramana Runkana. 2024. Agentic retrieval-augmented generation for time series analysis. *arXiv preprint arXiv:2408.14484*.
- Vipula Rawte, Swagata Chakraborty, Agnibh Pathak, Anubhav Sarkar, SM_Towhidul Islam Tonmoy, Aman Chadha, Amit Sheth, and Amitava Das. 2023. The troubling emergence of hallucination in large language models-an extensive definition, quantification, and prescriptive remediations. Association for Computational Linguistics.
- Stephen Robertson, Hugo Zaragoza, and 1 others. 2009. The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends® in Information Retrieval*, 3(4):333–389.

- Jingqing Ruan, Yihong Chen, Bin Zhang, Zhiwei Xu, Tianpeng Bao, Hangyu Mao, Ziyue Li, Xingyu Zeng, Rui Zhao, and 1 others. 2023. Tptu: Task planning and tool usage of large language model-based ai agents. In *NeurIPS 2023 Foundation Models for Decision Making Workshop*.
- Aditi Singh. 2023. Exploring language models: A comprehensive survey and analysis. In *2023 International Conference on Research Methodologies in Knowledge Management, Artificial Intelligence and Telecommunication Engineering (RMKMATE)*, pages 1–4. IEEE.
- Junde Wu, Jiayuan Zhu, and Yuyuan Liu. 2025. [Agentic reasoning: Reasoning llms with tools for the deep research](#). *Preprint*, arXiv:2502.04644.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025a. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*.
- Chenxiao Yang, Nathan Srebro, David McAllester, and Zhiyuan Li. 2025b. Pencil: Long thoughts with short memory. *arXiv preprint arXiv:2503.14337*.
- 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.
- Bin Zhang, Hangyu Mao, Jingqing Ruan, Ying Wen, Yang Li, Shao Zhang, Zhiwei Xu, Dapeng Li, Ziyue Li, Rui Zhao, and 1 others. Controlling large language model-based agents for large-scale decision-making: An actor-critic approach. In *ICLR 2024 Workshop on Large Language Model (LLM) Agents*.
- Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, and 1 others. 2023. Siren’s song in the ai ocean: a survey on hallucination in large language models. *arXiv preprint arXiv:2309.01219*.
- Penghao Zhao, Hailin Zhang, Qinhan Yu, Zhengren Wang, Yunteng Geng, Fangcheng Fu, Ling Yang, Wentao Zhang, Jie Jiang, and Bin Cui. 2024. Retrieval-augmented generation for ai-generated content: A survey. *arXiv preprint arXiv:2402.19473*.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, and 1 others. 2023. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 1(2).
- Yuxiang Zheng, Dayuan Fu, Xiangkun Hu, Xiaojie Cai, Lyumanshan Ye, Pengrui Lu, and Pengfei Liu. 2025. [Deepresearcher: Scaling deep research via reinforcement learning in real-world environments](#). *Preprint*, arXiv:2504.03160.
- Yizhang Zhu, Shiyin Du, Boyan Li, Yuyu Luo, and Nan Tang. 2024. Are large language models good statisticians? *arXiv preprint arXiv:2406.07815*.

A Query Drift

To illustrate the query drift phenomenon discussed in Section 4.3, we provide a concrete example here:

Original User Query: "I have been away from my previous company for 3 years. Can I still recover my housing provident fund contributions?"

Model-Reformulated Query: "Regulations on housing provident fund withdrawal time limits after resignation"

Analysis: The original query seeks information about recovering or reclaiming housing provident fund contributions from a previous employer after a 3-year gap. However, the model’s autonomous reformulation shifts the focus to withdrawal regulations and time limits, fundamentally changing the intent from recovery/reclaim to withdrawal procedures. This query drift leads to retrieval of irrelevant chunks about standard withdrawal processes rather than the specific legal provisions for recovering contributions from former employers, ultimately resulting in an incorrect or incomplete answer.

This example demonstrates how autonomous query reformulation in agentic systems can deviate from the user’s original intent, particularly affecting the retrieval of relevant evidence and degrading overall system performance.

B Retrieval Laziness

To validate the retrieval laziness phenomenon discussed in Section 4.3, we conducted a controlled experiment to measure how context length affects the model’s tendency to initiate follow-up retrieval calls.

Experimental Setup: We manually injected redundant irrelevant chunks into the initial retrieval results to reach specific context lengths while ensuring that the retrieved chunks contained only partial golden chunks—insufficient to fully answer the user’s question. We then measured the probability that the model would make a subsequent chunk_search call.

Context Length	Follow-up Retrieval Probability
3k tokens	95%
6k tokens	90%
9k tokens	50%
12k tokens	25%

Table 6: Retrieval laziness validation: probability of initiating follow-up retrieval calls at different context lengths.

Results: As shown in Table 6, the model’s propensity to continue searching decreases dramatically as context length increases. With shorter contexts, the model correctly identifies information gaps and performs follow-up searches in the majority of cases. However, this rate drops significantly when the context becomes longer, demonstrating severe retrieval laziness in extended contexts.

This experiment confirms our hypothesis that longer working contexts exacerbate retrieval laziness, leading to premature termination of the search process even when critical information is missing.

C Chunk Delete Tool Failure in Combined Strategy

To illustrate the chunk delete tool failure discussed in Section 4.4, we provide a concrete case study demonstrating how longer working contexts impair the tool’s effectiveness.

Query: "What was the area of flower cultivation in Kunming in the year when the first 'China-Kunming Dounan Flower Exhibition' was held?"

Context: The combined strategy (Token-Constrained Top- K_{max} + Iterative) retrieved 10 chunks in the first round, creating a working context of 14,507 tokens. Among these chunks, one contained the crucial information that the first exhibition was held in 2023.

Chunk Delete Tool Behavior: When presented with this extended context, the chunk delete tool incorrectly removed 9 out of 10 chunks, including the chunk containing the correct answer about the 2023 exhibition.

Resulting Error: Without the correct chunk, the model provided an incorrect answer based on incomplete information.

Analysis: This demonstrates how extended contexts cause cognitive overload in the chunk delete tool, leading to the removal of genuinely useful information critical for accurate answers.