

A Comprehensive Approach to Summarization: The Role of Retrieval-Augmented Generation

Vanita Mhaske¹, Rhishikesh Kadam²

¹Assistant Professor, PVG's College of Science and Commerce, Pune.

²Individual Researcher, Pune.

Abstract

There are many types of summaries that can be created, depending on the nature of the input documents—whether they are related to law, medicine, or other fields. It is important to understand the subject matter first, as different documents require different handling approaches. Highlighting key points is crucial for focusing on specific sentences. Depending on the topic and the desired output, various summarization models can be used. In this context, we particularly focus on the RAG system and how it can be beneficial for achieving better results.

Keywords: Retrieval methods, Generative methods, Dynamic knowledge integration, Fluent language generation, RAG (Retrieval-Augmented Generation)

1. Introduction

Before the advent of advanced methods like RAG, various approaches were employed to address language-related problems. These can broadly be categorized into two types:

1. **Extractive (Retrieval of Information Without Altering Sentence Structure)**
2. **Abstractive (Retrieval of Information With Sentence Reformation)**

These methods laid the foundation for advanced frameworks like RAG, which combine retrieval with generation for more dynamic and accurate results.

1. Retrieval-Based Methods

These methods retrieve relevant data while maintaining the original wording and structure of the sentences. Common examples include extractive summarization, where key phrases or sentences are directly picked from the source without modification [1].

a. Traditional Information Retrieval (IR) Methods

TF-IDF (Term Frequency-Inverse Document Frequency):

Ranked documents based on term frequency-inverse document frequency scores. Efficient but lacked semantic understanding. **TF-IDF** is a statistical measure used to evaluate the importance of a word in a document relative to a collection (or corpus) of documents [2]. It is widely used in text mining and information retrieval tasks, such as document ranking and keyword extraction [3].

BM25 (Best Matching 25):

An improvement over TF-IDF, incorporating term saturation and document length normalization [4]. Widely used in search engines for its relevance-based ranking. **BM25** is an advanced information retrieval algorithm that builds on the foundation of **TF-IDF**, addressing some of its key limitations [5]. It is widely regarded as a robust and effective ranking function for search engines and text retrieval system

b. Neural Retrieval Models

Dense Representations:

Models like **DPR (Dense Passage Retrieval)** used embeddings (e.g., from BERT) to represent queries and documents in a shared semantic space. Dense Representations are advanced techniques in information retrieval and natural language processing that encode text (queries and documents) into fixed-length vectors in a shared semantic space [6]. Unlike traditional methods that rely on lexical matching, dense representations focus on capturing the semantic meaning of text, enabling more effective and context-aware retrieval [7].

Siamese Networks:

Used for encoding queries and documents into vector spaces for similarity comparisons. Siamese Networks are a type of neural network architecture designed to learn the similarity between two inputs by mapping them into a shared vector space [8]. This architecture has been widely used for tasks like **semantic textual similarity**, **duplicate detection**, and **retrieval-based applications**.

Dual-Encoder Models:

Learned better embeddings for retrieving semantically relevant documents. **Dual-Encoder Models** are a type of architecture used for learning efficient, semantic representations of text for tasks like **document retrieval**, **question answering**, and **semantic search** [9]. These models process queries and documents separately but map both to a shared embedding space where semantically similar items are closer together.

2. Generative Models

These approaches retrieve information but modify the sentence structure to ensure better readability, coherence, or contextual relevance. This is commonly seen in abstractive summarization, where the essence of the information is retained, but the phrasing is reformulated to provide a concise and fluent output.

a. Early Generative Models

RNNs and LSTMs:

Used for sequence generation but struggled with long-range dependencies. **Recurrent Neural Networks (RNNs)** and **Long Short-Term Memory networks (LSTMs)** are types of neural networks designed to handle sequential data, where the order of information is important, such as in **language modeling**, **text generation**, and **speech recognition** [10]. These models are designed to process inputs sequentially and maintain a form of memory to capture the relationships between elements in the sequence.

Seq2Seq Models:

Encoder-decoder frameworks (e.g., for machine translation or text summarization). **Seq2Seq (Sequence-to-Sequence) Models** are a type of neural network architecture used for tasks that involve transforming one sequence of data into another, such as **machine translation**, **text summarization**, and **speech recognition** [11]. These models consist of two main components: an **encoder** and a **decoder**, which work together to process and generate sequences of different lengths.

b. Transformer-Based Models

BERT (Bidirectional Encoder Representations from Transformers):

BERT (Bidirectional Encoder Representations from Transformers), introduced in 2018 by Google, revolutionized the field of natural language processing (NLP) by significantly improving the performance of a wide range of language understanding tasks. BERT is based on the **Transformer** architecture, which uses self-attention mechanisms to process the input text in parallel, making it faster

and more efficient compared to previous models like RNNs and LSTMs. Focused on masked language modeling for understanding context but wasn't generative.

GPT (Generative Pre-trained Transformer):

GPT developed by OpenAI, is a family of autoregressive language models designed for natural language generation tasks. Unlike models like BERT, which are primarily focused on language understanding, GPT is built to generate coherent, contextually appropriate text based on an initial prompt or input. The model's ability to generate text is one of its standout features, making it particularly valuable for tasks that require creative writing, dialogue generation, and other forms of open-ended text generation. Enabled text generation but relied entirely on its pre-trained knowledge.

T5 (Text-to-Text Transfer Transformer):

T5 developed by Google Research in 2019, is a versatile model designed to handle a wide range of natural language processing (NLP) tasks. T5 is based on the **Transformer** architecture, and its key innovation is the conversion of all NLP tasks into a unified **text-to-text** format. This means that regardless of the specific task (e.g., translation, summarization, question answering), both the input and output are treated as text sequences. This approach simplifies the model's design, making it more flexible and easier to fine-tune for different tasks. Converted all NLP tasks into a text-to-text format but still suffered from hallucination and outdated knowledge.

BART:

BART is a transformer-based model introduced by Facebook AI in 2019. It combines the best of both **autoregressive models** (like GPT) and **autoencoding models** (like BERT) to create a versatile model for both understanding and generating text. BART was specifically designed for **sequence-to-sequence tasks**, such as text summarization, translation, and dialogue generation. Its unique approach to pre-training allows it to excel in tasks that require both text generation and understanding, such as text summarization and conversation generation. Combined denoising pre-training with text generation, improving on summarization and dialogue tasks.

3. Hybrid Approaches

Before RAG, some systems tried combining retrieval and generation, but they lacked the tight integration seen in RAG.

a. Retrieve-Then-Read

Systems first retrieved relevant documents, which were then summarized or synthesized by a generative model. Examples:

- **Open-Domain Question Answering:** Systems like DrQA used retrieval methods (e.g., TF-IDF) to find passages, followed by a reader (e.g., BiDAF) for extractive answers.
- **Generative QA Models:** Early versions used retrieval to provide context for language models like GPT.

b. Fusion-Based Models

Fusion-in-Decoder (FiD) merged retrieved documents directly into the decoding process for generation.

Graph-Based Approaches Leveraged knowledge graphs to add structured information to generative models.

Limitations:

Early retrieval and generative models faced several challenges. They often relied on simplistic techniques like **TF-IDF**, which lacked deep contextual understanding, resulting in less accurate retrieval and limited generation quality. The integration of context was minimal, leading to disjointed outputs.

Generative models also struggled with **hallucination**, where they produced incorrect or nonsensical information due to reliance on static training data. Additionally, these models had difficulty handling **dynamic** or **domain-specific knowledge** without retraining, which hindered their real-time applicability. Furthermore, **retrieval-based methods** couldn't synthesize new information, limiting their effectiveness in tasks requiring contextual generation or creativity

2. Conclusion

While both retrieval and generative methods were powerful, they each had limitations in handling tasks that required dynamic knowledge integration and fluent, coherent language generation. **RAG** (Retrieval-Augmented Generation) effectively addressed these challenges by combining retrieval with generation. It enabled **dynamic knowledge updates** from external databases, allowing models to stay current with new information. Additionally, it provided **grounded generation**, reducing hallucination and ensuring more accurate outputs. RAG's design also supported **scalability**, making it capable of handling complex or domain-specific tasks. This innovation filled the gaps in earlier approaches, paving the way for more reliable and accurate NLP systems.

3. References

1. Jain, D., Borah, M. D., & Biswas, A. (2021). Summarization of legal documents: Where are we now and the way forward. *Computer Science Review*, 40, 100388.
2. Havrlant, L., & Kreinovich, V. (2017). A simple probabilistic explanation of term frequency-inverse document frequency (tf-idf) heuristic (and variations motivated by this explanation). *International Journal of General Systems*, 46(1), 27-36.
3. Christian, H., Agus, M. P., & Suhartono, D. (2016). Single document automatic text summarization using term frequency-inverse document frequency (TF-IDF). *ComTech: Computer, Mathematics and Engineering Applications*, 7(4), 285-294.
4. Jian, F., Huang, J. X., Zhao, J., & He, T. (2018, June). A new term frequency normalization model for probabilistic information retrieval. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval* (pp. 1237-1240).
5. Xu, R. (2014). *POS weighted TF-IDF algorithm and its application for an MOOC search engine*. 2014 International Conference on Audio, Language and Image Processing, 868-873.
6. Althammer, S., Hofstätter, S., Sertkan, M., Verberne, S., & Hanbury, A. (2022, April). PARM: A paragraph aggregation retrieval model for dense document-to-document retrieval. In *European Conference on Information Retrieval* (pp. 19-34). Cham: Springer International Publishing.
7. Xiao, Q., Li, S., & Chen, L. (2023). Topic-DPR: Topic-based Prompts for Dense Passage Retrieval. *arXiv preprint arXiv:2310.06626*.
8. Chicco, D. (2021). Siamese neural networks: An overview. *Artificial neural networks*, 73-94.
9. Wang, J., Huang, J. X., & Sheng, J. (2024). An efficient long-text semantic retrieval approach via utilizing presentation learning on short-text. *Complex & Intelligent Systems*, 10(1), 963-979.
10. Sak, H., Senior, A. W., & Beaufays, F. (2014). Long short-term memory recurrent neural network architectures for large scale acoustic modeling.
11. Keneshloo, Y., Shi, T., Ramakrishnan, N., & Reddy, C. K. (2019). Deep reinforcement learning for sequence-to-sequence models. *IEEE transactions on neural networks and learning systems*, 31(7), 2469-2489.

12. Gayval, B., & Mhaske, V. (2023). Evaluation of Descriptive Probability Approach, Cosi Pretrained mo. *Journal of Scientific Research*, 67(2).