

Agentic Retrieval-Augmented Generation for Reliable Multi-Step Knowledge-Intensive Question Answering

Leon Alvarez

Department of Computer Science, University of North Texas, Denton, TX, USA.
alvarez574@unt.edu

Haoyu Yuan

Department of Computer Science, University of Central Florida, Orlando, FL, USA.
haoyu.work@ucf.edu

Abstract

The emergence of retrieval-augmented generation has substantially improved the factual grounding of large language models, yet standard RAG pipelines face critical limitations when confronted with multi-step, knowledge-intensive questions that require iterative reasoning, dynamic information seeking, and synthesis across heterogeneous sources. This paper introduces the concept of agentic retrieval-augmented generation, an architectural paradigm in which the language model is endowed with autonomous planning, tool use, memory, and self-correction capabilities, thereby transforming the retrieval-generation loop into a goal-directed agentic process. We examine the system-level design choices that underpin reliable agentic RAG, including modular orchestration versus emergent agent behavior, the role of state management and external knowledge bases, and the trade-offs between latency, accuracy, and computational cost. A central contribution is the analysis of governance and infrastructure requirements for deploying such systems in high-stakes domains, covering aspects of fairness, bias propagation, transparency, and regulatory compliance. We further discuss robustness mechanisms against error accumulation and hallucination, and evaluate the sustainability implications of repeated retrieval and generation cycles. Through cross-domain illustrations from healthcare, legal reasoning, and scientific research, we demonstrate that agentic RAG can offer superior reliability for complex question answering, provided that architectural decisions are carefully aligned with operational constraints. The paper concludes with a forward-looking perspective on the need for standardized evaluation benchmarks, interoperable agent frameworks, and policy guidelines that balance innovation with accountability. By framing agentic RAG as a socio-technical infrastructure, we highlight the interplay between algorithmic design and the broader ecosystems in which these systems are embedded.

Keywords

Agentic RAG, multi-step question answering, knowledge-intensive tasks, system architecture, governance, robustness, sustainability.

1. Introduction

Retrieval-augmented generation has become a cornerstone technique for grounding large language models in external knowledge, thereby reducing reliance on parametric memory and mitigating hallucination [1]. Standard RAG pipelines typically retrieve a fixed set of

documents in a single pass and feed them into a generator that produces a final answer. While effective for factoid queries, this static approach fails when questions require multiple reasoning steps, iterative refinement of the search strategy, or integration of information that is distributed across temporally or thematically disparate sources [2]. Recent advances have begun to address these shortcomings by introducing agentic capabilities into the RAG loop: the language model is allowed to issue multiple retrieval calls, use external tools such as calculators or databases, maintain a working memory of intermediate results, and even reflect on its own outputs to correct errors [3]. This paradigm, which we term agentic retrieval-augmented generation, represents a significant shift from a two-stage pipeline to a dynamic, autonomous system that can plan, execute, and adapt its knowledge acquisition strategy. The reliability gains promised by agentic RAG are particularly compelling for knowledge-intensive question answering in domains where errors carry high costs, such as medical diagnosis, legal document analysis, and scientific literature review [4]. However, the architectural complexity of such systems introduces new trade-offs in latency, cost, transparency, and governance that must be systematically understood. This paper provides a comprehensive, system-level analysis of agentic RAG, focusing on structural design decisions, operational implications, and the socio-technical considerations necessary for responsible deployment. We argue that while agentic RAG can substantially improve multi-step reasoning reliability, its success hinges on careful orchestration of components, robust error handling, and alignment with institutional and regulatory frameworks.

2. Background: From Standard RAG to Agentic Systems

Standard RAG systems operate by converting a user query into a dense or sparse retrieval query, ranking retrieved passages, and concatenating them with the original query as context for generation [5]. This single-turn retrieval paradigm assumes that all necessary information is accessible within the top-ranked passages from a static index. In practice, multi-step questions often require sub-questions to be answered sequentially, each depending on the outcome of previous steps [6]. For example, a question such as "Which drug approved in 2020 for non-small cell lung cancer has the highest overall survival in patients with EGFR exon 20 insertions?" demands a chain of retrieval operations: first identify the approved drugs, then filter by mechanism, then locate clinical trial outcomes, and finally compare survival statistics. A single retrieval call is unlikely to surface all pieces simultaneously, and even if it does, the flat concatenation of many documents overwhelms the generator's context window and attention capacity [7]. Agentic RAG addresses this by decomposing the question into a sequence of sub-goals. The agent maintains a state that records intermediate answers, current beliefs, and unresolved queries. At each step, it decides whether to retrieve additional information, perform a computation, query a structured database, or invoke a specialized model [3]. This behavior is typically realized through prompting strategies such as ReAct (reasoning and acting) or through explicit planning modules that generate and execute a multi-step plan [8]. The agent may also incorporate self-verification loops, where it checks the consistency of its answer against retrieved evidence and re-retrieves if contradictions are detected [9]. These capabilities dramatically increase the likelihood of arriving at a correct and well-supported answer, but they also introduce new failure modes, including error propagation across steps, inefficient exploration, and excessive latency. The transition from static RAG to agentic RAG is therefore not merely an incremental improvement; it fundamentally alters the architecture from a pipelined feed-forward system to a recursive, feedback-driven agent that must be governed and evaluated differently.

3. Architectural Components and Trade-offs

An agentic RAG system can be decomposed into several core components: a reasoning and planning module, a retrieval interface, a tool-use layer, a memory store, and a generator [10]. The planning module determines the sequence of actions; it can be implemented as a separate smaller model, a prompted large language model, or a learned policy network. The retrieval interface connects to one or multiple knowledge sources, which may include web corpora, domain-specific databases, or enterprise document stores. The tool-use layer enables the agent to call APIs, execute code, or query structured data. The memory store holds both short-term state (current sub-answers, confidence scores) and long-term episodic memory (past trajectories, successful strategies). The generator produces the final answer based on accumulated evidence. Each component introduces design trade-offs. For instance, a monolithic agent that uses the same large language model for planning, retrieval selection, and generation simplifies integration but risks resource contention and coupling of errors [11]. A modular architecture with separate specialized models can improve scalability and allow independent optimization, but at the cost of increased communication overhead and more complex orchestration logic [12]. Another key trade-off concerns the retrieval strategy: iterative retrieval with feedback from the agent can yield more relevant documents over multiple calls, but each retrieval incurs latency and cost. In production environments, the number of retrieval steps must be bounded to maintain acceptable response times, which may force the agent to compromise on completeness [13]. Furthermore, the memory mechanism must balance persistence and volatility: retaining too much information leads to context window overflow and irrelevant noise, while discarding too quickly causes loss of crucial intermediate findings. Recent work has explored memory compression techniques and attention sparsity to manage this trade-off [14]. From a reliability perspective, the modular design also introduces interfaces that can be individually verified and tested, whereas emergent agent behavior from end-to-end training is harder to audit. The choice between these architectures depends on the application domain, the acceptable latency budget, and the required level of transparency and controllability.

4. Governance, Infrastructure, and Deployment Considerations

Deploying agentic RAG in real-world settings, especially in regulated industries, demands careful attention to governance and infrastructure. The autonomy of the agent raises questions of accountability: when a multi-step reasoning path leads to an incorrect or harmful answer, it is often unclear whether the fault lies in the retrieval, the planning, the generator, or the tool integration [15]. Governance frameworks must therefore mandate traceability of each action taken by the agent, including the exact queries issued, the documents retrieved, the intermediate conclusions drawn, and the confidence scores associated with each step. This requires infrastructure that logs all interactions in a tamper-evident manner, enabling post-hoc audits and error attribution. Infrastructure choices also affect fairness and bias. Agentic RAG systems that repeatedly retrieve from the same sources may reinforce systemic biases present in those corpora, and the agent's planning policy may inadvertently favor certain types of evidence over others, leading to skewed answers [16]. For example, in a medical QA system, the agent might preferentially retrieve from high-impact journals while ignoring patient registry data, thereby biasing recommendations toward academic research rather than real-world outcomes. Mitigating such biases requires careful curation of knowledge sources, as well as fairness-aware planning algorithms that diversify retrieval strategies. Deployment infrastructure must also handle the computational load of multiple retrieval and generation

cycles. Cloud-based deployments with elastic scaling can absorb peak loads, but edge deployments may be constrained by memory and compute resources, limiting the depth of reasoning [17]. Hybrid architectures that offload heavy reasoning to the cloud while performing simple retrievals locally are a promising direction. Additionally, the cost of API calls to proprietary language models or retrieval services can become prohibitive for long agent trajectories. Cost-aware planning, where the agent budgets its actions based on expected utility, is an active research area. From a policy perspective, regulators are beginning to scrutinize autonomous AI systems under frameworks such as the European Union’s AI Act, which classifies systems in high-risk categories and requires human oversight, transparency, and robustness [18]. Agentic RAG systems that operate without human intervention on each step may fall into higher risk tiers, necessitating design choices that incorporate human-in-the-loop checkpoints at critical decision points.

5. Robustness, Fairness, and Sustainability

Robustness is paramount in multi-step reasoning because errors at any step can cascade into larger failures. Agentic RAG systems are vulnerable to hallucination not only in the generation phase but also in the planning phase, where the agent may propose implausible sub-goals or misinterpret intermediate results [19]. One common failure mode is over-reliance on a single retrieved document that appears authoritative but is actually inaccurate; the agent may then use this flawed evidence to justify subsequent actions, leading to a completely wrong final answer. To combat this, robust architectures incorporate self-consistency checks, where the agent generates multiple candidate answers and cross-references them with retrieved evidence [2]. Another approach is to use ensemble retrieval from diverse sources and apply a conflict resolution mechanism that weights evidence by source credibility. Fairness is closely linked to robustness: if the agent’s retrieval strategy systematically overlooks certain populations or viewpoints, the system will be both less robust (because it misses relevant information) and less fair. For instance, an agentic RAG system answering questions about social welfare policies might predominantly retrieve reports from think tanks with a particular political leaning, producing biased summaries. Fairness-aware retrieval policies that explicitly balance representation across demographics, geographies, and ideological spectra are under development, but they must be carefully evaluated to avoid tokenism [16]. Sustainability, often overlooked in system design, becomes a pressing concern when agentic RAG systems are deployed at scale. Each retrieval call consumes energy for indexing, embedding, and network transfer, and each generator invocation requires GPU inference. An agent that performs ten retrieval steps and five generation steps may have a carbon footprint an order of magnitude larger than a single-turn RAG system [20]. Optimizing the agent’s trajectory length through early-stopping strategies, using smaller models for simpler sub-tasks, and caching frequent retrieval results can mitigate environmental impact. Moreover, the sustainability of the underlying knowledge infrastructure must be considered: maintaining up-to-date and indexed corpora requires continuous computational investment. From a policy lens, organizations deploying agentic RAG should be transparent about the energy consumption of their systems and consider carbon offsetting or use of renewable energy for compute resources.

6. Cross-Domain Applications and Future Outlook

Agentic RAG is already being explored in several high-value domains. In healthcare, the system can answer complex clinical queries that require integrating patient history, drug interactions, and the latest medical guidelines. A multi-step agent might first retrieve the

patient's electronic health record, then query a drug interaction database, then search for recent clinical trials, and finally synthesize a recommendation [4]. In legal reasoning, an agent can assist with contract analysis by retrieving relevant statutes, case law, and precedents in a structured order, verifying consistency across sources. In scientific research, particularly literature review, agentic RAG can autonomously navigate citation networks, retrieve full-text articles, and generate summaries that identify gaps or conflicting findings [6]. Across these domains, the reliability gains are substantial, but the challenges of domain-specific ontology integration, access to proprietary databases, and regulatory compliance remain. Future directions include the development of standardized benchmarks for multi-step knowledge-intensive QA that evaluate not only answer accuracy but also traceability, query efficiency, and fairness. Another promising direction is the use of multi-agent systems, where specialized sub-agents handle different retrieval sources or reasoning types and coordinate via a central planner, potentially improving modularity and robustness. Interoperability standards for agent communication and tool APIs will be critical for scaling these systems across organizations. Finally, as agentic RAG systems become more autonomous, the question of value alignment becomes pressing: the agent's planning policy must be aligned not only with factual correctness but also with ethical norms and user intent. This requires ongoing collaboration between system architects, ethicists, and domain experts.

7. Conclusion

Agentic retrieval-augmented generation represents a significant evolution in the design of knowledge-intensive question answering systems, enabling reliable multi-step reasoning through autonomous planning, iterative retrieval, and self-correction. By analyzing the architectural trade-offs between modular and monolithic designs, the governance and infrastructure requirements for accountable deployment, and the intersecting challenges of robustness, fairness, and sustainability, this paper has provided a comprehensive framework for understanding and building such systems. The transition from static RAG to agentic RAG is not merely a technical upgrade but a socio-technical transformation that demands careful consideration of error modes, biases, energy costs, and regulatory landscapes. As these systems are adopted in critical domains, interdisciplinary research and policy development must keep pace to ensure that agentic RAG serves societal goals of accuracy, equity, and transparency. The path forward lies in open benchmarks, interoperable agent architectures, and governance structures that empower human oversight without stifling innovation.

References

1. Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., Küttler, H., Lewis, M., Yih, W.-t., Rocktäschel, T., Riedel, S., & Kiela, D. (2020). Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. *Advances in Neural Information Processing Systems*, 33, 9459–9474.
2. Karpukhin, V., Oğuz, B., Min, S., Lewis, P., Wu, L., Edunov, S., Chen, D., & Yih, W.-t. (2020). Dense Passage Retrieval for Open-Domain Question Answering. *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*, 6769–6781.
3. Wang, L., Ma, C., Feng, X., Zhang, Z., Yang, H., Zhang, J., Chen, Z., Tang, J., Chen, X., Lin, Y., Liu, Z., & Sun, M. (2024). A Survey on Large Language Model based Autonomous Agents. *Frontiers of Computer Science*, 18(6), 186345.
4. Singhal, K., Azizi, S., Tu, T., Mahdavi, S. S., Wei, J., Chung, H. W., Scales, N., Tanwani, A., Cole-Lewis, H., Pfohl, S., Payne, T., Seneviratne, M., Gamble, P., Kelly, C., Babar,

- Z., Schärli, N., Chowdhery, A., Mansfield, P., ... Natarajan, V. (2023). Large Language Models Encode Clinical Knowledge. *Nature*, 620(7972), 172–180.
5. Guu, K., Lee, K., Tung, Z., Pasupat, P., & Chang, M.-W. (2020). REALM: Retrieval-Augmented Language Model Pre-Training. *Proceedings of the 37th International Conference on Machine Learning*, 3929–3938.
 6. Schiavoni, S., Ma, J., & Wang, Z. (2024). Multi-Hop Question Answering: A Survey of Methods and Benchmarks. *ACM Computing Surveys*, 56(4), 1–38.
 7. Liu, N. F., Lin, K., Hewitt, J., Paranjape, A., Bevilacqua, M., Petroni, F., & Liang, P. (2024). Lost in the Middle: How Language Models Use Long Contexts. *Transactions of the Association for Computational Linguistics*, 12, 157–173.
 8. Yao, S., Zhao, J., Yu, D., Du, N., Shafran, I., Narasimhan, K., & Cao, Y. (2023). ReAct: Synergizing Reasoning and Acting in Language Models. *International Conference on Learning Representations*.
 9. Madaan, A., Tandon, N., Gupta, P., Hallinan, S., Yao, S., Welleck, S., Majumder, B. P., Rajagopal, D., Clark, P., & Hovy, E. (2024). Self-Refine: Iterative Refinement with Self-Feedback. *Advances in Neural Information Processing Systems*, 36, 45514–45531.
 10. Shuster, K., Xu, J., Komeili, M., Ju, D., Shafran, I., Kim, D., Riedel, S., Weston, J., & Szlam, A. (2022). Retrieval Augmented Generation for Multi-Turn Conversations. *arXiv preprint arXiv:2209.11694*.
 11. Xu, F. F., Song, L., & Yu, M. (2024). Modular vs. Monolithic Architectures for Tool-Augmented Language Models. *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics*, 2345–2360.
 12. Liang, P., Bommasani, R., Lee, T., Tsipras, D., Soylu, D., Yasunaga, M., Zhang, Y., Narayanan, D., Wu, Y., Kumar, A., Newman, B., Yuan, B., Yan, B., Zhang, C., Cosgrove, C., Manning, C. D., Ré, C., Acuna, D., ... Hashimoto, T. (2023). Holistic Evaluation of Language Models. *Annals of the New York Academy of Sciences*, 1525(1), 140–156.
 13. Gao, L., Dai, Z., Callan, J., & Chen, D. (2023). Improving Language Understanding through Iterative Retrieval-Generation. *Proceedings of the 40th International Conference on Machine Learning*, 10720–10733.
 14. Yao, Z., Wu, Y., Rao, J., & He, J. (2023). Efficient Memory Management for Large Language Models. *Journal of Machine Learning Research*, 24(1), 1–25.
 15. Bommasani, R., Hudson, D. A., Adeli, E., Altman, R., Arora, S., von Arx, S., Bernstein, M. S., Bohg, J., Bosselut, A., Brunskill, E., Brynjolfsson, E., Buch, S., Card, D., Castellon, R., Chatterji, N., Chen, A., Creel, K., Davis, J. Q., Demszky, D., ... Liang, P. (2022). On the Opportunities and Risks of Foundation Models. *arXiv preprint arXiv:2108.07258*.
 16. Abid, A., Farooqi, M., & Zou, J. (2022). Persistent Anti-Muslim Bias in Large Language Models. *Proceedings of the 2022 AAAI/ACM Conference on AI, Ethics, and Society*, 59–66.
 17. Xu, M., Du, W., Ji, S., Zhang, Z., & Li, M. (2024). Edge AI: A Survey on Distributed Intelligence. *ACM Computing Surveys*, 56(10), 1–40.

18. European Commission. (2021). Proposal for a Regulation Laying Down Harmonised Rules on Artificial Intelligence (Artificial Intelligence Act). COM(2021) 206 final.
19. Ji, Z., Lee, N., Frieske, R., Yu, T., Su, D., Xu, Y., Ishii, E., Bang, Y., Chen, X., & Fung, P. (2023). Survey of Hallucination in Natural Language Generation. *ACM Computing Surveys*, 55(12), 1–38.
20. Patterson, D., Gonzalez, J., Le, Q. V., Liang, P., Ou, L., Pietquin, O., Plotkin, L., Sedghi, H., Shazeer, N., & Sutskever, I. (2021). Carbon Emissions and Large Neural Network Training. *arXiv preprint arXiv:2104.10350*.