

# Summary of Papers

Boudjenane Zoubida Asmaa

February 8, 2025

paper of : **Retrieval-Augmented Generation for AI-Generated Content: A Survey** (Penghao Zhao et al.)

**Comprehensive RAG overview :** This paper provides a broad review of Retrieval-Augmented Generation (RAG), covering key areas such as existing retrieval methods categorized into:

- 1)-Sparse Retriever
- 2)-dense retriever
- 3)-others : (AST) abstract syntax trees , k-hop neighbor searches , Entity Recognition (NER)

and generation introduce 4 typical generators that are frequently used in RAG : transformer model , LSTM , Diffusion model , GAN .

**foundational paradigms of RAG** categorize RAG foundations into 4 classes

- 1)-Query-based RAG
- 2)-Latent Representation-based RAG
- 3)-Logit-based RAG
- 4)-Speculative RAG

**Enhancements in RAG:** Innovations and improvements across different retrieval and generation methodologies. including :

- input enhancements
- retriever enhancement
- generator enhancement
- result enhancement
- RAG Pipeline Enhancement)

**Applications:** Examines how RAG is used across different modalities and AI-generated content (AIGC) methods

- 1)-rag for text : Question Answering , Fact Verification, Commonsense Reasoning, Human-Machine Conversation , Neural Machine Translation, Event Extraction, Summarization
- 2)-rag for code : code generation , code summarization , code completion , automatic program repair , text-to -sql
- 3)-rag for knowledge : Knowledge Base Question Answering, Knowledge-augmented Open-domain Question Answering , Table for Question Answering
- 4)-rag for image : Image Generation , ImageCaptioning , 5)-rag for video : Video Captioning , Video QA and Dialogue
- 6)-rag for audio: Audio Generation , Audio Captioning 7)-rag for 3d: Text-to-3D
- 8)-rag for science : drug discovery , Biomedical Informatics Enhancement , Math Applications

**Challenges and Future Directions:** Identifies limitations and proposes research opportunities for enhancing RAG models.

### **Challenges**

- Noises in Retrieval Results
- Extra Overhead
- Increased System Complexity
- Lengthy Context
- Novel Design of Augmentation Methodologies

**Potential Future Directions**

- Flexible RAG Pipelines
- Broader Applications
- Efficient Deployment and Processing
- Incorporating Long-tail and Real-time Knowledge
- Combined with Other Techniques:

paper of **Retrieval-Augmented Generation for Large Language Models: A Survey (Yunfan Gao et al.)**

This paper provides a detailed and structured survey of Retrieval-Augmented Generation (RAG), mapping its evolution, core technologies, evaluation methods, and future directions. 1)-It focuses on how RAG integrates with Large Language Models (LLMs) and categorizes research into three paradigms:

**Comprehensive Review of RAG:**

-Naive RAG – Basic retrieval-based models.

-Advanced RAG – More optimized retrieval techniques and better generation models.

-Modular RAG – A more structured approach, integrating retrieval, generation, and augmentation as distinct but interdependent components.

-RAG vs Fine-tuning

**-Analysis of Core Components:**

**Retrieval:** Retrieval Source, Indexing methods, query optimization, and embedding strategies.

**Generator:** How LLMs process retrieved documents to generate responses adjusting the retrieved content including (Context Curation, reranking, Context Selection/Compression) and adjusting the LLM (LLM Fine-tuning)

**Augmentation:** Techniques that improve response quality, including re-ranking and reinforcement learning.

**Evaluation Framework for RAG:**

Covers 26 downstream tasks and 50 datasets. Summarizes current benchmarks, evaluation metrics, and tools used to assess RAG performance.

**Challenges and Future Directions:**

Identifies current limitations in retrieval accuracy, model efficiency, and scalability. Proposes future research areas to enhance RAG models, including better retrieval integration, multimodal capabilities, and low-resource language adaptation.

our paper :an overview of RAG with innovative aspects:

## **Introduction**

**Background and Related Work** : Existing RAG surveys and their focus ,  
Gap analysis

## **Foundations of Retrieval-Augmented Generation (RAG)**

How RAG Works: key components:

### **1) Indexing** (Data Preparation for Retrieval)

Chunking Methods: Fixed-length, semantic, hierarchical, and adaptive chunking.

Embedding Models: Lexical (BM25, TF-IDF) vs. Neural embeddings (BERT, SBERT, FAISS).

Indexing Techniques: Flat search, HNSW, IVF-PQ, and hybrid indexing for speed and accuracy.

### **2) Retrieval** (Finding Relevant Information)

Retrieval Models: Sparse (BM25), Dense Retrieval (DPR, Contriever) – Improves document retrieval quality.

Hybrid Retrieval (Dense + Sparse) – Combines BM25 with neural retrievers for better accuracy.

Memory-Augmented RAG – Stores past retrievals for future use, improving contextual continuity. Query Optimization: Neural query expansion, relevance feedback, and self-improving retrieval

### **3) Generator:** Uses the retrieved content to generate more accurate and context-aware text.

Comparison with traditional NLP models (GPT, T5, BERT) and why RAG is better for knowledge-intensive tasks.

## **Efficient RAG for Low-Resource Languages**

Knowledge distillation to train smaller, faster RAG models.

Fine-tuning RAG on domain-specific datasets for specialized applications (law, healthcare, etc.).

How these innovations reduce computational costs, making RAG more accessible

## **Applications of RAG in various NLP tasks:**

Question answering

Summarization

Code generation

Conversational AI

## **Challenges and Limitations of RAG**

Retrieval quality issues – How to ensure retrieved documents are relevant and factually correct.

Computational overhead – RAG is more expensive than standard generative models.

Bias in retrieval sources – If retrieval data is biased, the generated output will be as well.

Security risks – How adversarial attacks can manipulate retrieval results.

## **Arabic AI Generation and the Role of RAG**

### **1. The State of Arabic NLP**

Why Arabic AI generation is underdeveloped:

Morphological complexity – Arabic is root-based, making tokenization harder.

Dialects vs. Standard Arabic – No single Arabic dataset covers all dialects.

Lack of high-quality training data – Arabic corpora are smaller and less diverse than English ones.

### **Existing Arabic AI Models and Their Limitations**

AraBERT – Good for classification but not generative tasks.

AraGPT – Lacks high-quality retrieval mechanisms.

Arabic-T5 –

### **. How RAG Can Improve Arabic AI**

Future Research Directions Developing open-source Arabic RAG datasets. Building retrieval modules for Arabic-specific knowledge bases.

Exploring reinforcement learning for better retrieval in Arabic.

Creating multilingual RAG models with Arabic support.

### **Future Research Directions**

Developing open-source Arabic RAG datasets.

Building retrieval modules for Arabic-specific knowledge bases.

Exploring reinforcement learning for better retrieval in Arabic.

Creating multilingual RAG models with Arabic support.