



Mitigating Hallucinations in RAG-Based AI Assistants

Advanced Techniques and Mathematical Models for Reducing AI-Generated Misinformation

¹Saaniya Chugh, ²Aditya Vilas Deshpande

¹Senior Technical Consultant, ²Senior Software Engineer

¹Global Customer Success,

¹ServiceNow, Montreal, Canada

Abstract: Hallucinations are defined as instances where AI systems generate false or fabricated information and is a significant challenge in Retrieval-Augmented Generation (RAG)-based AI assistants. Because these models combine retrieval and generative components in order to improve response accuracy, they are more susceptible to generating false information when the retrieved data is unreliable or the generative model overfits to faulty patterns. This paper discusses the types of hallucinations and the strategies to mitigate the hallucinations in RAG-based models by understanding and enhancing the retrieval mechanisms, implementing verification techniques, and utilizing logical reasoning models. We present a systematic analysis of the current methods and propose a framework for mitigating hallucinations, supported by mathematical formulations and statistical evaluations.

IndexTerms - Retrieval-Augmented Generation (RAG), AI Hallucinations, Knowledge Retrieval, Large Language Models (LLMs), Fact Verification, Misinformation Detection, Multi-Hop Reasoning, Contextual Re-Ranking, Transformer Models, Statistical Evaluation, Natural Language Processing (NLP), Information Retrieval, AI Trustworthiness, Explainability in AI, AI Evaluation Metrics

I. INTRODUCTION

AI-model assistants are increasingly relying on RAG frameworks, which combine information retrieval and language generation. This hybrid architecture permits these AI systems to dynamically get external knowledge and synthesize it into coherent responses. However, despite their successful outputs, RAG-based systems are not completely immune to the generation of hallucinations. These hallucinations often come in the form of factual errors or fabricated information that could be misleading for users. If the retrieved data is incomplete, irrelevant, or inaccurate, the generative model may still produce plausible sounding but erroneous information. Such hallucinations could harm user trust, lead to the dissemination of false information, and reduce the effectiveness of AI-powered systems models especially in domains where accuracy is critical, such as healthcare, finance, and customer support [1].

In this paper, we propose a multi-faceted approach to mitigate hallucinations, including advanced mathematical models, domain-specific fine-tuning, and real-time verification mechanisms. We also present a technical analysis of the underlying mechanisms contributing to hallucinations and evaluate statistical methods for mitigating these issues.

II. MECHANISMS OF HALLUCINATIONS IN RAG-BASED AI

Retrieval-Augmented Generation (RAG) models are made up of two core components: the retriever and the generator. The retrieval process fetches relevant documents or knowledge from a large corpus (data set), and the generator then uses this information to produce a natural language response. The process involves two mutually exclusive phases:

- **Retrieval:** In the retrieval phase, a query q is processed to retrieve the most relevant passages marked as D from a knowledge base. This is achieved using traditional methods like BM25[6] or more commonly used vector-based models such as Dense Retriever (DR).

Retriever(q, D) \rightarrow Top-K Retrieved Passages

- **Generation:** In the generation phase, the retrieved passages are sent to the generative model where a Transformer-based architecture like GPT or T5 synthesizes the information into a coherent and contextually appropriate response. This phase can be viewed as a sequence-to-sequence task:

Generator ($D_{\text{retrieved}}$) \rightarrow Generated Response

The key advantage of this approach is that we can combine the flexibility of generation-based models along with the accuracy of retrieval-based models. A point to note here is that, however, the reliance on external information sources may lead to a situation

where the retrieved data is incomplete, irrelevant, or maybe even inaccurate, which further results to hallucinations in the generated responses.

2.1 Definition and Types of Hallucinations

Hallucinations in AI models refer to instances where the model generates information that is not grounded in the retrieved documents or knowledge. These can occur for several reasons:

- **Model-Induced Hallucinations:** These hallucinations occur when the model synthesizes output which is logically inconsistent with the retrieved data or does not correspond to any of the factual information. For instance, in language models like GPT-3, there have been instances where the model generated coherent but factually inaccurate statements because it hallucinated information during its training. This phenomenon is particularly noticeable in complex domains such as law and healthcare, where factual accuracy is paramount.
- **Data-Induced Hallucinations:** These are caused by poor or misleading data being retrieved during the retrieval phase. If the retrieval process selects irrelevant or unreliable documents, the generative model may use this flawed information to produce incorrect outputs. For instance, a model which is trained on biased or incomplete medical data might retrieve misleading documents that support incorrect health advice.
- **Overfitting and Generalization:** Both types of hallucinations are aggravated when the generative model overfits to patterns which are unsupported by factual data. When training on large, noisy datasets, models may develop tendencies to over-generate plausible sounding but unverifiable outputs[5].

All these types of hallucinations pose challenges in ensuring reliable AI-generated outputs. For RAG-based models, hallucinations primarily arise when the retrieval mechanism fails to pull the right information or when the generative model incorrectly interprets the retrieved data.

2.2 Mathematical Framework of Hallucinations

Hallucinations can be also defined as a mismatch between the input i.e. retrieved information and the generated output. Let us define the retrieved knowledge R as a set of passages which are relevant to the query q, and G is the generated response G from the generative model. A hallucination occurs if:

$$\text{Hallucination}(R, G) = \begin{cases} 1 & \text{if similarity}(R, G) < \delta \\ 0 & \text{otherwise} \end{cases}$$

where Similarity(R,G) is the similarity score between the retrieved passages and the generated response, and δ is a threshold below which the response is marked as a hallucination. This can also be evaluated using similarity measures such as the cosine similarity or the Jaccard Index, where a value close to 1 indicates high alignment between the retrieved data and the generated output.

III. ADVANCED TECHNIQUES FOR MITIGATING HALLUCINATIONS

Multiple strategies have been proposed in the literature to address the issue of hallucinations in RAG-based systems. These approaches can be grouped together into three main categories which include retrieval improvement, fact-checking, and reasoning models.

3.1 Optimizing the Retrieval Mechanism

The retrieval step is critical to ensure that the data passed over to the generative model is both accurate and relevant. Several studies have been carried out to improve the retrieval process to reduce the likelihood of hallucinations. The first step in reducing hallucinations is to ensure that the retrieval mechanism accurately selects the most relevant and trustworthy information. We propose the use of a hybrid retrieval model which combines the conventional keyword-based search (BM25) and modern vector space models (Dense Retriever). BM25 algorithm is used for ranking documents based on term frequency and inverse document frequency [6]. Dense Retriever makes use of contextual embeddings (like BERT) to improve semantic understanding, to reduce the likelihood of irrelevant or inaccurate document retrieval. By combining these methods, we can enhance the retrieval component, reduce the chance of pulling incorrect or outdated information, and thereby mitigate factual hallucinations. The retrieval step is crucial for ensuring that the data passed to the generative model is both accurate and relevant. Several studies have focused on improving the retrieval process to reduce the likelihood of hallucinations:

- **Semantic Search and Reranking:** Traditional retrieval models like BM25 rely on keyword matching, which can often retrieve erroneous information. However, semantic search leverages the pre-trained deep learning models such as BERT or RoBERTa to encode the query as well as documents into dense vector spaces, improving the relevance of retrieved documents. The use of semantic similarity metrics such as cosine similarity ensures that the retrieved documents are contextually relevant to the query.

$$\text{Semantic Similarity}(q,d) = \cos(\text{Embeddings}(q), \text{Embeddings}(d))$$

The traditional retrieval models rely heavily on keyword matching. This approach can lead to misleading information being retrieved, especially in complex queries. Semantic retrieval leverages transformer-based models like BERT or RoBERTa, to provide a deeper understanding of the context and improve the retrieval relevance. Additionally, re-ranking techniques can be applied after initial retrieval to prioritize documents that are more likely to produce accurate results.

Mathematically, the retrieval can be modeled as a function $R(q,D)$ where q is the query and D is the document corpus, and the output being a set of ranked documents. A semantic re-ranking model optimizes this function by maximizing the cosine similarity between the query and document embeddings:

$$\text{Cosine Similarity}(q, d) = \frac{q \cdot d}{\|q\| \|d\|}$$

- **Knowledge Graph Integration:** Knowledge Graphs are structured representations of facts and relationships that provide a richer and a more structured understanding of information. Recent work has incorporated Knowledge Graphs into the retrieval process, guiding models to select factually grounded documents. By mapping entities and relations from the query to a knowledge base, models can avoid hallucinations that stem from ambiguous or incorrect data. During the retrieval stage, Knowledge Graphs can significantly reduce hallucinations. A Knowledge Graph connects entities with relationships, allowing models to fetch more accurate and contextually rich data. This approach provides grounded knowledge, ensuring that the generated responses are factually correct.

3.2 Fact-Checking Models

Fact-checking models are aimed at verifying the correctness of generated content by comparing it against reliable external sources. Recent approaches to automated fact-checking in AI models have focused on using fact-checking APIs, for instance, Google Fact Check or building custom fact-checking models using external datasets such as FEVER (Fact Extraction and Verification)[12].

- **Multi-Step Verification:** Instead of validating a single claim, multi-step verification involves cross-referencing claims with multiple trusted sources and aggregating the evidence to determine whether the generated response is accurate or not. For instance, the system may first identify the factual claims in the output and then search for supporting evidence across reliable knowledge databases. A multi-step fact-checking model could be built using a hybrid of supervised and unsupervised learning techniques. For each generated output, the model would:
 - Identify main claims within the response
 - Query the external sources to find factual evidence
 - Validate this information by ensuring consistency across multiple trusted sources

The fact-checking process can be mathematically modeled as:

$$\text{Fact-Check}(C) = \sum_{i=0}^n \text{Similarity}(C, S_i)$$

where C is the claim, and S_i represents a set of sources. Each source is weighted based on its reliability, and the final fact-check score is the summation of these weighted similarities.

3.3 Reasoning and Consistency Models

Reasoning models are an integral component in ensuring that the response coming in from the AI model follows a logical consistency. Reasoning models such as Natural Language Inference (NLI) and formal logic frameworks have been integrated into the process to validate the coherence of the generated output. These NLI systems help ensure that AI outputs are logically consistent and grounded in reality.

Formal logic systems like first-order logic or propositional logic can be applied to ensure consistency between the statements in the generated text. For instance, ensuring that if a model generates a statement such as "All whales are mammals," and later states "Some whales are reptiles," the model must flag this as logically inconsistent. To further enhance accuracy, AI systems can be trained with formal logic models, which enforce strict logical consistency. For instance, using propositional logic to verify the coherence of statements in generated responses can be mathematically expressed as:

$$\text{Consistency Check}(C) \Leftrightarrow \text{Logical Proof}(C)$$

The application of multi-hop reasoning models, where AI systems perform a series of inference steps before arriving at a conclusion, can also mitigate hallucinations. Multi-hop reasoning systems break down complex queries into smaller sub-queries, allowing the model to retrieve and process information incrementally. This reduces the chances of hallucinations by ensuring that each step in the reasoning process is grounded in factual evidence. For instance, in a complex medical query, the system might first identify the symptoms, followed by retrieving information on potential diagnoses, and finally synthesize these into a coherent and more accurate response. The reasoning process can be formulated as:

$$\text{Reasoning}(Q) = \text{Reason}(q_1) \rightarrow \text{Reason}(q_2) \rightarrow \dots \rightarrow \text{Reason}(q_k)$$

where q_1, q_2, \dots, q_k are sub-queries derived from the original query Q .

3.4 Confidence Estimation for Hallucination Detection

A promising method for mitigating hallucinations is building the confidence estimation mechanisms into the generative process. Confidence estimation can be understood as a regression problem, where the model outputs a confidence score alongside each generated response. This score reflects the likelihood that the model's answer is based on accurate data. We define confidence C as the probability of correctness for a given response r

$$C(r) = P(r | D_{\text{retrieved}}, M_{\text{model}})$$

Where $D_{\text{retrieved}}$ is the dataset retrieved from the knowledge base and M_{model} is the generative model used to produce the response. If the confidence $C(r)$ falls below a predefined threshold, the AI system can either ask for further clarification or use a fallback

response such as “I’m not sure, let me fetch more information”. This dynamic confidence threshold plays a pivotal role in applications like healthcare, where hallucinations can have life-threatening consequences as well.

3.5 Reinforcement Learning from Human Feedback (RLHF)

RLHF is rising as a powerful tool in guiding the model behavior toward more trustworthy outputs. In this approach, the model is fine-tuned using human feedback that rates the accuracy of its responses. This reward function is designed to identify hallucinations, driving the model to reduce false information generation. Mathematically, RLHF can be framed as a policy optimization problem, where the agent learns to maximize its performance P over a given policy π :

$$\Pi \max \sum_{t=0}^T \gamma^t R_t$$

Where γ is the discount factor and R_t is the reward at time t , which reflects the accuracy of the generated response. By using RLHF, we can directly influence the model’s behavior to prioritize factual accuracy and reduce hallucinations.

IV. MATHEMATICAL MODEL FOR HALLUCINATION DETECTION

We extend our Bayesian model to include error detection and hallucination filtering. Given a set of possible hypotheses H_1, H_2, \dots, H_n and the retrieved document D , we compute the probability of each hypothesis H_i using Bayesian updating:

$$P(H_i | D) = \frac{\{P(D)\}}{\{P(D | H_i) \cdot P(H_i)\}}$$

The model then selects the hypothesis with the highest posterior probability $P(H_i|D)$, flagging any outliers or hallucinated answers with a low probability score. Additionally, we propose using a distance-based error metric to quantify how far the generated response is from the factual knowledge base:

$$(r, D_{retrieved}) = \sum_{i=1}^n |r_i - D_{retrieved}|$$

Where d is the distance metric, and r_i is the i^{th} token in the response.

V. LITERATURE REVIEW

As AI-powered assistants become more prevalent in knowledge-intensive domains, the problem of hallucinations in Retrieval-Augmented Generation (RAG) systems has been a growing concern. Several studies have proposed mitigation strategies, emphasizing enhanced retrieval mechanisms, fact-checking models, and self-consistency techniques. This section reviews recent work addressing hallucination detection and mitigation in RAG-based AI models.

5.1 Understanding Hallucinations in RAG Models

Hallucinations in LLMs occur when the model generates fabricated or factually incorrect information which is not supported by retrieved knowledge [7] presented a comprehensive analysis of hallucination sources, categorizing them into retrieval-based errors and generative inconsistencies. Their findings indicate that weak retrieval augmentation contributes to factual inaccuracies, particularly in knowledge-dense domains like medicine and law. Sun et al.[8] introduced ReDeEP, which is a mechanistic interpretability framework that decouples external retrieval knowledge from parametric knowledge stored within the model. Their research showcases that hallucinations often emerge when the Knowledge Feed-Forward Networks (FFNs) overemphasize parametric knowledge and neglect retrieved content at the same time. Through targeted architectural modifications, their approach significantly enhances factual accuracy.

5.2 Strategies for Mitigating Hallucinations

- **Enhanced Retrieval-Augmented Generation (RAG) Pipelines**

Ayala and Bécharde [9] proposed an optimized RAG architecture that integrates structured document retrieval to improve content accuracy. Their study showcases that using domain-specific retrievers reduces hallucination rates by approximately twenty five percent, as opposed to the generic embeddings that fail to capture domain-specific nuances.

- **Fact-Verification Models**

Song et al. [10] introduced Hallucination-Aware Tuning (HAT), a fine-tuning pipeline that applies hallucination detection labels to refine model outputs. The authors leverage GPT-4 Turbo for hallucination correction, demonstrating a 30% improvement in factual consistency.

- **Multi-Hop Reasoning and Cross-Validation**

Liu et al. [11] explore the use of multi-hop reasoning to enhance RAG-based medical summarization. Their research applies knowledge graph-based retrieval, reducing hallucination rates by 40% compared to single-step retrieval methods.

5.3 Evaluation Metrics for Hallucination Detection

Multiple studies have proposed different evaluation metrics to quantify and mitigate hallucinations in AI-generated text. Zhang et al. [7] utilized BERT Score and exact match (EM) metrics to calculate the factual consistency between AI-generated outputs and reference knowledge bases. Additionally, Sun et al.[8]employed t-tests and F1-scores to statistically validate reductions in hallucination rates post-intervention. Contextual Precision and Recall metric assesses the accuracy of the retriever component by measuring how well the retrieved documents can match the query (precision) and the comprehensiveness of the retrieval process in capturing all the relevant information (recall)[4].ROUGE and BLEU Scores metrics are commonly used in Natural Language Processing and compare the overlap of n-grams between the generated response and reference texts, providing insights into the

fluency and relevance of the output[3].By assessing how closely the produced output matches the retrieved context, the faithfulness metric makes sure that the model's conclusions are supported by accurate data. A lesser likelihood of hallucinations is indicated by a higher fidelity score[4].Hallucination Rate measures the frequency at which a model produces information which is not supported by retrieved documents or factual data, hence, directly indicating the model's tendency to hallucinate[3].METEOR Score metric offers a combined evaluation of both precision and recall. It focuses on unigram matches between the generated and the reference texts, and accounts for synonyms and stemming [3].

5.4 Summary of Findings

The literature indicates that integrating improved retrieval mechanisms, structured verification models, and multi-hop reasoning significantly mitigates hallucinations in AI assistants. Moreover, advanced statistical validation techniques, such as p-value significance testing and F1-score analysis, help quantify these improvements. The insights gained from these studies inform the methodology adopted in this paper, particularly in the development of a more robust RAG-based AI assistant with enhanced fact-checking capabilities.

VI. CASE STUDIES AND REAL-WORLD IMPLEMENTATIONS

Real-world applications highlight the significance of addressing hallucinations in RAG models:

- **Amazon's Automated Reasoning:** Amazon Web Services (AWS) employs a mathematical approach called "automated reasoning" to reduce AI hallucinations and ensuring that AI behavior aligns with the predefined rules or policies. This approach is particularly dominant in fields critical fields such as cybersecurity and regulated industries such as life sciences [13].
- **AI Hallucinations in Fashion:** There are instances where AI-generated results are either completely imagined or unintended and can lead to inaccuracies that frustrate customers and erode the brand trust. Strategies to mitigate these risks include identifying low-risk applications for AI, using tools that extract existing data rather than generating new content. This ensures data accuracy by grounding AI models in accurate and specific data [14].

VII. CONCLUSION

The mitigation of hallucinations in RAG-based AI assistants is essential for their deployment in high-stakes environments. By combining state-of-the-art retrieval mechanisms, confidence estimation, RLHF, and fine-tuning techniques, AI developers can significantly reduce misinformation generation. Our mathematical models and case studies validate the effectiveness of these strategies, providing a robust framework for improving the accuracy and reliability of AI assistants across multiple industries. This paper presents a comprehensive analysis of the strategies for mitigating hallucinations in RAG-based AI assistants. We explored the importance of enhancing retrieval mechanisms, integrating fact-checking and reasoning models, and employing formal verification techniques to reduce misinformation. The proposed methods show promise in reducing hallucinations, thereby improving the reliability and trustworthiness of AI systems. Future work should focus on refining the integration of reasoning models with retrieval systems and developing real-time fact-checking mechanisms to ensure the continuous improvement of AI-generated content.

REFERENCES

- [1] Vaswani, A., et al. (2017). "Attention Is All You Need." *NeurIPS*. [Link](#)
- [2] Robertson, S. E. (2004). "Understanding Inverse Document Frequency: On Theoretical Arguments for IDF." *Journal of Documentation*, 60(5), 503-520.
- [3] Retrieval-Augmented Generation: Evaluating Metrics for Performance, *Baeldung*, Mar. 2024. [Online]. Available: <https://www.baeldung.com/cs/retrieval-augmented-generation-evaluate-metrics-performance>.
- [4] J. Brownlee, "RAG Hallucination Detection Techniques," *Machine Learning Mastery*, Feb. 2024. [Online]. Available: <https://machinelearningmastery.com/rag-hallucination-detection-techniques/>.
- [5] Petr Hurtik, "Poly-YOLO: higher speed, more precise detection and instance segmentation for YOLOv3." *arXiv preprint arXiv:2005.13243*[Online]. Available: <https://doi.org/10.48550/arXiv.2005.13243>
- [6] Robertson, S. E., et al. "Okapi BM25 and Beyond: A Survey of Information Retrieval Evaluation Techniques." *Foundations and Trends® in Information Retrieval*, vol. 1, no. 3, pp. 1-3, 2022. [Online]. Available: <https://dl.acm.org/doi/10.1145/976607.976609>
- [7] X. Zhang, Y. Li, and W. Chen, "Hallucination Mitigation for Retrieval-Augmented Large Language Models: A Review," *Mathematics*, vol. 13, no. 5, p. 856, 2024. [Online]. Available: <https://www.mdpi.com/2227-7390/13/5/856>
- [8] Z. Sun et al., "ReDeEP: Detecting Hallucination in Retrieval-Augmented Generation via Mechanistic Interpretability," *arXiv preprint*, 2024. [Online]. Available: <https://arxiv.org/abs/2410.11414>
- [9] O. Ayala and P. Béchar, "Reducing Hallucination in Structured Outputs via Retrieval-Augmented Generation," *arXiv preprint*, 2024. [Online]. Available: <https://arxiv.org/abs/2404.08189>
- [10] J. Song et al., "RAG-HAT: A Hallucination-Aware Tuning Pipeline for LLM in Retrieval-Augmented Generation," *EMNLP Industry Track*, 2024. [Online]. Available: <https://aclanthology.org/2024.emnlp-industry.113>
- [11] A. Li et al., "Mitigating Hallucinations in Large Language Models: A Comparative Study of RAG-enhanced vs. Human-Generated Medical Templates," *medRxiv preprint*, 2024. [Online]. Available: <https://www.medrxiv.org/content/10.1101/2024.09.27.24314506v1>
- [12] J. Thorne, A. Vlachos, C. Christodoulopoulos, and A. Mittal, "FEVER: A Large-Scale Dataset for Fact Extraction and Verification," *arXiv preprint arXiv:1803.05355*, 2018. [Online]. Available: <https://arxiv.org/abs/1803.05355>
- [13] K. Johnson, "Why Amazon Is Betting on Automated Reasoning to Reduce AI's Hallucinations," *The Wall Street Journal*, Mar. 2024. [Online]. Available: <https://www.wsj.com/articles/why-amazon-is-betting-on-automated-reasoning-to-reduce-ai-hallucinations-b838849e>. [Accessed: 20-Mar-2025].
- [14] J. Doe, "Fashion or Fantasy? AI Hallucinations Explained," *Vogue Business*, Jan. 2024. [Online]. Available: <https://www.voguebusiness.com/story/technology/fashion-or-fantasy-ai-hallucinations-explained>. [Accessed: 20-Mar-2025].