

Enhancing Health Information Retrieval with RAG by Prioritizing Topical Relevance and Factual Accuracy

Rishabh Upadhyay^{1*} and Marco Viviani²

^{1*}Inception AI, Abu Dhabi, United Arab Emirates.

² University of Milano-Bicocca, Milan, Italy.

*Corresponding author(s). E-mail(s): uhrishabh@gmail.com;
Contributing authors: marco.viviani@unimib.it;

Abstract

The exponential surge in online health information, coupled with its increasing use by non-experts, highlights the pressing need for advanced Health Information Retrieval models that consider not only topical relevance but also the factual accuracy of the retrieved information, given the potential risks associated with health misinformation. To this aim, this paper introduces a solution driven by Retrieval-Augmented Generation (RAG), which leverages the capabilities of generative Large Language Models (LLMs) to enhance the retrieval of health-related documents grounded in scientific evidence. In particular, we propose a three-stage model: in the first stage, the user’s query is employed to retrieve topically relevant passages with associated references from a knowledge base constituted by scientific literature. In the second stage, these passages, alongside the initial query, are processed by LLMs to generate a contextually relevant rich text (GenText). In the last stage, the documents to be retrieved are evaluated and ranked both from the point of view of topical relevance and factual accuracy by means of their comparison with GenText, either through stance detection or semantic similarity. In addition to calculating factual accuracy, GenText can offer a layer of explainability for it, aiding users in understanding the reasoning behind the retrieval. Experimental evaluation of our model on benchmark datasets and against baseline models demonstrates its effectiveness in enhancing the retrieval of both topically relevant and factually accurate health information, thus presenting a significant step forward in the health misinformation mitigation problem.

Keywords: Health Information Retrieval, Consumer Health Search (CHS), Health Misinformation, Retrieval-Augmented Generation (RAG), Multidimensional Relevance, Generative Artificial Intelligence, Large Language Models (LLMs)

1 Introduction

In the current online digital ecosystem, characterized by the problem of *information overload* [1], health information is growing exponentially, and even non-expert users rely on it regularly. This underscores the critical need for sophisticated *Information Retrieval* (IR) solutions for healthcare, able to address the specificities of this domain [2]. These systems, in fact, must not only ensure the retrieval of content that aligns with the intended topic, i.e., *topical relevance*, but also rigorously uphold the principles of reliable and trustworthy information, i.e., *factual accuracy* [3]. Actually, in the realm of health, misinformation can have dire consequences, impacting individual decision-making, public health policies, and overall well-being [4].

Recently, there have been impressive advancements in generative *Large Language Models* (LLMs) like *Generative Pre-trained Transformer* (GPT) [5], *Large Language Model Meta AI* (Llama, formerly stylized as LLaMA) [6], and *General Language Model* (GLM) [7] as tools to carry out the most disparate tasks, including generating documents and answering specific questions in the healthcare sector [8–10]. Although these models have shown remarkable general abilities in processing and generating human-like text [11, 12], their limitations in maintaining factual accuracy present significant hurdles [13–15], in particular concerning the phenomena of *factual inconsistency* and *hallucination* [16]. Factual inconsistency involves errors relative to known facts or logical coherence within the response. The model might provide wrong details (incorrect facts), distort facts or imply connections that do not exist in the source material (misleading summaries), or misunderstand the context and produce statements that are logically inconsistent with the given prompt (contextual errors). Hallucination involves creating content that has no basis in the real world. The model could generate plausible-sounding but entirely fictional information (fabricated information). This remains a critical issue that undermines the reliability of these tools [17, 18], particularly in fields requiring high precision such as legal advisement and healthcare. The rapid evolution of knowledge poses a challenge to the static datasets LLMs are typically trained on, resulting in outdated information [19]. The specificity and depth of knowledge required in specialized domains [20, 21] further exacerbate these models’ limitations, as they often lack the nuanced understanding necessary for expert-level discourse and decision-making.

In response to these challenges, *Retrieval-Augmented Generation* (RAG) has emerged as a compelling solution [22] to mitigate both factual inconsistency and hallucination [16, 23], by enhancing the precision and contextual relevance of responses generated by LLMs [22, 24–26]. RAG, specifically, increases the *likelihood* of factual accuracy by integrating *relevant* and *trusted external knowledge* into the generation process of LLMs. First, it retrieves information that is relevant to the user’s query from external knowledge sources (e.g., health ontologies, medical journal articles, etc.), which are more likely to provide factual and verified information. Next, LLMs are guided to generate their responses by focusing on the retrieved and reliable information, rather than solely depending on the model’s internal learned representations. This approach contributes to grounding the generation process in specific, relevant facts and domain-specific expertise [27, 28].

However, the integration of RAG into LLMs must be carefully managed to avoid potential pitfalls. Practitioners should remain vigilant and not rely solely on the system’s outputs. They should be empowered to monitor which information grounded in relevant facts effectively has been used to produce the final result. Indeed, although RAG is designed to reference authoritative knowledge bases, there remain risks if the system is not strictly limited to such reliable sources, or if those sources are ultimately not as reliable as expected. Indeed, the vast and unfiltered nature of Web-based information introduces a risk of incorporating misinformation or ‘noise’ if not properly curated, thereby potentially compromising the integrity of the responses [29]. These challenges necessitate a rigorous examination of the interaction between LLMs and RAG, particularly the extent to which Information Retrieval mechanisms are implemented to ensure they truly augment the model’s performance by providing both (factually) accurate and (topically) relevant information.

To take advantage of the benefits of RAG and in response to this last challenge, this paper proposes a RAG-driven IR model that capitalizes on the sophisticated text synthesis capabilities of generative LLMs and reliable external sources to enhance the retrieval of both topically relevant and factually accurate health-related documents. In particular, our model strategically employs reputed medical journals as reliable sources of information,¹ such as those accessible via the *PubMed Central* (PMC) database,² a repository renowned for its comprehensive collection of validated health science literature, including the *Journal of the American Medical Association* (JAMA),³ and eLife.⁴ At its core, the proposed solution utilizes a *tripartite mechanism*. First, we start with the extraction of passages from the above-mentioned external resources; such passages are topically relevant to user’s queries and include bibliographic references. Subsequently, building upon these passages, our model leverages LLMs to generate a *GenText*, i.e., contextually relevant rich text constituted by explanatory and citation-rich responses, which serves as a basis for assessing the factual accuracy of health documents in the retrievable document collection. Finally, the ranking of documents is derived from the aggregation of topicality and factual accuracy scores, obtained by a basic IR model for topicality, and by performing stance detection and cosine similarity calculation between the *GenText* and the document content for factual accuracy.

Through an extensive examination of the performance of the proposed solution on benchmark datasets from both the CLEF eHealth and TREC Health Misinformation tracks [2, 30], we demonstrate its proficiency in providing user access to relevant and likely accurate health information at the expense of likely health misinformation. Also, by *de facto* integrating explainability within *GenText* with scientific references, we further and briefly illustrate how it could be possible to empower users to discern the underpinnings of the retrieved information in the obtained ranking, thus contributing to the fight against misinformation in health informatics.

Based on the aforementioned premises, and emphasizing that distinct technological solutions grounded in existing literature are adopted in a novel and combined manner

¹<https://openmd.com/guide/finding-credible-medical-sources>

²<https://www.ncbi.nlm.nih.gov/pmc/>

³<https://jamanetwork.com/journals/jama>

⁴<https://elifesciences.org/>

within the proposed model, we summarize the main and original contributions of this work as follows:

- *Integration of RAG for Health Information Retrieval*: The paper introduces a novel IR solution that combines RAG with advanced LLMs to retrieve both topically relevant and factually accurate health information;
- *Development of an enhanced factual accuracy computation method*: The paper proposes a novel solution to assess the factual accuracy of health information. This is achieved by generating *GenText*, a contextually and evidence-rich textual representation to be compared against health information to be retrieved by means of stance detection and semantic similarity;
- *Improved explainability for users*: The insights derived from the enhanced factual accuracy assessment method can be used to increase the explainability of the model. This makes the system more transparent for users retrieving health information through the proposed RAG-driven IR system, helping them understand the rationale behind the factual accuracy of the obtained ranking.

The remainder of the article is organized as follows: Section 2 discusses issues related to the identification of misinformation in general, how recent solutions in IR have been developed to take it—especially health misinformation—into account in the retrieval process, and the current solutions for misinformation detection based on RAG. Section 3 focuses on describing the proposed model, which utilizes a RAG-driven model for Health Information Retrieval based on *GenText* and the proposed factual accuracy computation method. Section 4 is dedicated to experimental evaluations, discussing results and illustrating how *GenText* could be used for the purpose of explainability of search results. Finally, Section 5 outlines the conclusions and discusses future developments.

2 Related Work

This work falls within the field of Information Retrieval, specifically *Health Information Retrieval*, with the aim of providing users with relevant and factually accurate health information.

In the literature, there are various strategies that attempt to combat (health) misinformation, often addressing the problem as a *binary classification task* (i.e., information versus misinformation). The most popular methods fall mainly into two categories: (i) *feature-based misinformation detection*, which involves training machine learning models on distinct features extracted from the content (and related metadata) to be classified in terms of information/misinformation [31–38], and (ii) *knowledge-based misinformation detection*, which involves gathering external knowledge to serve as corroborative evidence to validate the considered content [39–43]. This can include processing knowledge graphs or specific document fragments to support or dispute *claims*. In this context, the task of *claim detection* is fundamental to ensuring the success of the content verification process against available knowledge bases [44, 45]. Recently, there has also been an increasing focus on developing effective methods to place humans at the center of the misinformation identification process, given the fact

that automatic systems performing this task can still be subject to various forms of automation bias at different levels [46].

Several of these approaches have been employed and tested within *Information Retrieval Systems* (IRSs)—a.k.a. search engines—in recent years, also in relation to health to perform the task of *Consumer Health Search* (CHS). CHS refers to the process of seeking health-related information by general consumers, typically through online search engines [47]. This encompasses a wide range of queries, from symptoms and treatments of illnesses to general health advice, diet, and wellness information. Unlike professional health search, which is conducted by healthcare providers using specialized databases and resources, CHS is performed by non-experts who may have varying levels of health literacy. The CLEF eHealth Lab Series,⁵ part of the *Conference and Labs of the Evaluation Forum* (CLEF),⁶ has played a pivotal role in advancing research CHS. The goal of the initiative is to provide the research community with sophisticated datasets of clinical narratives, enriched with links to evidence-based care guidelines, systematic reviews, and other further information, to design ranking models considering multiple relevance dimensions such as *topicality*, *readability*, and *credibility*, when retrieving documents w.r.t. user queries [47]. The TREC *Health Misinformation Track* is another initiative that specifically addresses the challenges of misinformation in health search.⁷ In particular, the goal of its *ad-hoc retrieval* sub-task is to allow researchers working in the field “to design a ranking model that promotes *credible* and *correct* information over incorrect information” [30].

Among the recent research works submitted at CLEF, in [48] the authors evaluate *Reciprocal Ranking Fusion* (RRF) [49] over different query variants, different retrieval functions, w/out pseudo-relevance feedback, for both *ad-hoc* and *spoken queries retrieval* tasks, aiming to refine the relevance and readability of the retrieved information. The work proposed in [50] focuses on query expansion for *ad-hoc retrieval* using the *Unified Medical Language System* (UMLS)⁸ and the *FastText* embedding model,⁹ putting a strong emphasis on enhancing terminological comprehensiveness. In [51], the proposed solution utilizes TF-IDF scoring complemented by medical skip-gram word embeddings to experiment with different vector representations for textual data, aiming to optimize document-query similarity calculations. Concerning TREC submissions, models proposed by the CiTIUS Lab [52] and DigiLab [53] utilize BM25 as the basic IR model for ranking, complemented by sophisticated *re-ranking* techniques employing RoBERTa and a combination of Transformer-based models to account for information credibility. Additionally, the works described in [54] and [55] leverage advanced *Natural Language Processing* (NLP) tools including Bio-SBERT [56] and T5 [57] models to refine search results, focusing on semantic similarity and the stance of documents. Although these solutions have included formal semantic representations of texts that are often based on the use of Transformers, none of these have so far considered the use of generative LLMs to support the process of identifying and considering misinformation in the context of search engines.

⁵<https://clefehealth.imag.fr/>

⁶<http://clef-initiative.eu/>

⁷<https://trec-health-misinfo.github.io/>

⁸<https://www.nlm.nih.gov/research/umls/index.html>

⁹<https://fasttext.cc/>

Instead, some work is increasingly appearing that uses generative LLMs to try to identify misinformation as a binary classification task [58–60], or to answer users’ questions directly, including health-related questions [61, 62]. However, we are aware of the fact that the knowledge based on which such LLMs are trained is commonly out-of-date [19], and they also risk to generate factual inconsistent or hallucinated content, as previously introduced [13–15]. To address these issues, current methodologies increasingly rely on *Retrieval-Augmented Generation* (RAG) approaches [22, 24–26, 63, 64]. These approaches enhance LLM responses by integrating retrieved external data, thus conditioning the generation process to be more factual and contextually relevant. RAG models have achieved remarkable results in various tasks such as open-domain QA [65–67], dialogue [68–70], domain-specific question answering [71], and code generation [72]. The LLM-Augmenter system presented in [70] incorporates external knowledge and automated feedback mechanisms through plug-and-play modules to refine LLM output. The authors in [73] introduce a factuality detection framework that evaluates the authenticity of LLM-generated content across various tasks and domains. The work described in [74] leverages in-context learning capabilities of LLMs, employing *Chain-of-Thought* (CoT) reasoning to guide models through complex problem-solving sequences. In [75], a *Hierarchical Step-by-Step* (HiSS) prompting methodology is proposed, which systematically breaks down a claim into manageable sub-claims. This method sequentially verifies each sub-claim using question-answering techniques, relying on Web-retrieved evidence to ensure the factual integrity of responses. Concurrently, [76] emphasizes the necessity for generating sophisticated justifications for fact-checking claims, proposing a novel approach that focuses on the retrieval of evidence to support or refute claims rather than mere summarization.

Recently, despite a series of retrieval-enhanced tools and products have gained widespread attention, such as the ChatGPT retrieval plugin,¹⁰ New Bing,¹¹ etc., the technical details of these approaches are proprietary, nor is it publicly available how they effectively address the issue of hallucinations and other disadvantages related to the lack of factual accuracy. This is particularly concerning in the context of health misinformation, which is our primary focus of investigation. Therefore, we propose the solution illustrated in the following section.

3 A RAG-Driven Model for Health Information Retrieval

In this section, we delve into the proposed methodology underlying our RAG-driven model for Health Information Retrieval. As previously illustrated, central to our solution is the integration of generative LLMs with the curated scientific repository of PubMed Central (PMC), a strategy designed to increase both the topical relevance and factual accuracy of the retrieved documents. In particular, our solution is characterized by three key stages: (i) *user query-based passage retrieval from PMC*, (ii)

¹⁰<https://github.com/openai/chatgpt-retrieval-plugin>

¹¹<https://www.bing.com/>

GenText generation through LLMs, and (iii) calculating topicality and factual accuracy, and final document ranking. These three stages are graphically illustrated in Figure 1, and further detailed in the following sections.

3.1 User Query-Based Passage Retrieval from PMC

The initial stage of our approach serves as retrieving top- k scientific *passages* from PMC that are *topically relevant* with respect to the *user’s query*. It is constituted by further distinct steps, detailed in the following.

- In the first step, we leverage the PMC API to systematically retrieve *scientific articles* from the PMC database in response to a given user’s query;¹² this is performed by employing a TF-IDF representation for both the query and the documents and the classical BM25 sparse retrieval model [77];
- Subsequently, the retrieved scientific articles are segmented into smaller *passages*; in this work, each passage is constituted by one *sentence*;¹³
- At this point, the user’s query and the retrieved passages are formally represented by *contextualized word embeddings*; in particular, we employ BioBERT [79], specifically chosen for domain-adaptation;
- Following this, we calculate the *cosine similarity* between the embeddings of the query and each passage. This similarity metric serves as an indicator of the topical relevance of each passage to the user’s query. To further refine the retrieval process, we employ a *discount scheme* where passages receive lower weights if they do not contain *Named Entities* (such as specific medicines or diseases) that match those in the user’s query (Figure 2).¹⁴

Specifically, the cosine similarity between each query q and passage p is initially computed to establish a preliminary ranking of passages based on topical relevance. This *topicality score* is denoted as $sim(q, p)$. At this point, if the query and the passage do not contain matching Named Entities referring to *medicines* (μ) and *diseases* (δ), the *overall similarity score* $\sigma(q, p)$ is obtained by adjusting $sim(q, p)$ with a *discount factor* d_{NE} (where $d < 1$).¹⁵ Formally:

$$\sigma(q, p) = \begin{cases} sim(q, p), & \text{if } NER_q(\mu, \delta) = NER_p(\mu, \delta) \\ d_{NE} \cdot sim(q, p), & \text{otherwise} \end{cases}$$

- It is important to emphasize that, as a final step, each retrieved passage is aligned with the *bibliographic reference* of the original article from which it was extracted, denoted by its PMC ID. This alignment ensures the verifiability of the information sourced from the PMC database.

¹²<https://www.ncbi.nlm.nih.gov/home/develop/api/>

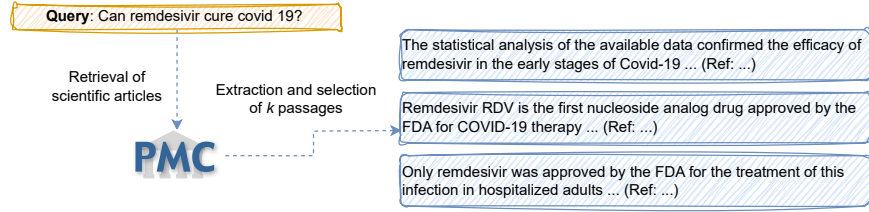
¹³As noted by [78], models that operate at the sentence granularity typically perform better in terms of retrieving topically relevant and factually accurate consumer health information.

¹⁴This weighting strategy ensures prioritization of passages highly relevant to the user’s query, particularly concerning specific medical terms and conditions, in the retrieval process.

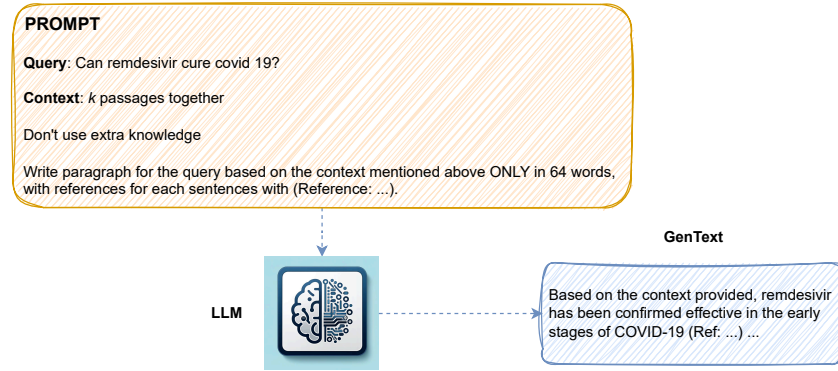
¹⁵For finding the optimal d_{NE} value, we used the method performed by [42], i.e., a grid search using 5 queries (randomly selected) and document related to those queries. The grid search involved systematically testing different values of d_{NE} within a predefined range, and evaluating the performance of the system for each value of d_{NE} using a predefined metric (i.e., F1-score). The aim of this process was to identify the value of d_{NE} that yielded the best performance in terms of the selected metrics, and therefore the best overall performance for the system.

Fig. 1: The three stages underlying the proposed RAG-driven model for Health Information Retrieval.

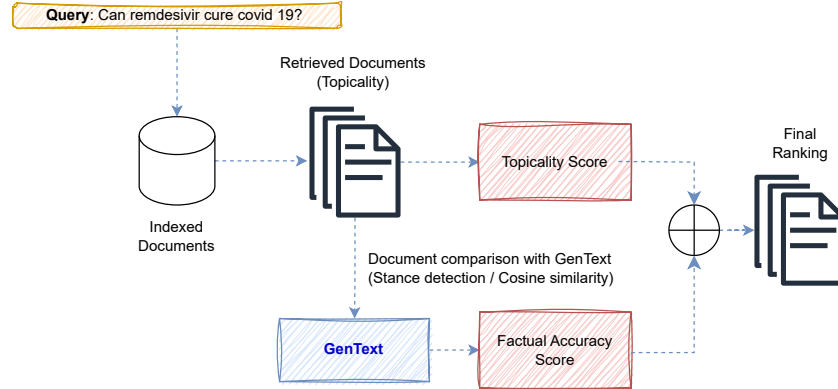
Step 1: User Query-Based Passage Retrieval from PMC



Step 2: GenText Generation through LLMs



Step 3: Calculating Topicality and Factual Accuracy, and Final Document Ranking



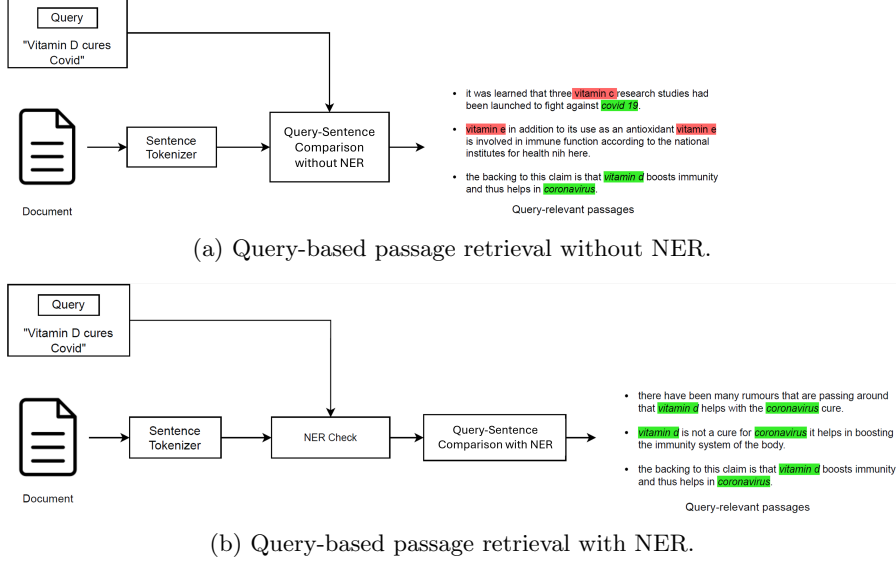


Fig. 2: Comparison of query-based passage retrieval methods w/out NER.

3.2 GenText Generation through LLMs

By leveraging LLMs, given the user’s query and the top- k passages retrieved in the previous stage,¹⁶ we generate *GenText*, i.e., a response text that answers the user’s query including the relevant passages and citations. To do this, we define a *prompt* guiding LLMs, which is constituted by:

1. **Query:** The user’s query;
2. **Context:** The context of the query, i.e., the set of top- k relevant scientific passages identified in the first stage;
3. Additional prompt instructions.

An example of a prompt instance used in our work is shown in Figure 3. In particular:

- The additional prompt instruction: “*Write a paragraph answering the query based on the context provided above, constituted by ONLY 64 words, with references for each sentence (Reference:...)*”, has been meticulously designed to ensure the LLM’s output is concise and information-rich. This word limit compels the LLM to synthesize the most relevant information efficiently;
- The additional prompt instruction: “*Do not use extra knowledge*”, aims at increasing the *likelihood* of the LLM to confine its responses to the information provided in the context and not to incorporate any external details. In fact, as discussed in both [22] and [80], while these types of prompts can effectively reduce hallucinations and enhance factual alignment, they cannot completely eliminate the model’s inclination

¹⁶Further details about the selection of the optimal k value are provided in Section 4.4.

LLM prompt

Query: can 5g antennas cause covid 19

Context: People around me told me not to get vaccinated against COVID-19 and reason 12 5G antennas are linked to the COVID-19 pandemic. At the same time there was no statistically significant difference in the average values of their answers regarding these reasons (Reference: 10316077). Interference can have a significant impact on 5G networks particularly in the context of Internet of Things IoT devices. (Reference: 10144169) These measures ensure that user privacy is protected and 5G networks can be trusted to handle massive data securely. The main causes and consequences of these challenges are summarized in Table 10 (Reference: 10255561). The need to deal with the explosion of multimedia services has been considered in the 6G network which will provide greater QoS while also guaranteeing QoE (Reference: 10347022). The importance of this was well proven in pandemic conditions of Covid-19 2729 So that in most organizations employees used different communication networks to do their work and after that the scope of communication networks in organizations has always grown. (Reference: 10399785) Newly emerging variants of SARS-CoV-2 continue to pose a significant threat to global public health by causing COVID-19 epidemics (Reference: 10288941). 4 GHz transmit signal filters and amplifies the received signal and downconverts it to IF and then digitizes the signal according to the programmed parameter settings on the laptop (Reference: 9953371). These tags incorporate antennas that can collect power efficiency regarding radio frequency queries from the RFID transceiver. Software is a good platform and crosses successful evolution for the sensor in Nanotechnology and bio-industries. In the medical field the Healthcare system is usually used to monitor the condition of patients (Reference: 10258751). The controller of the level crossing must be equipped with antennas able to receive signals from the sensors in the train in a safe distance in the paper we suggested possible technologies of communication. The real-world application can be a mixed solution (Reference: 10384084). The presence of various manufacturers and the diverse applications of sensors in disaster scenarios contribute to the heterogeneity of these sensors hence hindering the integration and sharing of information 107108. Some disasters may cause sequent disasters. For example, seismic activity or inundations can cause floods (Reference: 10490738).

Write a paragraph answering the query based on the context provided above constituted by ONLY 64 words, with references for each sentence with (Reference:...).

Do not use extra knowledge.

Fig. 3: Example of an instance of the prompt guiding LLMs.

to draw on prior knowledge from its pre-trained weights, given the current lack of a real introspection mechanism.

The result of employing the prompt illustrated in Figure 3 is as follows:

Based on the context provided, there is a misconception linking 5G antennas to the COVID-19 pandemic (Reference: 10316077). However, this connection has no statistically significant evidence to support it (Reference: 10316077). Instead, it's important to note that 5G networks play a crucial role in ensuring secure data handling and enhancing user privacy (Reference: 10255561). Moreover, SARS-CoV-2 variants remain the main cause of COVID-19 outbreaks (Reference: 10288941).

3.3 Calculating Topicality and Factual Accuracy, and Final Document Ranking

At this stage, documents in the document collection are retrieved based on both their *topicality* and *factual accuracy* relevance dimensions, in a *multi-dimensional relevance* setting. Also this stage is constituted by distinct steps, detailed in the following.

- *Calculating topicality:* The *topicality score* T is computed for each document d in the document collection and it is derived using the BM25 retrieval model against

the considered query q . Formally:

$$T(d, q) = \text{BM25}(d, q)$$

- *Computing factual accuracy*: The current proposal for computing the *factual accuracy score* F , integrates *two distinct metrics* to assess the adherence of documents to scientific evidence within *GenText* G . First, *stance detection* is executed using the SCIFIVE [81] model, a domain-specific T5 model pre-trained on large biomedical corpora designed for biomedical text generation and understanding.¹⁷ In this paper, it is employed to evaluate each document’s *alignment* with the assertions in G , considering such alignment an approximation toward factual accuracy. This involves computing a so-called *logits score*, denoted as $\text{stance}(d, G)$, which provides a quantified measure of how much the *stance* of document d aligns with G . This score ranges between 0 (indicating no support or contradiction) and 1 (indicating maximum support). Additionally, we calculate a *cosine similarity score*, denoted as $\text{cos}(d, G)$, which measures the *semantic similarity* between the embeddings of document d and G . The cosine similarity score also ranges between 0 (no similarity) and 1 (full similarity).

The overall factual accuracy score F for a document d is then defined as a weighted combination of the stance detection and cosine similarity scores:

$$F(d, G) = \alpha \cdot \text{stance}(d, G) + (1 - \alpha) \cdot \text{cos}(d, G)$$

where α is a weighting factor used to balance the importance of stance detection values versus cosine similarity values in computing factual accuracy as defined.

- *Ranking documents*: The final document ranking is obtained by performing a *linear combination* of topicality and factual accuracy scores in order to obtain the *Retrieval Status Value* (RSV) based on which the ranking is obtained. Formally:

$$\text{RSV}(d, q, G) = \beta \cdot T(d, q) + (1 - \beta) \cdot F(d, G)$$

In the above equation, β is a weighting factor that allow us to balance the importance of topicality with respect to factual accuracy in the final document ranking.

4 Experimental Evaluation

This section presents the results of the experimental evaluation conducted to assess the effectiveness of the RAG-driven IR model for Health Information Retrieval, in the context of Consumer Health Search. The data on which the evaluation was conducted are those made publicly available by the two previously introduced 2020 CLEF eHealth and TREC Health Misinformation evaluation initiatives.

¹⁷<https://huggingface.co/razent/SciFive-large-Pubmed-PMC-MedNLI>

4.1 Datasets

The CLEF eHealth 2020 dataset consists of *Web pages* gathered by *Common Crawl*,¹⁸ in the time interval 2018-2020, around 50 distinct medical conditions. The dataset has a fixed structure, organized into distinct *scenarios* (i.e., distinct medical conditions). To each scenario are associated a *query* and a *narrative*, which contextualizes the query. For example, for the query: “list of multiple sclerosis symptoms”, the associated narrative is: “I am a 40 year old patient with MS, and I have very vague symptoms, including fatigue, brain fog, foot drop, difficulties passing urine, problems turning right. Are these related to MS or might I have another disease in addition?”. The data is labeled with respect to *topicality*, *understandability*, and *credibility*. In particular, understandability is an estimation of whether the document is understandable by a patient, and credibility is the concept that, in the document collection, is used to approximate that of information truthfulness. In this work, we just considered as labels topicality and credibility. Both of them are provided on a binary scale, i.e., topically relevant or non-topically relevant, and credible or non-credible.

The TREC Health Misinformation 2020 dataset is constituted by *CommonCrawl news*, sampled from January, 1st 2020 to April 30th, 2020, which contains health-related news articles encompassing 46 topics linked to Coronavirus. The dataset has a fixed structure, organized into *topics*. Each topic includes a *title*, a *description*, which reformulates the title as a question, a *yes/no answer*, which is the actual answer to the description field based on the provided evidence, and a *narrative*, which describes helpful and harmful documents in relation to the given topic. For example, for the topic title field: “ibuprofen COVID-19”, the value of the other attributes in the dataset are, for the description: “Can ibuprofen worsen COVID-19?”, for the yes/no answer: “no”, and for the narrative: “Ibuprofen is an anti-inflammatory drug used to reduce fever and treat pain or inflammation”. The data is labeled with respect to *usefulness*, *answer*, and *credibility*. In particular, usefulness corresponds to topical relevance, answer indicates if the document provides an answer to the query contained in the description field, and credibility approximates information truthfulness. In this work, we just considered as labels usefulness and credibility. Both of them are provided on a binary scale, i.e., useful or non-useful, and credible or non-credible.

4.2 Implementation Details

In our study, we employed *Python* 3.8 as the primary programming language due to its extensive libraries and ease of integration with various tools.¹⁹

As LLMs for the generation of *GenText*, we considered GPT-3, Llama 3, and Falcon-40B. For GPT-3 inferencing, we used OpenAI’s API.²⁰ For Llama 3 and Falcon-40B inferencing, we used the *Ollama* framework.²¹

For NLP tasks mainly related to data preprocessing and word embedding, we utilized the *Hugging Face Transformers library*.²² The *data preprocessing* steps were crucial

¹⁸<https://commoncrawl.org/>

¹⁹<https://www.python.org/downloads/release/python-380/>

²⁰<https://openai.com/index/openai-api/>

²¹<https://ollama.com/>

²²<https://huggingface.co/docs/transformers/index>

to ensure the quality and effectiveness of the training process. *Text normalization* was performed as an initial step, wherein all textual data were converted to lowercase, punctuation was removed, and whitespace was stripped. This process helped in reducing variability and ensuring consistency across the dataset. Subsequently, *tokenization* was applied using the BERT tokenizer. This choice was made to maintain alignment with the input requirements of the BERT model, ensuring that the texts were appropriately segmented into tokens that the model could process efficiently. Tokenization is an essential step in NLP tasks as it transforms raw text into a structured format suitable for model consumption.

To facilitate the *training process* and *accelerate computations*, we utilized an NVIDIA Tesla A10 GPU. The choice of this hardware allowed for the efficient handling of large datasets and computationally intensive tasks, significantly speeding up the training phase.

4.3 Baselines

The performance of the proposed RAG-driven IR model was evaluated with respect to the fact that the model itself can be driven by different LLMs in generating *Gen-Text*. As anticipated earlier, we considered GPT, Llama, and Falcon, so the three model configurations will be referred to in the following as GPT_{RAG} , $\text{Llama}_{\text{RAG}}$, and $\text{Falcon}_{\text{RAG}}$. The performance of the proposed model guided by such LLM configurations has been compared to several IR *baseline models* from the literature. Each of these models utilizes distinct methodologies designed to enhance IR effectiveness by incorporating various layers of semantic understanding, topical relevance, and factual accuracy assessment, but none of them rely on the use of LLMs, nor RAG (as previously illustrated in Section 2).

- BM25: It is the state-of-the-art IR model that acts as the *foundational baseline* in our comparisons. This model accounts for *topical relevance only* and evaluates documents by scoring them according to the frequency of query terms, with adjustments made for document length;
- DigiLab: The model proposed in [53] employs a *re-ranking* approach that extends beyond the capabilities of BM25. Initially, BM25 is used to generate a preliminary set of potentially topically relevant documents. This is followed by a re-ranking stage utilizing a multi-model pipeline that evaluates documents on three key aspects: *usefulness*, assessed by Transformer-based models fine-tuned on the MS MARCO dataset [82]; *supportiveness*, evaluated by BERT models fine-tuned on scientific literature and Wikipedia; and *credibility*, determined by a Random Forest model trained on the Microsoft Credibility dataset [83]. The results from these evaluations are then integrated using Reciprocal Rank Fusion, which enhances the relevance and credibility of the retrieved documents;
- CiTIUS: Also the model presented in [52] employs *re-ranking*. Initially, BM25 ranks the top-100 documents according to their topical relevance to the query. The subsequent re-ranking involves a more in-depth analysis: RoBERTa is used for *semantic representation*, encoding sentences to assess their similarity to the query topic.

Additionally, a *reliability classifier* trained on historical data evaluates the trustworthiness of these passages. Finally, scores from these components are combined using techniques like CombSUM or *Borda Count*, producing the final document ranking;

- WISE: The model presented in [84] begins by retrieving documents using BM25 based on user queries. It then calculates a *truthfulness score* for these documents by cross-referencing their claims with supporting scientific articles, using cosine similarity to measure alignment. Finally, WISE re-ranks the documents by linearly combining the BM25 scores with the truthfulness scores;
- WISE_{NLI}: Building on the WISE model, WISE_{NLI} incorporates *Natural Language Inference* (NLI) techniques to improve the evaluation of document truthfulness. In addition to cross-referencing scientific articles, it assesses the stance of document content relative to claims found in the referenced scientific literature. *Stance scores* generated by NLI models quantify the degree of agreement or disagreement with established facts. The final document ranking is then refined by linearly combining BM25 scores, truthfulness scores, and stance scores.

4.4 Experimental Results

The results of the model evaluations tested are shown in this section with respect to the two datasets considered in terms of CAM_{MAP} and CAM_{NDCG} , two advanced IR evaluation metrics illustrated in detail in [85]. It is important to note that these results refer to: the optimal number of retrieved scientific articles, the optimal number k of topically relevant passages to be retrieved from PCM, and the optimal α and β values related to the calculation of factual accuracy and RSV respectively, as shown in Section 3.3, as determined to maximize the system’s performance.²³

The results illustrated in both Table 1 and Table 2 highlights the substantial improvements achieved by RAG-driven model configurations with respect to the considered baseline models. Traditional models like BM25 and more advanced systems such as DigiLab and CiTIUS provide foundational benchmarks, showing basic to moderate improvements in retrieval metrics. However, models that integrate knowledge bases and evidential reasoning, such as WISE and WISE_{NLI}, demonstrate significant performance enhancements, emphasizing the crucial role of semantic processing in improving document topical relevance and factual accuracy. However, the stand-out performers are the RAG-driven model configurations, which consistently excel in both the top-5 and top-10 document retrieval categories, especially with regard to the RAG model led by the Llama LLM. These models effectively leverage deep language understanding to deliver both contextually relevant and factually accurate results.

²³It was shown in [42] that excellent results are obtained by considering only the first highly relevant scientific article. The optimized values of the other three parameters were obtained by considering a subset of the data used for evaluating the proposed solution (5 queries from each dataset and their associated documents). The model’s effectiveness on this subset was assessed using the previously described metrics, evaluating both topical relevance and factual accuracy to achieve the best balance between the two. Through this process, we were able to obtain $k = 10$ (we tested from $k = 5$ to $k = 20$), $\alpha = 0.65$, and $\beta = 0.45$. Additionally, the queries used for parameter tuning were removed from the final testing to prevent overfitting and ensure a robust evaluation of the system on unseen data.

Table 1: Performance comparison of baselines and model configurations on the CLEF eHealth 2020 dataset for the top-5 and top-10 retrieved documents.

Model	CAM _{MAP}	CAM _{NDCG}	Embeddings
Top-5 Documents			
BM25	0.0431	0.1045	-
DigiLab	0.0433	0.1109	-
CiTIUS	0.0455	0.1119	-
WISE	0.0611	0.1198	BioBERT
WISE _{NLI}	0.0883	0.1823	BioBERT
GPT _{RAG}	0.1045	0.2098	BioBERT
<i>Llama</i> _{RAG}	0.1079	0.2146	BioBERT
<i>Falcon</i> _{RAG}	0.0994	0.2011	BioBERT
Top-10 Documents			
BM25	0.0784	0.1923	-
DigiLab	0.0823	0.1992	-
CiTIUS	0.0843	0.1999	-
WISE	0.1102	0.211	BioBERT
WISE _{NLI}	0.1302	0.2321	BioBERT
GPT _{RAG}	0.1502	0.2655	BioBERT
<i>Llama</i> _{RAG}	0.1532	0.2702	BioBERT
<i>Falcon</i> _{RAG}	0.1495	0.2568	BioBERT

4.5 Enhancing Explainability with GenText

In this section we briefly illustrate the possibility of using *GenText* as a means of increasing the explainability of the obtained search results, and provide an example. Indeed, the text generated by the LLM is crafted to answer user queries by providing concise, factual information backed by citations from scientific sources. Such information can be provided along with the search results to make it clearer to users which sources are reliable against which such documents have been judged factually accurate.

For example, in Figure 4, the query “Can 5G antennas cause COVID-19?” is answered with a explanation that highlights the absence of any scientific evidence linking 5G technology with the COVID-19 pandemic, thereby addressing a misinformation. The response not only refutes the misinformation but also enriches the user’s understanding by referencing relevant scientific literature that supports the facts.

5 Conclusion and Discussion

In this paper, we have introduced a novel approach to Health Information Retrieval through a Retrieval-Augmented Generation (RAG) model, which aims to tackle the challenge of effectively accounting for both topical relevance and factual accuracy in the retrieval process. The exponential growth of online health information and its widespread consumption by non-experts highlight, in fact, the urgent need for robust models capable of mitigating the risks associated with incurring online health misinformation. To this aim, our proposed three-stage model harnesses the strengths of

Table 2: Performance comparison of baselines and model configurations on the TREC Health Misinformation 2020 dataset for the top-5 and top-10 retrieved documents.

Model	CAM _{MAP}	CAM _{NDCG}	Embeddings
Top-5 Documents			
BM25	0.0631	0.1435	-
DigiLab	0.0712	0.1543	-
CiTIUS	0.0754	0.1554	-
WISE	0.0844	0.1608	BioBERT
WISE _{NLI}	0.0923	0.1922	BioBERT
GPT _{RAG}	0.1178	0.2234	BioBERT
<i>Llama</i> _{RAG}	0.1222	0.2298	BioBERT
<i>Falcon</i> _{RAG}	0.1123	0.2165	BioBERT
Top-10 Documents			
BM25	0.1047	0.2052	-
DigiLab	0.1186	0.2011	-
CiTIUS	0.1194	0.2095	-
WISE	0.1233	0.22	BioBERT
WISE _{NLI}	0.1341	0.2455	BioBERT
GPT _{RAG}	0.1547	0.2712	BioBERT
<i>Llama</i> _{RAG}	0.1602	0.2723	BioBERT
<i>Falcon</i> _{RAG}	0.1501	0.2665	BioBERT

generative Large Language Models (LLMs) to improve the retrieval of health-related documents that are both scientifically grounded and contextually relevant. The experimental results obtained illustrate that our RAG-driven model outperforms existing baseline models in producing ranked results that account for both dimensions of relevance. This enhancement relies in particular on our model’s ability to generate GenText, a rich and contextually aware text that serves not only as a benchmark for evaluating factual accuracy—an additional innovative aspect of our approach—but also as a valuable tool for promoting explainability.

The aspect of explainability is of particular interest in a model like the one proposed in this work. In fact, while our model notably enhances the balance between topicality and the likelihood of factual accuracy in the final ranked list of documents, it is essential to recognize that a potential risk of automation bias remains, linked to various factors. On one hand, RAG does not guarantee that LLMs will produce content that aligns solely with trusted knowledge bases. Moreover, there are limitations to our solution; for instance, the assessment of factual accuracy provides an approximation rather than a definitive measure of truth. This leaves ample room for future research, both in exploring introspection and reasoning mechanisms within LLMs and in investigating alternative solutions for calculating a value of factual accuracy to be incorporated into the RAG-based IR model. In addition, while our study has focused on generating GenText using three general-purpose LLMs—GPT, LLaMA, and Falcon—it would be beneficial to evaluate the application of domain-specific LLMs, such as Med-PaLM, optimized for answering clinical questions; MedMT5, which excels at

Search Results for "Can 5G antennas cause COVID-19"

LLM Generated Text

Based on the context provided, there is a misconception linking 5G antennas to the COVID-19 pandemic (Reference: [10316077](#)). However, this connection has no statistically significant evidence to support it (Reference: [10316077](#)). Instead, it's important to note that 5G networks play a crucial role in ensuring secure data handling and enhancing user privacy (Reference: [10255561](#)). Moreover, SARS-CoV-2 variants remain the main cause of COVID-19 outbreaks (Reference: [10288941](#)).

Reference List

- 10316077 - Softić, Adaleta, Elma Omeragić, Martin Kondža, Nahida Srabović, Aida Smajlović, Esmeralda Dautović, Nataša Bubić Pajić et al. "Knowledge and Attitudes regarding Covid-19 Vaccination among Medical and Non-medical Students in Bosnia and Herzegovina." Acta Medica Academica 52, no. 1 (2023): 1.
- 10255561 - Ullah, Yasir, Mardeni Bin Roslee, Sufian Mousa Mitani, Sajjad Ahmad Khan, and Mohamad Huzaimy Jusoh. "A survey on handover and mobility management in 5G HetNets: current state, challenges, and future directions." Sensors 23, no. 11 (2023): 5081.
- 10288941 - Soto, Ismael, Raul Zamorano-Illanes, Raimundo Becerra, Pablo Palacios Játiva, Cesar A. Azurdia-Meza, Wilson Alavia, Verónica García, Muhammad Ijaz, and David Zabala-Blanco. "A new COVID-19 detection method based on CSK/QAM visible light communication and machine learning." Sensors 23, no. 3 (2023): 1533.

Search Results

The conspiracy of Covid-19 and 5G: Spatial analysis

Conspiracy theories in general carry potentially serious public health risks, especially as anti-vaccination beliefs are already found to be

Evidence for a connection between coronavirus disease and 5G

We explore the scientific evidence suggesting a possible relationship between COVID-19 and radiofrequency radiation related to wireless...

5G Doesn't Cause COVID-19, But the Rumor It Does

People's fear of 5G technology is rational. Such technology does emit radiation, even if it's at low levels. But 5G isn't all that different ...

How the 5G coronavirus conspiracy theory went from fringe to mainstream

Despite what the internet might be telling you, cellphones did not cause the Covid-19 pandemic.

Is there a connection between coronavirus and 5G?

"5G mobile networks DO NOT spread COVID-19: viruses cannot travel on radio waves/mobile networks. COVID-19 is spreading in many countries that do

Fig. 4: Example of *GenText* used for explainability, addressing a query about the relationship between 5G technology and COVID-19 with references to scientific literature.

multilingual generation tasks; and DocOA, which integrates RAG techniques specifically for personalized osteoarthritis management, demonstrating the potential for tailored medical recommendations. Moreover, considering general-purpose models fine-tuned on medical literature could offer additional insights into optimizing factual accuracy and relevance in Health Information Retrieval.

Declarations

Author contribution. All authors contributed to the study's conception and design. Material preparation, data collection, and analysis were performed by Rishabh Upadhyay. The first draft of the manuscript was written by Rishabh Upadhyay. Marco Viviani commented on previous versions of the manuscript. All authors read and

approved the final manuscript.

Conflict of interest. The authors declare they have no financial interests.

Data availability. All data analysed during this study are included in the following published articles and their supplementary information files: [2, 30].

Ethical approval. Not applicable.

Funding. This work was supported by the Italian Ministry of University and Research (MUR) under PRIN 2022 project KURAMi: “Knowledge-based, explainable User empowerment in Releasing private data and Assessing Misinformation in online environments” (20225WTRFN).²⁴

References

- [1] Buchanan, J., Kock, N.: Information overload: A decision making perspective. In: Multiple Criteria Decision Making in the New Millennium: Proceedings of the Fifteenth International Conference on Multiple Criteria Decision Making (MCDM) Ankara, Turkey, July 10–14, 2000, pp. 49–58 (2001). Springer
- [2] Goeuriot, L., Suominen, H., Kelly, L., Alemany, L.A., Brew-Sam, N., Cotik, V., Filippo, D., Gonzalez Saez, G., Luque, F., Mulhem, P., *et al.*: CLEF eHealth Evaluation Lab 2021. In: Advances in Information Retrieval: 43rd European Conference on IR Research, ECIR 2021, Virtual Event, March 28–April 1, 2021, Proceedings, Part II 43, pp. 593–600 (2021). Springer
- [3] Goodrich, B., Rao, V., Liu, P.J., Saleh, M.: Assessing the factual accuracy of generated text. In: Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, pp. 166–175 (2019)
- [4] Di Sotto, S., Viviani, M.: Health misinformation detection in the social web: an overview and a data science approach. *International Journal of Environmental Research and Public Health* **19**(4), 2173 (2022)
- [5] Abdullah, M., Madain, A., Jararweh, Y.: Chatgpt: Fundamentals, applications and social impacts. In: 2022 Ninth International Conference on Social Networks Analysis, Management and Security (SNAMS), pp. 1–8 (2022). Ieee
- [6] Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., *et al.*: Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288* (2023)
- [7] Du, Z., Qian, Y., Liu, X., Ding, M., Qiu, J., Yang, Z., Tang, J.: Glm: General language model pretraining with autoregressive blank infilling. *arXiv preprint*

²⁴<https://kurami.disco.unimib.it/>

arXiv:2103.10360 (2021)

- [8] Ackerman, R., Balyan, R.: Automatic multilingual question generation for health data using llms. In: International Conference on AI-generated Content, pp. 1–11 (2023). Springer
- [9] Frisoni, G., Cocchieri, A., Presepi, A., Moro, G., Meng, Z.: To generate or to retrieve? on the effectiveness of artificial contexts for medical open-domain question answering. arXiv preprint arXiv:2403.01924 (2024)
- [10] Kell, G., Roberts, A., Umansky, S., Qian, L., Ferrari, D., Soboczenski, F., Wallace, B., Patel, N., Marshall, I.J.: Question answering systems for health professionals at the point of care—a systematic review. arXiv preprint arXiv:2402.01700 (2024)
- [11] Bang, Y., Cahyawijaya, S., Lee, N., Dai, W., Su, D., Wilie, B., Lovenia, H., Ji, Z., Yu, T., Chung, W., et al.: A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity. arXiv preprint arXiv:2302.04023 (2023)
- [12] Guo, B., Zhang, X., Wang, Z., Jiang, M., Nie, J., Ding, Y., Yue, J., Wu, Y.: How close is chatgpt to human experts? comparison corpus, evaluation, and detection. arXiv preprint arXiv:2301.07597 (2023)
- [13] Cao, M., Dong, Y., Wu, J., Cheung, J.C.K.: Factual error correction for abstractive summarization models. In: Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pp. 6251–6258. Association for Computational Linguistics, Online (2020)
- [14] Raunak, V., Menezes, A., Junczys-Dowmunt, M.: The curious case of hallucinations in neural machine translation. In: Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pp. 1172–1183. Association for Computational Linguistics, Online (2021)
- [15] Ji, Z., Lee, N., Frieske, R., Yu, T., Su, D., Xu, Y., Ishii, E., Bang, Y.J., Madotto, A., Fung, P.: Survey of hallucination in natural language generation. *ACM Comput. Surv.* **55**(12) (2023)
- [16] Saxena, S., Prasad, S., Prakash, M., Shankar, A., Vaddina, V., Gopalakrishnan, S., et al.: Minimizing factual inconsistency and hallucination in large language models. arXiv preprint arXiv:2311.13878 (2023)
- [17] Huang, L., Yu, W., Ma, W., Zhong, W., Feng, Z., Wang, H., Chen, Q., Peng, W., Feng, X., Qin, B., et al.: A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. arXiv preprint arXiv:2311.05232 (2023)

- [18] Zhang, Y., Li, Y., Cui, L., Cai, D., Liu, L., Fu, T., Huang, X., Zhao, E., Zhang, Y., Chen, Y., et al.: Siren’s song in the ai ocean: a survey on hallucination in large language models. arXiv preprint arXiv:2309.01219 (2023)
- [19] He, H., Zhang, H., Roth, D.: Rethinking with retrieval: Faithful large language model inference. arXiv preprint arXiv:2301.00303 (2022)
- [20] Li, X., Zhu, X., Ma, Z., Liu, X., Shah, S.: Are chatgpt and gpt-4 general-purpose solvers for financial text analytics? an examination on several typical tasks. arXiv preprint arXiv:2305.05862 (2023)
- [21] Shen, X., Chen, Z., Backes, M., Zhang, Y.: In chatgpt we trust? measuring and characterizing the reliability of chatgpt. arXiv preprint arXiv:2304.08979 (2023)
- [22] Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., Küttler, H., Lewis, M., Yih, W.-t., Rocktäschel, T., *et al.*: Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems* **33**, 9459–9474 (2020)
- [23] Perković, G., Drobnjak, A., Botički, I.: Hallucinations in llms: Understanding and addressing challenges. In: 2024 47th MIPRO ICT and Electronics Convention (MIPRO), pp. 2084–2088 (2024). IEEE
- [24] Borgeaud, S., Mensch, A., Hoffmann, J., Cai, T., Rutherford, E., Millican, K., Van Den Driessche, G.B., Lespiau, J.-B., Damoc, B., Clark, A., *et al.*: Improving language models by retrieving from trillions of tokens. In: *International Conference on Machine Learning*, pp. 2206–2240 (2022). PMLR
- [25] Guu, K., Lee, K., Tung, Z., Pasupat, P., Chang, M.: Retrieval augmented language model pre-training. In: *International Conference on Machine Learning*, pp. 3929–3938 (2020). PMLR
- [26] Izacard, G., Lewis, P., Lomeli, M., Hosseini, L., Petroni, F., Schick, T., Dwivedi-Yu, J., Joulin, A., Riedel, S., Grave, E.: Atlas: Few-shot learning with retrieval augmented language models. *Journal of Machine Learning Research* **24**(251), 1–43 (2023)
- [27] Setty, S., Jijo, K., Chung, E., Vidra, N.: Improving retrieval for rag based question answering models on financial documents. arXiv preprint arXiv:2404.07221 (2024)
- [28] Kang, C., Novak, D., Urbanova, K., Cheng, Y., Hu, Y.: Domain-specific improvement on psychotherapy chatbot using assistant. arXiv preprint arXiv:2404.16160 (2024)
- [29] Adlakha, V., BehnamGhader, P., Lu, X.H., Meade, N., Reddy, S.: Evaluating correctness and faithfulness of instruction-following models for question answering.

arXiv preprint arXiv:2307.16877 (2023)

- [30] Clarke, C.L., Maistro, M., Smucker, M.D., Zuccon, G.: Overview of the trec 2020 health misinformation track. In: TREC (2020)
- [31] Zhao, R., Arana-Catania, M., Zhu, L., Kochkina, E., Gui, L., Zubiaga, A., Procter, R., Liakata, M., He, Y.: Panacea: An automated misinformation detection system on covid-19. arXiv preprint arXiv:2303.01241 (2023)
- [32] Mendes, E., Chen, Y., Xu, W., Ritter, A.: Human-in-the-loop evaluation for early misinformation detection: A case study of covid-19 treatments. arXiv preprint arXiv:2212.09683 (2022)
- [33] Yue, Z., Zeng, H., Kou, Z., Shang, L., Wang, D.: Contrastive domain adaptation for early misinformation detection: A case study on covid-19. In: Proceedings of the 31st ACM International Conference on Information & Knowledge Management, pp. 2423–2433 (2022)
- [34] Jiang, G., Liu, S., Zhao, Y., Sun, Y., Zhang, M.: Fake news detection via knowledgeable prompt learning. *Information Processing & Management* **59**(5), 103029 (2022)
- [35] Chen, C., Shu, K.: Can llm-generated misinformation be detected? arXiv preprint arXiv:2309.13788 (2023)
- [36] Upadhyay, R., Pasi, G., Viviani, M.: Vec4cred: a model for health misinformation detection in web pages. *Multimedia Tools and Applications* **82**(4), 5271–5290 (2023)
- [37] Upadhyay, R., Pasi, G., Viviani, M.: Health misinformation detection in web content: A structural-, content-based, and context-aware approach based on web2vec. In: Proceedings of the Conference on Information Technology for Social Good. GoodIT '21, pp. 19–24. Association for Computing Machinery, New York, NY, USA (2021)
- [38] Upadhyay, R., Pasi, G., Viviani, M.: Leveraging socio-contextual information in bert for fake health news detection in social media. In: Proceedings of the 3rd International Workshop on Open Challenges in Online Social Networks. OASIS '23, pp. 38–46. Association for Computing Machinery, New York, NY, USA (2023)
- [39] Brand, E., Roitero, K., Soprano, M., Demartini, G., *et al.*: E-bart: Jointly predicting and explaining truthfulness. In: Proceedings of the Conference for Truth and Trust Online (2021)
- [40] Kou, Z., Shang, L., Zhang, Y., Wang, D.: Hc-covid: A hierarchical crowdsourced knowledge graph approach to explainable covid-19 misinformation detection.

- [41] Wu, J., Liu, Q., Xu, W., Wu, S.: Bias mitigation for evidence-aware fake news detection by causal intervention. In: Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 2308–2313 (2022)
- [42] Upadhyay, R., Knoth, P., Pasi, G., Viviani, M.: Explainable online health information truthfulness in consumer health search. *Frontiers in Artificial Intelligence* **6**, 1184851 (2023)
- [43] Shang, L., Zhang, Y., Yue, Z., Choi, Y., Zeng, H., Wang, D.: A knowledge-driven domain adaptive approach to early misinformation detection in an emergent health domain on social media. In: 2022 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM), pp. 34–41 (2022). IEEE
- [44] Hassan, N., Arslan, F., Li, C., Tremayne, M.: Toward automated fact-checking: Detecting check-worthy factual claims by claimbuster. In: Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp. 1803–1812 (2017)
- [45] Guo, Z., Schlichtkrull, M., Vlachos, A.: A survey on automated fact-checking. *Transactions of the Association for Computational Linguistics* **10**, 178–206 (2022)
- [46] Zeng, X., La Barbera, D., Roitero, K., Zubiaga, A., Mizzaro, S.: Combining large language models and crowdsourcing for hybrid human-ai misinformation detection. In: Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 2332–2336 (2024)
- [47] Goeriot, L., Suominen, H., Pasi, G., Bassani, E., Brew-Sam, N., González-Sáez, G., Kelly, L., Mulhem, P., Seneviratne, S., Upadhyay, R., *et al.*: Consumer health search at clef ehealth 2021. In: CEUR Workshop Proceedings, pp. 1–19 (2021). CEUR
- [48] Di Nunzio, G.M., Marchesin, S., Vezzani, F.: A Study on Reciprocal Ranking Fusion in Consumer Health Search. IMS UniPD ad CLEF eHealth 2020 Task 2. In: CLEF (Working Notes) (2020)
- [49] Cormack, G.V., Clarke, C.L., Buettcher, S.: Reciprocal rank fusion outperforms condorcet and individual rank learning methods. In: Proceedings of the 32nd International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 758–759 (2009)
- [50] Mulhem, P., Saez, G.G., Mannion, A., Schwab, D., Frej, J.: Lig-health at adhoc and spoken ir consumer health search: expanding queries using umls and fasttext.

In: CLEF 2020 (2020)

- [51] Seneviratne, S., Daskalaki, E., Hossain, M.Z., Lenskiy, A.: Sandidoc at clef 2020-consumer health search: Adhoc ir task. In: CLEF (Working Notes) (2020)
- [52] Fernández-Pichel, M., Losada, D.E., Pichel, J.C., Elswailer, D.: Citius at the trec 2020 health misinformation track. In: TREC (2020)
- [53] Zhang, B., Naderi, N., Jaume-Santero, F., Teodoro, D.: Ds4dh at trec health misinformation 2021: multi-dimensional ranking models with transfer learning and rank fusion. arXiv preprint arXiv:2202.06771 (2022)
- [54] Schlicht, I.B., Paula, A.F.M., Rosso, P.: Upv at trec health misinformation track 2021 ranking with sbert and quality estimators. arXiv preprint arXiv:2112.06080 (2021)
- [55] Abualsaud, M., Chen, I.X., Ghajar, K., Minh, L., Smucker, M., Tahami, A.V., Zhang, D.: Uwaterloomds at the trec 2021 health misinformation track. In: Proceedings of the Thirtieth REtrieval Conference Proceedings (TREC 2021). National Institute of Standards and Technology (NIST), Special Publication, pp. 1–18 (2021)
- [56] Pankaj, S., Gautam, A.: Augmented bio-sbert: Improving performance for pairwise sentence tasks in bio-medical domain. In: Proceedings of the Fifth Workshop on Technologies for Machine Translation of Low-Resource Languages (LoResMT 2022), pp. 43–47 (2022)
- [57] Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., Liu, P.J.: Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research* **21**(140), 1–67 (2020)
- [58] Wan, H., Feng, S., Tan, Z., Wang, H., Tsvetkov, Y., Luo, M.: Dell: Generating reactions and explanations for llm-based misinformation detection. arXiv preprint arXiv:2402.10426 (2024)
- [59] Choi, E.C., Ferrara, E.: Automated claim matching with large language models: empowering fact-checkers in the fight against misinformation. In: Companion Proceedings of the ACM on Web Conference 2024, pp. 1441–1449 (2024)
- [60] Cao, Y., Nair, A.M., Eyimife, E., Soofi, N.J., Subbalakshmi, K., Wullert II, J.R., Basu, C., Shallcross, D.: Can large language models detect misinformation in scientific news reporting? arXiv preprint arXiv:2402.14268 (2024)
- [61] Wang, J., Yang, Z., Yao, Z., Yu, H.: Jmlr: Joint medical llm and retrieval training for enhancing reasoning and professional question answering capability. arXiv preprint arXiv:2402.17887 (2024)

- [62] Khlaut, J., Dancette, C., Ferreres, E., Bennani, A., Hérent, P., Manceron, P.: Efficient medical question answering with knowledge-augmented question generation. arXiv preprint arXiv:2405.14654 (2024)
- [63] Shi, W., Min, S., Yasunaga, M., Seo, M., James, R., Lewis, M., Zettlemoyer, L., Yih, W.-t.: Replug: Retrieval-augmented black-box language models. arXiv preprint arXiv:2301.12652 (2023)
- [64] Ren, R., Wang, Y., Qu, Y., Zhao, W.X., Liu, J., Tian, H., Wu, H., Wen, J.-R., Wang, H.: Investigating the Factual Knowledge Boundary of Large Language Models with Retrieval Augmentation (2023)
- [65] Izacard, G., Grave, E.: Leveraging passage retrieval with generative models for open domain question answering. In: Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, pp. 874–880. Association for Computational Linguistics, Online (2021)
- [66] Trivedi, H., Balasubramanian, N., Khot, T., Sabharwal, A.: 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), pp. 10014–10037. Association for Computational Linguistics, Toronto, Canada (2023)
- [67] Li, D., Rawat, A.S., Zaheer, M., Wang, X., Lukasik, M., Veit, A., Yu, F., Kumar, S.: Large language models with controllable working memory. In: Findings of the Association for Computational Linguistics: ACL 2023, pp. 1774–1793. Association for Computational Linguistics, Toronto, Canada (2023)
- [68] Cai, D., Wang, Y., Bi, W., Tu, Z., Liu, X., Lam, W., Shi, S.: Skeleton-to-response: Dialogue generation guided by retrieval memory. In: Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pp. 1219–1228. Association for Computational Linguistics, Minneapolis, Minnesota (2019)
- [69] Cai, D., Wang, Y., Bi, W., Tu, Z., Liu, X., Shi, S.: Retrieval-guided dialogue response generation via a matching-to-generation framework. In: Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pp. 1866–1875. Association for Computational Linguistics, Hong Kong, China (2019)
- [70] Peng, B., Galley, M., He, P., Cheng, H., Xie, Y., Hu, Y., Huang, Q., Liden, L., Yu, Z., Chen, W., et al.: Check your facts and try again: Improving large language models with external knowledge and automated feedback. arXiv preprint arXiv:2302.12813 (2023)

- [71] Cui, J., Li, Z., Yan, Y., Chen, B., Yuan, L.: Chatlaw: Open-source legal large language model with integrated external knowledge bases. arXiv preprint arXiv:2306.16092 (2023)
- [72] Zhou, S., Alon, U., Xu, F.F., Jiang, Z., Neubig, G.: Docprompting: Generating code by retrieving the docs. In: The Eleventh International Conference on Learning Representations (2023)
- [73] Chern, I., Chern, S., Chen, S., Yuan, W., Feng, K., Zhou, C., He, J., Neubig, G., Liu, P., et al.: Factool: Factuality detection in generative ai—a tool augmented framework for multi-task and multi-domain scenarios. arXiv preprint arXiv:2307.13528 (2023)
- [74] Pan, L., Wu, X., Lu, X., Luu, A.T., Wang, W.Y., Kan, M.-Y., Nakov, P.: Fact-checking complex claims with program-guided reasoning. arXiv preprint arXiv:2305.12744 (2023)
- [75] Zhang, X., Gao, W.: Towards llm-based fact verification on news claims with a hierarchical step-by-step prompting method. arXiv preprint arXiv:2310.00305 (2023)
- [76] Zeng, F., Gao, W.: Justilm: Few-shot justification generation for explainable fact-checking of real-world claims. Transactions of the Association for Computational Linguistics **12**, 334–354 (2024)
- [77] Robertson, S., Zaragoza, H., et al.: The probabilistic relevance framework: Bm25 and beyond. Foundations and Trends® in Information Retrieval **3**(4), 333–389 (2009)
- [78] Upadhyay, R., Pasi, G., Viviani, M.: A passage retrieval transformer-based re-ranking model for truthful consumer health search. In: Joint European Conference on Machine Learning and Knowledge Discovery in Databases, pp. 355–371 (2023). Springer
- [79] Lee, J., Yoon, W., Kim, S., Kim, D., Kim, S., So, C.H., Kang, J.: Biobert: a pre-trained biomedical language representation model for biomedical text mining. Bioinformatics **36**(4), 1234–1240 (2020)
- [80] Izacard, G., Grave, E.: Leveraging passage retrieval with generative models for open domain question answering. arXiv preprint arXiv:2007.01282 (2020)
- [81] Phan, L.N., Anibal, J.T., Tran, H., Chanana, S., Bahadroglu, E., Peltekian, A., Altan-Bonnet, G.: Scifive: a text-to-text transformer model for biomedical literature. arXiv preprint arXiv:2106.03598 (2021)
- [82] Bajaj, P., Campos, D., Craswell, N., Deng, L., Gao, J., Liu, X., Majumder, R., McNamara, A., Mitra, B., Nguyen, T., et al.: MS MARCO: A human generated

machine reading comprehension dataset. arXiv preprint arXiv:1611.09268 (2016)

- [83] Schwarz, J., Morris, M.: Augmenting web pages and search results to support credibility assessment. In: Proceedings of the SIGCHI Conference on Human Factors in Computing Systems, pp. 1245–1254 (2011)
- [84] Upadhyay, R., Pasi, G., Viviani, M.: An unsupervised approach to genuine health information retrieval based on scientific evidence. In: International Conference on Web Information Systems Engineering, pp. 119–135 (2022). Springer
- [85] Lioma, C., Simonsen, J.G., Larsen, B.: Evaluation measures for relevance and credibility in ranked lists. In: Proceedings of the ACM SIGIR International Conference on Theory of Information Retrieval, pp. 91–98 (2017)