# Adaptive Memory Retrieval Augmentation with Self-Checks (AMRAG)

## Abstract

This paper presents the framework of Adaptive Memory Retrieval Augmentation with Self-Checks, AMRAG, for enhancing the accuracy and reliability of a Retrieval-Augmented Generation, RAG system. It incorporates dynamic query refinement and web search integration for context augmentation, while at the same time diminishing hallucinations and increasing relevance by including self-check mechanisms for retrieved documents. In this paper, it is illustrated that the AMRAG framework outperforms the traditional RAG framework with a series of experiments across multiple NLP tasks, therefore providing a way to integrate external knowledge into large language models more robustly. This study adds up to a growing body of literature on RAGs, which has an adaptive and more reliable approach and has the potential to transform the way in which retrieval-based generation tasks are handled within different domains.

## 1. Introduction

Retrieval-Augmented Generation (RAG) has emerged as a powerful method for enhancing the performance of large language models (LLMs) by incorporating external knowledge into the response generation process. While RAG systems have shown remarkable improvements over traditional generation methods, they still face significant challenges, particularly in terms of retrieval accuracy and the generation of hallucinations—false or irrelevant information.

Existing RAG implementations often rely on static retrieval processes that do not adapt to the complexity or ambiguity of queries, leading to suboptimal performance. Moreover, the lack of robust self-check mechanisms within these systems further exacerbates the issue of hallucinations. To address these challenges, we propose the Adaptive Memory Retrieval Augmentation with Self-Checks (AMRAG) framework. AMRAG introduces dynamic query refinement and self-verification processes, enabling more precise and contextually relevant retrieval while significantly reducing the incidence of hallucinations.

## 2. Background and Related Work

### 2.1 Retrieval-Augmented Generation (RAG)

Large-scale PLMs such as BERT Devlin et al. 2019, GPT-3 Brown et al. 2020, and T5 Raffel et al. 2020 have revolutionized the NLP domain. Specifically, they demonstrate impressive performance on a wide range of downstream tasks by exploiting large amounts of knowledge captured in pre-training. However, one major drawback of techniques developed along this line is that they are fixed at deployment time; thus, they cannot use real-time information or perform knowledge-intensive tasks effectively.

To address this limitation, Retrieval-Augmented Generation (RAG) was introduced as a framework that combines the generative capabilities of PLMs with the precision of information retrieval systems. RAG allows models to fetch relevant documents from external sources during the generation process, thus enhancing their ability to provide accurate and up-to-date responses (Lewis et al., 2020). Despite the success of RAG systems, challenges such as retrieval inaccuracy and the generation of hallucinated content—irrelevant or incorrect information—remain prevalent (Ji et al., 2023).

### 2.2 Limitations of Existing RAG Systems

**Self-RAG**, proposed by Asai et al. (2023), introduces a self-reflection mechanism to decide when retrieval is necessary, reducing irrelevant document usage and potential hallucinations. However, its effectiveness depends heavily on the quality of initial retrievals. If the initial retrieval fails, the system may proceed with insufficient context, leading to suboptimal responses.

Additionally, Self-RAG may overlook useful information by not performing multiple retrievals for complex queries.

**CRAG** (Corrective Retrieval-Augmented Generation) by Yan et al. (2024) enhances retrieval robustness by performing corrective web searches when initial retrievals are inadequate. However, CRAG does not store retrieved information, requiring repeated web searches for similar queries. This process can be resource-intensive, increasing latency and computational demands, particularly in scenarios with frequent, similar queries.

Another advancement is the RQ-RAG framework, which introduces query refinement strategies to handle complex and ambiguous queries more effectively (Chan et al., 2024). RQ-RAG allows the model to decompose and rewrite queries dynamically, improving the relevance of retrieved documents. However, while this approach improves retrieval, it does not fully address the issue of hallucinations, as it focuses primarily on query optimization rather than on verifying the generated content.

### 2.3 Hallucination Detection and Context Augmentation

In recent advancements in Retrieval-Augmented Generation (RAG), the challenge of hallucinations—where a model generates incorrect or irrelevant information—has gained significant attention. Traditional RAG systems often struggle to maintain the accuracy of generated content, especially when the retrieval process returns documents of varying relevance. To address this, several methods have been proposed that aim to refine the retrieval process and enhance the accuracy of the generated content.

One notable approach is the **Corrective Retrieval-Augmented Generation (CRAG)** framework, proposed by Yan et al. (2024). CRAG introduces a novel method for hallucination detection by employing a lightweight LLM to evaluate the quality of the retrieved documents before they are used in generation. This secondary LLM acts as a critic, detecting potential inaccuracies or irrelevant information that could lead to hallucinations in the generated text. Additionally, CRAG incorporates

large-scale web searches to augment the context when the initial retrieval from a static corpus does not provide sufficient information. This ensures a broader and more reliable base of knowledge for generating accurate responses.

### 2.4 The Need for Adaptive Retrieval and Self-Checks

The persistent challenges in RAG systems highlight the need for more adaptive retrieval mechanisms and robust self-check processes. Adaptive retrieval involves dynamically refining queries and using multiple retrieval attempts to ensure that the most relevant documents are obtained. This is particularly important in handling complex queries where a single retrieval attempt may not suffice.

Moreover, incorporating self-check mechanisms within the generation process can significantly reduce the likelihood of hallucinations. Self-checks involve verifying the generated content against the retrieved documents to ensure consistency and accuracy. This approach not only improves the reliability of the output but also enhances the model's ability to handle a broader range of queries effectively.

The Adaptive Memory Retrieval Augmentation with Self-Checks (AMRAG) framework proposed in this paper builds on these advancements by integrating dynamic query refinement with self-verification processes. AMRAG is designed to adapt to query complexity and to verify the generated content in real-time, thus addressing the key limitations of existing RAG systems.

# 4. Methodology/Framework

The AMRAG framework integrates several novel components designed to enhance the retrieval and generation processes in RAG systems. The workflow of AMRAG is illustrated in Figure 1, and its key components are described below.

### 4.1 Query Analysis and Decomposition

The query analysis and decomposition process is a critical component of the AMRAG framework, designed to enhance the system's ability to accurately retrieve relevant information from

complex or ambiguous queries. This process begins as soon as the user inputs a query into the system.

### 4.1.1 Decomposition with GPT-4o-mini

If the initial analysis indicates that the query is complex or contains multiple facets, the AMRAG framework proceeds to decompose the query into simpler sub-queries. This decomposition is performed using GPT-4o-mini, a lightweight language model specifically fine-tuned for this purpose.

GPT-4o-mini takes the complex query as input and generates several sub-queries, each targeting a specific aspect of the original query. These sub-queries are designed to be more precise and narrowly focused, making it easier for the retrieval system to identify and extract relevant information from the document corpus.

**For instance,** using the earlier example, GPT-4o-mini might break down the query into the following sub-queries:

- "What are the benefits of using renewable energy?"
- "What are the challenges of using renewable energy?"
- "How does renewable energy impact urban areas?"

By breaking down the query in this manner, AMRAG ensures that each sub-query can be treated independently during the retrieval process, leading to a more comprehensive and accurate aggregation of information in the final response.

### 4.1.2 Integration into Retrieval Process

Once the sub-queries are generated, each is processed separately through the retrieval system. The documents retrieved in response to each sub-query are then aggregated and synthesized to form a complete response to the original complex query.

This decomposition process not only improves the accuracy of the retrieved information but also helps in reducing potential hallucinations by focusing the retrieval on specific, well-defined aspects of the query. The use of GPT-4o-mini ensures that this decomposition is done efficiently, making it suitable for real-time applications within the AMRAG framework.

### 4.2 Retrieval and Relevance Checking

The decomposed query is used to retrieve documents from a pre-built vector database. The retrieval process is iterative, with the system employing query rewrites if the initial retrieval does not yield relevant results. If relevant documents are still not found, **AMRAG uses external web searches** to augment the available data, ensuring the most comprehensive information is retrieved.

Retrieved documents are then evaluated for relevance. If a document is deemed relevant, it proceeds to the next stage; if not, the query is rewritten and the retrieval process is repeated. This cycle continues until either relevant document are found or a maximum number of attempts is reached.

### 4.3 Generation and Self-Check Mechanisms

Once relevant documents are identified, they are used to generate the QA context needed for answering the query. The system then generates a response based on this context. Inspired by the self-reflective approach introduced by Asai et al. (2023) in their SELF-RAG framework, AMRAG incorporates a self-check mechanism that critically evaluates the generated content for potential inaccuracies or hallucinations. This mechanism involves a reflection process, where the model reviews and critiques its own output to ensure consistency and factual accuracy. If a hallucination or inconsistency is detected, the system loops back to the query refinement phase to correct the process. If no issues are found, the response is finalized and delivered to the user

### 4.4 Final Output

The final response, verified and refined through multiple stages of self-checking and query adjustment, is then presented to the user. This ensures that the generated answer is both accurate and contextually appropriate.

# 5. Results

## 5.1 Testing Methodology

To evaluate the effectiveness of the Adaptive Memory Retrieval Augmentation with Self-Checks (AMRAG) framework compared to a traditional Retrieval-Augmented Generation (RAG) system, we conducted a comprehensive testing process. The evaluation was designed to measure the quality of the responses generated by both systems across a variety of natural language processing tasks.

For the testing, we utilized the **RAGTruth** dataset (Niu et al., 2023), which comprises a collection of questions paired with their corresponding context passages. These passages were converted into a vector database, enabling efficient retrieval during the evaluation. The dataset included question-and-answer (Q/A) pairs, where the questions served as queries to test both the AMRAG and traditional RAG systems. This setup allowed us to assess the systems' ability to retrieve relevant information and generate accurate, contextually appropriate responses.

### 5.1.1 Binary Classification by a Large Language Model (LLM)

The primary method of evaluation involved using a large language model (LLM) to perform binary classification on the responses generated by both the simple RAG and AMRAG systems. For each query, both RAG and AMRAG generated a response. These responses were then fed into the LLM GPT-4, which was tasked with determining which response was superior.

The LLM was not only required to choose between "Response 1" (from RAG) and "Response 2" (from AMRAG) but also to provide a reason for its choice. This approach ensured that the evaluation was not only quantitative (i.e., choosing the better response) but also qualitative, as the reasons provided by the LLM helped to understand the nuances behind its decisions.

## 5.2 Performance Metrics

The primary metric used to assess the performance was the percentage of instances where AMRAG's response was judged to be better than the simple RAG system's response. This metric provides a direct comparison of the two systems' effectiveness.

## 5.3 Comparative Analysis

The results of the testing showed that AMRAG outperformed the simple RAG system in 88% of the cases. Specifically, in the majority of queries, the LLM classified AMRAG's response as better due to various factors such as increased relevance, reduced hallucination, and more coherent context integration.

### 5.3.1 Reasons for Superior Performance of AMRAG

The LLM's reasoning provided valuable insights into why AMRAG's responses were often superior. Some of the key reasons highlighted by the LLM include:

- **Improved Relevance:** AMRAG's dynamic query refinement allowed for more precise retrieval of relevant documents, which in turn led to responses that were more aligned with the user's query.
- **Reduction in Hallucinations:** The self-check mechanisms integrated within AMRAG effectively minimized the occurrence of hallucinations. The LLM noted that AMRAG's responses were more factual and contained fewer speculative elements compared to the simple RAG system.
- **Contextual Coherence:** The AMRAG framework was better at maintaining the coherence of the context throughout the response. This was particularly evident in complex queries where maintaining a consistent narrative was challenging.

## 5.4 Statistical Significance

To ensure the robustness of the results, we performed statistical analysis on the classification outcomes. The high percentage (88%) of cases

where AMRAG outperformed the simple RAG system was found to be statistically significant, confirming that the observed improvements were not due to random chance.

## 5.5 Discussion

The results of this testing underscore the advantages of the AMRAG framework over traditional RAG systems. The combination of adaptive query refinement and self-check mechanisms not only enhances the relevance and accuracy of the responses but also contributes to a more reliable and trustworthy output.

The reasons provided by the LLM for selecting AMRAG's responses over those of the simple RAG system further validate the design choices made in developing AMRAG. By addressing the common pitfalls of traditional RAG systems—such as retrieval inaccuracies and hallucinations—AMRAG sets a new standard for retrieval-augmented generation in NLP tasks.

# 6. Future Work

One of the most promising areas for future work is the development and integration of a Self-Memory Storage System (SMSS) within the AMRAG framework. The concept of SMSS draws inspiration from the framework proposed by Cheng et al. (2023) in their work on retrieval-augmented text generation with self-memory, where the model iteratively generates and selects its own outputs as memory for subsequent rounds of generation. By adapting this approach, SMSS would allow AMRAG to retain useful documents from previous queries, enabling more efficient and accurate response generation over time

## 6.1 Enhancing Retrieval Efficiency with Self-Memory Storage System (SMSS)

One of the most promising areas for future work is the development and integration of a Self-Memory Storage System (SMSS) within the AMRAG framework. The SMSS would serve as a dynamic repository for storing documents and information retrieved from web searches that have been successfully utilized in generating accurate responses. This addition aims to improve both the efficiency and accuracy of the retrieval process.

### 6.1.1 Concept and Purpose

The SMSS is envisioned as an adaptive memory module that retains valuable information from previous queries. By storing these documents, the system could bypass redundant web searches for similar future queries, significantly reducing retrieval times and computational overhead. This stored information could also contribute to more accurate and contextually relevant responses, leveraging past retrievals to inform future outputs.

### 6.1.2 Anticipated Impact

The implementation of SMSS could lead to significant improvements in retrieval efficiency and response accuracy, particularly for queries that share similarities with previously encountered ones. Over time, the SMSS would enable AMRAG to build a more personalized, adaptive memory, enhancing its performance across a broader range of NLP tasks.

## 6.2 Comprehensive Testing and Evaluation

While the initial results have demonstrated the potential of the AMRAG framework, more extensive testing is required to fully understand its capabilities and limitations. Future work should involve rigorous testing across a broader range of datasets and natural language processing tasks to evaluate the framework's performance in diverse contexts. Specifically, the following areas should be explored:

- **Robustness Against Diverse Query Types:** Testing AMRAG's ability to handle a wide variety of queries, including those that are ambiguous, multi-faceted, or domain-specific.
- **Effectiveness of Self-Check Mechanisms:** Evaluating the impact of the self-check mechanisms on reducing hallucinations, particularly in complex and open-ended queries.
- **Comparison with State-of-the-Art Models:** Conducting comparative studies against other advanced retrieval-augmented

generation models, such as CRAG and SELF-RAG, to benchmark AMRAG's performance.

### 6.3 Evaluation Metrics

To comprehensively evaluate the performance of the AMRAG framework, a variety of metrics will be employed, covering different aspects of system functionality:

- **Relevance and Accuracy:** Metrics such as Precision@K, Recall@K, Exact Match (EM), and F1 Score will be used to measure the accuracy and relevance of the retrieved and generated content.
- **Robustness:** Hallucination Rate and Error Rate will help assess the framework's ability to generate accurate and factually correct responses, while Consistency Score will evaluate stability.
- **Efficiency:** Query Latency, Retrieval Time, and Response Generation Time will be critical in evaluating the real-time applicability and efficiency of the system.
- **User Experience:** User Satisfaction Score and Engagement Rate will be utilized to gauge the overall effectiveness and appeal of the system from the user's perspective.
- **Comparative Analysis:** Improvement Over Baseline (IOB) and A/B Testing Results will be employed to compare AMRAG's performance against existing RAG models and frameworks.

# 7. Conclusion

This paper introduced the Adaptive Memory Retrieval Augmentation with Self-Checks (AMRAG) framework, a novel approach to enhancing the accuracy and reliability of retrieval-augmented generation systems. By integrating adaptive query refinement and self-verification mechanisms, AMRAG addresses the key challenges of traditional RAG systems, such as retrieval inaccuracy and hallucinations.

The results of our experiments demonstrate that AMRAG offers significant improvements in both retrieval precision and answer relevance. Future work could explore the application of AMRAG to more diverse datasets and its integration with real-time, dynamic environments.

# References:

1. Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I. and Amodei, D. (2020) 'Language Models are Few-Shot Learners', *arXiv preprint arXiv:2005.14165*.
2. Cheng, X., Luo, D., Chen, X., Liu, L., Zhao, D., & Yan, R. (2023). **Lift Yourself Up: Retrieval-augmented Text Generation with Self-Memory**. Proceedings of the 37th Conference on Neural Information Processing Systems (NeurIPS 2023).
3. Chan, C.-M., Xu, C., Yuan, R., Luo, H., Xue, W., Guo, Y. and Fu, J. (2024) 'RQ-RAG: Learning to Refine Queries for Retrieval Augmented Generation', *arXiv preprint arXiv:2404.00610*.
4. Devlin, J., Chang, M.-W., Lee, K. and Toutanova, K. (2019) 'BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding', *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 1, pp. 4171-4186.
5. Gao, L., Ma, X., Lin, J. and Callan, J. (2022) 'Precise Zero-Shot Dense Retrieval without Relevance Labels', *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics*, pp. 1706-1717.
6. Ji, Z., Lee, N., Frieske, R., Yu, T., Su, D., Xu, Y., Ishii, E., Bang, Y. and Madotto, A. (2023) 'Survey of Hallucination in Natural Language Generation', *ACM Computing Surveys*, 55(12), pp. 1-38.
7. Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., Küttler, H., Lewis, M., Yih, W.-T., Rocktäschel, T., Riedel, S. and Kiela, D. (2020) 'Retrieval-

augmented generation for knowledge-intensive NLP tasks', *Advances in Neural Information Processing Systems*, 33, pp. 9459-9474.

8. Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W. and Liu, P. (2020) 'Exploring the limits of transfer learning with a unified text-to-text transformer', *Journal of Machine Learning Research*, 21(1), pp. 5485-5550.

9. Asai, A., Wu, Z., Wang, Y., Sil, A., & Hajishirzi, H. (2023). **SELF-RAG: Learning to Retrieve, Generate, and Critique Through Self-Reflection**. Preprint. Available at arXiv:2310.11511.

10. Niu, C., Wu, Y., Zhu, J., Xu, S., Shum, K., Zhong, R., Song, J., & Zhang, T. (2023). **RAGTruth: A Hallucination Corpus for Developing Trustworthy Retrieval-Augmented Language Models**.