

---

# A Survey of Large Language Models and Temporal Knowledge Graphs in Retrieval-Augmented Generation and Temporal Reasoning

---

[www.surveyx.cn](http://www.surveyx.cn)

## Abstract

Large Language Models (LLMs), Temporal Knowledge Graphs (TKGs), and Retrieval-Augmented Generation (RAG) systems represent cutting-edge advancements in computational linguistics and artificial intelligence. These technologies collectively enhance information retrieval and question-answering systems by incorporating time-sensitive context and reasoning capabilities. This survey explores the integration of LLMs with TKGs and RAG, emphasizing their transformative potential in managing domain-specific and temporally complex queries. The study highlights significant improvements in retrieval accuracy and reasoning efficiency, as demonstrated by frameworks like Topo-RAG and enhancements in Elasticsearch-based RAG systems. Despite these advancements, challenges persist in scalability, retrieval quality, and bias mitigation. The survey identifies the need for continued research into optimizing retrieval methodologies, integrating structured knowledge, and addressing ethical considerations. Future directions include refining temporal reasoning algorithms, developing comprehensive benchmarks, and exploring novel applications across diverse domains. By addressing these challenges, the integration of LLMs, TKGs, and RAG systems can further advance the capabilities of information retrieval technologies, enabling more precise and contextually aware responses to user queries in dynamic environments.

## 1 Introduction

### 1.1 Significance of Large Language Models

Large Language Models (LLMs) are integral to modern computational linguistics and artificial intelligence, significantly advancing natural language processing (NLP) capabilities for applications like automated text generation, language translation, and sentiment analysis. Models such as OpenAI's GPT series and Google's BERT have shown remarkable performance improvements, driven by large datasets and enhanced computational power. Despite these advancements, challenges persist, including the tendency to generate inaccurate outputs in specialized contexts and issues related to hallucinations. Researchers are increasingly integrating knowledge graphs to provide structured, contextually relevant information that enhances the reliability of LLM outputs, paving the way for more trustworthy AI applications [1, 2, 3, 4]. By leveraging vast datasets, LLMs create coherent, contextually relevant text, impacting fields like healthcare and information retrieval.

The significance of LLMs is further emphasized by their potential to democratize access to specialized knowledge, such as medical information, effectively addressing disparities in healthcare delivery [5]. However, their reliance on static training data often results in outdated responses, particularly in dynamic sectors, and their limited ability to perform multi-hop reasoning within extensive textual contexts due to pre-defined context lengths presents additional challenges [6].

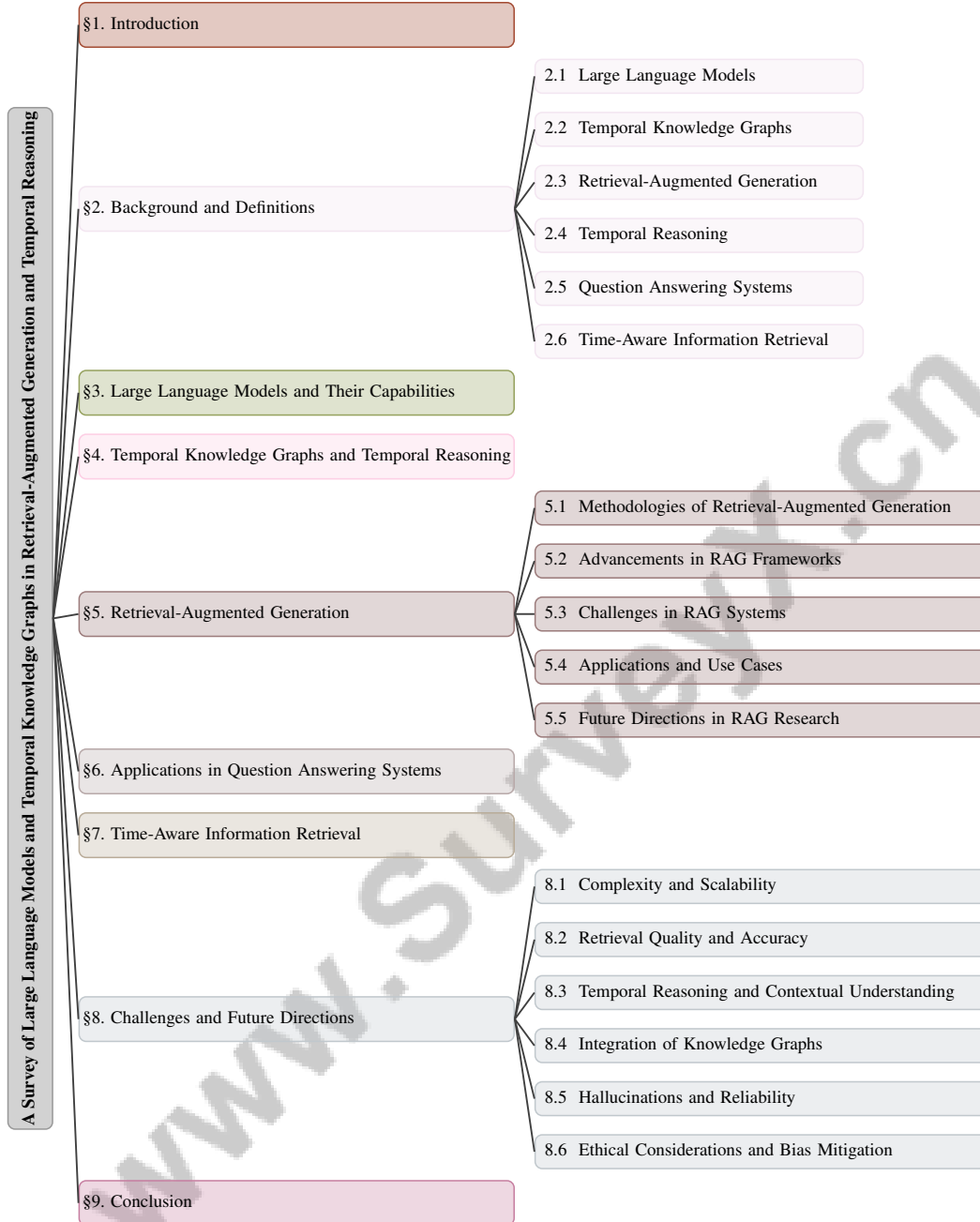


Figure 1: chapter structure

To mitigate these challenges, Retrieval-Augmented Generation (RAG) systems have been proposed, effectively integrating retrieval mechanisms with generative language models to enhance output accuracy and reduce hallucinations, especially in specialized domains. RAG techniques improve response quality by incorporating up-to-date information and streamline knowledge-intensive tasks, despite concerns regarding their complex implementation and response times [7, 8]. By integrating domain-specific and time-sensitive data, RAG systems enhance the factual accuracy and contextual relevance of LLM outputs, optimizing knowledge utilization during inference.

The ongoing development of LLMs in specialized fields underscores their potential for advanced data analysis and knowledge discovery [5]. Future research should focus on the synergies between LLMs and dynamic knowledge bases, such as knowledge graphs, to address existing limitations in knowledge representation and NLP [9].

---

## 1.2 Integration with Temporal Knowledge Graphs

Integrating Large Language Models (LLMs) with Temporal Knowledge Graphs (TKGs) is a crucial advancement for enhancing temporal reasoning and information retrieval capabilities. TKGs offer a structured framework for time-sensitive data, addressing LLM limitations such as hallucinations and static knowledge reliance by grounding outputs in accurate, current information. This integration is vital for overcoming challenges related to temporal concepts and logic, essential for tasks like historical event analysis and time-sensitive information retrieval [10].

Employing Knowledge Graph-based Retrieval-Augmented Generation (RAG) methodologies enhances LLM performance by enabling dynamic integration of external knowledge, thereby improving accuracy and contextual relevance in generated responses. For example, the Topology-aware Retrieval-augmented Generation (Topo-RAG) framework incorporates topological relationships from knowledge graphs to enhance text generation [11]. This ensures that LLM outputs are grounded in structured knowledge, thereby improving their temporal reasoning capabilities.

Frameworks like SynapticRAG, which integrates synaptic dynamics into RAG, enhance memory retrieval by considering temporal relevance [12]. The DRAGON-AI framework exemplifies the potential of integrating LLMs and RAG for automating ontology term generation, definitions, and relationships, further illustrating the benefits of LLM and TKG integration [13].

Incorporating TKGs into LLMs not only bolsters temporal reasoning but also significantly enhances performance in time-aware information retrieval tasks. This integration paves the way for advanced, reliable question-answering systems capable of providing precise, contextually aware responses in temporally complex environments. As research progresses, exploring synergies between LLMs and dynamic knowledge bases will be essential for overcoming existing limitations in knowledge representation and NLP [1].

## 1.3 Relevance in Question Answering and Information Retrieval

The integration of Large Language Models (LLMs) with Temporal Knowledge Graphs (TKGs) has significantly advanced question answering (QA) systems and information retrieval, particularly for domain-specific and time-sensitive queries. Retrieval-Augmented Generation (RAG) plays a crucial role in this integration by mitigating information hallucination, where LLMs may generate fluent yet factually incorrect responses. By incorporating domain-specific knowledge through RAG, LLMs can notably enhance the factual accuracy of their outputs. Studies underscore RAG's effectiveness in reducing hallucinations and introducing new knowledge, although its capacity to assist in deeper reasoning tasks remains limited. Furthermore, flawed information during the retrieval phase can compromise the reliability of generated content. Recent advancements, such as the Credibility-aware Generation (CAG) framework, aim to enhance model performance by enabling credibility assessment of retrieved information, thus improving resilience against noisy data [14, 15].

Despite progress in knowledge-augmented methods like RAG and Generation-Augmented Generation (GAG), significant challenges persist in their application to QA tasks, particularly with complex, long-context queries and ensuring robust responses to private enterprise documents. The effectiveness of RAG is influenced by the quality of text retrieval, often hindered by inaccuracies in parsing professional documents in formats like PDFs. Innovations such as enhanced PDF structure recognition and synthetic question generation have shown promise in improving retrieval precision and answer depth. However, achieving comprehensive understanding and synthesis from large, diverse document sets remains a critical challenge for optimizing these systems in real-world applications [16, 8, 17]. Innovative frameworks like THREAD utilize a new knowledge granularity called 'logic unit' to organize information into structured, interconnected units, enhancing RAG's capacity to process complex queries.

In Open-Domain Question-Answering (ODQA), LLMs integrated with RAG methodologies exhibit improved performance, particularly in managing long contexts and reasoning abilities [18]. Benchmark systems comparing RAG and fine-tuning techniques have yielded valuable insights for optimizing LLM-based knowledge systems [19].

Adaptive Contrastive Decoding (ACD) offers a promising technique, outperforming baseline methods in handling noisy contexts while maintaining robust performance in accurate contexts. The inte-

---

gration of retriever and memory-adaptive note-enhanced RAG systems addresses factual errors and hallucinated outputs, further enhancing LLM reliability in open-domain QA tasks [20].

The combined integration of LLMs with TKGs, enhanced by RAG techniques and innovative methodologies, has significantly improved the performance of QA systems and information retrieval processes. This synergy enables sourcing highly relevant text chunks, enhancing accuracy and contextual relevance of responses to user queries. Notably, advancements in RAG, such as sophisticated chunking techniques, query expansion, and modules like the Query Rewriter and Knowledge Filter, address challenges like irrelevant knowledge and ambiguous queries, refining LLM inputs. These enhancements streamline retrieval processes and bolster response reliability across diverse datasets, ultimately leading to a more effective user experience [21, 22].

## 1.4 Structure of the Survey

This survey provides a comprehensive examination of the integration of Large Language Models (LLMs) with Temporal Knowledge Graphs (TKGs) and their applications in Retrieval-Augmented Generation (RAG) and temporal reasoning. The introduction underscores the critical role of LLMs and TKGs, emphasizing their combined potential to enhance retrieval-augmented generation capabilities and improve temporal reasoning. This integration facilitates accurate, contextually relevant responses and addresses challenges related to data retrieval and user intent interpretation, leading to reliable applications across fields requiring real-time information processing and domain-specific knowledge [22, 23, 24, 25].

Subsequent sections delve into core concepts and definitions, providing foundational knowledge on LLMs, TKGs, RAG, temporal reasoning, QA systems, and time-aware information retrieval. The survey continues by analyzing the development, strengths, and limitations of LLMs and their applications in temporal reasoning, followed by an exploration of TKGs, challenges in temporal reasoning, and future research directions.

Further sections focus on the application of LLMs, TKGs, and RAG in QA systems, discussing evaluation, educational applications, enhancing factual accuracy, and innovations in complex question answering. The importance of incorporating temporal context in information retrieval systems is addressed, along with challenges and future directions in time-aware retrieval systems.

The survey concludes by identifying pressing challenges in integrating LLMs with TKGs and RAG frameworks, highlighting critical future research directions such as addressing complexity and scalability in LLM applications, improving retrieval quality for enhanced response accuracy, developing robust temporal reasoning capabilities, optimizing knowledge graph integration, mitigating hallucinations in outputs, and considering ethical implications for technology deployment across various domains [23, 21, 22, 25, 26]. The following sections are organized as shown in Figure 1.

## 2 Background and Definitions

### 2.1 Large Language Models

Large Language Models (LLMs) are advanced systems that generate human-like text using extensive datasets and sophisticated neural architectures. They operate on probabilistic principles, learning patterns from training data to produce coherent outputs. However, factual accuracy and reliability remain challenges, particularly in specialized domains due to static training data, which can lead to inaccuracies and hallucinations [1, 11, 27]. Despite these issues, LLMs excel in formal language understanding, although their generative capabilities are less robust.

In information retrieval and question answering, LLMs show promise, especially when integrated with structured knowledge like knowledge graphs. This integration aids in tackling complex queries by elucidating implicit relationships, proving advantageous in tasks like temporal event forecasting and reasoning [13]. LLMs process partial term inputs by retrieving relevant information from ontologies and knowledge sources, as demonstrated in systems like DRAGON-AI [13]. Temporal Complex Events (TCEs) benchmarks assess LLMs' proficiency in managing temporal dynamics and understanding extensive text, facilitating performance comparisons among models [28, 3].

Integrating LLMs with Retrieval-Augmented Generation (RAG) methodologies enhances their capabilities, particularly in multilingual contexts where diverse global knowledge presents challenges

---

[29]. Though effective benchmarks for multilingual RAG are lacking, this research avenue is crucial for maximizing LLMs’ potential across various linguistic environments. LLMs are foundational to this survey, offering essential capabilities for advanced natural language processing tasks. Their integration with external knowledge sources and specialized frameworks broadens their applicability in dynamic and complex domains, necessitating continuous performance evaluations, especially for predicting future events based on historical data [1].

## 2.2 Temporal Knowledge Graphs

Temporal Knowledge Graphs (TKGs) extend traditional knowledge graphs by incorporating temporal information, enabling the modeling of dynamic relationships over time. They capture temporal dynamics absent in static representations, fostering a nuanced understanding of time-sensitive information [30]. TKGs address challenges in temporal reasoning, notably the lack of structured temporal representation, which impedes effective temporal task resolution [31].

TKGs are vital for temporal knowledge forecasting, predicting future events based on historical data, and enhancing Temporal Knowledge Graph Reasoning (TKGR) by leveraging temporal information to derive new knowledge about complex relations [32, 31]. Despite their advantages, TKGs face challenges related to complexity and incompleteness, alongside LLMs’ tendencies to generate hallucinations [33]. The development of knowledge graph-enhanced large language models (KGLLMs) aims to integrate structured temporal data with language models, enhancing reasoning capabilities and factual accuracy [34].

Efforts like the SDAAP dataset, an open-source textual knowledge dataset, improve spectral analysis through annotated literature and knowledge instruction data, enhancing TKG applications for better temporal reasoning [35]. Bridging static and dynamic knowledge representations, TKGs are pivotal in advancing temporal reasoning and information retrieval, enabling precise, contextually aware responses to temporally complex queries.

## 2.3 Retrieval-Augmented Generation

Retrieval-Augmented Generation (RAG) enhances LLM capabilities by integrating retrieval mechanisms with generative processes, allowing access to up-to-date information from external sources. This integration addresses static dataset limitations and constrained context windows, improving factual accuracy and contextual relevance in knowledge-intensive tasks where hallucinations and inaccuracies are prevalent [12].

Advanced retrieval techniques, such as Graph Neural Network-enhanced Retrieval (GNN-Ret), optimize retrieval by constructing graphs of related passages, enhancing retrieval quality [36]. The Chain-of-Verification Retrieval Augmented Generation (CoV-RAG) framework introduces a verification module that iteratively refines knowledge retrieval and response generation, enhancing robustness [37]. Dynamic-Relevant Retrieval-Augmented Generation (DR-RAG) improves document retrieval accuracy and efficiency in multi-hop question-answering systems, demonstrating RAG’s capability to manage complex queries effectively [38].

Despite advancements, RAG systems face challenges retrieving irrelevant or noisy information, affecting LLM output quality. Optimizing retrieval and refining content before integration into prompts is crucial [39]. Benchmarks like PaperQA, a RAG agent for answering scientific questions, enhance LLM output quality and reliability [40]. RAG bridges the gap between static knowledge and dynamic information retrieval, enabling accurate, contextually aware responses in complex information retrieval and question-answering tasks. Developing standardized frameworks and efficient knowledge caching strategies will further enhance RAG systems, addressing modern information retrieval demands [41].

## 2.4 Temporal Reasoning

Temporal reasoning is crucial for processing time-sensitive queries, enabling systems to understand chronological sequences and causal relationships, particularly in scenarios involving TKGs and LLMs. Challenges in temporal reasoning arise from current methodologies’ limitations in interpreting temporal data [31]. LLMs often struggle with reasoning about temporal relationships, especially in tabular data, limiting real-world applications [42].

---

This issue is exacerbated by the complexity of Temporal Complex Events (TCEs), necessitating precise event-timestamp correlation and understanding of chronological and causal connections [28]. Benchmarks evaluate models’ understanding and reasoning about time-related events, providing a structured approach to assessing temporal reasoning capabilities [43]. Empirical studies emphasize context’s importance in temporal reasoning, revealing LLMs’ struggles with processing time-related information [10]. Enhancing LLMs’ temporal reasoning abilities is critical for dynamic domains, where accurate, contextually aware responses are paramount.

## 2.5 Question Answering Systems

Question Answering (QA) systems deliver precise, contextually relevant answers by leveraging extensive datasets and sophisticated algorithms. The integration of LLMs and TKGs has significantly enhanced open-domain and domain-specific contexts. RAG frameworks improve QA systems’ ability to manage complex queries by combining retrieval mechanisms with LLMs, enhancing response accuracy and relevance. Modules like the Query Rewriter and Knowledge Filter address query ambiguity and irrelevant knowledge, while multi-query generation overcomes information plateaus. Frameworks like DR-RAG optimize retrieval by improving document relevance and efficiency, enhancing performance [16, 8, 21, 38].

In open-domain QA tasks, efficiently processing and retrieving information from lengthy contexts is challenging due to input length constraints. RAG frameworks address these issues by embedding retrieval mechanisms, allowing LLMs to dynamically access relevant external knowledge, enhancing contextual relevance and factual accuracy. This is critical for Knowledge Base Question Answering (KBQA), which answers natural language questions based on facts stored in knowledge bases [27]. Adaptive contrastive decoding techniques maintain performance even with noisy contextual data.

Domain-specific QA systems benefit from LLMs and RAG integration, requiring precise, contextually tailored responses for specialized fields. For example, existing QA systems for Adobe products struggle to retrieve relevant information due to a lack of domain-specific training and frequent updates [18]. LLMs and RAG frameworks enhance these systems’ adaptability to dynamic domain-specific knowledge.

Deploying QA applications over proprietary documents introduces challenges like data security and response reliability, crucial for enterprise environments where data sensitivity and accuracy are paramount. Evolving named entities pose challenges for RAG frameworks, emphasizing advanced entity resolution mechanisms for accurately linking linguistic expressions to dynamic entities, improving retrieval efficiency and reducing hallucinations [7, 8, 44].

Incorporating knowledge graph-based retrieval techniques enhances QA systems by emphasizing domain-specific information, improving answer accuracy and quality. The HybridRAG approach integrates knowledge graphs with vector retrieval methods for complex information extraction from unstructured financial documents, outperforming traditional methods. The SubgraphRAG framework optimizes structured knowledge retrieval to support LLMs, minimizing hallucinations and grounding responses in relevant data. Customer service applications using knowledge graphs maintain historical issue structures, enhancing retrieval accuracy and reducing resolution times. These advancements illustrate how integrating knowledge graphs into QA systems leads to precise, contextually relevant outputs across domains [45, 46, 47]. Continuous development and refinement of QA systems, driven by LLMs, TKGs, and RAG integration, will advance natural language processing and information retrieval capabilities across diverse domains.

## 2.6 Time-Aware Information Retrieval

Time-Aware Information Retrieval (TAIR) emphasizes incorporating temporal context to enhance the relevance and accuracy of retrieved information, crucial in dynamic domains where time significantly influences relevance [10]. Integrating temporal knowledge within retrieval systems addresses traditional models’ limitations, which often overlook time’s significance in shaping context and meaning [31].

TAIR’s relevance lies in leveraging TKGs and LLMs to enhance retrieval systems’ temporal reasoning capabilities. By incorporating structured temporal data, TAIR systems handle time-sensitive queries, providing contextually appropriate and temporally relevant responses [32]. This is important for

---

applications requiring understanding of temporal sequences and causal relationships, like historical event analysis and future event forecasting [28].

A key challenge in TAIR systems is integrating temporal data without compromising efficiency and accuracy. Advanced methodologies, such as graph neural networks and dynamic retrieval frameworks, optimize retrieval, ensuring effective utilization of temporal context without introducing noise or irrelevant information. These approaches illustrate TAIR’s potential to enhance LLMs and TKGs, improving their ability to deliver accurate, contextually relevant responses to complex queries involving temporal elements. TAIR optimizes retrieval by refining information sourcing and integration, addressing challenges like ambiguous queries and irrelevant knowledge. Through advanced techniques like sophisticated chunking, query expansion, and enhanced filtering, TAIR ensures LLMs and TKGs better understand and respond to time-sensitive inquiries, leading to higher reliability and user satisfaction [21, 22].

Developing benchmarks for evaluating temporal reasoning capabilities underscores TAIR’s importance in advancing information retrieval technologies. By assessing models’ performance in handling time-related queries, researchers identify improvement areas and drive innovation [43]. As demand for time-sensitive information grows, integrating TAIR methodologies within LLMs and TKGs will shape the future of information retrieval and question-answering systems.

### 3 Large Language Models and Their Capabilities

#### 3.1 Development and Evolution of LLMs

The progression of Large Language Models (LLMs) is marked by significant strides in computational linguistics and AI, aimed at refining human-like text processing and generation. This evolution encompasses phases like pre-training, domain adaptation, reinforcement learning from human feedback, and retrieval-augmented generation [1]. Initially, LLMs focused on pre-training large networks on extensive datasets to capture linguistic patterns, forming a basis for coherent text generation. Domain adaptation further honed LLMs for specific fields, enhancing specialized performance [1].

Reinforcement learning from human feedback introduced adaptability, allowing models to refine responses based on user interactions, thus improving contextual accuracy in dynamic settings [1]. Retrieval-Augmented Generation (RAG) marks a pivotal advance, addressing static dataset limitations by integrating external knowledge, thereby enhancing factual accuracy and contextual relevance [25, 12]. Figure 2 illustrates the development and evolution of LLMs, highlighting key training phases, advanced techniques, and future directions. It showcases the progression from foundational training methods to advanced retrieval-augmented strategies and lifelong learning, emphasizing ongoing research into dynamic editing and retrieval improvements. Emerging frameworks like GraphRAG leverage graph databases to optimize retrieval and generation, highlighting graph-based indexing’s role in complex information tasks [48].

Lifelong learning strategies further contribute to LLM evolution, structuring understanding and categorizing learning scenarios [49]. Future research will likely explore dynamic in-context editing and improved retrieval strategies to enhance multi-hop reasoning and external information integration. Recent approaches have shown improved performance in open-domain questions by effectively combining retrieved passages with LLMs, surpassing traditional methods [6, 24].

#### 3.2 Strengths of LLMs

LLMs have revolutionized natural language processing, significantly enhancing task efficiency and accuracy. Their integration with knowledge graphs (KGs) improves interpretability and factual accuracy in applications like question answering and recommendation systems [33, 50]. The combination with Retrieval-Augmented Generation (RAG) technology further customizes and enhances AI response accuracy, addressing limitations of rule-based systems and enabling dynamic, context-aware interactions [20]. Tailored retrieval mechanisms enhance multi-modal reasoning, allowing LLMs to excel in integrating diverse data types [51].

In domain-specific applications, LLMs adapt to rapid changes, delivering timely, relevant responses crucial for user satisfaction [18]. Automation capabilities, exemplified by systems like DRAGON-AI, illustrate efficiency gains by reducing manual workloads [13]. LLMs manage extensive narratives

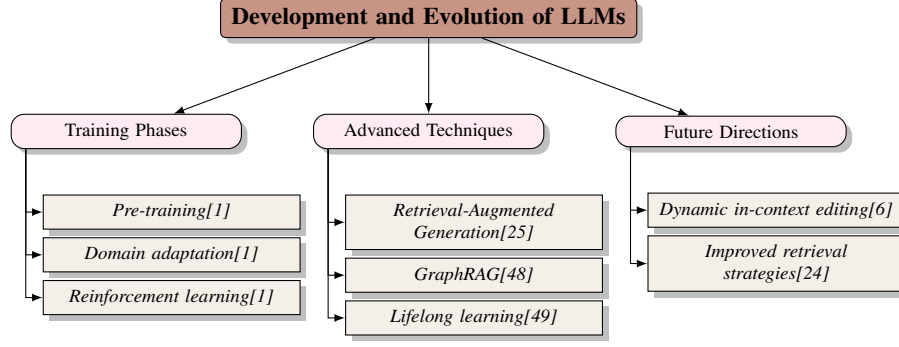


Figure 2: This figure illustrates the development and evolution of Large Language Models (LLMs), highlighting key training phases, advanced techniques, and future directions. It showcases the progression from foundational training methods to advanced retrieval-augmented strategies and lifelong learning, emphasizing ongoing research into dynamic editing and retrieval improvements.

effectively, with retrieval methods playing a critical role in handling complex information [28]. The ARM-RAG framework enhances performance without extensive retraining, reducing computational costs [52].

In education, LLMs surpass human-generated content in clarity and usability, as seen in models like ChatGLM [53]. Systems like PaperQA achieve expert-level accuracy with fewer errors [40]. These strengths position LLMs as essential tools in advancing natural language processing and information retrieval across various applications.

### 3.3 Limitations and Challenges

LLMs face significant challenges, including hallucinations where models produce fluent but factually incorrect outputs due to static datasets lacking current information [5]. Creating domain-specific expertise without extensive resources results in unreliable outputs in specialized fields [5]. Training data biases complicate ethical deployment, affecting fairness and accuracy [1].

RAG methodologies, while beneficial, introduce challenges like scattered relevant information across documents, complicating retrieval [52]. LLM performance varies across languages, with high-resource languages outperforming low-resource ones, necessitating improved multilingual benchmarks [29]. In reasoning tasks, LLMs struggle with complex scenarios due to difficulties in retrieving relevant reasoning chains, leading to suboptimal performance [52]. Context window limits restrict LLMs' ability to utilize all relevant information, affecting reasoning and accuracy, especially in domain-specific contexts [5].

Ongoing research is crucial to improve LLM factual accuracy, adaptability, and reasoning. Advancements like RAG show promise by integrating external knowledge and assisting multi-hop reasoning, though challenges persist, including preprocessing to filter irrelevant information and inherent constraints in deep reasoning tasks. Innovative approaches like dynamic in-context editing and frameworks like FRAMES aim to bridge these gaps, facilitating complex information integration [54, 6, 55, 15]. Developing specialized benchmarks and methodologies will be vital for evaluating and improving LLM performance, ensuring reliability and relevance in complex environments.

### 3.4 Applications in Temporal Reasoning

LLMs are increasingly recognized for their potential in temporal reasoning, where interpreting time-related information is crucial for managing complex queries. Recent advancements focus on enhancing LLMs' temporal reasoning capabilities through sophisticated strategies and external knowledge sources. The Dynamic In-Context Editing (DICE) method exemplifies these efforts by enabling dynamic information editing and retrieval for effective multi-hop reasoning [6].

Multi-modal reasoning frameworks utilizing diverse examples improve LLM accuracy and adaptability in temporal tasks [51]. Frameworks like THINKON-GR1 combine knowledge graphs and text-based retrieval for deep, faithful reasoning, enhancing temporal query management [56]. The



---

development of frameworks like ToG-2, evaluated on multiple benchmarks, demonstrates LLM effectiveness in complex reasoning tasks [56].

As research progresses, ongoing development of innovative frameworks and methodologies will be essential for effective applications in dynamic environments. Recent studies highlight challenges in processing tabular data and performing multi-hop reasoning within extensive contexts. Enhancements to datasets like TempTabQA and approaches like C.L.E.A.R show promise in improving LLM performance. Techniques like RAG and the Retrieve-Plan-Generation (RPG) framework illustrate potential for integrating external knowledge to bolster reasoning, addressing limitations related to noise and context length. Continuous updates and adaptations will be critical for maintaining LLM effectiveness in real-world applications [57, 15, 42, 6, 58]. Integrating dynamic retrieval and reasoning strategies will further enhance LLMs’ effectiveness in managing temporal information, ensuring relevance and reliability in addressing temporally complex queries.

## 4 Temporal Knowledge Graphs and Temporal Reasoning

### 4.1 Structure and Function of Temporal Knowledge Graphs

Temporal Knowledge Graphs (TKGs) are sophisticated data structures designed to model dynamic relationships and events over time, providing a robust framework for temporal reasoning and information retrieval. Their architecture seamlessly integrates temporal data with language models and retrieval-augmented generation systems, optimizing the processing of time-sensitive information [30]. A pivotal component of this architecture is the G-Indexing stage, which constructs and indexes graph databases to organize temporal data efficiently, thus facilitating effective retrieval and reasoning [48].

As illustrated in Figure 3, the structure and integration methods of TKGs are depicted, highlighting key components such as G-Indexing, G-Retrieval, and G-Generation. This figure also presents the integration of Graph Neural Networks (GNNs) and TPP-LLM-Embedding models, alongside challenges and advancements in temporal reasoning, retrieval performance, and predictive accuracy.

Sophisticated retrieval methods bolster the operational mechanisms of TKGs. The G-Retrieval method, for instance, extracts contextually relevant and temporally accurate information from user queries, crucial for generating precise responses in time-sensitive domains [48]. The subsequent G-Generation stage synthesizes responses using retrieved graph data, producing coherent outputs that incorporate temporal dynamics [59].

The integration of Graph Neural Networks (GNNs) with TKGs, as exemplified by the G-Retriever framework, enhances the understanding of complex temporal relationships within textual graphs. This approach leverages the strengths of GNNs, LLMs, and RAG to facilitate effective question answering, showcasing the potential of retrieval-augmented methodologies in improving graph comprehension and reasoning capabilities [60].

Moreover, the TPP-LLM-Embedding model integrates temporal and event-type representations into language models, distinguishing it from conventional models and improving retrieval performance [61]. This model underscores the significance of incorporating temporal sequences and event types into TKG architecture to enhance the handling of complex temporal queries.

The architecture of TKGs is further enriched by knowledge-infused prompts, as demonstrated by the sLA-tKGF framework, which utilizes historical TKG data, web information, and pre-trained language model descriptions to improve forecasting capabilities. This approach highlights the importance of integrating diverse data sources into TKGs to enhance predictive accuracy and reasoning capabilities [62].

Thus, the structure and function of TKGs are characterized by their ability to merge temporal data with advanced retrieval and generation methodologies, enabling effective processing of time-sensitive information and supporting sophisticated temporal reasoning tasks. As research advances, refining indexing, retrieval, and generation techniques will continue to expand TKGs’ potential in dynamic and temporally complex domains [33].

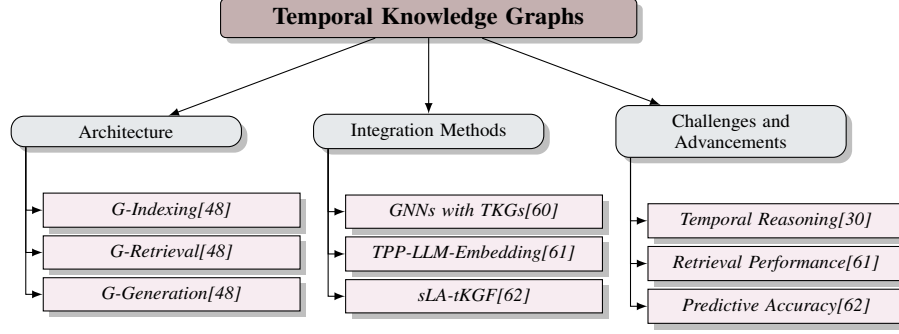


Figure 3: This figure illustrates the structure and integration methods of Temporal Knowledge Graphs (TKGs), highlighting key components such as G-Indexing, G-Retrieval, and G-Generation. It also presents the integration of Graph Neural Networks (GNNs) and TPP-LLM-Embedding models, alongside challenges and advancements in temporal reasoning, retrieval performance, and predictive accuracy.

## 4.2 Challenges in Temporal Reasoning

Temporal reasoning within Temporal Knowledge Graphs (TKGs) faces numerous challenges, complicating the integration of temporal data with Large Language Models (LLMs). A primary issue is the focus of existing benchmarks on surface-level understanding, which neglects the complexities of fine-grained lexical semantics and the nuanced interpretation of temporal data [3]. This limitation hampers the development of robust temporal reasoning capabilities, as models are inadequately trained to manage the intricacies of temporal relationships and events.

Inaccurate factual recall and hallucinations can arise when LLMs interact with TKGs, often exacerbated by data biases and potential future data leakage, compromising the reliability of temporal reasoning systems [63]. Aligning user search intents with retrieved data complicates temporal reasoning, necessitating accurate interpretation and response to complex temporal queries.

Integrating diverse contextual information during training significantly enhances the robustness and accuracy of temporal question-answering systems [10]. However, managing and synthesizing vast amounts of contextual data presents challenges. Additionally, issues in temporal reasoning parallel those in machine learning concerning limited labeled data, where the scarcity of high-quality temporal data impedes the development of effective reasoning models [64].

To address these challenges, developing benchmarks that evaluate nuanced understanding of temporal semantics and methodologies to mitigate biases and hallucinations is essential. Systematic refinements in training and evaluation methodologies, including enhancements like the TempTabQA dataset for tabular temporal question answering and the DPrompt tuning approach for efficient preprocessing, can significantly improve the temporal reasoning capabilities of LLMs. These advancements enable LLMs to provide more accurate and contextually relevant responses in dynamic and temporally intricate domains while addressing outdated benchmarks and the limitations of Retrieval-Augmented Generation (RAG) in facilitating deeper reasoning [57, 42, 15].

## 4.3 Applications and Case Studies

Temporal Knowledge Graphs (TKGs) have significant applications across various domains, particularly in enhancing temporal reasoning capabilities. In educational contexts, LLMs are utilized to generate high-quality quiz questions, demonstrating their practical utility in creating educational content, as highlighted by structured surveys with experts emphasizing the importance of LLMs in education [65].

In medicine, TKGs improve the accuracy and contextual relevance of models used for medication consultations. The MedicineQA benchmark enhances the evaluation of LLMs in the medical domain, contributing to the development of more accurate models for medication-related queries [66]. This illustrates the potential of TKGs to support complex reasoning tasks in specialized fields.

---

The integration of TKGs with Retrieval-Augmented Generation (RAG) models has been explored for domain-specific queries. An end-to-end system combining RAG with curated datasets enhances the factual accuracy of LLM responses, showcasing the effectiveness of TKGs in providing domain-specific knowledge [67]. This is particularly beneficial in fields requiring precise and contextually aware information retrieval.

Benchmarks like TRAM offer a comprehensive evaluation framework for temporal reasoning tasks, addressing existing gaps and providing a unified platform for assessing the temporal reasoning capabilities of LLMs [43]. These benchmarks are crucial for advancing the understanding and application of TKGs in temporal reasoning, enabling more accurate predictions and analyses of temporal events.

The KGI system exemplifies successful TKG integration with retrieval-augmented generation models, achieving state-of-the-art performance in knowledge-intensive tasks. This integration highlights the potential of TKGs to enhance LLM capabilities across various applications, particularly in tasks requiring the synthesis of complex information [68].

Frameworks like DRAGON-AI, assessed by ontology editors, further illustrate TKG applications in generating definitions and relationships across diverse ontologies [13]. This demonstrates the versatility of TKGs in supporting ontology generation and management, contributing to the advancement of knowledge representation and reasoning.

#### 4.4 Future Directions in Temporal Knowledge Graphs

Future developments in Temporal Knowledge Graphs (TKGs) aim to tackle critical research challenges and opportunities, enhancing their utility in temporal reasoning and information retrieval. Refining temporal reasoning methodologies to improve interpretability and accuracy is a promising direction. This includes developing advanced algorithms that better understand and process temporal semantics, enabling TKGs to provide contextually aware and precise responses to time-sensitive queries [3].

Another significant area of research is integrating TKGs with machine learning models to enhance reasoning capabilities. Exploring neural networks and deep learning techniques can improve the efficiency and effectiveness of temporal data processing within TKGs. By leveraging machine learning strengths, TKGs can achieve more robust temporal reasoning, allowing them to handle complex temporal queries with greater accuracy [31].

Developing comprehensive benchmarks for evaluating TKG performance in temporal reasoning tasks is also crucial. These benchmarks can provide a standardized framework for assessing TKG capabilities, identifying areas for improvement, and guiding the development of more effective temporal reasoning models [43]. Establishing clear evaluation criteria enables researchers to systematically enhance TKG performance across various applications.

Exploring novel applications for TKGs in diverse domains presents significant potential for future research. Investigating TKG use in healthcare, finance, and education, where processing and reasoning about temporal information is critical, can unlock new opportunities for leveraging temporal data innovatively, contributing to knowledge representation and reasoning advancement [66].

Finally, addressing data quality and completeness challenges in TKGs is essential. Ensuring the accuracy and reliability of temporal data is crucial for effective TKG functioning, particularly in applications requiring precise and contextually relevant information. Future research should focus on developing methodologies to improve data quality and completeness, enhancing the overall utility of TKGs in temporal reasoning and information retrieval [67].

## 5 Retrieval-Augmented Generation

Retrieval-Augmented Generation (RAG) represents a pivotal advancement in natural language processing by combining information retrieval with generative modeling. This section delves into RAG methodologies that enhance Large Language Models (LLMs) by integrating external knowledge, thereby improving the accuracy and contextual relevance of generated responses in open-domain and domain-specific tasks. Table 2 presents a detailed summary of the diverse methodologies and advancements in Retrieval-Augmented Generation (RAG), illustrating the integration of retrieval

Category	Feature	Method
<b>Methodologies of Retrieval-Augmented Generation</b>	Memory and Reasoning Enhancement	ARM-RAG[52]
<b>Advancements in RAG Frameworks</b>	Integrated Improvement Contextual Retrieval	RAG[25], RLLM[69] GNN-Ret[36]
<b>Challenges in RAG Systems</b>	Uncertainty and Context Handling	E2E-AFG[70], CAR[71], DR-RAG[38], CoV-RAG[37], ERAgent[41], FILCO[72]
	Retrieval and Performance Optimization	ToG-2[56]
<b>Applications and Use Cases</b>	Domain Knowledge Utilization	HIRO[73], SBRAG[74], RAG-KG[46], C-RAG[75], AIT[20], CRAG[76]
	Information Retrieval Techniques	PolyRAG[77]
<b>Future Directions in RAG Research</b>	Data Efficiency Enhancements	RGC[78], PR3[16]

Table 1: This table provides a comprehensive overview of the methodologies, advancements, challenges, applications, and future directions in Retrieval-Augmented Generation (RAG) systems. It categorizes various features and methods, highlighting key innovations and challenges faced within the RAG framework. The table serves as a summary of current research trends and potential future developments in enhancing the capabilities of Large Language Models (LLMs) through RAG.

Category	Feature	Method
<b>Methodologies of Retrieval-Augmented Generation</b>	Memory and Reasoning Enhancement	ARM-RAG[52]
<b>Advancements in RAG Frameworks</b>	Integrated Improvement Contextual Retrieval	RAG[25], RLLM[69] GNN-Ret[36]
<b>Challenges in RAG Systems</b>	Uncertainty and Context Handling	E2E-AFG[70], CAR[71], DR-RAG[38], CoV-RAG[37], ERAgent[41], FILCO[72]
	Retrieval and Performance Optimization	ToG-2[56]
<b>Applications and Use Cases</b>	Domain Knowledge Utilization	HIRO[73], SBRAG[74], RAG-KG[46], C-RAG[75], AIT[20], CRAG[76]
	Information Retrieval Techniques	PolyRAG[77]
<b>Future Directions in RAG Research</b>	Data Efficiency Enhancements	RGC[78], PR3[16]

Table 2: This table provides a comprehensive overview of the methodologies, advancements, challenges, applications, and future directions in Retrieval-Augmented Generation (RAG) systems. It categorizes various features and methods, highlighting key innovations and challenges faced within the RAG framework. The table serves as a summary of current research trends and potential future developments in enhancing the capabilities of Large Language Models (LLMs) through RAG.

mechanisms with generative processes to enhance Large Language Models (LLMs). Additionally, Table 4 provides a detailed comparison of various methodologies within Retrieval-Augmented Generation (RAG), showcasing the distinct features and optimization strategies of ERAgent, ARM-RAG, and Dynamic In-Context Editing (DICE) in the context of enhancing Large Language Models (LLMs). The following subsection outlines specific methodologies employed in RAG, highlighting their role in advancing LLM capabilities.

## 5.1 Methodologies of Retrieval-Augmented Generation

RAG methodologies significantly enhance LLMs by integrating retrieval mechanisms with generative processes. By merging retrieved passages with LLMs, these methods facilitate access to external knowledge, improving output accuracy and relevance. Techniques such as iterative response generation and feedback loops have demonstrated over a 10% improvement in answer accuracy, as measured by metrics like the Ragas score [2, 24]. The RAG process typically involves data collection, preprocessing, vector embedding creation, content retrieval, context augmentation, and response generation, enabling LLMs to utilize up-to-date, domain-specific information effectively.

The ERAgent framework exemplifies advancements in RAG, enhancing accuracy, efficiency, and personalization through advanced modules [41]. This modular approach optimizes retrieval and generation processes, adapting to diverse user needs. Similarly, the ARM-RAG system combines RAG with an auxiliary memory mechanism to retrieve reasoning chains from past experiences, improving LLM reasoning capabilities in complex scenarios [52].

Dynamic retrieval strategies, such as the Dynamic In-Context Editing (DICE) method, enable LLMs to engage in multi-hop reasoning within extensive textual contexts, effectively gathering and integrating relevant data [17, 21, 25, 6, 38]. DICE’s ability to decompose questions into sub-questions and retrieve necessary information distinguishes it from traditional approaches, emphasizing adaptability in retrieval processes.

Hybrid retrieval-augmented generation methods integrate knowledge graphs and document retrieval, facilitating complex reasoning by leveraging diverse data sources. This integration enhances LLMs' reasoning capabilities in information retrieval tasks. Additionally, benchmarks like the LitQA dataset, which sources questions from recent literature, provide a standardized framework for evaluating RAG systems in knowledge-intensive tasks [40].

RAG methodologies are essential for advancing natural language processing, particularly in complex tasks such as information retrieval and question answering. By integrating external knowledge bases and employing advanced retrieval techniques, RAG significantly enhances the accuracy and contextual relevance of LLM-generated responses. Innovations like the Dynamic-Relevant Retrieval-Augmented Generation (DR-RAG) framework exemplify how strategic document retrieval optimizes efficiency while maintaining high response quality. Refining text chunk retrieval and employing sophisticated algorithms address common challenges like scalability and relevance, positioning RAG as a transformative approach in the NLP landscape [8, 22, 38]. As research progresses, these methodologies will continue to evolve, offering new avenues for integrating and utilizing external knowledge across diverse domains.

## 5.2 Advancements in RAG Frameworks

Method Name	Retrieval Mechanisms	Integration Techniques	Evaluation Frameworks
RGC[78]	Dynamic Caching System	Prefix-aware Replacement	Average Tft Throughput
RLLM[69]	Hierarchical Fm-Index	Jointly Optimized Processes	Accuracy, F1 Score
GNN-Ret[36]	Graph Neural Networks	Semantic Distances Integration	F1 Score
ERAGent[41]	Knowledge Retriever	Modular System Integration	Exact Match
RAG[25]	Retriever TO Identify	Retrieval And Generation	Assessing The System

Table 3: Overview of recent advancements in retrieval-augmented generation (RAG) frameworks, detailing the retrieval mechanisms, integration techniques, and evaluation frameworks employed by various methods. The table highlights the innovative strategies implemented by each method to enhance the performance and accuracy of large language models (LLMs).

Recent advancements in RAG frameworks have significantly enhanced LLM capabilities by integrating sophisticated retrieval mechanisms with generative processes. A notable development is RAGCache, which caches intermediate states in a knowledge tree structure sensitive to document retrieval order, employing a prefix-aware replacement policy to improve retrieval efficiency and accuracy [78]. This framework exemplifies the potential for optimizing knowledge management within RAG systems.

The RetroLLM framework further advances RAG by eliminating the need for separate embedding models, enabling joint optimization of retrieval and generation tasks, thereby enhancing performance across applications [69]. The GNN-Ret framework, leveraging graph neural networks to assess relatedness between passages, demonstrates superior performance in retrieval tasks, particularly in question-answering scenarios [36]. The integration of graph-based methodologies into RAG frameworks underscores the potential for improving retrieval precision and supporting contextually aware LLM responses.

ERAGent combines enhanced question rewriting, selective knowledge retrieval, and personalized response generation, significantly improving RAG system adaptability and efficiency [41]. Comprehensive benchmarks, as introduced by [7], provide structured frameworks for evaluating RAG practices, facilitating systematic improvements in methodologies. These benchmarks establish clear evaluation criteria, guiding the development of more effective RAG systems capable of delivering accurate and contextually relevant responses across various domains. Table 3 provides a comprehensive summary of the key retrieval mechanisms, integration techniques, and evaluation frameworks utilized in recent RAG frameworks, illustrating the diverse approaches to optimizing large language models.

Advancements in RAG frameworks reflect ongoing efforts to refine retrieval and generation processes, ensuring LLMs provide accurate, contextually relevant, and efficient responses across diverse tasks. As research progresses, these methodologies will continue to evolve, offering new avenues for integrating and utilizing external knowledge sources [25].

### 5.3 Challenges in RAG Systems

The implementation of RAG systems faces several challenges that affect their effectiveness across applications. A primary obstacle is the difficulty in quantifying and managing uncertainty associated with retrieval, which undermines the trustworthiness of generated responses [71]. This uncertainty is compounded by distracting content in retrieved passages that mislead generation models and degrade response accuracy [72].

As illustrated in Figure 4, these challenges can be categorized into three primary areas: uncertainty management, retrieval effectiveness, and system optimization. Each category highlights specific issues such as quantifying uncertainty, managing noise in knowledge base documents, and improving caching efficiency, drawing on insights from recent research.

RAG systems often struggle to filter out irrelevant content, leading to hallucinatory outputs and reduced reliability [70]. Challenges include the absence of relevance labels for document chunking and extracting relevant information from lengthy documents [16]. Additionally, the retriever’s inability to return relevant results for vague or incomplete queries, coupled with language models’ tendency to generate hallucinations, complicates RAG system effectiveness [37].

Current RAG methods often lack depth and completeness in retrieved information necessary for complex reasoning tasks [56]. Low relevance of dynamic-relevant documents can result in incomplete or inaccurate answers in multi-hop QA scenarios [38]. Furthermore, the inability to handle multi-faceted semantic information and retain previously retrieved knowledge during long-term interactions presents significant challenges [41].

Existing RAG systems do not effectively cache and share intermediate states of retrieved knowledge across multiple requests, leading to redundant computations and increased latency [78]. The lack of comprehensive evaluation results and the inability to capture the complexity of real-world queries and knowledge integration in existing benchmarks further hinder optimization efforts [79].

To address these challenges, ongoing research and development are essential for optimizing RAG systems. Enhancing retrieval mechanisms to ensure accurate, contextually relevant responses is crucial. Recent advancements, like the DR-RAG framework, illustrate how integrating external knowledge and innovative retrieval strategies can significantly improve answer accuracy and efficiency. Addressing issues such as scalability, bias, and redundancy in information retrieval will enhance the robustness and applicability of RAG models in knowledge-intensive tasks [22, 21, 38, 8, 64].

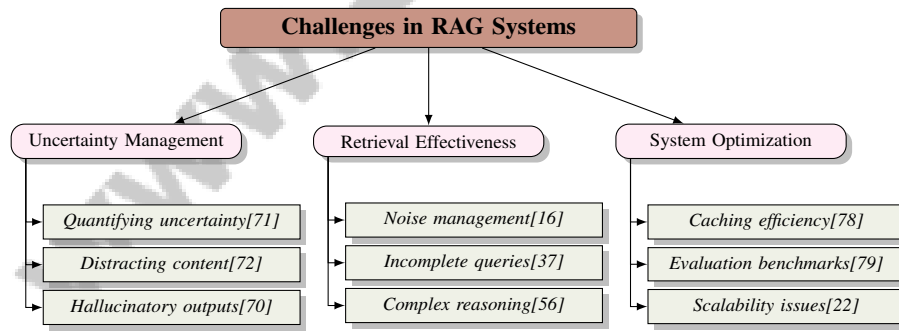


Figure 4: This figure illustrates the primary challenges faced by Retrieval-Augmented Generation (RAG) systems, categorized into uncertainty management, retrieval effectiveness, and system optimization. Each category highlights specific issues such as quantifying uncertainty, managing noise, and improving caching efficiency, drawing on insights from recent research.

### 5.4 Applications and Use Cases

RAG systems exhibit remarkable versatility across domains, significantly enhancing LLM capabilities in complex information retrieval and generation tasks. In customer service, RAG methodologies improve automated question answering by integrating domain-specific knowledge, resulting in accurate, contextually relevant responses [46]. This application highlights RAG’s potential to enhance customer interactions through precise and timely answers.

---

In education, RAG systems support adaptive AI tutors that leverage course-specific materials to generate personalized responses, enriching the learning experience for students [20]. This adaptability underscores RAG's role in tailoring educational content to specific learning needs.

The medical field benefits significantly from RAG systems, particularly in enhancing medical reasoning. The Self-BioRAG framework has shown substantial performance improvements over state-of-the-art models on major medical question-answering benchmarks, integrating relevant medical knowledge to support healthcare professionals [74].

Moreover, the integration of RAG with advanced retrieval methodologies, such as the PR3 methodology, enhances retrieval precision and recall, allowing for comprehensive information synthesis from knowledge bases [16]. This approach is particularly beneficial in fields requiring precise and contextually aware information retrieval.

The HIRO framework has demonstrated significant improvements in managing context for LLMs, achieving notable performance increases on datasets like NarrativeQA, thereby optimizing RAG methodologies [73]. Similarly, the PolyRAG system uses structured querying approaches to provide accurate answers by combining ontology precision with the recall of knowledge graphs and raw text [77].

RAG systems are also applied in various NLP tasks, including language modeling, machine translation, text summarization, question answering, and dialogue systems [5]. The integration of evaluators within RAG systems enhances semantic relevance and logical consistency by comparing retrieved chunks with external recommendations, thereby improving output reliability [75]. Frameworks like CRAG systematically enhance output reliability by addressing inaccuracies in retrieved documents [76].

The RAGCache framework significantly reduces computation time and resource usage through efficient knowledge caching, leading to faster response times and higher throughput [78]. This efficiency is crucial for applications requiring rapid information retrieval and response generation.

The diverse applications and use cases of RAG systems underscore their pivotal role in enhancing LLM capabilities. By integrating external data sources, RAG systems improve accuracy, transparency, and contextuality in responses, crucial for sectors requiring domain-specific knowledge and real-time information retrieval. The categorization of user queries into explicit, implicit, interpretable, and hidden rationale types allows for a tailored approach to leveraging external data, addressing complexities in deploying data-augmented LLMs effectively [23, 25]. RAG systems continue to expand their applicability and effectiveness in various settings, enhancing accuracy and contextual relevance in complex information retrieval and question-answering tasks.

## 5.5 Future Directions in RAG Research

The future development of RAG systems presents numerous opportunities for enhancing capabilities and applicability across diverse domains. One promising direction involves improving caching strategies to adapt to varying retrieval patterns, as exemplified by the RAGCache framework, which could enhance scalability for larger datasets [78]. This enhancement would allow RAG systems to manage extensive information more efficiently, optimizing retrieval processes.

Another critical area for future research is refining knowledge selection strategies, essential for ensuring the relevance and accuracy of retrieved information [79]. Developing comprehensive evaluation methodologies will enable systematic assessment of these strategies, guiding the creation of more precise and contextually aware RAG systems.

Exploring multi-hop iterative searches offers another avenue for enhancing retrieval capabilities, allowing for in-depth exploration of knowledge bases [16]. This approach could significantly improve RAG systems' ability to handle complex queries, enhancing their utility in knowledge-intensive tasks.

Future research should also focus on optimizing metadata generation processes to support effective retrieval and generation tasks [16]. Improving the quality and relevance of metadata will enhance RAG systems' accuracy in information retrieval, ensuring generated responses are precise and contextually relevant.

Integrating hybrid approaches that combine preloading with selective retrieval can enhance RAG systems' adaptability and efficiency by leveraging advanced modules such as the Query Rewriter and

Knowledge Filter. These modules improve knowledge retrieval and response generation accuracy, addressing challenges like irrelevant knowledge and redundancy. Incorporating techniques like multi-query retrieval can optimize the retrieval phase, leading to more relevant outputs and better resource utilization in knowledge-intensive tasks [80, 8, 81, 21]. These methods could result in more responsive and context-aware retrieval processes, enabling RAG systems to meet user needs across various applications.

Finally, developing accurate benchmark datasets and task-specific tuning for LLMs is essential for advancing RAG systems. Comprehensive benchmarks encompassing various application scenarios and refined evaluation metrics are crucial for providing insights into current methodologies in RAG frameworks. These benchmarks facilitate identifying limitations in existing systems, such as challenges posed by irrelevant knowledge retrieval and the need for multi-perspective views in knowledge-dense domains like law and medicine. By systematically evaluating RAG techniques—such as enhancements to the Query Rewriter and the introduction of credibility-aware generation frameworks—researchers can better understand the effectiveness of different approaches, guiding future advancements and improving the precision and reliability of RAG systems [14, 82, 21, 38].

Focusing on these research directions will significantly enhance RAG systems’ capabilities, leading to more accurate and contextually relevant responses. This includes implementing innovative retrieval methodologies, such as the DR-RAG framework, which optimizes document retrieval efficiency and improves answer accuracy by minimizing LLM access frequency. Additionally, refining the RAG process through advanced text chunking, query expansion, metadata incorporation, and re-ranking algorithms will address current limitations, resulting in more effective and reliable applications across various domains, including complex question-answering tasks [22, 38].

Feature	ERAGent	ARM-RAG	Dynamic In-Context Editing (DICE)
<b>Integration Mechanism</b>	Modular Approach	Auxiliary Memory	Multi-hop Reasoning
<b>Optimization Focus</b>	Personalization	Reasoning Chains	Data Integration
<b>Unique Feature</b>	Advanced Modules	Complex Scenarios	Sub-question Decomposition

Table 4: Table comparing the features of three Retrieval-Augmented Generation (RAG) methodologies: ERAGent, ARM-RAG, and Dynamic In-Context Editing (DICE). Each method is evaluated based on its integration mechanism, optimization focus, and unique feature, highlighting the diverse approaches to enhancing Large Language Models (LLMs) through retrieval-augmented processes.

## 6 Applications in Question Answering Systems

Integrating advanced methodologies into question answering systems has significantly enhanced performance and reliability. This section delves into the applications of Large Language Models (LLMs) across various question answering contexts, emphasizing their transformative potential in diverse domains. The following subsection evaluates LLMs in open-domain question answering, highlighting recent advancements that enhance their effectiveness and contextual relevance in generating accurate responses.

### 6.1 Evaluation of LLMs in Open-Domain QA

LLMs in open-domain question answering (QA) have benefited from Retrieval-Augmented Generation (RAG) methodologies, which improve response accuracy and contextual relevance. Recent studies highlight RAG frameworks’ effectiveness in enhancing LLM performance across various QA scenarios. For instance, the RetroLLM framework demonstrated superior performance on datasets such as NQ, TriviaQA, HotpotQA, PopQA, and 2WIKI, showcasing its ability to leverage retrieval processes effectively [69]. Similarly, the ES-RAG model, with an accuracy of 68.29

In multi-hop question answering, the StrategyQA dataset, containing 2,780 queries, was used to assess meta-prompting optimized RAG systems, revealing notable improvements over baseline models [88]. The FRAMES dataset, with its 824 questions, provided a diverse range of topics for evaluating LLM performance in complex reasoning tasks [55]. These evaluations underscore RAG frameworks’ critical role in advancing LLM performance in open-domain QA tasks.



Benchmark	Size	Domain	Task Format	Metric
ES-RAG[83]	100,000	Question Answering	Answering Questions	Accuracy, F1 Score
FRAMES[55]	824	Question Answering	Multi-hop Reasoning	Accuracy
TRAM[43]	526,668	Temporal Reasoning	Multiple Choice	Accuracy, F1-score
Complex-TR[84]	10,800	Temporal Question Answering	Multi-Hop Question Answering	Set Acc., Ans. F1
CurriculumQA[80]	200	Education	Question Answering	Hit Rate, ROUGE
SRAG[85]	45,000	Question Answering	Question Answering	Accuracy, Recall@1
RAG-Ophthalmology[86]	70,000	Ophthalmology	Long-form Question Answering	Accuracy, Evidence Attribution
CAT[87]	12,000	Conversational AI	Question Answering	Recall, F2

Table 5: The table provides a comprehensive overview of representative benchmarks utilized in evaluating large language models (LLMs) in open-domain question answering (QA) tasks. It details various datasets, their sizes, domains, task formats, and evaluation metrics, highlighting the diversity and complexity of the benchmarks used to assess LLM performance. This information is crucial for understanding the scope and challenges in optimizing retrieval-augmented generation methodologies across different QA scenarios.

As illustrated in Figure 5, the figure categorizes the key frameworks and techniques evaluated in open-domain QA using LLMs. It highlights the methodologies into RAG frameworks, multi-hop QA datasets, and enhanced techniques, underscoring significant contributions from recent studies. Enhanced techniques, including advanced text chunk retrieval and effective content restructuring, enable LLMs to deliver more accurate and contextually relevant answers, meeting complex user demands and reducing hallucination risks. Innovations like the Refiner framework exemplify how restructuring retrieved content can lead to substantial gains in answer accuracy and efficiency [89, 22, 15]. Table 5 presents a detailed summary of key benchmarks employed in the evaluation of large language models for open-domain question answering, illustrating the breadth of datasets and metrics critical for assessing the efficacy of retrieval-augmented generation frameworks.

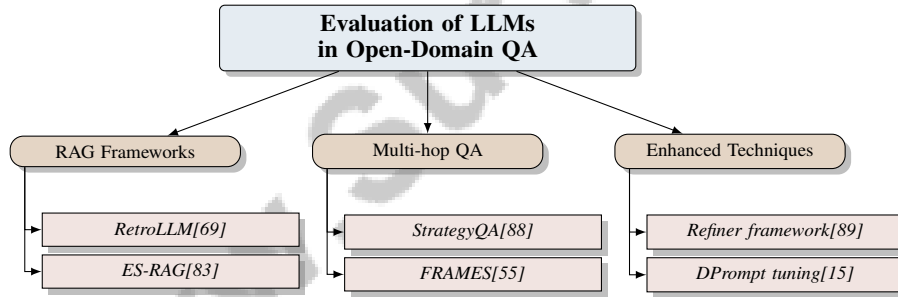


Figure 5: This figure illustrates the key frameworks and techniques evaluated in open-domain question answering (QA) using large language models (LLMs). It categorizes the methodologies into RAG frameworks, multi-hop QA datasets, and enhanced techniques, highlighting significant contributions from recent studies.

## 6.2 Educational Applications and Quiz Generation

LLMs show significant promise in educational contexts, particularly in generating contextually relevant quiz questions tailored to specific courses. This capability is enhanced through the integration of LLMs with Temporal Knowledge Graphs (TKGs), improving the accuracy and alignment of educational content with course material [65]. Techniques such as the Time-aware Retrieve-Rewrite-Retrieve-Rerank framework enable LLMs to manage complex temporal constraints, enhancing their ability to produce contextually appropriate educational materials [90, 42, 91, 24].

By leveraging Teaching Knowledge Graphs (TKGs), LLMs can generate questions that are not only factually correct but also pedagogically effective, enhancing student comprehension [53, 90, 22]. This integration significantly enhances the efficiency of question generation and supports the development of effective educational assessment systems. The adaptability of LLMs allows for ongoing updates and refinements of educational content, crucial for maintaining quiz relevance as new information emerges. By integrating external document retrieval techniques such as RAG, LLMs can enhance

---

their knowledge base and improve answer accuracy, ensuring educational materials are current and well-aligned with the latest data [24, 89].

### 6.3 Enhancing Factual Accuracy and Mitigating Hallucinations

Enhancing factual accuracy and mitigating hallucinations in LLMs are critical challenges addressed through various RAG methodologies. The integration of structured knowledge from Knowledge Graphs (KGs) significantly improves the factual accuracy of LLM outputs. For instance, the GNN-Ret framework retrieves supporting passages by considering their interrelatedness, thereby enhancing accuracy in question-answering tasks [36]. Similarly, the CoV-RAG framework effectively mitigates hallucination issues during both retrieval and generation stages, leading to improved accuracy [37].

The Refiner framework enhances downstream LLM performance by providing structured and contextually relevant outputs, demonstrating resilience against noisy document inputs [89]. RAG methodologies such as RetroLLM reduce token consumption and improve evidence accuracy by eliminating separate retrievers, thus enhancing overall system efficiency [69]. The integration of topological relationships in frameworks like Topo-RAG further enhances factual accuracy in generated texts [11]. The DR-RAG framework effectively retrieves both static and dynamic relevant documents, improving overall QA system performance [38]. These strategies underscore ongoing efforts to enhance factual accuracy and mitigate hallucinations in LLM outputs. By continuously refining RAG methodologies and establishing comprehensive benchmarks, researchers can significantly improve the reliability and precision of LLM-generated responses [14, 25, 22, 15].

### 6.4 Advancements in Contextual and Personalized Responses

Recent advancements in generating contextual and personalized responses have significantly enhanced LLM capabilities in adapting to user-specific needs. Frameworks such as Think-then-Act have optimized retrieval-augmented generation by refining query assessment and model capability evaluation [92]. The integration of relational databases with LLMs exemplifies advancements in this domain, allowing more precise and personalized responses by leveraging structured data sources [93]. These advancements highlight ongoing efforts to refine LLM contextual understanding and personalization capabilities. Employing advanced techniques such as RAG combines generative capabilities with precise information retrieval, addressing challenges like query ambiguity and irrelevant knowledge through innovative modules such as Query Rewriters and Knowledge Filters [24, 23, 21, 22, 25].

### 6.5 Domain-Specific Implementations

Domain-specific implementations of QA systems leveraging LLMs and TKGs have achieved significant advancements in addressing unique challenges in specialized fields. These implementations utilize the structured and temporal characteristics of TKGs to enhance LLM retrieval and reasoning capabilities. By integrating advanced methodologies such as a Time-aware Retrieve-Rewrite-Retrieve-Rerank framework, these approaches improve response accuracy and relevance in domain-specific contexts [26, 21, 91, 22].

In the financial sector, the FinanceBench dataset provides a robust benchmark for evaluating RAG models' performance in processing domain-specific financial information [19]. The RAGLAB framework offers an efficient, user-friendly library for fair comparisons of RAG algorithms, serving as a vital tool within the NLP community [94]. Datasets containing text from academic papers, financial reports, and other professional documents are crucial for assessing RAG systems' performance in academic and professional sectors [17]. Exploring multilingual capabilities is also critical, as demonstrated by benchmarks that provide comprehensive evaluation frameworks for multilingual Retrieval-Augmented Language Models (RALMs) [29]. While long-context LLMs have shown superior performance compared to RAG in certain scenarios, the lower computational cost of RAG makes it viable for specific use cases [95]. The domain-specific implementations of QA systems utilizing LLMs and TKGs continue to expand the frontiers of information retrieval and reasoning in specialized fields [22, 21].

---

## 6.6 Innovations in Multi-Hop and Complex Question Answering

Innovations in multi-hop and complex question answering have been significantly driven by advancements in frameworks that integrate RAG methodologies with LLMs. The GenGround framework exemplifies this progress by effectively leveraging both the inherent knowledge of LLMs and evidence from retrieved documents, outperforming existing methods in multi-hop question answering tasks [96]. Further advancements are seen in systems like ERASE, which enhance language models' accuracy in dynamic contexts by improving knowledge base updates [9]. These innovations highlight continuous advancements aimed at enhancing LLM capabilities in addressing intricate queries. By employing RAG and the proposed Refiner framework, LLMs can effectively integrate and synthesize information from diverse sources, adaptively extract relevant content, restructure it for clarity, and utilize multi-view retrieval methods to improve interpretability and precision in knowledge-dense domains [23, 82, 89].

## 7 Time-Aware Information Retrieval

### 7.1 Significance of Temporal Context in Information Retrieval

Incorporating temporal context in information retrieval systems is crucial for enhancing the relevance and accuracy of retrieved data, particularly in domains where temporal factors significantly influence data interpretation. Temporal context enables systems to prioritize information based on time-specific relevance, improving search precision in dynamic environments [10, 31]. This capability is vital for applications like historical data analysis and real-time event monitoring, where accurate interpretation and decision-making hinge on the temporal aspect of information.

As illustrated in Figure 6, the hierarchical structure of temporal context in information retrieval underscores its importance, the role of Temporal Knowledge Graphs (TKGs), and advanced retrieval methodologies. TKGs play a key role in time-aware retrieval, offering a structured framework for capturing temporal relationships and events. By integrating TKGs with Large Language Models (LLMs), retrieval systems enhance reasoning capabilities, processing time-sensitive queries with greater accuracy and contextual relevance [32, 28]. This integration is particularly beneficial for applications requiring an understanding of temporal sequences and causal relationships, such as forecasting future events based on historical data.

Advanced retrieval methodologies optimize the integration of temporal data without sacrificing efficiency. Techniques involving graph neural networks and dynamic retrieval frameworks enhance the handling of complex, time-sensitive queries. The TimeR 4 framework exemplifies the potential of time-aware systems to improve response accuracy and contextual relevance in temporally complex domains. By effectively integrating temporal knowledge from TKGs and employing strategies like retrieve-rewrite-rerank modules, these systems mitigate temporal hallucination and enhance semantic understanding of time constraints. Robust training with diverse contextual information significantly improves the performance of temporal question-answering systems, ensuring that LLMs deliver precise and contextually appropriate answers in scenarios requiring temporal reasoning [91, 6, 21, 10].

Establishing benchmarks for evaluating temporal reasoning capabilities underscores the growing importance of temporal context in advancing information retrieval technologies. Systematic assessment of model performance in handling time-related queries helps identify areas for improvement and drive innovation in the field [43]. As demand for time-sensitive information rises, integrating temporal context within retrieval systems will shape the future of information retrieval and question-answering systems.

### 7.2 Challenges in Time-Aware Information Retrieval

Time-aware information retrieval (TAIR) systems face challenges that hinder their ability to effectively process temporally relevant information. Key challenges include accurately understanding complex temporal semantics, as existing methods relying on pre-trained embeddings or graph neural networks often fall short. Integrating temporal knowledge from TKGs into LLMs is complicated by issues such as temporal hallucination and precise retrieval of semantically relevant facts. Frameworks like TimeR 4 aim to enhance temporal reasoning capabilities through innovative retrieval and reranking strategies, yet persistent obstacles limit TAIR systems' overall effectiveness [61, 91, 21, 97].

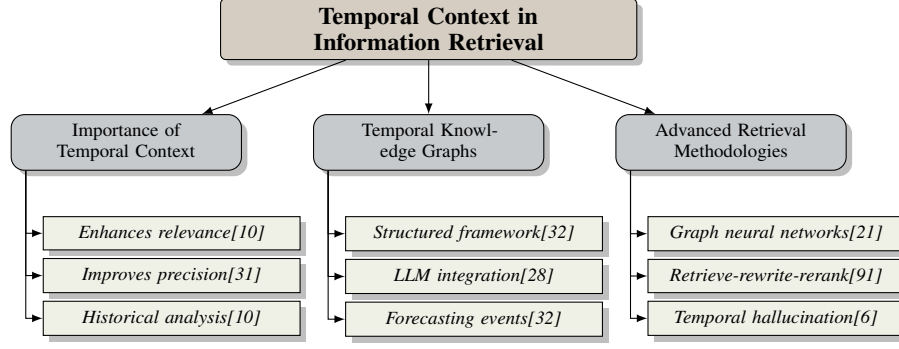


Figure 6: This figure illustrates the hierarchical structure of temporal context in information retrieval, highlighting its importance, the role of Temporal Knowledge Graphs, and advanced retrieval methodologies.

A primary challenge is integrating temporal data into existing frameworks without compromising efficiency and accuracy. This requires sophisticated methodologies to manage complex temporal sequences and causal relationships. Additionally, the scarcity of high-quality temporal data hampers robust TAIR development, as comprehensive datasets are crucial for training models capable of reasoning about time-sensitive queries [64].

The dynamic nature of temporal information presents challenges, as systems must adapt to changes in temporal data. Developing adaptive strategies that integrate timely data while minimizing noise is essential for enhancing retrieval accuracy. This approach is critical for addressing challenges in complex question-answering tasks, where timing and relevance significantly impact retrieval-augmented generation (RAG) systems' performance [98, 99, 81].

Moreover, aligning user search intents with retrieved data remains challenging, as ensuring that retrieved information accurately reflects users' temporal intents requires models to interpret complex temporal queries effectively [63]. Addressing these challenges necessitates ongoing research to enhance temporal reasoning capabilities. Advancing methodologies for integrating temporal data can significantly improve the accuracy and contextual relevance of Temporal Question Answering (TQA) systems, enabling more precise responses to time-sensitive queries. Findings demonstrate that diverse context training and frameworks like TimeR 4 effectively reduce temporal hallucination and improve temporal reasoning in LLMs. Comprehensive benchmarks such as TESRBench and RAG technique enhancements support robust TQA system development, capable of tackling complex temporal information retrieval challenges [61, 91, 21, 10].

### 7.3 Future Directions in Time-Aware Retrieval Systems

Future developments in time-aware information retrieval (TAIR) systems present opportunities for enhancing capabilities across various domains. Refining temporal reasoning methodologies to improve interpretability and accuracy is a promising direction, involving advanced algorithms that better understand temporal semantics, enhancing systems' ability to provide contextually aware and precise responses to time-sensitive queries [31].

Integrating machine learning models with TAIR systems to bolster retrieval capabilities is another significant research area. Exploring neural networks and deep learning techniques can improve temporal data processing efficiency and effectiveness. Leveraging machine learning strengths, TAIR systems can achieve more robust temporal reasoning, handling complex temporal queries with greater accuracy [36].

Developing comprehensive benchmarks for evaluating TAIR systems' performance in temporal reasoning tasks is crucial. These benchmarks establish a standardized framework for assessing capabilities, identifying improvement areas, and guiding more effective temporal reasoning models' development [43]. Clear evaluation criteria enable systematic TAIR performance enhancement across applications.

---

Exploring novel TAIR system applications in diverse domains presents significant potential for future research. Investigating TAIR use in areas like healthcare, finance, and education, where processing temporal information is critical, can unlock opportunities for leveraging temporal data innovatively, contributing to advancements in knowledge representation and reasoning [32].

Addressing data quality and completeness challenges in TAIR systems is essential. Ensuring temporal data accuracy and reliability is crucial for effective functionality, particularly in applications requiring precise and contextually relevant information. Future research should focus on improving data quality and completeness methodologies, enhancing TAIR systems' utility in temporal reasoning and information retrieval [64].

## 8 Challenges and Future Directions

The integration of Large Language Models (LLMs) with Retrieval-Augmented Generation (RAG) systems and Temporal Knowledge Graphs (TKGs) presents complex challenges, particularly concerning complexity and scalability. As illustrated in Figure 7, this figure highlights key areas of concern, including retrieval quality and accuracy, temporal reasoning, knowledge graph integration, hallucinations, and ethical considerations, outlining the challenges and potential research directions in each domain. The static nature of existing methods restricts LLM adaptability in dynamic environments, compounded by computational demands that risk overfitting [52]. Scalability is further hindered by resource-intensive processes like keyword extraction and graph construction, requiring multiple LLM queries and leading to redundant retrievals that affect performance [98, 21, 100, 81, 64].

The complexity of ontology management, as seen in frameworks like DRAGON-AI, requires collaboration among domain experts to ensure the accuracy and reliability of AI-generated content [13, 101, 102]. The opacity of LLM processes complicates understanding information sourcing and citation, affecting transparency and reliability.

Elasticsearch integration enhances RAG system efficiency but underscores significant computational demands for optimizing retrieval [83, 22, 81, 21]. Ongoing research is vital to develop efficient preprocessing techniques and scalable methodologies for LLM integration with TKGs and RAG systems to maintain adaptability in evolving information landscapes.

### 8.1 Complexity and Scalability

Integrating LLMs with TKGs and RAG systems involves challenges of complexity and scalability. The static nature of current methods limits LLM adaptability, reducing effectiveness in dynamic contexts [52]. This is compounded by the computational costs of adding parameters, risking overfitting.

Scalability challenges stem from resource-intensive processes like keyword extraction and graph construction, necessitating numerous LLM queries and limiting graph-based retrieval efficiency. Uniform recall in RAG systems often results in redundant retrievals, necessitating strategies that prioritize relevant information [98, 21, 100, 81, 64].

Ontology construction and maintenance, as shown by DRAGON-AI, require expert collaboration to ensure AI-generated content accuracy [13, 101, 102]. The opacity of LLM processes complicates understanding information sourcing and citation, impacting reliability.

Elasticsearch integration improves RAG system efficiency but highlights significant computational demands for optimizing retrieval [83, 22, 81, 21]. Addressing these challenges requires ongoing research to develop efficient preprocessing techniques and scalable methodologies for LLM integration with TKGs and RAG systems.

### 8.2 Retrieval Quality and Accuracy

Ensuring retrieval quality and accuracy in RAG systems is challenging when integrating LLMs with diverse knowledge sources. Relying on single sources like Wikipedia limits information scope [68], compounded by difficulties in processing complex documents, affecting fields like legal analysis [54].

As illustrated in Figure 8, the key challenges, current methods, and future directions in enhancing retrieval quality and accuracy in Retrieval-Augmented Generation (RAG) systems are depicted. This figure highlights the reliance on single sources, issues with processing complex documents, and

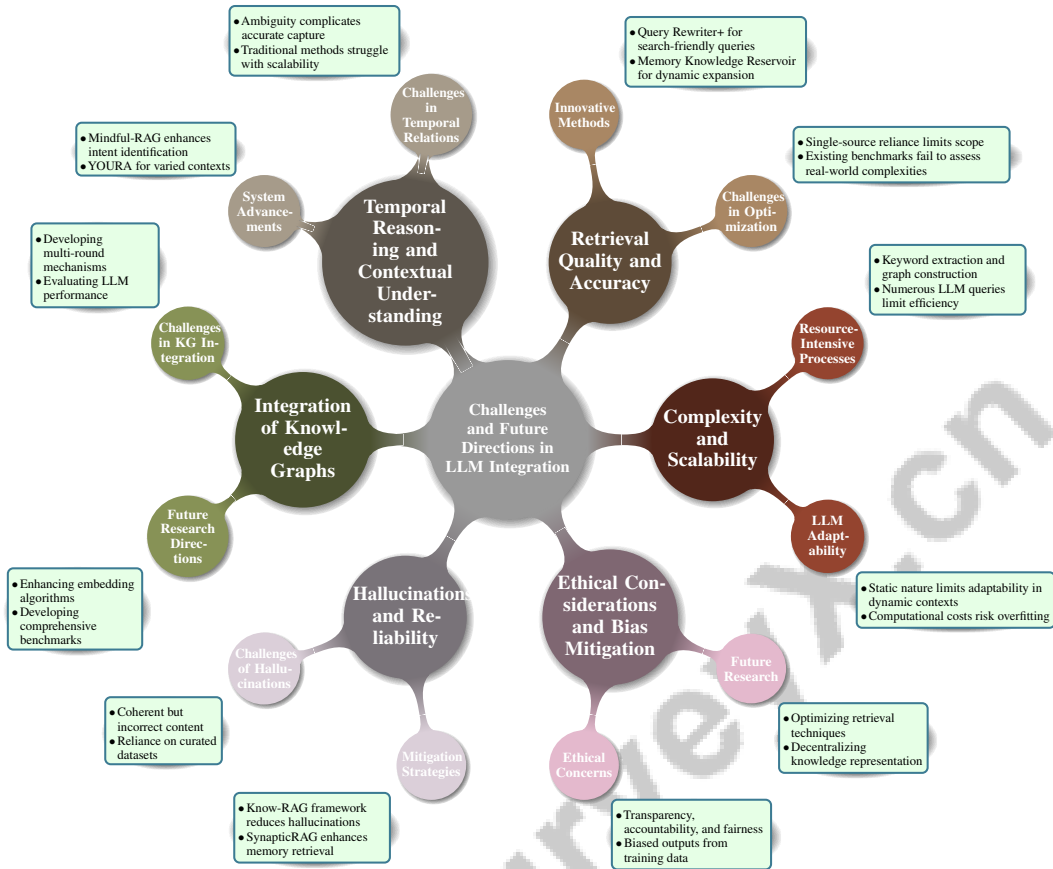


Figure 7: This figure illustrates the challenges and future directions in integrating Large Language Models (LLMs) with Retrieval-Augmented Generation (RAG) systems and Temporal Knowledge Graphs (TKGs). It highlights key areas such as complexity and scalability, retrieval quality and accuracy, temporal reasoning, knowledge graph integration, hallucinations, and ethical considerations, outlining the challenges and potential research directions in each domain.

inefficiencies in current optimization techniques. It further explores challenges related to retrieval quality and fusion methods while pointing towards innovative solutions such as advanced retrieval methodologies and the development of new evaluation metrics.

RAG systems lack systematic optimization, leading to inefficiencies [7]. Existing benchmarks often fail to assess real-world complexities, limiting retrieval quality and accuracy [56]. Challenges such as retrieval accuracy, data biases, and coherent integration persist [8].

Current RAG methods struggle with retrieval quality, processing efficiency, and fusion method interpretability, crucial for accurate outputs [5]. Despite advancements in frameworks like ERAGent, challenges remain in handling granular queries and ensuring relevance and accuracy [41].

Ongoing research should develop advanced retrieval methodologies and evaluation metrics to enhance RAG systems' accuracy and reliability. Innovations like Query Rewriter+ for generating search-friendly queries, the Knowledge Filter for refining information, and the Memory Knowledge Reservoir for dynamic expansion show promise in improving question-answering capabilities and response times [7, 21].

### 8.3 Temporal Reasoning and Contextual Understanding

Temporal reasoning and contextual understanding face challenges that impede LLM performance with TKGs. Ambiguity in temporal relations complicates accurate capture [103], and traditional methods struggle with scalability, requiring validation on diverse datasets [24].

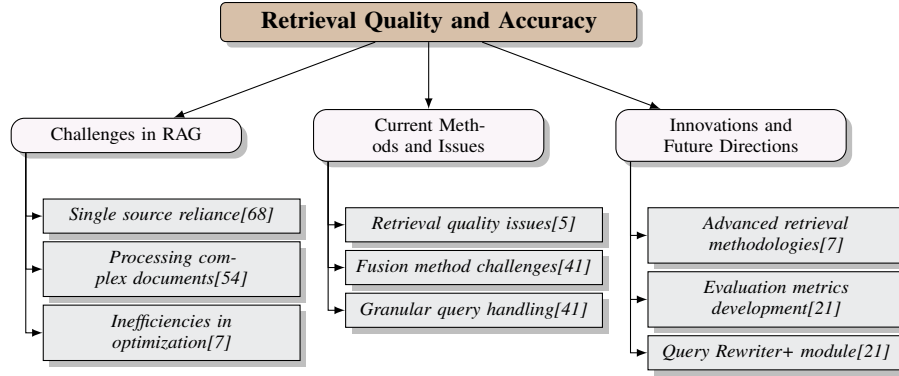


Figure 8: This figure illustrates the key challenges, current methods, and future directions in enhancing retrieval quality and accuracy in Retrieval-Augmented Generation (RAG) systems. It highlights the reliance on single sources, issues with processing complex documents, and inefficiencies in current optimization techniques. The diagram further explores current method challenges like retrieval quality and fusion method issues while pointing towards innovative solutions such as advanced retrieval methodologies and the development of new evaluation metrics.

LLM capabilities in temporal reasoning are often evaluated using benchmarks reliant on multiple-choice questions, skewing performance evaluations [43]. This underscores the need for robust frameworks to assess LLM temporal reasoning capabilities [84]. Existing approaches face challenges compared to uncompressed contexts, indicating a need for improved methodologies [104].

Systems like Mindful-RAG, enhancing intent identification and contextual alignment, show potential advancements [105]. However, challenges persist in managing noise when graph density is high, complicating retrieval accuracy [106].

Inconsistencies in retrieved contexts, especially from financial documents, complicate temporal reasoning [47]. The reliance on generated temporal graph quality and potential reasoning errors highlight the need for reliable methods [30].

Future research should enhance systems like YOURA in varied contexts and integrate them with other strategies to improve performance [107]. Improving retrieval processes to identify relevant facts needing updates and enhancing LLM capabilities for multi-hop reasoning are crucial areas [9]. The Daily Oracle benchmark emphasizes the need for continuous updates to maintain predictive accuracy, addressing LLM performance degradation [57].

#### 8.4 Integration of Knowledge Graphs

Integrating Knowledge Graphs (KGs) with LLMs enhances reliability and contextual understanding. KGs provide structured knowledge, addressing static dataset limitations and improving response accuracy in specialized domains like biomedical research [55].

A primary challenge in KG integration is developing multi-round mechanisms for LLMs to autonomously query knowledge bases, necessitating comprehensive evaluation frameworks for assessing LLM performance across diverse applications [55]. Exploring hybrid approaches combining RAG with fine-tuning is crucial for optimizing integration processes.

In educational contexts, integrating KGs with LLMs presents challenges related to adapting LLMs to domain-specific knowledge and effectively retrieving relevant information. Understanding strategies for embedding KGs into LLMs and evaluating enhancement techniques is essential for improving output accuracy and contextual relevance [24, 22, 105, 47]. Ensuring LLM-generated results' reliability is critical, highlighting the need for robust integration strategies.

Future research should prioritize enhancing embedding algorithms and integrating KGs to improve retrieval capabilities, particularly in knowledge-dense domains like law and medicine. Innovative frameworks like HybridRAG and Mindful-RAG, leveraging vector databases and KGs, can refine query intent recognition and context gathering, enhancing accuracy and relevance [82, 7, 21, 105, 47]. Developing comprehensive benchmarks for evaluating conversational agents and retrieval tasks is

---

essential for advancing KG integration with LLMs, providing insights into challenges and guiding effective strategies.

The integration of KGs with LLMs represents a significant research frontier, promising to enhance LLM capabilities by addressing hallucinations and domain-specific knowledge needs. Leveraging KGs' structured information can improve factual accuracy and interpretability, while LLMs can contribute to KG construction and validation. This synergy facilitates reliable AI applications and paves the way for innovative techniques in Natural Language Processing (NLP) [33, 90, 24]. Addressing challenges and exploring innovative integration strategies can enhance LLM output accuracy, reliability, and contextual understanding, expanding their applicability and effectiveness.

## 8.5 Hallucinations and Reliability

Hallucinations in LLMs pose a critical challenge, where models generate coherent but factually incorrect content. This issue is pronounced in RAG systems, where extensive contexts can introduce errors, affecting response reliability [20]. The reliance on curated datasets, which may not cover all inquiries, exacerbates this challenge [108].

Efforts to mitigate hallucinations include methodologies to improve LLM output reliability. The Know-RAG framework shows accuracy improvements and reduced hallucinations across domains [102]. SynapticRAG enhances memory retrieval accuracy and context awareness, simulating human-like memory processes [12]. Methods grounding LLM responses in relevant information also reduce hallucinations [18].

Challenges persist, as CoV-RAG relies heavily on query quality; poorly formulated queries can lead to suboptimal retrieval [37]. The correctness of scientific literature can propagate inaccuracies [40]. Linguistic inequalities in RALMs underscore the need for strategies to improve multilingual knowledge extraction [29].

Integrating KGs with RAG systems shows potential in providing structured, real-time information, reducing hallucinations and enhancing adaptability [34]. Ensuring retrieved chunks contain correct answers with specified confidence levels remains a challenge, as highlighted by approaches guaranteeing LLM output trustworthiness [71].

Future research should refine retrieval techniques and explore methods to abstract problem types, enhancing reasoning chain relevance [52]. Addressing data contamination and hallucination challenges in LLMs, alongside ethical implications in RAG deployments, remains crucial. Developing effective retrieval strategies and comprehensive evaluation frameworks can improve LLM output effectiveness and trustworthiness, ensuring relevance and accuracy in complex environments.

## 8.6 Ethical Considerations and Bias Mitigation

The integration of LLMs with RAG systems and TKGs necessitates examining ethical considerations and bias mitigation strategies. Ensuring transparency, accountability, and fairness is paramount, especially in sensitive domains [1]. A significant concern is biased outputs from inherent biases in training data and methodologies [19]. Addressing these biases requires transparent methodologies and robust accountability measures to ensure equitable LLM outputs.

The ethical implications of deploying LLMs, particularly in applications like ChatGPT, highlight the need for privacy and bias considerations. Some studies use synthetically generated data to safeguard user privacy, minimizing sensitive information exposure [63]. This underscores the importance of protecting user data while leveraging LLM capabilities.

Future research should optimize retrieval techniques and enhance graph data integration with language models to reduce biases and improve LLM output accuracy [17]. Decentralizing knowledge representation and developing methods to inject new facts into models are crucial for addressing knowledge uncertainty and ensuring LLMs remain unbiased [76].

Creating lifelong personal models for users emphasizes user-specific memory models' importance in enhancing personalization. This must be balanced with ethical considerations to prevent reinforcing existing biases and ensure fair information for all users [38]. Strategies for bias mitigation should consider ethical implications of using small-scale language models, as these may inherently possess limitations affecting outputs.



---

## 9 Conclusion

The convergence of Large Language Models (LLMs) with Temporal Knowledge Graphs (TKGs) and Retrieval-Augmented Generation (RAG) systems represents a transformative leap in the domains of information retrieval and reasoning. This survey elucidates the remarkable advancements achieved by integrating these technologies, notably enhancing retrieval precision, reasoning accuracy, and operational efficiency. Frameworks such as Topo-RAG exemplify how leveraging topological relationships can significantly refine text generation, addressing the constraints faced by earlier RAG methodologies. The symbiotic relationship between LLMs and Knowledge Graphs (KGs) has notably improved factual reasoning, establishing a robust framework for informed decision-making by anchoring outputs in timely and pertinent data.

The integration of Elasticsearch within the RAG framework has notably amplified the capabilities of question-answering systems, enhancing both accuracy and retrieval efficiency beyond conventional methods. This underscores the critical importance of sophisticated retrieval mechanisms in maximizing the potential of LLM and RAG systems.

Experimental evaluations of RAG-based scholarly QA systems demonstrate commendable performance in generating precise responses, yet highlight the need for continued enhancements in stability and efficiency. The potential for scaling inference computations to achieve near-linear improvements in RAG performance underscores the necessity of refining inference strategies in future research endeavors. Key insights reveal the broad applicability of RAG, the imperative to mitigate biases, and the ongoing need to explore retrieval mechanisms and model scalability.

Future research should aim to broaden benchmarks to include a wider array of formal languages and investigate the potential of code-pretrained models in Knowledge Base Question Answering (KBQA) tasks. By addressing existing challenges and exploring innovative research pathways, these technologies can be further honed to deliver more precise, contextually relevant, and efficient responses across diverse applications.

The amalgamation of LLMs, TKGs, and RAG systems holds substantial promise for propelling the frontiers of information retrieval and reasoning. The significant performance gains achieved through dynamic unsupervised learning techniques highlight the potential of integrating LLMs and TKGs to enhance system capabilities. Future investigations should continue to explore hybrid models that amalgamate diverse retrieval strategies and develop more efficient retrieval algorithms.

---

## References

- [1] Milad Moradi, Ke Yan, David Colwell, Matthias Samwald, and Rhona Asgari. Exploring the landscape of large language models: Foundations, techniques, and challenges, 2024.
- [2] Kazi Ahmed Asif Fuad and Lizhong Chen. Llm-ref: Enhancing reference handling in technical writing with large language models, 2024.
- [3] Jinyang Wu, Feihu Che, Xinxin Zheng, Shuai Zhang, Ruihan Jin, Shuai Nie, Pengpeng Shao, and Jianhua Tao. Can large language models understand uncommon meanings of common words?, 2024.
- [4] `<div style="text-align: center;"`.
- [5] Shangyu Wu, Ying Xiong, Yufei Cui, Haolun Wu, Can Chen, Ye Yuan, Lianming Huang, Xue Liu, Tei-Wei Kuo, Nan Guan, and Chun Jason Xue. Retrieval-augmented generation for natural language processing: A survey, 2024.
- [6] Weizhi Fei, Xueyan Niu, Guoqing Xie, Yanhua Zhang, Bo Bai, Lei Deng, and Wei Han. Retrieval meets reasoning: Dynamic in-context editing for long-text understanding, 2024.
- [7] Xiaohua Wang, Zhenghua Wang, Xuan Gao, Feiran Zhang, Yixin Wu, Zhibo Xu, Tianyuan Shi, Zhengyuan Wang, Shizheng Li, Qi Qian, Ruicheng Yin, Changze Lv, Xiaoqing Zheng, and Xuanjing Huang. Searching for best practices in retrieval-augmented generation, 2024.
- [8] Shailja Gupta, Rajesh Ranjan, and Surya Narayan Singh. A comprehensive survey of retrieval-augmented generation (rag): Evolution, current landscape and future directions. *arXiv preprint arXiv:2410.12837*, 2024.
- [9] Belinda Z. Li, Emmy Liu, Alexis Ross, Abbas Zeitoun, Graham Neubig, and Jacob Andreas. Language modeling with editable external knowledge, 2024.
- [10] Dan Schumacher, Fatemeh Haji, Tara Grey, Niharika Bandlamudi, Nupoor Karnik, Gagana Uday Kumar, Jason Cho-Yu Chiang, Paul Rad, Nishant Vishwamitra, and Anthony Rios. Context matters: An empirical study of the impact of contextual information in temporal question answering systems, 2024.
- [11] Yu Wang, Nedim Lipka, Ruiyi Zhang, Alexa Siu, Yuying Zhao, Bo Ni, Xin Wang, Ryan Rossi, and Tyler Derr. Augmenting textual generation via topology aware retrieval, 2024.
- [12] Yuki Hou, Haruki Tamoto, and Homei Miyashita. Integrating temporal representations for dynamic memory retrieval and management in large language models, 2024.
- [13] Sabrina Toro, Anna V Anagnostopoulos, Sue Bello, Kai Blumberg, Rhiannon Cameron, Leigh Carmody, Alexander D Diehl, Damion Dooley, William Duncan, Petra Fey, Pascale Gaudet, Nomi L Harris, Marcin Joachimiak, Leila Kiani, Tiago Lubiana, Monica C Munoz-Torres, Shawn O’Neil, David Osumi-Sutherland, Aleix Puig, Justin P Reese, Leonore Reiser, Sofia Robb, Troy Ruemping, James Seager, Eric Sid, Ray Stefancsik, Magalie Weber, Valerie Wood, Melissa A Haendel, and Christopher J Mungall. Dynamic retrieval augmented generation of ontologies using artificial intelligence (dragon-ai), 2024.
- [14] Ruotong Pan, Boxi Cao, Hongyu Lin, Xianpei Han, Jia Zheng, Sirui Wang, Xunliang Cai, and Le Sun. Not all contexts are equal: Teaching llms credibility-aware generation, 2024.
- [15] Jingyu Liu, Jiaen Lin, and Yong Liu. How much can rag help the reasoning of llm?, 2024.
- [16] Laurent Mombaerts, Terry Ding, Adi Banerjee, Florian Felice, Jonathan Taws, and Tarik Borogovac. Meta knowledge for retrieval augmented large language models, 2024.
- [17] Demiao Lin. Revolutionizing retrieval-augmented generation with enhanced pdf structure recognition, 2024.
- [18] Sanat Sharma, David Seunghyun Yoon, Franck Dernoncourt, Dewang Sultania, Karishma Bagga, Mengjiao Zhang, Trung Bui, and Varun Kotte. Retrieval augmented generation for domain-specific question answering, 2024.

- 
- [19] Zooney Nguyen, Anthony Annunziata, Vinh Luong, Sang Dinh, Quynh Le, Anh Hai Ha, Chanh Le, Hong An Phan, Shruti Raghavan, and Christopher Nguyen. Enhancing qa with domain-specific fine-tuning and iterative reasoning: A comparative study, 2024.
  - [20] Chenxi Dong, Yimin Yuan, Kan Chen, Shupeu Cheng, and Chujie Wen. How to build an adaptive ai tutor for any course using knowledge graph-enhanced retrieval-augmented generation (kg-rag), 2025.
  - [21] Yunxiao Shi, Xing Zi, Zijing Shi, Haimin Zhang, Qiang Wu, and Min Xu. Enhancing retrieval and managing retrieval: A four-module synergy for improved quality and efficiency in rag systems, 2024.
  - [22] Spurthi Setty, Harsh Thakkar, Alyssa Lee, Eden Chung, and Natan Vidra. Improving retrieval for rag based question answering models on financial documents, 2024.
  - [23] Siyun Zhao, Yuqing Yang, Zilong Wang, Zhiyuan He, Luna K. Qiu, and Lili Qiu. Retrieval augmented generation (rag) and beyond: A comprehensive survey on how to make your llms use external data more wisely, 2024.
  - [24] Ye Liu, Semih Yavuz, Rui Meng, Meghana Moorthy, Shafiq Joty, Caiming Xiong, and Yingbo Zhou. Modeling uncertainty and using post-fusion as fallback improves retrieval augmented generation with llms, 2024.
  - [25] Ayman Asad Khan, Md Toufique Hasan, Kai Kristian Kemell, Jussi Rasku, and Pekka Abrahamsson. Developing retrieval augmented generation (rag) based llm systems from pdfs: An experience report, 2024.
  - [26] Masoomali Fatehkia, Ji Kim Lucas, and Sanjay Chawla. T-rag: Lessons from the llm trenches, 2024.
  - [27] Jinxin Liu, Shulin Cao, Jiaxin Shi, Tingjian Zhang, Lunyiu Nie, Linmei Hu, Lei Hou, and Juanzi Li. How proficient are large language models in formal languages? an in-depth insight for knowledge base question answering, 2024.
  - [28] Zhihan Zhang, Yixin Cao, Chenchen Ye, Yunshan Ma, Lizi Liao, and Tat-Seng Chua. Analyzing temporal complex events with large language models? a benchmark towards temporal, long context understanding, 2024.
  - [29] Suhang Wu, Jialong Tang, Baosong Yang, Ante Wang, Kaidi Jia, Jiawei Yu, Junfeng Yao, and Jinsong Su. Not all languages are equal: Insights into multilingual retrieval-augmented generation, 2024.
  - [30] Siheng Xiong, Ali Payani, Ramana Kompella, and Faramarz Fekri. Large language models can learn temporal reasoning. *arXiv preprint arXiv:2401.06853*, 2024.
  - [31] Jiapu Wang, Kai Sun, Linhao Luo, Wei Wei, Yongli Hu, Alan Wee-Chung Liew, Shirui Pan, and Baocai Yin. Large language models-guided dynamic adaptation for temporal knowledge graph reasoning, 2024.
  - [32] Ruotong Liao, Xu Jia, Yangzhe Li, Yunpu Ma, and Volker Tresp. Gentkg: Generative forecasting on temporal knowledge graph with large language models, 2024.
  - [33] Shirui Pan, Linhao Luo, Yufei Wang, Chen Chen, Jiapu Wang, and Xindong Wu. Unifying large language models and knowledge graphs: A roadmap. *IEEE Transactions on Knowledge and Data Engineering*, 2024.
  - [34] Linyao Yang, Hongyang Chen, Zhao Li, Xiao Ding, and Xindong Wu. Give us the facts: Enhancing large language models with knowledge graphs for fact-aware language modeling, 2024.
  - [35] Jiheng Liang, Ziru Yu, Zujie Xie, and Xiangyang Yu. A quick, trustworthy spectral knowledge qa system leveraging retrieval-augmented generation on llm, 2024.
  - [36] Zijian Li, Qingyan Guo, Jiawei Shao, Lei Song, Jiang Bian, Jun Zhang, and Rui Wang. Graph neural network enhanced retrieval for question answering of llms, 2024.

- 
- [37] Bolei He, Nuo Chen, Xinran He, Lingyong Yan, Zhenkai Wei, Jinchang Luo, and Zhen-Hua Ling. Retrieving, rethinking and revising: The chain-of-verification can improve retrieval augmented generation, 2024.
- [38] Zijian Hei, Weiling Liu, Wenjie Ou, Juyi Qiao, Junming Jiao, Guowen Song, Ting Tian, and Yi Lin. Dr-rag: Applying dynamic document relevance to retrieval-augmented generation for question-answering, 2024.
- [39] Tristan Kenneweg, Philip Kenneweg, and Barbara Hammer. Retrieval augmented generation systems: Automatic dataset creation, evaluation and boolean agent setup, 2024.
- [40] Jakub Lála, Odhran O’Donoghue, Aleksandar Shtedritski, Sam Cox, Samuel G. Rodrigues, and Andrew D. White. Paperqa: Retrieval-augmented generative agent for scientific research, 2023.
- [41] Yunxiao Shi, Xing Zi, Zijing Shi, Haimin Zhang, Qiang Wu, and Min Xu. Eragent: Enhancing retrieval-augmented language models with improved accuracy, efficiency, and personalization, 2024.
- [42] Irwin Deng, Kushagra Dixit, Vivek Gupta, and Dan Roth. Enhancing temporal understanding in llms for semi-structured tables, 2024.
- [43] Yuqing Wang and Yun Zhao. Tram: Benchmarking temporal reasoning for large language models, 2024.
- [44] Jinyoung Kim, Dayoon Ko, and Gunhee Kim. Dynamicer: Resolving emerging mentions to dynamic entities for rag, 2024.
- [45] Mufei Li, Siqi Miao, and Pan Li. Simple is effective: The roles of graphs and large language models in knowledge-graph-based retrieval-augmented generation, 2025.
- [46] Zhentao Xu, Mark Jerome Cruz, Matthew Guevara, Tie Wang, Manasi Deshpande, Xiaofeng Wang, and Zheng Li. Retrieval-augmented generation with knowledge graphs for customer service question answering, 2024.
- [47] Bhaskarjit Sarmah, Benika Hall, Rohan Rao, Sunil Patel, Stefano Pasquali, and Dhagash Mehta. Hybridrag: Integrating knowledge graphs and vector retrieval augmented generation for efficient information extraction, 2024.
- [48] Boci Peng, Yun Zhu, Yongchao Liu, Xiaohe Bo, Haizhou Shi, Chuntao Hong, Yan Zhang, and Siliang Tang. Graph retrieval-augmented generation: A survey, 2024.
- [49] Junhao Zheng, Shengjie Qiu, Chengming Shi, and Qianli Ma. Towards lifelong learning of large language models: A survey, 2024.
- [50] Boci Peng, Yun Zhu, Yongchao Liu, Xiaohe Bo, Haizhou Shi, Chuntao Hong, Yan Zhang, and Siliang Tang. Graph retrieval-augmented generation: A survey. *arXiv preprint arXiv:2408.08921*, 2024.
- [51] Bingshuai Liu, Chenyang Lyu, Zijun Min, Zhanyu Wang, Jinsong Su, and Longyue Wang. Retrieval-augmented multi-modal chain-of-thoughts reasoning for large language models, 2024.
- [52] Eric Melz. Enhancing llm intelligence with arm-rag: Auxiliary rationale memory for retrieval augmented generation, 2023.
- [53] Yanxin Chen and Ling He. Research on the application of large language models in automatic question generation: A case study of chatglm in the context of high school information technology curriculum, 2024.
- [54] Hung Phan, Anurag Acharya, Rounak Meyur, Sarthak Chaturvedi, Shivam Sharma, Mike Parker, Dan Nally, Ali Jannesari, Karl Pazdernik, Mahantesh Halappanavar, Sai Munikoti, and Sameera Horawalavithana. Examining long-context large language models for environmental review document comprehension, 2024.

- 
- [55] Satyapriya Krishna, Kalpesh Krishna, Anhad Mohananey, Steven Schwarcz, Adam Stambler, Shyam Upadhyay, and Manaal Faruqui. Fact, fetch, and reason: A unified evaluation of retrieval-augmented generation, 2024.
  - [56] Shengjie Ma, Chengjin Xu, Xuhui Jiang, Muzhi Li, Huaren Qu, Cehao Yang, Jiaxin Mao, and Jian Guo. Think-on-graph 2.0: Deep and faithful large language model reasoning with knowledge-guided retrieval augmented generation, 2025.
  - [57] Hui Dai, Ryan Teehan, and Mengye Ren. Are llms prescient? a continuous evaluation using daily news as the oracle, 2024.
  - [58] Yuanjie Lyu, Zihan Niu, Zheyong Xie, Chao Zhang, Tong Xu, Yang Wang, and Enhong Chen. Retrieve-plan-generation: An iterative planning and answering framework for knowledge-intensive llm generation, 2024.
  - [59] Giorgio Roffo. Exploring advanced large language models with llmsuite, 2024.
  - [60] Xiaoxin He, Yijun Tian, Yifei Sun, Nitesh V. Chawla, Thomas Laurent, Yann LeCun, Xavier Bresson, and Bryan Hooi. G-retriever: Retrieval-augmented generation for textual graph understanding and question answering, 2024.
  - [61] Zefang Liu and Yin Zhu Quan. Retrieval of temporal event sequences from textual descriptions, 2025.
  - [62] al-augmented generation meets da.
  - [63] Geethan Sannidhi, Sagar Srinivas Sakhinana, and Venkataramana Runkana. Retrieval-augmented generation meets data-driven tabula rasa approach for temporal knowledge graph forecasting, 2024.
  - [64] Marc Pickett, Jeremy Hartman, Ayan Kumar Bhowmick, Raquib ul Alam, and Aditya Vempaty. Better rag using relevant information gain, 2025.
  - [65] Dominic Lohr, Marc Berges, Abhishek Chugh, Michael Kohlhase, and Dennis Müller. Leveraging large language models to generate course-specific semantically annotated learning objects, 2024.
  - [66] Zhongzhen Huang, Kui Xue, Yongqi Fan, Linjie Mu, Ruoyu Liu, Tong Ruan, Shaoting Zhang, and Xiaofan Zhang. Tool calling: Enhancing medication consultation via retrieval-augmented large language models, 2024.
  - [67] Jiarui Li, Ye Yuan, and Zehua Zhang. Enhancing llm factual accuracy with rag to counter hallucinations: A case study on domain-specific queries in private knowledge-bases, 2024.
  - [68] Md Faisal Mahbub Chowdhury, Michael Glass, Gaetano Rossiello, Alfio Gliozzo, and Nandana Mihindukulasooriya. Kgi: An integrated framework for knowledge intensive language tasks, 2022.
  - [69] Xiaoxi Li, Jiajie Jin, Yujia Zhou, Yongkang Wu, Zhonghua Li, Qi Ye, and Zhicheng Dou. Retrollm: Empowering large language models to retrieve fine-grained evidence within generation, 2024.
  - [70] Yun Jiang, Zilong Xie, Wei Zhang, Yun Fang, and Shuai Pan. E2e-afg: An end-to-end model with adaptive filtering for retrieval-augmented generation, 2024.
  - [71] Pouria Rouzrokh, Shahriar Faghani, Cooper U. Gamble, Moein Shariatnia, and Bradley J. Erickson. Conflare: Conformal large language model retrieval, 2024.
  - [72] Zhiruo Wang, Jun Araki, Zhengbao Jiang, Md Rizwan Parvez, and Graham Neubig. Learning to filter context for retrieval-augmented generation, 2023.
  - [73] Krish Goel and Mahek Chandak. Hiro: Hierarchical information retrieval optimization, 2024.
  - [74] Minbyul Jeong, Jiwoong Sohn, Mujeen Sung, and Jaewoo Kang. Improving medical reasoning through retrieval and self-reflection with retrieval-augmented large language models, 2024.

- 
- [75] Joel Suro. Semantic tokens in retrieval augmented generation, 2024.
- [76] Shi-Qi Yan, Jia-Chen Gu, Yun Zhu, and Zhen-Hua Ling. Corrective retrieval augmented generation, 2024.
- [77] Rubing Chen, Xulu Zhang, Jiaxin Wu, Wenqi Fan, Xiao-Yong Wei, and Qing Li. Multi-level querying using a knowledge pyramid, 2024.
- [78] Chao Jin, Zili Zhang, Xuanlin Jiang, Fangyue Liu, Xin Liu, Xuanzhe Liu, and Xin Jin. Ragcache: Efficient knowledge caching for retrieval-augmented generation, 2024.
- [79] Shuo Yu, Mingyue Cheng, Jiqian Yang, Jie Ouyang, Yucong Luo, Chenyi Lei, Qi Liu, and Enhong Chen. Multi-source knowledge pruning for retrieval-augmented generation: A benchmark and empirical study, 2025.
- [80] Anum Afzal, Juraj Vladika, Gentrit Fazlija, Andrei Staradubets, and Florian Matthes. Towards optimizing a retrieval augmented generation using large language model on academic data, 2024.
- [81] Florin Cuconasu, Giovanni Trappolini, Federico Siciliano, Simone Filice, Cesare Campagnano, Yoelle Maarek, Nicola Tonellotto, and Fabrizio Silvestri. The power of noise: Redefining retrieval for rag systems, 2024.
- [82] Guanhua Chen, Wenhan Yu, and Lei Sha. Unlocking multi-view insights in knowledge-dense retrieval-augmented generation, 2024.
- [83] Jiajing Chen, Runyuan Bao, Hongye Zheng, Zhen Qi, Jianjun Wei, and Jiacheng Hu. Optimizing retrieval-augmented generation with elasticsearch for enhanced question-answering systems, 2024.
- [84] Qingyu Tan, Hwee Tou Ng, and Lidong Bing. Towards robust temporal reasoning of large language models via a multi-hop qa dataset and pseudo-instruction tuning, 2024.
- [85] Heydar Soudani, Evangelos Kanoulas, and Faegheh Hasibi. Fine tuning vs. retrieval augmented generation for less popular knowledge, 2024.
- [86] Aidan Gilson, Xuguang Ai, Thilaka Arunachalam, Ziyu Chen, Ki Xiong Cheong, Amisha Dave, Cameron Duic, Mercy Kibe, Annette Kaminaka, Minali Prasad, Fares Siddig, Maxwell Singer, Wendy Wong, Qiao Jin, Tiarnan D. L. Keenan, Xia Hu, Emily Y. Chew, Zhiyong Lu, Hua Xu, Ron A. Adelman, Yih-Chung Tham, and Qingyu Chen. Enhancing large language models with domain-specific retrieval augment generation: A case study on long-form consumer health question answering in ophthalmology, 2024.
- [87] Nick Alonso, Tomás Figliolia, Anthony Ndirango, and Beren Millidge. Toward conversational agents with context and time sensitive long-term memory, 2024.
- [88] João Rodrigues and António Branco. Meta-prompting optimized retrieval-augmented generation, 2024.
- [89] Zhonghao Li, Xuming Hu, Aiwei Liu, Kening Zheng, Sirui Huang, and Hui Xiong. Refiner: Restructure retrieval content efficiently to advance question-answering capabilities, 2024.
- [90] Tuan Bui, Oanh Tran, Phuong Nguyen, Bao Ho, Long Nguyen, Thang Bui, and Tho Quan. Cross-data knowledge graph construction for llm-enabled educational question-answering system: A case study at hcmut, 2024.
- [91] Timer 4: Time-aware retrieval-au.
- [92] Yige Shen, Hao Jiang, Hua Qu, and Jihong Zhao. Think-then-act: A dual-angle evaluated retrieval-augmented generation, 2024.
- [93] Zongyue Qin, Chen Luo, Zhengyang Wang, Haoming Jiang, and Yizhou Sun. Relational database augmented large language model, 2024.

- 
- [94] Xuanwang Zhang, Yunze Song, Yidong Wang, Shuyun Tang, Xinfeng Li, Zhengran Zeng, Zhen Wu, Wei Ye, Wenyuan Xu, Yue Zhang, Xinyu Dai, Shikun Zhang, and Qingsong Wen. Raglab: A modular and research-oriented unified framework for retrieval-augmented generation, 2024.
- [95] Zhuowan Li, Cheng Li, Mingyang Zhang, Qiaozhu Mei, and Michael Bendersky. Retrieval augmented generation or long-context llms? a comprehensive study and hybrid approach, 2024.
- [96] Zhengliang Shi, Weiwei Sun, Shen Gao, Pengjie Ren, Zhumin Chen, and Zhaochun Ren. Generate-then-ground in retrieval-augmented generation for multi-hop question answering, 2024.
- [97] Brian J Chan, Chao-Ting Chen, Jui-Hung Cheng, and Hen-Hsen Huang. Don’t do rag: When cache-augmented generation is all you need for knowledge tasks, 2025.
- [98] Ruobing Wang, Daren Zha, Shi Yu, Qingfei Zhao, Yuxuan Chen, Yixuan Wang, Shuo Wang, Yukun Yan, Zhenghao Liu, Xu Han, Zhiyuan Liu, and Maosong Sun. Retriever-and-memory: Towards adaptive note-enhanced retrieval-augmented generation, 2024.
- [99] Yizheng Huang and Jimmy Huang. A survey on retrieval-augmented text generation for large language models. *arXiv preprint arXiv:2404.10981*, 2024.
- [100] Chunjing Gan, Dan Yang, Binbin Hu, Hanxiao Zhang, Siyuan Li, Ziqi Liu, Yue Shen, Lin Ju, Zhiqiang Zhang, Jinjie Gu, Lei Liang, and Jun Zhou. Similarity is not all you need: Endowing retrieval augmented generation with multi layered thoughts, 2024.
- [101] Hamed Babaei Giglou, Tilahun Abedissa Taffa, Rana Abdullah, Aida Usmanova, Ricardo Usbeck, Jennifer D’Souza, and Sören Auer. Scholarly question answering using large language models in the nfdi4datascience gateway, 2024.
- [102] Know-rag: An adaptive approach f.
- [103] Event temporal relation extracti.
- [104] Sourav Verma. Contextual compression in retrieval-augmented generation for large language models: A survey, 2024.
- [105] Garima Agrawal, Tharindu Kumarage, Zeyad Alghamdi, and Huan Liu. Mindful-rag: A study of points of failure in retrieval augmented generation, 2024.
- [106] Tiezheng Guo, Chen Wang, Yanyi Liu, Jiawei Tang, Pan Li, Sai Xu, Qingwen Yang, Xianlin Gao, Zhi Li, and Yingyou Wen. Leveraging inter-chunk interactions for enhanced retrieval in large language model-based question answering, 2024.
- [107] Yun Joon Soh, Hanxian Huang, Yuandong Tian, and Jishen Zhao. You only use reactive attention slice for long context retrieval, 2024.
- [108] Liang Zhang, Katherine Jijo, Spurthi Setty, Eden Chung, Fatima Javid, Natan Vidra, and Tommy Clifford. Enhancing large language model performance to answer questions and extract information more accurately, 2024.

---

**Disclaimer:**

SurveyX is an AI-powered system designed to automate the generation of surveys. While it aims to produce high-quality, coherent, and comprehensive surveys with accurate citations, the final output is derived from the AI's synthesis of pre-processed materials, which may contain limitations or inaccuracies. As such, the generated content should not be used for academic publication or formal submissions and must be independently reviewed and verified. The developers of SurveyX do not assume responsibility for any errors or consequences arising from the use of the generated surveys.

www.SurveyX.cn