

RES-RAG: Reinforced Evidence Set Retrieval for Retrieval-Augmented Generation

Anonymous ACL submission

Abstract

Retrieval-augmented generation (RAG) is widely used to incorporate external knowledge into large language models (LLMs), effectively improving answer quality in question-answering (QA) tasks. Since external knowledge may far exceed the context window of LLMs, retrieving the key evidence they need within a limited context is challenging. Existing heuristic retrieval methods often add more evidence in context to improve answer quality, but some non-essential evidence may contribute little, potentially leading to over-retrieval and increased noise. To overcome these challenges, we propose RES-RAG, which, for the first time, models evidence retrieval as an action decision problem under context constraints. Compared to existing methods such as top- k and truncation, RES-RAG can adaptively retrieve evidence that contributes to answer quality. Specifically, we construct a lightweight retrieval strategy model. We use reinforcement learning (RL) to optimize the retrieval process, making evidence retrieval quality-oriented. Furthermore, we design a joint quality-cost function for the optimization signal in RL. We use the F1 score as a positive reward and penalize over-retrieval. This achieves a balance between answer quality and retrieval cost. Experiments across multiple open-domain QA benchmarks demonstrate that RES-RAG achieves higher answer quality at lower evidence cost, reflecting a superior quality-cost trade-off. Code and data can be found at <https://anonymous.4open.science/r/RES-RAG-0106>.

1 Introduction

Retrieval-Augmented Generation (RAG) has been emerged as a popular paradigm for mitigating the hallucination and knowledge timeliness problems of large language models (LLMs) (Shuster et al., 2021; Asai et al., 2024; Béchar and Ayala, 2024; Huang et al., 2025; Tonmoy et al., 2024; Izacard

Q: Where does the Civil Liberties Act place the blame for the internment of U.S. citizens?

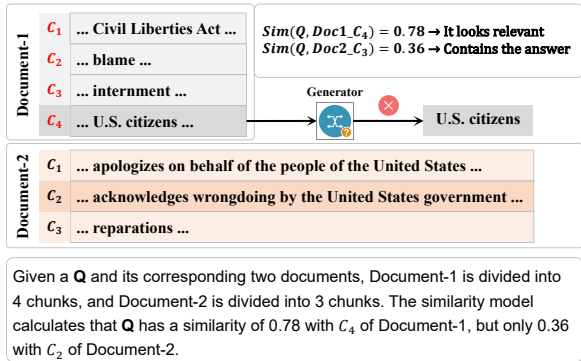


Figure 1: Given a query, similarity-based retrieval methods tend to select Document-1 because its text blocks match the query terms. However, Document-2 contains the key answer to the query.

and Grave, 2020; Guu et al., 2020; Izacard et al., 2023). Its core idea is to retrieve external evidence from updatable corpora and combine it with the model’s parameterized knowledge to answer given open questions (Karpukhin et al., 2020; Izacard et al., 2023; Kim and Lee, 2024; Gao et al., 2023c; Vu et al., 2023; Lewis et al., 2020). However, as longer and more complex contexts are introduced, the model’s context window needs to expand longer and longer. As a result, the computational and latency costs of long-context reasoning typically increase significantly with input size, making evidence retrieval under budget constraints a key trade-off between performance and cost (Lin et al., 2025; Taguchi et al., 2025; Lahmy and Yozevitch, 2025).

Existing research typically focuses on optimizing each component (Gao et al., 2023b; Formal et al., 2021; Ren et al., 2021; Pradeep et al., 2020; Press et al., 2022). For example, stronger retrievers (Izacard et al., 2021; Xiong et al., 2020; Ni et al., 2021; Xiao et al., 2022a), more robust rerankers, or more reliable generators (Yu et al., 2021; Zemlyan-

067 [skiy et al., 2021](#); [Glass et al., 2022](#); [Su et al., 2022a](#)).
068 In these pipelined designs, a common assumption
069 is that the context can be simplified to a set of
070 text fragments ([Chen et al., 2020](#); [Li et al., 2021](#)).
071 Furthermore, evidence retrieval relies on heuristic
072 truncation, fixed windows, or similarity-based
073 top- k strategies. This approach presents two main
074 challenges when dealing with long, complex con-
075 texts or scattered evidence: (1) it introduces a large
076 amount of irrelevant noise that interferes with the
077 generator’s reasoning; (2) it may miss crucial evi-
078 dence within the context window, leading to in-
079 complete answers ([Karpukhin et al., 2020](#); [Krishna
080 et al., 2021](#); [Dai et al., 2025](#)). These challenges
081 limit existing paradigms’ ability to approach the
082 optimal efficiency frontier along the answer quality–
083 evidence cost dimension.

084 To address the above challenges, we propose
085 RES-RAG (Reinforced Evidence Set Retrieval for
086 Retrieval-Augmented Generation). Given a query
087 and its corresponding set of context documents, we
088 train a lightweight model as the evidence retriever
089 that outputs the set of retrieved chunk indices for
090 each document. This design explicitly models evi-
091 dence retrieval as a structured prediction task, en-
092 abling the system to produce controllable results
093 under context constraints. Meanwhile, structured
094 output has been shown to significantly improve the
095 model’s controllability and usability ([Scholak et al.,
096 2021](#); [Muralidharan and Thomas, 2024](#); [Beurer-
097 Kellner et al., 2024](#); [Park et al., 2024](#)).

098 First, we place evidence retrieval within a
099 budget-aware reinforcement learning (RL) opti-
100 mization framework. We call the maximum con-
101 text capacity available for retrieving evidence the
102 budget, and the context actually consumed by the
103 evidence the cost. The set of indices output by the
104 retrieval model is treated as actions. This process
105 offers two key advantages: (1) the retrieval results
106 are presented in an indexed format, reducing the
107 complexity of the action space; (2) the generator
108 receives a compact set of retrieved evidence, re-
109 ducing the burden of long contextual noise and
110 increasing the effective evidence density.

111 Secondly, to improve the RL optimization pro-
112 cess, we design a quality-cost joint reward: the
113 reward signal is measured by the F1 score; redun-
114 dant evidence and over-retrieval are penalized by an
115 over-budget penalty. This reward characterizes the
116 two-dimensional trade-off between answer qual-
117 ity and evidence cost on the same scale, thereby
118 driving the model to converge at the minimally

sufficient evidence granularity. Optimizing evi- 119
dence retrieval via RL under joint quality–cost 120
constraints provides a practical way to improve the 121
quality–cost trade-off in the RAG pipeline ([Guo 122
et al., 2025](#); [Li et al., 2025b](#)). 123

124 Finally, we systematically evaluate RES-RAG’s
125 performance across several popular QA datasets,
126 demonstrating that it can improve QA task quality
127 at lower cost.

128 The main contributions of this paper are as fol-
129 lows:

- We propose an optimization framework centered 130
on evidence retrieval that optimizes re- 131
trieval strategies via RL. 132
- We design a joint quality-cost function to im- 133
prove the trade-off between answer quality 134
and evidence cost. 135
- We evaluate the system across multiple open- 136
domain question-and-answer benchmarks and 137
report end-to-end metrics and costs. 138

139 In summary, our goal is to propose a budget-
140 aware evidence-retrieval mechanism for QA tasks,
141 enabling LLMs to select better factual evidence
142 within a limited budget and to optimize the quality
143 and cost of answer generation.

2 Related Work 144

2.1 RAG for LLM Inference 145

146 A typical RAG pipeline usually includes steps such
147 as query rewriting or expansion, candidate retrieval,
148 reranking, and fusing query and evidence to gener-
149 ate the answer ([Ma et al., 2023](#); [Wang et al.,
150 2024b](#)). In QA scenarios, RAG performance de-
151 pends not only on retrieval recall but also heavily
152 on the evidence organization method and the gener-
153 ator’s ability to locate evidence within the con-
154 text. Therefore, the coupling effect of retrieval —
155 evidence—generation has been widely discussed
156 ([Barnett et al., 2024](#); [Siriwardhana et al., 2023](#);
157 [Besta et al., 2024](#)). Although existing work has
158 made significant progress in retrieval training ([San-
159 thanam et al., 2021](#); [Shi et al., 2023](#)), reranker de-
160 sign ([Zhang et al., 2023](#); [Edge et al., 2024](#)), and gener-
161 ator fusion strategies ([de Jong et al., 2022](#); [Aki-
162 moto et al., 2024](#)), most methods still simplify evi-
163 dence retrieval to global top- k or similarity thresh-
164 old filtering ([Saxena et al., 2025](#); [Xin et al., 2025](#)),
165 assuming that these heuristic retrieval strategies
166 can adequately support downstream short-answer

generation. This default assumption is prone to failure in situations involving long contexts and scattered evidence. In particular, when complementary information across chunks, or the elimination of interfering segments, is required, relevance ranking with a fixed k often fails to align with the final metric (Hei et al., 2024; Mortaheb et al., 2025).

2.2 Evidence Retrieval and Reranking

To improve the evidence quality of RAG, one type of method sorts candidate passages using a dense retriever, then truncates the top- k chunks as the evidence (Pradeep et al., 2023; Zheng et al., 2024). The advantage of this type of method is its relatively simple training and deployment. Still, its retrieval strategy is usually fixed-length, either globally or weakly adaptive, making it difficult to consider the contribution of different passages to the answer (Xu et al., 2023). Another type of research focuses on more refined evidence retrieval, such as learning a gating mechanism on the candidate set to decide which chunks to retain (Wang et al., 2023; Sheng et al., 2025; Xie et al., 2024). These selectors are often trained with supervisory signals to reduce noise and increase evidence density (Shen et al., 2022). However, supervised retrieval is often limited by labeling costs or pseudo-label bias, and a gap remains between local relevance and final answerability under long context conditions, making it difficult for the retriever to generalize stably (Su et al., 2022b; Asai et al., 2024; Christmann et al., 2024). Furthermore, previous work has pointed out that adding more seemingly relevant fragments reduces the generator’s attention allocation to key evidence, thereby impairing the stability of answer generation (Günther et al., 2024; Ma et al., 2022; Wen et al., 2004). Therefore, evidence retrieval is also a problem of achieving high-density, low-interference evidence assembly under a limited context budget (Dai et al., 2025; Wang et al., 2024d).

2.3 RL for Retrieval Policies

To address the target-mismatch problem of inconsistency between relevance scores and answer quality, some works use reinforcement learning to treat retrieval as a decision-making process and train the retrieval model using the generator’s answer quality as the reward signal (Kulkarni et al., 2024; Wang et al., 2024c). In this paradigm, the policy model can leverage end-to-end feedback to learn which combinations of evidence are more likely to yield a correct answer, thereby surpassing heuris-

tics based on fixed thresholds or static top- k rules (Java et al., 2025; Tang et al., 2025). Existing RL-related work covers multiple levels, including using RL to learn query rewriting and retrieval action sequences (Khattab et al., 2022; Li et al., 2025a; Xu et al., 2025), passage retrieval/stopping strategies (Zhu et al., 2025), and treating the reranker as a learnable ranking strategy (Zhang et al., 2025; Gao et al., 2023a). However, these methods often employ a global retrieval action space, making it difficult to directly characterize the structured constraints on how context contributes evidence and the extent to which each piece of evidence contributes (Liu and Lapata, 2019; Xiao et al., 2022b).

We propose RES-RAG to enable the retrieval strategy to learn how to retrieve the sufficient and necessary set of evidence in complex, long-context scenarios.

3 Method

We propose RES-RAG, a budget-aware framework for evidence retrieval and answer generation in complex long-context. Given a query and its corresponding set of HTML documents, RES-RAG models evidence retrieval as a structured decision. The retriever policy selects a compact set of chunk indices under a chunk budget, and a frozen generator \mathcal{G} produces the answer conditioned on the retrieved evidence and the JSON-LD¹ annotation. We optimize the retriever via RL, balancing rewards to encourage minimal sufficient evidence and high-quality answers.

3.1 Task Definition

Each instance consists of a query q , a set of M HTML documents:

$$\mathcal{D} = \{d_1, \dots, d_M\}, \quad (1)$$

and a gold answer a^* . Each document d_m ($1 \leq m \leq M$) includes structured annotation in JSON-LD format. Our goal is to learn a retriever policy π_θ that outputs a structured action:

$$\hat{S} = \{\hat{S}_1, \dots, \hat{S}_M\}, \quad (2)$$

where \hat{S}_m denotes the retrieved chunk indices from d_m . In parallel, we extract JSON-LD annotation from each document via HTML parsing²:

$$J_m = \text{JSONLD}(d_m), \quad 1 \leq m \leq M, \quad (3)$$

¹<https://json-ld.org/>

²<https://www.crummy.com/software/BeautifulSoup/>

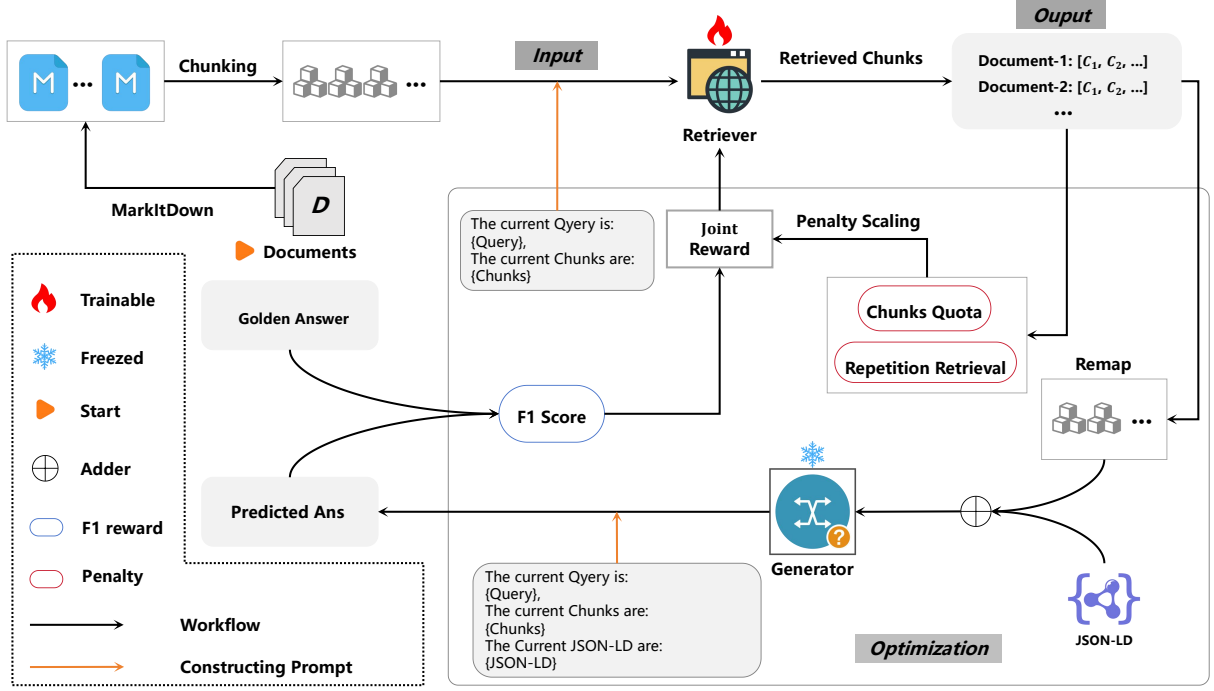


Figure 2: RES-RAG framework. Our training uses the GRPO algorithm. The optimization objective of RES-RAG is the Retriever, and the joint reward is calculated using the predicted answer generated by the Generator and a penalty.

and denote the aggregated annotation as $J = \{J_1, \dots, J_M\}$. Given \hat{S} , retrieved indices are remapped to evidence text and, together with J , are fed into a frozen generator \mathcal{G} to produce an answer \hat{a} .

3.2 Document Processing and Chunking

We convert each d_m into a Markdown string using `Markdown`³: $x_m = \text{Markdown}(d_m)$. We then split x_m into N_m chunks: $\mathcal{X}_m = \{x_{m,1}, x_{m,2}, \dots, x_{m,N_m}\}$. Each chunk is constrained by a maximum token length L :

$$\text{len}(\text{Tok}_{\mathcal{G}}(x_{m,n})) \leq L, \quad 1 \leq n \leq N_m, \quad (4)$$

where $\text{Tok}_{\mathcal{G}}(\cdot)$ uses the same tokenizer as the retriever, ensuring the chunk length matches the actual context usage at inference time. Algorithm A provides the overall chunking procedure. For convenience, we denote the set of all chunk lists as $\mathcal{X} = \{\mathcal{X}_1, \dots, \mathcal{X}_M\}$.

3.3 Components of RES-RAG and Reward Design

Structured Evidence Retriever The Retriever takes (q, \mathcal{X}) as input and outputs a structured action \hat{S} specifying which chunks to include as evidence.

To make evidence decisions explicit and auditable, we serialize \hat{S} as a JSON object:

$$\hat{S} = \{ \text{"Document_1"} : [s_{1,1}, \dots, s_{1,K_1}], \dots, \text{"Document_M"} : [s_{M,1}, \dots, s_{M,K_M}] \}, \quad (5)$$

where each $s_{m,k}$ is an integer chunk index satisfying $s_{m,k} \in \{1, \dots, N_m\}$, and we allow $\hat{S}_m = \emptyset$, this means no evidence be retrieved from d_m . Within each d_m , we encourage π_{θ} to output indices in descending order of their estimated relevance to the q . We parameterize the retriever as a conditional generative policy: $\hat{S} \sim \pi_{\theta}(\cdot | q, \mathcal{X})$, implemented by an instruction-tuned LLM that directly outputs the JSON action. To improve parsability and reduce invalid actions, we produce constraints during decoding (Appendix B).

Generator and Evidence Remapping Given a sampled action \hat{S} , we remap indices to evidence text by selecting corresponding chunks and serializing them in a fixed order:

$$E(\hat{S}) = \text{Concat}_{m=1}^M \text{Concat}_{n \in \hat{S}_m} x_{m,n}, \quad (6)$$

where `Concat` follows the index order in each list \hat{S}_m . We then use \mathcal{G} to produce an answer:

$$\hat{a} = \mathcal{G}(q, E(\hat{S}), J). \quad (7)$$

³<https://github.com/microsoft/markitdown>

We serialize (E, J) into \mathcal{G} prompt using a fixed template (Appendix C).

Chunk-level Budget Penalty We penalize over-retrieval at the chunk level to encourage compact evidence sets. For each d_m , define the retrieved-chunk count $C_m(\hat{S}) = |\hat{S}_m|$. Let N_m be the number of chunks in d_m after chunking. We define a segmented per-document penalty with a free region proportional to $\sqrt{N_m}$:

$$\phi_m(\hat{S}) = \begin{cases} 0, & \text{if } C_m(\hat{S}) \leq \sqrt{N_m}, \\ 0.5, & \text{if } \sqrt{N_m} < C_m(\hat{S}) \leq N_m, \\ 1, & \text{if } C_m(\hat{S}) > N_m. \end{cases} \quad (8)$$

The overall budget penalty averages across documents and scales by $\eta > 0$:

$$P_{\text{budget}}(\hat{S}) = \eta \cdot \frac{1}{M} \sum_{m=1}^M \phi_m(\hat{S}). \quad (9)$$

Answer Quality and Retrieval Consistency We measure answer quality using token-level F1 between \hat{a} and a^* : $\text{F1}_{\text{ans}}(\hat{a}, a^*) \in [0, 1]$, with standard normalization, including lowercasing, punctuation removal, and whitespace normalization. To stabilize learning when \mathcal{G} produces unusable answers, we use a fixed non-positive penalty $\gamma_{\text{bad}} < 0$:

$$U_{\text{ans}}(\hat{a}) = \begin{cases} \alpha_{\text{ans}} \cdot \text{F1}_{\text{ans}}(\hat{a}, a^*), & \hat{a} \text{ is not empty,} \\ \gamma_{\text{bad}}, & \text{otherwise,} \end{cases} \quad (10)$$

where $\alpha_{\text{ans}} > 0$. Beyond answer quality, we add an auxiliary retrieval-consistency signal against heuristic labels \tilde{S} (Sec . 3.4.1), which provides a denser learning signal than answer-only rewards. Although \hat{S}_m and \tilde{S}_m are output as index lists, each index uniquely identifies a chunk text $x_{m,n}$, so we compute consistency at the *chunk level*. Let $\text{set}(\cdot)$ remove duplicates from a list, and define the retrieved and labeled evidence-chunk sets for d_m as

$$\begin{aligned} \hat{\mathcal{E}}_m &= \{x_{m,n} \mid n \in \text{set}(\hat{S}_m)\}, \\ \tilde{\mathcal{E}}_m &= \{x_{m,n} \mid n \in \text{set}(\tilde{S}_m)\}. \end{aligned} \quad (11)$$

We then compute per-document chunk-level F1:

$$\text{F1}_m(\hat{\mathcal{E}}_m, \tilde{\mathcal{E}}_m) = \frac{2|\hat{\mathcal{E}}_m \cap \tilde{\mathcal{E}}_m|}{|\hat{\mathcal{E}}_m| + |\tilde{\mathcal{E}}_m|}, \quad (12)$$

with the convention that $\text{F1}_m = 1$ if both sets are empty, and 0 if only one is empty. The macro-

averaged retrieval consistency is

$$\text{F1}_{\text{ret}}(\hat{S}, \tilde{S}) = \frac{1}{M} \sum_{m=1}^M \text{F1}_m(\hat{\mathcal{E}}_m, \tilde{\mathcal{E}}_m). \quad (13)$$

Joint Reward Putting everything together, for valid actions, we define the reward:

$$\begin{aligned} R(\hat{S}) &= \text{clip}\left(U_{\text{ans}}(\hat{a}) \right. \\ &\quad \left. + \alpha_{\text{ret}} \cdot \text{F1}_{\text{ret}}(\hat{S}, \tilde{S}) \right. \\ &\quad \left. - P_{\text{budget}}(\hat{S}), R_{\text{min}}, R_{\text{max}}\right), \end{aligned} \quad (14)$$

where $\alpha_{\text{ret}} > 0$, and $\text{clip}(\cdot, R_{\text{min}}, R_{\text{max}})$ improves optimization stability. In our design, we apply an additional fixed negative reward for samples caused by external failures.

3.4 Training the Retriever via GRPO

We optimize only the Retriever π_θ using GRPO. The \mathcal{G} is frozen and serves as part of the environment that outputs \hat{a} based on \hat{S} , then outputs a scalar reward:

$$(q, \mathcal{X}) \xrightarrow{\pi_\theta} \hat{S} \xrightarrow{\text{remap}} (E, J) \xrightarrow{\mathcal{G}} \hat{a} \xrightarrow{\text{score}} R(\hat{S}). \quad (15)$$

3.4.1 SFT Warm-up

To stabilize structured generation and reduce RL exploration difficulty, we warm up the Retriever with supervised fine-tuning (SFT). We construct heuristic labels \tilde{S} using a similarity-based rule and fine-tune π_θ to imitate the labeled JSON actions.

Specifically, for each d_m , we compute a similarity score between the q and each chunk $x_{m,n}$ using a sentence embedding model⁴. We keep chunks with similarity greater than a threshold (0.4 in our design) as the top K_m indices to form the label list:

$$\tilde{S}_m = [\tilde{s}_{m,1}, \dots, \tilde{s}_{m,K_m}], \quad \tilde{s}_{m,k} \in \{1, \dots, N_m\}. \quad (16)$$

All indices in both \tilde{S} and \hat{S} are arranged by relevance to, in the same form as Equation 5. The resulting checkpoint initializes the RL policy, and a frozen copy serves as the reference policy, $\pi_{\theta_{\text{ref}}}$.

3.4.2 Group Sampling and Relative Advantage

For each training instance t , we sample a group of W candidate actions using the behavior policy $\pi_{\theta_{\text{old}}}$:

$$\{\hat{S}_{t,1}, \dots, \hat{S}_{t,W}\} \sim \pi_{\theta_{\text{old}}}(\cdot \mid q_t, \mathcal{X}_t). \quad (17)$$

⁴<https://huggingface.co/sentence-transformers/all-mpnet-base-v2>

Each candidate is evaluated by the \mathcal{G} , producing reward $\hat{R}_{t,w} = R(\hat{S}_{t,w})$. To mitigate reward-scale differences across prompts and reduce gradient variance, we compute group-relative advantages by normalizing rewards within each prompt group:

$$\bar{R}_t = \frac{1}{W} \sum_{w=1}^W \hat{R}_{t,w}, \quad \sigma_t = \sqrt{\frac{1}{W} \sum_{w=1}^W (\hat{R}_{t,w} - \bar{R}_t)^2}, \quad (18)$$

$$\text{ADV}_{t,w} = \frac{\hat{R}_{t,w} - \bar{R}_t}{\sigma_t + \epsilon_{\text{adv}}}, \quad (19)$$

where ϵ_{adv} is a constant for numerical stability.

3.4.3 Clipped Policy Optimization Objective

Let $\pi_{\theta_{\text{old}}}$ denote the behavior policy used to sample actions. We define the importance ratio

$$\rho_{t,w} = \frac{\pi_{\theta}(\hat{S}_{t,w} | q_t, \mathcal{X}_t)}{\pi_{\theta_{\text{old}}}(\hat{S}_{t,w} | q_t, \mathcal{X}_t)}. \quad (20)$$

We optimize a clipped policy-gradient objective:

$$\begin{aligned} \mathcal{L}(\theta) = \mathbb{E} \left[\min \left(\rho_{t,w} \text{ADV}_{t,w}, \right. \right. \\ \left. \left. \text{clip}(\rho_{t,w}, 1 - \epsilon, 1 + \epsilon) \text{ADV}_{t,w} \right) \right] \\ - \beta \cdot \mathbb{E} \left[\text{KL} \left(\pi_{\theta}(\cdot | q_t, \mathcal{X}_t) \parallel \pi_{\theta_{\text{ref}}}(\cdot | q_t, \mathcal{X}_t) \right) \right] \end{aligned} \quad (21)$$

where ϵ is the clipping range, and $\beta \geq 0$ controls KL regularization to the reference policy, and $\mathbb{E}(\cdot)$ is taken over actions $\hat{S}_{t,w}$ sampled from $\pi_{\theta_{\text{old}}}(\cdot | q_t, \mathcal{X}_t)$.

4 Experiments

We conduct experiments on the following questions:

RQ1: How can RES-RAG improve performance in terms of quality-cost trade-off?

RQ2: How does the design of the joint quality-cost reward function in RES-RAG affect the effectiveness of evidence retrieval?

RQ3: Can RES-RAG maintain universality across RAG systems with different document structures and generator configurations?

4.1 Experimental Settings

During the training phase, we run GRPO on two NVIDIA RTX 4090 GPUs, and during the inference phase, we deploy the trained model on one

NVIDIA RTX 4090 GPU. The retrieval model we used is Qwen⁵, a small model with only 4B parameters.

Datasets We conducted experiments with RES-RAG and various baseline models on six open-domain QA datasets, including: (1) two ambiguous datasets, ASQA (Stelmakh et al., 2022) and AbmigQA (Min et al., 2020); (2) three multi-hop QA datasets, Hotpot-QA (Yang et al., 2018), 2wikiMultiHop (Ho et al., 2020), and MusiQue (Trivedi et al., 2022); and (3) an open web QA dataset, Trivia-QA (Joshi et al., 2017).

Evaluation Metrics Our goal is for Retriever to retrieve as few query-relevant evidence chunks as possible from complex web environments. Therefore, we evaluate the short answer output by the generator \mathcal{G} as the final result. For ambiguous datasets, we report both Hit@1 and EM metrics; for other datasets, we report the EM metric. Hit@1 indicates that at least one golden answer perfectly matches the \mathcal{G} 's output.

Implementation Details For each query, we sample $W = 4$ candidate structured actions. We keep the KL coefficient at $\beta = 0.05$. More training parameters can be found in Appendix D.

4.2 Comparisons with Other Methods

We compare RES-RAG against widely used baselines, including lexical retrieval BM25 (Robertson et al., 2009), dense retrievers BGE (Xiao et al., 2024) and E5-Mistral (Wang et al., 2024a), and compression systems LongLLMLingua (Jiang et al., 2024) and JinaAI Reader (JinaAI, 2024). For retrieval baselines (BM25/BGE/E5-Mistral), we retrieve evidence chunks and pass them to the same Generator \mathcal{G} with the same chunk budget as in RES-RAG. For compression systems (LongLLMLingua/JinaAI Reader), we follow their standard pipelines to fit the same maximum prompt constraint. To be fair, we redeploy the above methods, and Table 1 reports detailed indicators.

Specifically, RES-RAG achieves a better quality-cost frontier: it improves answer correctness while using fewer evidence chunks, indicating that the learned retriever policy tends to select minimally sufficient evidence. This behavior is more pronounced on ambiguous data, where the ambiguity structure implies that blindly increasing retrieved

⁵<https://huggingface.co/Qwen/Qwen3-4B-Instruct-2507>

Table 1: Comparison results of the RES-RAG method with other baselines. Optimal and suboptimal results are indicated by bold and underlined, respectively. The symbol † indicates that our method achieved optimality.

Method	ASQA		HotpotQA	TriviaQA	AmbigQA		2Wiki	MuSiQue
	EM	Hit@1	EM	EM	EM	Hit@1	EM	EM
Llama-3.1-8B-Instruct-8k								
BM25	3.25	<u>48.75</u>	<u>23.93</u>	62.29	17.80	46.60	14.00	13.25
BGE	<u>3.50</u>	43.50	20.40	63.91	20.42	50.79	13.50	12.25
E5-Mistral	3.25	43.56	24.43	61.50	25.65	<u>51.83</u>	21.50	15.00
LongLLMLingua	2.50	43.32	15.87	66.25	<u>23.04</u>	46.60	26.25	14.25
JinaAI Reader	3.00	45.25	14.86	58.25	16.75	49.46	35.29	13.46
RES-RAG	4.75 †	52.75 †	23.68	<u>66.06</u>	25.65	57.59 †	<u>35.00</u>	<u>14.79</u>
Llama-3.1-70B-Instruct-8k								
BM25	1.75	51.50	<u>31.49</u>	73.92	23.56	56.54	45.25	18.50
BGE	2.25	57.25	27.96	74.00	22.51	47.12	39.75	16.75
E5-Mistral	<u>3.75</u>	55.50	<u>31.49</u>	78.25	<u>31.41</u>	66.49	49.00	<u>21.50</u>
LongLLMLingua	2.50	52.75	20.65	<u>77.48</u>	26.70	63.35	<u>49.75</u>	18.50
JinaAI Reader	2.25	<u>57.50</u>	18.89	61.25	31.94	<u>70.16</u>	29.34	14.75
RES-RAG	4.00 †	58.25 †	31.74 †	75.25	<u>31.41</u>	72.25 †	56.03 †	22.75 †

evidence often introduces distractors. RES-RAG mitigates this by selecting only the small subset of chunks that directly support at least one gold answer. In multi-hop, where reasoning requires composing evidence across documents, RES-RAG better allocates the budget, reducing redundant near-duplicate chunks. On Trivia-QA, which contains open-web pages with diverse structures and substantial noise, RES-RAG, rather than over-retrieving semantically similar but non-evidential chunks, learns evidence patterns aligned with the Generator’s success signal. Dense retrievers still rely on fixed retrieval rules, and compression systems do not explicitly optimize whether the generator can correctly select supporting discrete evidence. In summary, the results in Table 1 answer **RQ1**: RES-RAG improves the quality-cost trade-off by learning to retrieve compact, highly useful evidence with a joint quality-cost objective.

4.3 Ablation Experiments Analysis

We analyze how the joint quality–cost reward in RES-RAG shapes evidence retrieval behavior. Table 2 reports ablation results. Starting from the full RES-RAG objective, we consider: (1) **top- k -only**; (2) **SFT-only**; (3) **Joint quality–cost reward**.

A fixed top- k strategy cannot adapt evidence size to instance difficulty: increasing k improves recall

but often over-retrieves, while decreasing k under-retrieves on hard instances, yielding an inferior quality–cost frontier. SFT-only provides a strong warm start by imitating heuristic labels \tilde{S} , but it still does not optimize the utility–cost trade-off, so the retriever remains sensitive to heuristic biases and budget violations.

Optimizing utility only encourages over-retrieval since extra chunks are rarely penalized, increasing over-budget rates and weakening cost control. Overly emphasizing cost pushes the policy to retrieve too few chunks, reducing evidence cost but hurting answer correctness. These extremes confirm that both utility and cost are necessary for retrieval strategies. Compared to a linear penalty, the segmented penalty better distinguishes mild vs. severe budget violations, making the learning signal more robust and stabilizing the quality–cost frontier. Removing retrieval-consistency makes the reward sparse with respect to discrete index selection: many rollouts obtain similar answer rewards despite different evidence. The per-document F1 signal against \tilde{S} densifies credit assignment, improving retrieval quality and budgeted selection behavior, which then translates to better downstream answer quality under the same budget.

Overall, the ablations answer **RQ2**: the joint reward design is crucial for learning effective, budget-

RES-RAG	ASQA		HotpotQA	TriviaQA	AmbigQA		2Wiki	MuSiQue
	EM	Hit@1	EM	EM	EM	Hit@1	EM	EM
Llama-3.1-8B-Instruct-8k								
(top- k) <i>w/o</i> SFT	2.56	43.25	18.43	56.25	17.34	48.95	23.75	13.25
(Warm-up) <i>w</i> SFT	3.75	44.50	21.16	61.75	21.47	51.83	27.25	14.75
(GRPO) <i>w</i> Joint reward	4.75	52.75	23.68	66.06	25.65	57.59	35.00	14.79

Table 2: Ablation experiment results.

feasible evidence retrieval.

4.4 Generality Experiment Studies

We test whether RES-RAG maintains effectiveness across (1) document structures with different noise profiles and layouts, and (2) generator configurations with different decoding behaviors and context capacities.

We compare RES-RAG on datasets with different sources: Wikipedia-style versus open-web pages. Across these variations, RES-RAG maintains an advantage in quality–cost trade-off, suggesting that the learned policy is not overfitting to a specific surface structure but instead captures transferable evidence selection behaviors. Although RES-RAG optimizes the Retriever, the reward is induced by the Generator \mathcal{G} , raising the question of whether the learned retriever is tightly coupled to a particular generator. We therefore swap \mathcal{G} across different configurations while keeping the retriever fixed, and re-evaluate under the same chunk budget (Table 1). We observe that the benefits of RES-RAG persist across generator choices: absolute scores vary with generator capacity, indicating that the retriever policy learns evidence selection strategies that are broadly useful rather than exploiting idiosyncrasies of a single generator. The results in Table 1 answer **RQ3**: RES-RAG exhibits strong universality across both document structures and generator configurations, while retaining the core benefit of budget-aware, minimally sufficient evidence retrieval.

5 Conclusions

This paper presents RES-RAG, a budget-aware evidence retrieval framework for complex multi-document RAG scenarios. By formulating retrieval as a structured action that outputs per-document chunk indices, RES-RAG turns combinatorial evidence selection into a controllable decision problem that can be optimized for the retriever. We

Dataset	Method	EvidTok	
		Avg	Max
ASQA	Top- k	318k	7259k
	RES-RAG	1k	7k
HotpoQA	Top- k	223k	5138k
	RES-RAG	2k	6k
TriviaQA	Top- k	53k	504k
	RES-RAG	2k	20k
AmbigQA	Top- k	288k	15631k
	RES-RAG	1k	7k
2Wiki	Top- k	11k	247k
	RES-RAG	1k	8k
MuSiQue	Top- k	78k	2359k
	RES-RAG	1k	10k

Table 3: Quality–cost comparison under identical settings. EvidTok counts retrieved evidence tokens. Avg represents the average number of tokens, and Max represents the maximum.

further introduce a joint quality–cost reward that unifies answer utility and evidence cost into a single scalar signal, enabling reinforcement learning to converge toward minimally sufficient evidence under a fixed context budget. Experiments across multiple open-domain QA benchmarks demonstrate that RES-RAG achieves a better trade-off between answer quality and evidence cost, and that indexed, traceable evidence selection reduces long-context noise while improving evidence density for downstream generation.

Limitations

The optimization signal relies on short-answer F1 and related anomaly penalties. This may not faithfully capture correctness for questions requiring richer reasoning, multi-hop justification, or long-form generation, and may encourage reward-

571	specific behaviors. In addition, although structured indexing improves controllability, the approach still depends on valid, parseable outputs. In other domains or under adversarial noise, format-compliance failures may increase or decrease robustness.	
572		
573		
574		
575		
576		
577	References	
578	Kosuke Akimoto, Kunihiro Takeoka, and Masafumi Oyamada. 2024. Context quality matters in training fusion-in-decoder for extractive open-domain question answering. <i>arXiv preprint arXiv:2403.14197</i> .	
579		
580		
581		
582	Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2024. Self-rag: Learning to retrieve, generate, and critique through self-reflection. <i>ICLR</i> .	
583		
584		
585		
586	Scott Barnett, Stefanus Kurniawan, Srikanth Thudumu, Zach Brannelly, and Mohamed Abdelrazek. 2024. Seven failure points when engineering a retrieval augmented generation system. In <i>Proceedings of the IEEE/ACM 3rd International Conference on AI Engineering-Software Engineering for AI</i> , pages 194–199.	
587		
588		
589		
590		
591		
592		
593	Patrice Béchard and Orlando Marquez Ayala. 2024. Reducing hallucination in structured outputs via retrieval-augmented generation. <i>arXiv preprint arXiv:2404.08189</i> .	
594		
595		
596		
597	Maciej Besta, Ales Kubicek, Robert Gerstenberger, Marcin Chrapek, Roman Niggl, Patrik Okanovic, Yi Zhu, Patrick Iff, Michal Podstawski, Lucas Weitzendorf, and 1 others. 2024. Multi-head rag: Solving multi-aspect problems with llms. <i>arXiv preprint arXiv:2406.05085</i> .	
598		
599		
600		
601		
602		
603	Luca Beurer-Kellner, Marc Fischer, and Martin Vechev. 2024. Guiding llms the right way: Fast, non-invasive constrained generation. <i>arXiv preprint arXiv:2403.06988</i> .	
604		
605		
606		
607	Wenhu Chen, Hanwen Zha, Zhiyu Chen, Wenhan Xiong, Hong Wang, and William Wang. 2020. Hybridqa: A dataset of multi-hop question answering over tabular and textual data. <i>arXiv preprint arXiv:2004.07347</i> .	
608		
609		
610		
611	Philipp Christmann, Svitlana Vakulenko, Ionut Teodor Sorodoc, Bill Byrne, and Adrià de Gispert. 2024. Retrieving contextual information for long-form question answering using weak supervision. <i>arXiv preprint arXiv:2410.08623</i> .	
612		
613		
614		
615		
616	Yuqin Dai, Guoqing Wang, Yuan Wang, Kairan Dou, Kaichen Zhou, Zhanwei Zhang, Shuo Yang, Fei Tang, Jun Yin, Pengyu Zeng, and 1 others. 2025. Evinote-rag: Enhancing rag models via answer-supportive evidence notes. <i>arXiv preprint arXiv:2509.00877</i> .	
617		
618		
619		
620		
621	Michiel de Jong, Yury Zemlyanskiy, Joshua Ainslie, Nicholas FitzGerald, Sumit Sanghai, Fei Sha, and William Cohen. 2022. Fido: Fusion-in-decoder optimized for stronger performance and faster inference. <i>arXiv preprint arXiv:2212.08153</i> .	623 624 625
622		
	Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt, Dasha Metropolitan, Robert Osazuwa Ness, and Jonathan Larson. 2024. From local to global: A graph rag approach to query-focused summarization. <i>arXiv preprint arXiv:2404.16130</i> .	626 627 628 629 630 631
	Thibault Formal, Benjamin Piwowarski, and Stéphane Clinchant. 2021. Splade: Sparse lexical and expansion model for first stage ranking. In <i>Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , pages 2288–2292.	632 633 634 635 636 637
	Ge Gao, Jonathan D Chang, Claire Cardie, Kianté Brantley, and Thorsten Joachim. 2023a. Policy-gradient training of language models for ranking. <i>arXiv preprint arXiv:2310.04407</i> .	638 639 640 641
	Luyu Gao, Xueguang Ma, Jimmy Lin, and Jamie Callan. 2023b. Precise zero-shot dense retrieval without relevance labels. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1762–1777.	642 643 644 645 646
	Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yixin Dai, Jiawei Sun, Haofen Wang, and Haofen Wang. 2023c. Retrieval-augmented generation for large language models: A survey. <i>arXiv preprint arXiv:2312.10997</i> , 2(1).	647 648 649 650 651
	Michael Glass, Gaetano Rossiello, Md Faisal Mahbub Chowdhury, Ankita Rajaram Naik, Pengshan Cai, and Alfio Gliozzo. 2022. Re2g: Retrieve, rerank, generate. <i>arXiv preprint arXiv:2207.06300</i> .	652 653 654 655
	Michael Günther, Isabelle Mohr, Daniel James Williams, Bo Wang, and Han Xiao. 2024. Late chunking: contextual chunk embeddings using long-context embedding models. <i>arXiv preprint arXiv:2409.04701</i> .	656 657 658 659
	Yucan Guo, Miao Su, Saiping Guan, Zihao Sun, Xiaolong Jin, Jiafeng Guo, and Xueqi Cheng. 2025. Routerag: Efficient retrieval-augmented generation from text and graph via reinforcement learning. <i>arXiv preprint arXiv:2512.09487</i> .	660 661 662 663 664
	Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. 2020. Retrieval augmented language model pre-training. In <i>International conference on machine learning</i> , pages 3929–3938. PMLR.	665 666 667 668
	Zijian Hei, Weiling Liu, Wenjie Ou, Juyi Qiao, Junming Jiao, Guowen Song, Ting Tian, and Yi Lin. 2024. Dr-rag: Applying dynamic document relevance to retrieval-augmented generation for question-answering. <i>arXiv preprint arXiv:2406.07348</i> .	669 670 671 672 673
	Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. 2020. Constructing a multi-hop qa dataset for comprehensive evaluation of reasoning	674 675 676

677	steps. In <i>Proceedings of the 28th International Conference on Computational Linguistics</i> , pages 6609–6625.	Kiseung Kim and Jay-Yoon Lee. 2024. Re-rag: Improving open-domain qa performance and interpretability with relevance estimator in retrieval-augmented generation. <i>arXiv preprint arXiv:2406.05794</i> .	733
678			734
679			735
680	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. <i>ACM Transactions on Information Systems</i> , 43(2):1–55.		736
681		Kalpesh Krishna, Aurko Roy, and Mohit Iyyer. 2021. Hurdles to progress in long-form question answering. <i>arXiv preprint arXiv:2103.06332</i> .	737
682			738
683			739
684		Mandar Kulkarni, Praveen Tangarajan, Kyung Kim, and Anusua Trivedi. 2024. Reinforcement learning for optimizing rag for domain chatbots. <i>arXiv preprint arXiv:2401.06800</i> .	740
685			741
686			742
687	Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2021. Unsupervised dense information retrieval with contrastive learning. <i>arXiv preprint arXiv:2112.09118</i> .		743
688			744
689		Moshe Lahmy and Roi Yozevitch. 2025. Replace, don't expand: Mitigating context dilution in multi-hop rag via fixed-budget evidence assembly. <i>arXiv preprint arXiv:2512.10787</i> .	745
690			746
691			747
692	Gautier Izacard and Edouard Grave. 2020. Leveraging passage retrieval with generative models for open domain question answering. <i>arXiv preprint arXiv:2007.01282</i> .		748
693		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. <i>Advances in neural information processing systems</i> , 33:9459–9474.	749
694			750
695			751
696	Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. 2023. Atlas: Few-shot learning with retrieval augmented language models. <i>Journal of Machine Learning Research</i> , 24(251):1–43.		752
697			753
698			754
699			755
700		Junlong Li, Yiheng Xu, Lei Cui, and Furu Wei. 2021. Markuplm: Pre-training of text and markup language for visually-rich document understanding. <i>arXiv preprint arXiv:2110.08518</i> .	756
701			757
702	Abhinav Java, Srivathsan Koundinyan, Nagarajan Natarajan, and Amit Sharma. 2025. Frugalrag: Learning to retrieve and reason for multi-hop qa. <i>arXiv preprint arXiv:2507.07634</i> .		758
703			759
704			760
705			761
706	Huiqiang Jiang, Qianhui Wu, Xufang Luo, Dongsheng Li, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. 2024. Longllmlingua: Accelerating and enhancing llms in long context scenarios via prompt compression. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1658–1677.		762
707			763
708			764
709			765
710			766
711			767
712			768
713	JinaAI. 2024. Reader-lm: Small language models for cleaning and converting HTML to Markdown . [Online; accessed 2024-10-05].		769
714			770
715			771
716	Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. 2017. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1601–1611.		772
717			773
718			774
719			775
720			776
721			777
722	Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick SH Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for open-domain question answering. In <i>EMNLP (1)</i> , pages 6769–6781.		778
723			779
724			780
725			781
726			782
727	Omar Khattab, Keshav Santhanam, Xiang Lisa Li, David Hall, Percy Liang, Christopher Potts, and Matei Zaharia. 2022. Demonstrate-search-predict: Composing retrieval and language models for knowledge-intensive nlp. <i>arXiv preprint arXiv:2212.14024</i> .		783
728			784
729			785
730			786
731			787
732			788

786	Sewon Min, Julian Michael, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2020. Ambigqa: Answering ambiguous open-domain questions. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 5783–5797.	with selection in rag for sensitive domains. <i>arXiv preprint arXiv:2505.16014</i> .	841 842
792	Matin Mortaheb, Mohammad A Amir Khojastepour, Srimat T Chakradhar, and Sennur Ulukus. 2025. Re-ranking the context for multimodal retrieval augmented generation. <i>arXiv preprint arXiv:2501.04695</i> .	Torsten Scholak, Nathan Schucher, and Dzmitry Bahdanau. 2021. Picard: Parsing incrementally for constrained auto-regressive decoding from language models. <i>arXiv preprint arXiv:2109.05093</i> .	843 844 845 846
797	Jananee Muralidharan and Tiju Thomas. 2024. Deliberate problem-solving with a large language model as a brainstorm aid using a checklist for prompt generation. <i>The Journal of the Association of Physicians of India</i> , 72(5):89–90.	Xiaoyu Shen, Svitlana Vakulenko, Marco Del Tredici, Gianni Barlacchi, Bill Byrne, and Adrià de Gispert. 2022. Low-resource dense retrieval for open-domain question answering: A comprehensive survey. <i>arXiv preprint arXiv:2208.03197</i> .	847 848 849 850 851
801	Jianmo Ni, Chen Qu, Jing Lu, Zhuyun Dai, Gustavo Hernandez Abrego, Ji Ma, Vincent Y Zhao, Yi Luan, Keith B Hall, Ming-Wei Chang, and 1 others. 2021. Large dual encoders are generalizable retrievers. <i>arXiv preprint arXiv:2112.07899</i> .	Boheng Sheng, Jiacheng Yao, Meicong Zhang, and Guoxiu He. 2025. Dynamic chunking and selection for reading comprehension of ultra-long context in large language models. <i>arXiv preprint arXiv:2506.00773</i> .	852 853 854 855 856
802	Kanghee Park, Jiayu Wang, Taylor Berg-Kirkpatrick, Nadia Polikarpova, and Loris D’Antoni. 2024. Grammar-aligned decoding. <i>Advances in Neural Information Processing Systems</i> , 37:24547–24568.	Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Rich James, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. 2023. Replug: Retrieval-augmented black-box language models. <i>arXiv preprint arXiv:2301.12652</i> .	857 858 859 860 861
803	Ronak Pradeep, Xueguang Ma, Xinyu Zhang, Hang Cui, Ruizhou Xu, Rodrigo Nogueira, and Jimmy Lin. 2020. H2oloo at trec 2020: When all you got is a hammer... deep learning, health misinformation, and precision medicine. <i>Corpus</i> , 5(d3):d2.	Kurt Shuster, Spencer Poff, Moya Chen, Douwe Kiela, and Jason Weston. 2021. Retrieval augmentation reduces hallucination in conversation. <i>arXiv preprint arXiv:2104.07567</i> .	862 863 864 865
804	Ronak Pradeep, Sahel Sharifymoghaddam, and Jimmy Lin. 2023. Rankzephyr: Effective and robust zero-shot listwise reranking is a breeze! <i>arXiv preprint arXiv:2312.02724</i> .	Shamane Siriwardhana, Rivindu Weerasekera, Elliott Wen, Tharindu Kaluarachchi, Rajib Rana, and Suranga Nanayakkara. 2023. Improving the domain adaptation of retrieval augmented generation (rag) models for open domain question answering. <i>Transactions of the Association for Computational Linguistics</i> , 11:1–17.	866 867 868 869 870 871 872
805	Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A Smith, and Mike Lewis. 2022. Measuring and narrowing the compositionality gap in language models. <i>arXiv preprint arXiv:2210.03350</i> .	Ivan Stelmakh, Yi Luan, Bhuwan Dhingra, and Ming-Wei Chang. 2022. Asqa: Factoid questions meet long-form answers. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 8273–8288.	873 874 875 876 877
806	Ruiyang Ren, Yingqi Qu, Jing Liu, Wayne Xin Zhao, Qiaoqiao She, Hua Wu, Haifeng Wang, and Ji-Rong Wen. 2021. Rocketqav2: A joint training method for dense passage retrieval and passage re-ranking. <i>arXiv preprint arXiv:2110.07367</i> .	Dan Su, Xiaoguang Li, Jindi Zhang, Lifeng Shang, Xin Jiang, Qun Liu, and Pascale Fung. 2022a. Read before generate! faithful long form question answering with machine reading. <i>arXiv preprint arXiv:2203.00343</i> .	878 879 880 881 882
807	Stephen Robertson, Hugo Zaragoza, and 1 others. 2009. The probabilistic relevance framework: Bm25 and beyond. <i>Foundations and Trends® in Information Retrieval</i> , 3(4):333–389.	Hongjin Su, Jungo Kasai, Chen Henry Wu, Weijia Shi, Tianlu Wang, Jiayi Xin, Rui Zhang, Mari Ostendorf, Luke Zettlemoyer, Noah A Smith, and 1 others. 2022b. Selective annotation makes language models better few-shot learners. <i>arXiv preprint arXiv:2209.01975</i> .	883 884 885 886 887 888
808	Keshav Santhanam, Omar Khattab, Jon Saad-Falcon, Christopher Potts, and Matei Zaharia. 2021. Colbertv2: Effective and efficient retrieval via lightweight late interaction. <i>arXiv preprint arXiv:2112.01488</i> .	Chihiro Taguchi, Seiji Maekawa, and Nikita Bhutani. 2025. Efficient context selection for long-context qa: No tuning, no iteration, just adaptive- <i>k</i> . <i>arXiv preprint arXiv:2506.08479</i> .	889 890 891 892
809	Yash Saxena, Ankur Padia, Mandar S Chaudhary, Kalpa Gunaratna, Srinivasan Parthasarathy, and Manas Gaur. 2025. Ranking free rag: Replacing re-ranking	Xiaqiang Tang, Qiang Gao, Jian Li, Nan Du, Qi Li, and Sihong Xie. 2025. Mba-rag: a bandit approach for adaptive retrieval-augmented generation through	893 894 895

896	question complexity. In <i>Proceedings of the 31st International Conference on Computational Linguistics</i> , pages 3248–3254.	952
897		953
898		954
899	SMTI Tonmoy, SM Zaman, Vinija Jain, Anku Rani, Vipula Rawte, Aman Chadha, and Amitava Das. 2024. A comprehensive survey of hallucination mitigation techniques in large language models. <i>arXiv preprint arXiv:2401.01313</i> , 6.	955
900		956
901		957
902		
903		
904	Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022. Musique: Multi-hop questions via single-hop question composition. <i>Transactions of the Association for Computational Linguistics</i> , 10:539–554.	958
905		959
906		960
907		961
908		962
909	Tu Vu, Mohit Iyyer, Xuezhi Wang, Noah Constant, Jerry Wei, Jason Wei, Chris Tar, Yun-Hsuan Sung, Denny Zhou, Quoc Le, and 1 others. 2023. Freshllms: Refreshing large language models with search engine augmentation. <i>arXiv preprint arXiv:2310.03214</i> .	963
910		964
911		965
912		966
913		967
914	Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang, Rangan Majumder, and Furu Wei. 2024a. Improving text embeddings with large language models. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 11897–11916.	968
915		969
916		970
917		971
918		972
919		973
920	Xiaohua Wang, Zhenghua Wang, Xuan Gao, Feiran Zhang, Yixin Wu, Zhibo Xu, Tianyuan Shi, Zhengyuan Wang, Shizheng Li, Qi Qian, and 1 others. 2024b. Searching for best practices in retrieval-augmented generation. <i>arXiv preprint arXiv:2407.01219</i> .	974
921		975
922		976
923		977
924		978
925		979
926	Zheng Wang, Shu Teo, Jieer Ouyang, Yongjun Xu, and Wei Shi. 2024c. M-rag: Reinforcing large language model performance through retrieval-augmented generation with multiple partitions. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1966–1978.	980
927		981
928		982
929		983
930		984
931		985
932		986
933	Zhiruo Wang, Jun Araki, Zhengbao Jiang, Md Rizwan Parvez, and Graham Neubig. 2023. Learning to filter context for retrieval-augmented generation. <i>arXiv preprint arXiv:2311.08377</i> .	987
934		988
935		989
936		990
937	Zilong Wang, Zifeng Wang, Long Le, Huaixiu Steven Zheng, Swaroop Mishra, Vincent Perot, Yuwei Zhang, Anush Mattapalli, Ankur Taly, Jingbo Shang, and 1 others. 2024d. Speculative rag: Enhancing retrieval augmented generation through drafting. <i>arXiv preprint arXiv:2407.08223</i> .	991
938		992
939		993
940		994
941		995
942		996
943	Ji-Rong Wen, Ni Lao, and Wei-Ying Ma. 2004. Probabilistic model for contextual retrieval. In <i>Proceedings of the 27th annual international ACM SIGIR conference on Research and development in information retrieval</i> , pages 57–63.	997
944		998
945		999
946		1000
947		1001
948	Shitao Xiao, Zheng Liu, Yingxia Shao, and Zhao Cao. 2022a. Retromae: Pre-training retrieval-oriented language models via masked auto-encoder. <i>arXiv preprint arXiv:2205.12035</i> .	1002
949		1003
950		1004
951		1005
		1006
		1007
		1008
		1009
		1010
		1011
		1012
		1013
		1014
		1015
		1016
		1017
		1018
		1019
		1020
		1021
		1022
		1023
		1024
		1025
		1026
		1027
		1028
		1029
		1030
		1031
		1032
		1033
		1034
		1035
		1036
		1037
		1038
		1039
		1040
		1041
		1042
		1043
		1044
		1045
		1046
		1047
		1048
		1049
		1050
		1051
		1052
		1053
		1054
		1055
		1056
		1057
		1058
		1059
		1060
		1061
		1062
		1063
		1064
		1065
		1066
		1067
		1068
		1069
		1070
		1071
		1072
		1073
		1074
		1075
		1076
		1077
		1078
		1079
		1080
		1081
		1082
		1083
		1084
		1085
		1086
		1087
		1088
		1089
		1090
		1091
		1092
		1093
		1094
		1095
		1096
		1097
		1098
		1099
		1100
		1101
		1102
		1103
		1104
		1105
		1106
		1107
		1108
		1109
		1110
		1111
		1112
		1113
		1114
		1115
		1116
		1117
		1118
		1119
		1120
		1121
		1122
		1123
		1124
		1125
		1126
		1127
		1128
		1129
		1130
		1131
		1132
		1133
		1134
		1135
		1136
		1137
		1138
		1139
		1140
		1141
		1142
		1143
		1144
		1145
		1146
		1147
		1148
		1149
		1150
		1151
		1152
		1153
		1154
		1155
		1156
		1157
		1158
		1159
		1160
		1161
		1162
		1163
		1164
		1165
		1166
		1167
		1168
		1169
		1170
		1171
		1172
		1173
		1174
		1175
		1176
		1177
		1178
		1179
		1180
		1181
		1182
		1183
		1184
		1185
		1186
		1187
		1188
		1189
		1190
		1191
		1192
		1193
		1194
		1195
		1196
		1197
		1198
		1199
		1200
		1201
		1202
		1203
		1204
		1205
		1206
		1207
		1208
		1209
		1210
		1211
		1212
		1213
		1214
		1215
		1216
		1217
		1218
		1219
		1220
		1221
		1222
		1223
		1224
		1225
		1226
		1227
		1228
		1229
		1230
		1231
		1232
		1233
		1234
		1235
		1236
		1237
		1238
		1239
		1240
		1241
		1242
		1243
		1244
		1245
		1246
		1247
		1248
		1249
		1250
		1251
		1252
		1253
		1254
		1255
		1256
		1257
		1258
		1259
		1260
		1261
		1262
		1263
		1264
		1265
		1266
		1267
		1268
		1269
		1270
		1271
		1272
		1273
		1274
		1275
		1276
		1277
		1278
		1279
		1280
		1281
		1282
		1283
		1284
		1285
		1286
		1287
		1288
		1289
		1290
		1291
		1292
		1293
		1294
		1295
		1296
		1297
		1298
		1299
		1300
		1301
		1302
		1303
		1304
		1305
		1306
		1307
		1308
		1309
		1310
		1311
		1312
		1313
		1314
		1315
		1316
		1317
		1318
		1319
		1320
		1321
		1322
		1323
		1324
		1325
		1326
		1327
		1328
		1329
		1330
		1331
		1332
		1333
		1334
		1335
		1336
		1337
		1338
		1339
		1340
		1341
		1342
		1343
		1344
		1345
		1346
		1347
		1348
		1349
		1350
		1351
		1352
		1353
		1354
		1355
		1356
		1357
		1358
		1359
		1360
		1361
		1362
		1363
		1364
		1365
		1366
		1367
		1368
		1369

1009 Le Zhang, Bo Wang, Xipeng Qiu, Siva Reddy, and
 1010 Aishwarya Agrawal. 2025. Rearank: Reasoning
 1011 re-ranking agent via reinforcement learning. *arXiv*
 1012 *preprint arXiv:2505.20046*.

1013 Zhebin Zhang, Xinyu Zhang, Yuanhang Ren, Saijiang
 1014 Shi, Meng Han, Yongkang Wu, Ruofei Lai, and Zhao
 1015 Cao. 2023. Iag: Induction-augmented generation
 1016 framework for answering reasoning questions. *arXiv*
 1017 *preprint arXiv:2311.18397*.

1018 Yuanhang Zheng, Peng Li, Wei Liu, Yang Liu, Jian
 1019 Luan, and Bin Wang. 2024. Toolrerank: Adap-
 1020 tive and hierarchy-aware reranking for tool retrieval.
 1021 *arXiv preprint arXiv:2403.06551*.

1022 Siyuan Zhu, Chengdong Xu, Kaiqiang Ke, and Chao
 1023 Yu. 2025. Context-picker: Dynamic context selec-
 1024 tion using multi-stage reinforcement learning. *arXiv*
 1025 *preprint arXiv:2512.14465*.

1026 A Chunking Process

1027 We employ a structure-aware chunking pipeline
 1028 to transform each HTML webpage into an ordered
 1029 list of chunks for retrieval. Our design
 1030 goal is two-fold: (1) preserve evidence-bearing
 1031 structures (e.g., headings, tables, and lists) af-
 1032 ter conversion to Markdown; and (2) enforce a
 1033 strict length constraint so that each chunk can
 1034 be consumed reliably by the Retriever. We
 1035 use MARKITDOWN to convert HTML to Mark-
 1036 down, preserving document structure (headings,
 1037 lists, and tables) while removing many HTML-
 1038 specific artifacts. For each document d_m , we out-
 1039 put an ordered chunk sequence $\{c_i\}_{i=1}^{N_d}$, where
 1040 each chunk stores: $\text{type} \in \{\text{text}, \text{table}, \text{list}\}$,
 1041 a heading_path (e.g., H1→H2→H3), the chunk
 1042 token_len measured by the Retriever tokenizer,
 1043 and an offset indicating its position in the origi-
 1044 nal document.

1045 We first remove low-signal HTML elements,
 1046 normalize whitespace, and drop hyperlink targets
 1047 while keeping anchor text. The cleaned HTML
 1048 is then converted to Markdown using MARKIT-
 1049 DOWN. We treat each table/list block in Markdown
 1050 as an atomic evidence unit and extract it into a stan-
 1051 dalone chunk. For the remaining narrative text,
 1052 we chunk by heading boundaries: encountering a
 1053 heading closes the current text chunk and starts
 1054 a new one under the updated heading path. We
 1055 prepend the heading path to the chunk content to
 1056 retain discourse context.

1057 We set a maximum chunk length $T_{\max} = 20480$
 1058 tokens (Retriever tokenizer). If a text chunk ex-
 1059 ceeds T_{\max} , we split it into multiple sub-chunks
 1060 at sentence (or paragraph) boundaries to maintain

Algorithm 1 Structure-aware Chunking

Require: HTML webpage h , tokenizer τ , max to-
 kens T_{\max}

Ensure: Ordered chunks $\{c_i\}_{i=1}^N$

```

1:  $h^{\text{clean}} \leftarrow \text{CleanHTML}(h)$ 
2:  $m \leftarrow \text{MarkItDown}(h^{\text{clean}})$ 
3:  $\mathcal{B} \leftarrow \text{ParseBlocks}(m)$   $\triangleright$  linear blocks with
   type, heading_path, and offset
4:  $\mathcal{C} \leftarrow [], \text{buf} \leftarrow \emptyset, p \leftarrow \emptyset, \text{off} \leftarrow 0$ 
5: for all  $b \in \mathcal{B}$  in order do
6:   if  $b.\text{type} = \text{heading}$  then
7:     if  $\text{buf} \neq \emptyset$  then
8:        $\mathcal{C}$   $\leftarrow$ 
9:        $\mathcal{C} \parallel \text{MakeChunk}(\text{buf}, \text{text}, p, \text{off})$ 
10:       $\text{buf} \leftarrow \emptyset$ 
11:     end if
12:     $p \leftarrow b.\text{heading\_path}; \text{off} \leftarrow$ 
13:     $b.\text{offset}$ 
14:    else if  $b.\text{type} \in \{\text{table}, \text{list}\}$  then
15:       $\mathcal{C}$   $\leftarrow$ 
16:       $\mathcal{C} \parallel \text{MakeChunk}(b.\text{content}, b.\text{type}, p, b.\text{offset})$ 
17:    else
18:       $\text{buf} \leftarrow \text{buf} \parallel b.\text{content}$ 
19:    end if
20:  end for
21: if  $\text{buf} \neq \emptyset$  then
22:    $\mathcal{C} \leftarrow \mathcal{C} \parallel \text{MakeChunk}(\text{buf}, \text{text}, p, \text{off})$ 
23: end if
24:  $\mathcal{C} \leftarrow \text{SplitOversized}(\mathcal{C}, \tau, T_{\max})$ 
25:  $\mathcal{C} \leftarrow \text{Deduplicate}(\mathcal{C})$ 
26:  $\mathcal{C} \leftarrow \text{SortByOffset}(\mathcal{C})$ 
27:  $\{c_i\}_{i=1}^N \leftarrow \text{ReindexFromOne}(\mathcal{C})$ 
28: return  $\{c_i\}_{i=1}^N$ 

```

1061 readability. If a single table/list chunk exceeds
 1062 T_{\max} , we further split it by rows (for tables) or by
 1063 top-level items (for lists).

1064 Finally, we merge table/list chunks and text
 1065 chunks, sort them by offset, remove duplicates
 1066 by normalized hashing, and reindex them to obtain
 1067 $\{c_i\}_{i=1}^{N_d}$.

1068 B Decoding Constraints

1069 To prevent format violations and ensure that the
 1070 retriever policy $\pi_{\theta_{\text{ref}}}$ always outputs a *parsable*
 1071 JSON string, we apply guided decoding with a de-
 1072 coding constraint over the output space. The target
 1073 output follows a fixed schema: a single JSON ob-
 1074 ject that maps each document id $i \in \{1, \dots, m\}$
 1075 to an array of chunk indices. We performance

the following syntactic constraints during generation: (1) the JSON object contains *exactly* m document fields; (2) each field value is an array of integers (chunk indices); (3) each index is 1-based and bounded by the corresponding document length, i.e., $s \in \{1, \dots, M_i\}$ for document i with M_i chunks. With these constraints, every decoded action is valid by construction and can be directly parsed and evaluated by the downstream reward computation.

C Prompts

Retriever Prompt. The retriever is the *learned policy* optimized in RES-RAG. Given a query q and a set of chunked documents $\mathcal{X} = \{d_m\}_{m=1}^M$, it is instructed to select a *minimally sufficient* subset of evidence under a chunk budget, and to output the selection in a *structured* format. Specifically, the prompt uses a JSON object that maps each document to a list of 1-indexed chunk IDs (cf. decoding constraints in Appendix B), enabling deterministic parsing and budget accounting. The prompt template is shown in Figure 5.

Generator Prompt. The generator is a frozen LLM accessed via an online API. It is used only for reward computation: conditioned on the query q and the evidence selected by the retriever (and optional document metadata), it produces a short answer \hat{a} , which is then scored against references to obtain the answer-quality signal (F1) used in our joint reward. No parameters of the generator are updated during training. The prompt template is shown in Figure 6.

D Training parameters

We optimize the retriever with GRPO implemented in TRL. To stabilize policy updates, we regularize the retriever against a frozen reference policy via a KL penalty with coefficient β . During training, rollouts are sampled with temperature and nucleus sampling (top- p). Table 4 summarizes the key hyperparameters, grouped by context budget, optimization, sampling, and reward shaping.

E Inferring Example

During inference, our trained model π_θ receives a query and multiple document chunks, and outputs the indices of the chunks associated with the query in each document. We map these indexes back to the text form of the chunks, then feed the query

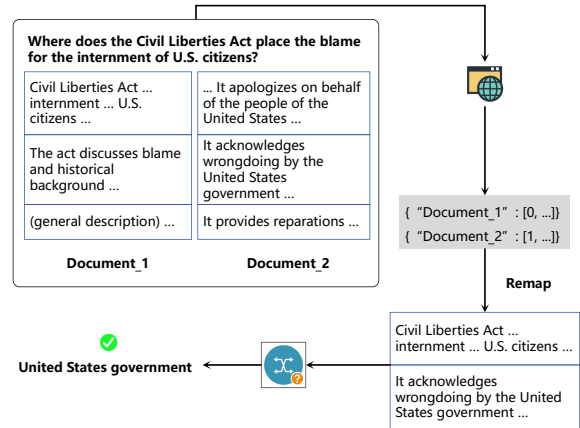


Figure 3: Inference example.

along with these texts and each document’s JSON-LD into the generator, producing an answer. For consistency, we use the example in Figure 1, and this process is illustrated in Figure 3.

F Cost Study

We report the retrieval cost of RES-RAG and a similarity-threshold heuristic in terms of (1) the number of selected chunks and (2) the resulting token budget of the retrieved evidence. Figure 4 summarizes the average and maximum cost per query on six QA benchmarks.

Overall, the similarity-threshold baseline tends to over-retrieve, especially on datasets with ambiguous or multi-hop questions. For example, on ASQA, Ambig, and HotpotQA, it selects substantially more chunks on average (e.g., 22.93 vs. 4.60 on ASQA and 20.82 vs. 4.32 on Ambig), which further translates into a much higher token budget (e.g., 318k vs. 7.45k average tokens on ASQA and 288k vs. 6.86k on Ambig). More importantly, the heuristic exhibits severe tail cost: the maximum number of retrieved tokens can reach millions across several datasets, indicating a high risk of extreme-context inputs.

In contrast, RES-RAG consistently controls both the mean and the tail of retrieval cost. This behavior is expected because the policy is optimized under an explicit quality–cost objective, which encourages early stopping and avoids including long or redundant chunks unless they are necessary. Notably, on Musique, the chunk-count gap is moderate (8.71 vs. 4.92), while the token gap remains large (78k vs. 7.30k on average), suggesting that token cost is dominated not only by how many chunks are selected but also by how long the se-

Table 4: Key hyperparameters used in RES-RAG retriever training.

Name	Explanation	Value
chunk-tokens	Maximum token budget for each chunk	8192
max-prompt-length	Maximum input length to the retriever	8192
max-completion-length	Maximum generation length of the retriever action	1024
batch-size	Per-device batch size	1
gradient-accumulation-steps	Accumulation steps per update	4
β	KL regularization coefficient	0.05
temperature	Sampling temperature for training rollouts	1.0
top- p	Nucleus sampling threshold for training rollouts	0.8
W	Number of rollouts per instance	4
lr	Initial learning rate	1e-5
α_{ans}	Weight of the F1 in the joint reward	2.0
α_{ret}	Weight of the retrieval-consistency term	1.0
R_{min}	Lower bound for reward clipping	-5.0
R_{max}	Upper bound for reward clipping	5.0
η	Penalty coefficient for budget violation / over-retrieval	2.0

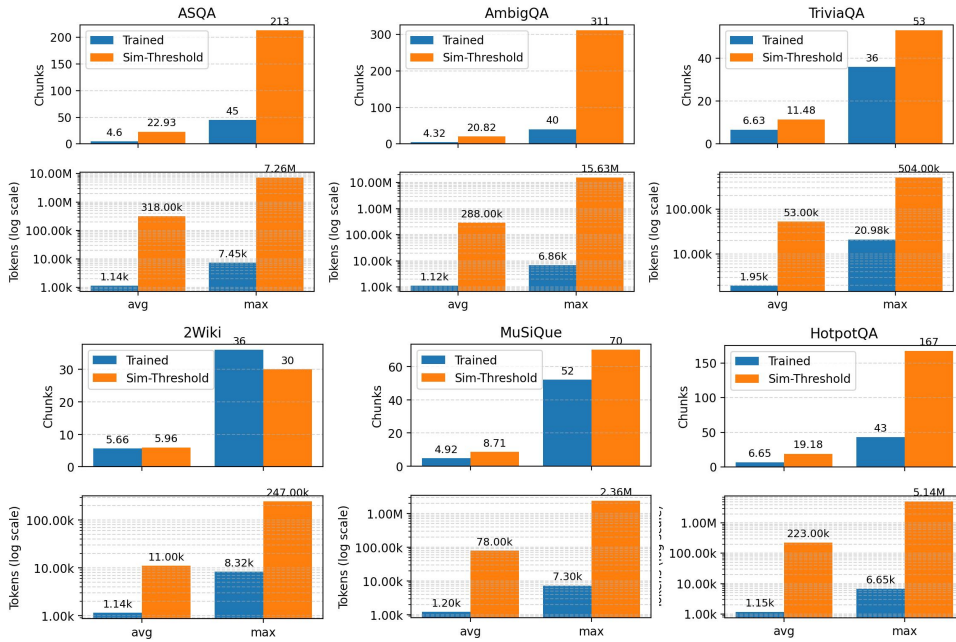


Figure 4: Comparison of costs between RES-RAG and top- k thresholding methods.

lected chunks are. On 2Wiki and TriviaQA, where evidence is typically more localized, the similarity-threshold baseline is less prone to over-retrieval, hence the average-cost gap is smaller; nevertheless, RES-RAG still reduces the risk of high-cost outliers.

These results highlight a practical advantage of RES-RAG: it provides a deployable mechanism to stabilize retrieval cost, reducing latency and the likelihood of exceeding the generator context window, while preserving answer quality.

system:

You are a Retrieval expert. Your sole responsibility is to receive a Query and Chunks from multiple documents, and then rigorously select the indexes(Integer) for the Chunks that are most relevant to the Query.

Input Instructions:

You will receive a Query string and a dictionary of Chunks from multiple documents, where the Chunks are in the following form:

```
{  
  "Document_i": [[chunk_0], ... [chunk_n]]  
}
```

Selection and Sorting Rules:

1. A Query typically corresponds to multiple documents, so you can select multiple Chunks from each document and output their indices (integers);
2. You need to carefully consider which Chunks in each document might contain the answer to the current Query. The Chunks you select will be input along with the Query to the downstream reader to generate the answer;

Output Format:

Regardless, you must adhere to the following output rules:

1. Your output must be a parsable JSON object;
2. Do not output any `` tags;
3. Do not output any introductory words;
4. Your output must be the indexes(Integer) of the selected Chunks in each document;
5. The Chunks you select must be sorted in descending order of similarity to the Query.

Output Example:

```
{  
  "Document_0": [13, 2, 45, 7],  
  "Document_1": [5, 17, 0, 3, 8, 6, 33],  
}
```

user:

Please select Chunks from each document that are relevant to the current Query.

The Query is: {question}

The document Chunks is as follows:

```
<JSON>  
{html}  
</JSON>
```

Figure 5: System and user prompt for the retriever model.

system:

You are a Reader expert. Your only task is to answer the given questions based only on the provided context paragraphs; output a list of short answers.

For each Query, you will receive multiple documents associated with it. Each document consists of multiple Chunks, all in Markdown format.

Answering Guidelines

1. Answer based on the context only; do not fabricate or introduce outside knowledge. If there is insufficient evidence, output [].
2. Short answers, using concise entity names/numbers/dates/phrases; no explanations, prefixes or suffixes (e.g., "The answer is").
3. When multiple answers are listed side by side (such as aliases, multiple entities, and multiple years), they are listed from high to low in order of confidence, and duplicates and noise are removed.
4. Normalization: Remove unnecessary spaces and punctuation; retain numbers/units of measure as is; format dates according to the original text; preserve capitalization of proper nouns.

Output format

Your output must be a parseable JSON array. Do not include any `` or any other descriptive terms other than JSON objects. Your output should look like this:

```
["ans1", ...]
```

user:

Please output your answer directly in the form of a parsable JSON array based on the following passages:

The current Query is: {question}.

The current Chunks are: {item}.

The current JSON-LD are: {JSON-LD}.

Figure 6: System and user prompt for the generator model.