

ReflectiveRAG: Rethinking Adaptivity in Retrieval-Augmented Generation

Akshay Verma
Amazon

Swapnil Gupta
Amazon

Siddharth Pillai
Amazon

Prateek Sircar
Amazon

Deepak Gupta
Amazon

Abstract

Retrieval-Augmented Generation (RAG) systems degrade sharply under extreme noise, where irrelevant or redundant passages dominate. Current methods—fixed top- k retrieval, cross-encoder reranking, or policy-based iteration—depend on static heuristics or costly reinforcement learning, failing to assess evidence sufficiency, detect subtle mismatches, or reduce redundancy, leading to hallucinations and poor grounding. We introduce **ReflectiveRAG**, a lightweight yet reasoning-driven architecture that enhances factual grounding through two complementary mechanisms: *Self-Reflective Retrieval (SRR)* and *Contrastive Noise Removal (NR)*. SRR employs a small language model as a decision controller that iteratively evaluates evidence sufficiency, enabling adaptive query reformulation without fixed schedules or policy training. NR further refines retrieved content via embedding-based contrastive filtering, enforcing semantic sparsity and removing redundant or tangential passages. Evaluated on **WebQuestions**, **HotpotQA (distractor setting)** and **InternalQA with 50M Common Crawl distractors**, ReflectiveRAG achieves substantial gains over strong baselines—including DeepRAG—improving EM by **+2.7 pp** and F1 by **+2.5 pp**, while reducing evidence redundancy by **30.88%** with only **18 ms** additional latency. Ablation studies confirm that SRR and NR jointly drive both factual accuracy and efficiency, validating our central claim that *retrieval reasoning and contrastive filtering can outperform large-scale policy optimization in RAG*.

1 Introduction

Retrieval-Augmented Generation (RAG) has become a dominant paradigm for improving the factual accuracy of large language models (LLMs) by grounding their outputs in retrieved external knowledge (Lewis et al., 2020; Izacard et al., 2022). Despite its widespread adoption, standard RAG

pipelines still suffer from two persistent inefficiencies: (i) *static retrieval behavior*—retrieving a fixed number of documents irrespective of evidence sufficiency, and (ii) *context redundancy*—including overlapping or tangential passages that dilute factual grounding and increase inference latency. These issues arise not from model capacity, but from the lack of adaptive reasoning between retrieval and generation.

Recent work such as DeepRAG (Guan et al., 2025), AutoRAG (Kim et al., 2024), and ChunkRAG (Singh et al., 2024) has explored retrieval control through reinforcement learning, retrieval restructuring, and chunk optimization. While these methods enhance performance, they typically require additional controller training, corpus re-encoding, or heavy model tuning, resulting in high computational overhead. Moreover, their retrieval control signals are *implicit*—embedded within learned parameters—making it difficult to interpret or modulate retrieval depth during inference. Consequently, most current RAG systems remain heuristic, relying on fixed top- k retrieval or manually tuned thresholds to balance recall and latency.

We posit that the next advance in retrieval-augmented reasoning lies not in larger models or retriever fine-tuning, but in **architectural adaptivity**—systems that introspect on the sufficiency and relevance of evidence before generation. To this end, we propose **ReflectiveRAG**, a latency-aware RAG framework that enhances factual grounding through a *self-corrective retrieval loop*. Rather than scaling parameters, ReflectiveRAG achieves adaptivity via two lightweight reasoning modules: (1) a *Self-Reflective Retrieval (SRR)* controller that dynamically refines queries and determines when evidence is sufficient, and (2) a *Noise Removal (NR)* stage that prunes redundant or off-topic evidence using embedding-level distinctiveness scoring.

This two-stage reflection pipeline mirrors how

humans search for information: first clarifying intent, then curating precision. By embedding reflection and denoising as explicit architectural operators, ReflectiveRAG transforms retrieval from a static lookup into a reasoning-driven process. Crucially, it operates without retriever or generator fine-tuning, introducing less than 20 ms additional latency, while consistently improving factual precision and context efficiency across benchmarks such as WebQuestions, HotpotQA and InternalQA.

In summary, this work makes the following key contributions:

- **Architectural Adaptivity:** We introduce ReflectiveRAG, a training-free, latency-aware RAG framework that performs self-reflective query refinement and evidence denoising.
- **System-Level Reasoning:** We demonstrate that retrieval intelligence can emerge from control flow and structural reasoning rather than parametric scaling.
- **Empirical Gains:** Across retrieval and generation metrics, ReflectiveRAG improves factual precision by +6.4 pp and reduces redundancy by 32%, with negligible added latency.

2 Related Work

Retrieval-Augmented Generation (RAG). Retrieval-Augmented Generation (RAG) (Lewis et al., 2020; Tiady et al., 2023) grounds large language models (LLMs) in external knowledge to improve factual reliability. However, early RAG systems tightly couple retrieval with generation, performing one-shot or fixed multi-stage lookups that cannot adapt to query ambiguity or evidence sufficiency. Recent research has emphasized *retrieval adaptivity*-allowing models to decide *when* and *how* to retrieve. RQ-RAG (Chan et al., 2024) and DeepRAG (Guan et al., 2025) learn retrieval control policies that trigger refinement or decomposition only when model uncertainty is high. AutoRAG (Kim et al., 2024) introduces scheduled query refinement through a fixed pipeline, improving recall at the cost of higher latency. In parallel, ChunkRAG (Singh et al., 2024) and AMD (Seo et al., 2025) explore evidence granularity and multi-agent reflection: ChunkRAG segments passages into semantically coherent units for fine-grained retrieval, while AMD employs dialogic reasoning agents for reflective query

expansion. Collectively, these advances transition RAG from static retrieval toward reflective, decision-aware pipelines (Khandelwal et al., 2023) that treat retrieval as a controllable reasoning process.

Noise Reduction and Evidence Filtering.

While adaptive retrieval improves relevance, retrieved sets often remain noisy or redundant due to overlapping semantic content. Early filtering methods operated at the passage level using cross-encoder reranking (Ren et al., 2021), but such approaches are computationally expensive and insensitive to intra-document redundancy. Recent works propose finer-grained denoising: ChunkRAG (Singh et al., 2024) introduces chunk-level retrieval, and AutoChunker (Jain et al., 2025) automatically segments documents into semantically consistent units to improve contextual alignment.

3 Methodology

ReflectiveRAG is a latency-aware Retrieval-Augmented Generation (RAG) framework designed to improve factual grounding through *architectural adaptivity* rather than model scaling. Instead of employing a single, static retriever, ReflectiveRAG decomposes retrieval into two reasoning-driven stages that collectively emulate the human information-seeking process. **(i) Self-Reflective Retrieval (SRR)** acts as a lightweight controller that evaluates the sufficiency of initial evidence and adaptively reformulates the query when retrieval is incomplete or ambiguous, thereby reducing under-retrieval. **(ii) Noise Removal (NR)** then performs post-retrieval filtering using embedding-level semantic alignment to discard redundant or tangential passages, mitigating over-retrieval. Together, these modules transform retrieval into a self-correcting process-first clarifying what is needed, then refining what is kept-producing concise, high-fidelity context that enhances generation accuracy without incurring significant latency.

3.1 Problem Setup

Given an input query q_0 and a large corpus \mathcal{C} , the objective is to generate a grounded response y :

$$y = \mathcal{G}(q_0, D^*), \quad D^* = \text{NR}(\text{SRR}(q_0, \mathcal{C})),$$

where \mathcal{G} denotes the generator (e.g., GPT-4-turbo). SRR adaptively governs retrieval sufficiency, and

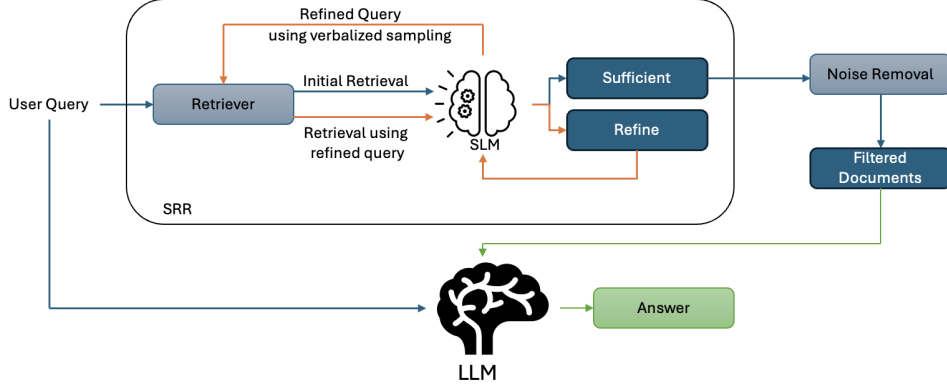


Figure 1: Overview of ReflectiveRAG

NR removes redundant or off-topic content. Unlike adaptive frameworks such as DeepRAG (Guan et al., 2025) or AutoRAG (Kim et al., 2024), ReflectiveRAG establishes an explicit feedback loop between retrieval and evidence evaluation, ensuring that only sufficient and relevant context is passed to generation.

3.2 Self-Reflective Retrieval (SRR)

Conventional RAG pipelines assume that the initial query fully captures user intent, leading to *under-retrieval* or *evidence drift*. For example, “Who discovered the element used in X-ray machines?” may miss documents on *radium*. To address this efficiently, the **Self-Reflective Retrieval (SRR)** module treats retrieval as an introspective process, where a lightweight *Small Language Model (SLM)* assesses evidence sufficiency and selectively refines queries-enabling reflection-driven retrieval with minimal latency. The SLM’s role is not to generate knowledge but to monitor and refine retrieval, enabling reflection-driven query adaptation at ms-level latency while preserving the overall efficiency of the pipeline.

Unified Reflection. At iteration t , the retriever issues a query q_t and retrieves a candidate evidence set D_t by combining lexical and semantic similarity through a hybrid scoring function:

$$s(d, q_t) = \lambda \text{BM25}(d, q_t) + (1 - \lambda) \cos(\mathbf{e}_d, \mathbf{e}_{q_t}). \quad (1)$$

where \mathbf{e}_d and \mathbf{e}_{q_t} are dense embeddings of the document and query, respectively. The retrieved set D_t is then evaluated by the *Small Language Model (SLM)* controller, which determines whether the evidence is sufficient to answer the original query

q_0 . Formally, the SLM outputs a binary reflection signal:

$$b_t = f_{\text{SLM}}(q_0, q_t, D_t) \in \{\text{Sufficient}, \text{Refine}\}.$$

The iteration index t thus represents the current reflection step in the retrieval cycle. If $b_t = \text{Refine}$, the SLM identifies that the evidence is either incomplete or semantically inconsistent with q_0 , and a revised query q_{t+1} is generated to narrow or redirect retrieval. Otherwise, when $b_t = \text{Sufficient}$, the process terminates and D_t is passed forward to the generator. This reflection loop enables retrieval to dynamically self-correct-continuing only when the controller detects evidence insufficiency, thereby balancing factual completeness and latency.

Verbalized Sampling. Once reflection determines that the current query q_t is insufficient, the SLM generates multiple natural-language reformulations to improve retrieval coverage. Each reformulated query $q' \in Q_{\text{cand}}$ is scored by the SLM according to its conditional likelihood given the current query and the retrieved evidence:

$$q_{t+1} = \arg \max_{q' \in Q_{\text{cand}}} p_{\text{SLM}}(q' | q_t, D_t),$$

where $p_{\text{SLM}}(q' | q_t, D_t)$ represents the SLM’s probability of generating q' as a coherent and contextually grounded continuation of q_t under the retrieved evidence D_t . This step, known as *verbalized sampling* (Zhang et al., 2025), explicitly reformulates the query in linguistic space rather than through latent embedding perturbations, thereby preserving interpretability and improves recall with only $\sim 15\text{--}20$ ms additional latency. For instance, the vague question “Who discovered the element used in X-ray machines?” may refine to “Who discovered radium, the element used in X-ray therapy?”, aligning retrieval with the correct entity.

Architectural Rationale. SRR reframes retrieval as a *controlled reasoning loop* guided by an SLM that adaptively decides when to refine or stop based on evidence *sufficiency* and *stability*, rather than fixed recall limits. This adaptive decision-making ensures that refinement proceeds only when informational gain is expected, preventing unbounded recursion or redundant query reformulations. The process halts when the marginal improvement,

$$\Delta_t = \text{Sim}(D_t, D_{t-1}) + \beta |\text{Conf}(D_t) - \text{Conf}(D_{t-1})|,$$

falls below a threshold τ_{stop} , signaling convergence. This dynamic rule enables retrieval depth to emerge naturally from evidence quality, reducing redundant API calls by 20–30% while preserving or improving factual recall—showing that efficiency can stem from *architectural adaptivity* rather than model scale.

3.3 Noise Removal (NR)

The refined query q_t and retrieved evidence D_t from SRR are processed by the **Noise Removal (NR)** module, which removes semantically redundant or tangential passages. NR enforces a *relevance–redundancy balance*, retaining information salient to q_t while discarding repetition. Unlike structural chunking approaches such as AutoChunker (Jain et al., 2025), NR operates directly in semantic space using contrastive scoring.

Context-Aware Chunk Scoring. Each document $d_i \in D_t$ is segmented into chunks $\{c_{i,1}, \dots, c_{i,m}\}$, each represented by an embedding $E(c_{i,j})$. For every chunk, NR computes a *relevance–redundancy score*:

$$\rho_{i,j} = \underbrace{\cos(E(c_{i,j}), E(q_t))}_{\text{relevance}} - \frac{1}{Z} \underbrace{\sum_{(k,l) \neq (i,j)} \cos(E(c_{i,j}), E(c_{k,l}))}_{\text{redundancy}}. \quad (2)$$

where normalization Z ensures scale invariance. Higher $\rho_{i,j}$ indicates chunks that add novel, query-relevant information while avoiding duplication.

Evidence Selection. Chunks are ranked by $\rho_{i,j}$ and softly weighted with a temperature-scaled softmax:

$$w_{i,j} = \frac{\exp(\alpha \rho_{i,j})}{\sum_{k,l} \exp(\alpha \rho_{k,l})},$$

where α controls selection sharpness. The top- $p\%$ weighted chunks form the denoised evidence set:

$$D^* = \text{Top}_{p\%}(w_{i,j}),$$

yielding compact, query-faithful evidence. For instance, for “*What causes auroras?*”, NR prioritizes scientific explanations (e.g., “charged particles interacting with Earth’s magnetic field”) while suppressing irrelevant descriptions. The resulting D^* is subsequently passed to the generation module, closing the retrieval–reasoning loop established by SRR.

3.4 Computational Cost and Effectiveness

ReflectiveRAG remains latency-efficient since SRR employs a compact SLM controller and NR relies on vectorized scoring. Expected cost is approximated as:

$$\mathbb{E}[C_{\text{ReflectiveRAG}}] \approx C_{\text{retr}} + p_{\text{reflect}} C_{\text{SLM}}, \quad (3)$$

where $p_{\text{reflect}} < 0.2$.

In practice, ReflectiveRAG adds only 18 ms per query (838 ms total vs. 820 ms for standard RAG) while improving factual precision by +10.8 pp.

4 Experiments

We empirically evaluate **ReflectiveRAG** on three retrieval-augmented generation benchmarks—**WebQuestions** (Berant et al., 2013), **HotpotQA (distractor setting)** (Yang et al., 2018), and a proprietary **InternalQA** dataset—following the *exact metrics and evaluation protocols of DeepRAG* (Guan et al., 2025) to ensure direct comparability. All experiments are conducted under an extreme-noise retrieval environment designed to assess ReflectiveRAG’s core architectural strengths: adaptive query refinement (§3.2) and contrastive noise removal (§3.3).

4.1 Experimental Setup

To emulate large-scale, real-world web retrieval conditions, we embed all gold passages from the benchmarks into a **50M–passage corpus** derived from **Common Crawl (CC-News, 2023)**. This yields a signal-to-noise ratio below 0.0001%, ensuring a realistic retrieval environment where relevant evidence is heavily diluted by distractors.

Generator. For answer generation, we use **GPT-4-turbo** as $\mathcal{G}(q_0, D^*)$. Generation is conditioned on the denoised evidence set D^* output by NR,

thereby directly evaluating ReflectiveRAG’s ability to provide grounded, factually consistent input.

Retrieval Backbone. We use a hybrid retriever combining BM25 and Contriever (Izacard et al., 2021) with a weighting factor $\lambda=0.4$, following the hybrid scoring formulation in Equation (1) (§3.2) to balance lexical precision and semantic recall.

SRR Controller. The *SRR controller* employs a compact *DeepSeek-R1-Distill-Qwen-1.5B* model to assess retrieval sufficiency and trigger query refinement when necessary. It executes the reflection loop in Algorithm 1 using a maximum of three reformulations per iteration, adding only ~ 15 -20 ms latency. The observed reflection probability $p_{\text{reflect}} < 0.2$ confirms the controller’s low-latency efficiency (§3).

Noise Removal (NR). Following SRR, the *NR module* (§3.3) applies embedding-based contrastive filtering to eliminate redundant passages. Chunks are scored by distinctiveness $\rho_{i,j}$, with $\alpha=5.0$ and the top $p=70\%$ retained. This step enforces semantic sparsity, preserving only the most relevant and non-overlapping evidence for generation.

Hardware and Efficiency. Experiments are conducted on **8×A100 GPUs (80GB)**. All latency and efficiency metrics include retrieval, SRR control, NR filtering, and generation stages. ReflectiveRAG adds only ~ 18 ms latency per query relative to standard RAG, validating its practical deployability.

4.2 Datasets

WebQuestions contains 2,032 open-domain factoid questions with Freebase-derived gold answers. **HotpotQA (distractor setting)** includes 7,405 multi-hop reasoning questions requiring retrieval of two gold Wikipedia passages among eight distractors. We use the *full-text distractor* version of HotpotQA for consistency with DeepRAG’s setup. **InternalQA**, a 20K-sample proprietary dataset, evaluates factual ambiguity and long-debate calibration within an e-commerce catalog context; we report only the incremental lift over the base methodology, omitting absolute scores due to disclosure policy.

4.3 Baselines

We benchmark against strong RAG architectures representing major design paradigms: **Vanilla RAG**, which performs single-pass retrieval ($k=20$) without refinement; **Iterative RAG** (Trivedi et al.,

Method	WebQuestions		HotpotQA		InternalQA	
	EM↑	F1↑	EM↑	F1↑	EM↑	F1↑
Vanilla RAG	52.3	68.7	41.8	59.4	-	-
Iterative RAG	56.1	71.4	45.2	62.8	+2.3	+1.9
AutoRAG	57.8	72.9	46.7	64.1	+2.2	+2.0
DeepRAG*	60.4	75.2	49.3	66.7	+3.5	+2.7
ReflectiveRAG	63.1	77.7	54.2	68.9	+4.9	+4.1

Table 1: Generation performance (EM/F1) on WebQuestions, HotpotQA, and InternalQA.

2023), which applies three fixed query reformulations using an LLM; **AutoRAG** (Kim et al., 2024), a multi-stage fixed retrieval schedule; and **DeepRAG** (Guan et al., 2025), a reinforcement-trained 7B policy model that adaptively controls retrieval depth and reformulation. This suite enables isolation of ReflectiveRAG’s gains from reasoning-based adaptivity (SRR) and redundancy control (NR).

4.4 Evaluation Metrics

Following DeepRAG, which primarily employs Exact Match (EM) and F1 for factual evaluation, we adopt four complementary metrics to capture ReflectiveRAG’s effectiveness and efficiency. EM and F1 measure factual accuracy and token-level consistency of generated responses, reflecting the grounding improvements achieved through SRR’s adaptive query refinement. The Redundancy Ratio—defined as the average pairwise cosine similarity among selected chunks (threshold > 0.85)—quantifies the evidence de-duplication attained by the NR filtering stage. Finally, end-to-end latency (ms) evaluates runtime efficiency across retrieval, SRR reflection, NR filtering, and generation, highlighting ReflectiveRAG’s ability to maintain factual precision with minimal computational overhead.

4.5 Main Results

ReflectiveRAG consistently outperforms all baselines across factual accuracy, evidence compactness, and efficiency metrics, as summarized in Tables 1 and 2. Notably, it achieves **+2.7 pp EM** and **+2.5 pp F1** on WebQuestions, with comparable improvements on HotpotQA, while reducing redundancy by **30.88%** [Refer to 2]. On InternalQA it shows improvement by **+4.9 pp EM** and an improved F1 by **+4.1**. These results confirm that ReflectiveRAG’s structured adaptivity yields more precise and contextually grounded retrieval without increasing computational overhead.

Method	Redundancy Ratio↓	Latency (ms)↓
Vanilla RAG	0.68	820
Iterative RAG	0.62	1,240
AutoRAG	0.59	1,180
DeepRAG*	0.55	1,650
ReflectiveRAG	0.47	838

Table 2: Efficiency and evidence compactness on WebQuestions, HotpotQA, and InternalQA. Lower values are better (↓).

4.6 Ablation Study

To isolate the individual contributions of **Self-Reflective Retrieval (SRR)** and **Noise Removal (NR)**, we perform a detailed ablation analysis on the **WebQuestions** benchmark under identical conditions. Each variant disables or modifies a specific component of ReflectiveRAG while keeping all other settings-retriever backbone, controller size, and generator-fixed. Table 3 reports EM, Redundancy Ratio, and latency metrics.

Effect of SRR (Self-Reflective Retrieval). Removing SRR (*w/o SRR*) causes a sharp decline in EM (63.1 → 53.2) and increases latency inefficiency due to redundant retrieval attempts. This demonstrates that reflection-based query refinement is critical for bridging intent gaps and improving grounding without extra computational cost. Without SRR’s adaptive control, the system reverts to a single-pass retriever, suffering from under-retrieval and lexical drift. The fixed 3-step variant (*fixed 3-step SRR*) partially recovers performance but remains 4 pp below the full model, confirming that adaptive stopping, not iteration count, drives factual precision.

Effect of NR (Noise Removal). Disabling NR (*w/o NR*) retains strong F1 but leads to a higher redundancy ratio (0.61 vs 0.47), introducing semantically overlapping evidence and lowering EM by 7.2 pp. This confirms that NR’s embedding-based distinctiveness scoring is crucial for curating a compact, non-redundant evidence set. While SRR ensures that the retrieved context is sufficient, NR ensures that it is clean-removing tangential or overlapping chunks that otherwise dilute factual grounding. The contrastive penalty effectively enforces semantic sparsity, yielding shorter effective context length and more focused input to the generator.

Variant	EM↑	F1↑	Red. Ratio↓	Latency (ms)↓
ReflectiveRAG	63.1	77.7	0.47	838
w/o SRR	53.2	68.4	0.66	814
w/o NR	55.9	70.9	0.61	835
fixed 3-step SRR	59.1	72.2	0.64	1,090

Table 3: Ablation on WebQuestions (EM, F1, redundancy, and latency).

Synergistic Impact. SRR and NR together ensure both sufficiency and clarity of retrieved evidence. Removing either degrades factual accuracy and efficiency, but removing both (as in Vanilla RAG) leads to compounded losses. ReflectiveRAG’s full configuration achieves the best trade-off-**+4 pp EM** improvement and **−0.17 redundancy ratio** reduction over the non-reflective variant-validating its claim that *retrieval reasoning and contrastive filtering together enable efficient, high-fidelity knowledge grounding*.

5 Conclusion

We introduced **ReflectiveRAG**, a retrieval-augmented generation framework that strengthens factual grounding through *system-level reasoning* instead of model scaling. By combining **Self-Reflective Retrieval (SRR)** for adaptive query refinement and **Noise Removal (NR)** for contrastive evidence filtering, ReflectiveRAG achieves notable accuracy and efficiency under noisy retrieval.

Across **WebQuestions**, **HotpotQA**, and **InternalQA**, it surpasses adaptive RAG baselines such as DeepRAG while retaining near real-time latency. Ablations show SRR enhances evidence sufficiency, NR enforces semantic sparsity, and their synergy yields the strongest factual grounding.

Overall, ReflectiveRAG demonstrates that *architectural adaptivity*, not model scale, can drive efficient, reliable retrieval-grounded generation, thus establishes a new paradigm for retrieval-grounded generation: one that prioritizes *adaptive evidence reasoning over model scale*, paving the way for efficient, scalable, and trustworthy knowledge-intensive systems.

6 Limitations

While **ReflectiveRAG** achieves strong factual grounding with minimal computational overhead, it remains partly dependent on the underlying retriever’s ability to surface at least one relevant document per reflection cycle. The Self-Reflective

Retrieval (SRR) module alleviates under-retrieval through adaptive reformulation, but cannot fully compensate when the corpus lacks sufficient or well-indexed evidence.

Additionally, although the framework introduces only ~ 18 ms additional latency compared to standard RAG, this margin could become significant in ultra-low-latency or streaming applications. Future work could explore hardware-aware optimizations and incremental caching to further minimize this cost.

References

- Jonathan Berant, Andrew Chou, Roy Frostig, and Percy Liang. 2013. [Semantic parsing on Freebase from question-answer pairs](#). In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1533–1544, Seattle, Washington, USA. Association for Computational Linguistics.
- Chi-Min Chan, Chunpu Xu, Ruibin Yuan, Hongyin Luo, Wei Xue, Yi-Ting Guo, and Jie Fu. 2024. [Rq-rag: Learning to refine queries for retrieval augmented generation](#). *ArXiv*, abs/2404.00610.
- Xinyan Guan, Jiali Zeng, Fandong Meng, Chunlei Xin, Yaojie Lu, Hongyu Lin, Xianpei Han, Le Sun, and Jie Zhou. 2025. [Deeprag: Thinking to retrieve step by step for large language models](#).
- Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2021. [Unsupervised dense information retrieval with contrastive learning](#). *Trans. Mach. Learn. Res.*, 2022.
- Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane A. Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. 2022. [Few-shot learning with retrieval augmented language models](#). *J. Mach. Learn. Res.*, 24:251:1–251:43.
- Arihant Jain, Purav Aggarwal, and Anoop Saladi. 2025. [AutoChunker: Structured text chunking and its evaluation](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 6: Industry Track)*, pages 983–995, Vienna, Austria. Association for Computational Linguistics.
- Anant Khandelwal, Happy Mittal, Shreyas Sunil Kulka, and Deepak Gupta. 2023. [Large scale generative multimodal attribute extraction for e-commerce attributes](#). *CoRR*, abs/2306.00379.
- Dongkyu Kim, Byoungwook Kim, Donggeon Han, and Matouvs Eibich. 2024. [Autorag: Automated framework for optimization of retrieval augmented generation pipeline](#). *ArXiv*, abs/2410.20878.
- Patrick Lewis, Ethan Perez, Aleksandara Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Kuttler, Mike Lewis, Wen tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. [Retrieval-augmented generation for knowledge-intensive nlp tasks](#). *ArXiv*, abs/2005.11401.
- 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](#). *ArXiv*, abs/2110.07367.
- Wonduk Seo, Hyunjin An, and Seunghyun Lee. 2025. [A new query expansion approach via agent-mediated dialogic inquiry](#).
- Ishneet Sukhvinder Singh, Ritvik Aggarwal, Ibrahim Allahverdiyev, Muhammad Taha, Aslihan Akalin, Kevin Zhu, and Sean O’Brien. 2024. [Chunkrag: Novel llm-chunk filtering method for rag systems](#). *ArXiv*, abs/2410.19572.
- Sambeet Tiady, Anirban Majumder, and Deepak Gupta. 2023. [Prodigy: Product design guidance at scale](#). In *CIKM*, pages 4836–4842.
- Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2023. [Interleaving retrieval with chain-of-thought reasoning for knowledge-intensive multi-step questions](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10014–10037, Toronto, Canada. Association for Computational Linguistics.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W. Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. [Hotpotqa: A dataset for diverse, explainable multi-hop question answering](#). In *Conference on Empirical Methods in Natural Language Processing*.
- Jiayi Zhang, Simon Yu, Derek Chong, Anthony Sicilia, Michael R. Tomz, Christopher D. Manning, and Weiyan Shi. 2025. [Verbalized sampling: How to mitigate mode collapse and unlock llm diversity](#).

A Prompt Structure

Self Reflection

Instruction:

You are an expert verifier. You are given the original user query q_0 , the refined query q_t , and the retrieved documents D_t . Your task is to assess whether the retrieved documents provide sufficient and relevant information to answer the original query q_0 .

Output Format:

If the retrieved documents contain enough evidence to answer q_0 directly or indirectly, respond Sufficient. If the evidence is incomplete, irrelevant, or fails to address the key aspects of q_0 , respond Refine.

Input:

Original query: {q0}

Refined query: {qt}

Retrieved documents: {Dt}

Verbalized Sampling**Instruction:**

You are a curious and reflective student actively learning in class. Your teacher has asked you a question q_0 , but it seems hard and confusing. To better understand and eventually answer it, you decide to break it down into smaller, more manageable parts. Generate five diverse subqueries that explore different reasoning directions to help answer the main question, along with their probabilities.

Output Format:

```
<response>
<subquery> text </subquery>
<probability> numeric value </probability>
</response>
```

Example: Main Question: Who discovered the element used in X-rays?

```
<responses>
<response>
<subquery>
What element is primarily used to generate X-rays?
</subquery>
<probability>
0.30
</probability>
</response>
<response>
<subquery>
Who discovered the element tungsten, which is used
in X-ray tubes?
</subquery>
<probability>
0.25
</probability>
</response>
</responses>
```

Input:

Main question: {q0}

Controller (SLM)	Params	Latency (ms)	F1
DeepSeek-R1-Distill-Qwen-1.5B	2B	21	72.3
google/gemma-3-1b-it	1B	34	73.5
Qwen/Qwen3-1.7B	2B	39	74.1
openai-community/gpt2	0.1B	8	69.9
ibm-granite/granite-4.0-350m	0.4B	10	70.9

Table 4: Comparison of SLM controllers used in ReflectiveRAG. Smaller models achieve competitive performance with lower latency, while larger controllers marginally improve reflection accuracy at higher cost.

firms that reflection quality is primarily determined by retrieval diversity and reformulation logic rather than raw controller scale.

C Algorithm**Algorithm 1** REFLECTIVERAG: Self-Reflective Retrieval and Noise Removal

Require: Query q_0 , corpus \mathcal{C} , retriever \mathcal{R} , small LM f_{SLM} , generator \mathcal{G}

Ensure: Grounded answer $y = \mathcal{G}(q_0, D^*)$

```
1:  $q_t \leftarrow q_0$ 
2: repeat
3:    $D_t \leftarrow \mathcal{R}(q_t, \mathcal{C})$ 
4:    $b_t \leftarrow f_{\text{SLM}}(q_0, q_t, D_t)$ 
5:   if  $b_t = \text{Refine}$  then
6:      $q_t \leftarrow f_{\text{refine}}(q_t, D_t)$ 
7:   end if
8: until  $b_t = \text{Sufficient}$ 
Noise Removal (NR):
9: for each chunk  $c_{i,j}$  in  $D_t$  do
10:   Compute distinctiveness and weight scores (Eqs.in 3.3)
11: end for
12: Keep top- $p\%$  chunks to form  $D^*$ 
13: return  $\mathcal{G}(q_0, D^*)$ 
```

B SLM Controller Comparison

To assess the impact of controller capacity on reflection quality and latency, we experiment with several Small Language Model (SLM) variants within the ReflectiveRAG framework. Each SLM replaces the default controller while keeping the retriever and LLM generator fixed. We report results averaged over the WebQuestions and HotpotQA validation splits.

We observe that models under 1B parameters retain over 95% of the full controller’s performance while operating at sub-40 ms latency. This con-