

VIETNAM NATIONAL UNIVERSITY, HO CHI MINH CITY
UNIVERSITY OF TECHNOLOGY
FACULTY OF COMPUTER SCIENCE AND ENGINEERING



Programming Intergration Project (CO3101)

Group 4:

“Research and build AI chatbots using Retrieval-Augmented Generation for a music-related website, with Speech-to-Text and Text-to-Speech integration”

Instructor(s): Nguyễn Quốc Minh

Students: Nguyễn Thiên Minh - 2312097
Huỳnh Đức Nhân - 2312420
Phạm Trần Minh Trí - 2313622

HO CHI MINH CITY, DECEMBER 2025



Contents

List of Figures	5
List of Tables	5
Member list & Workload	5
Source code	5
1 Introduction	6
1.1 Background and Motivation	6
1.2 Problem Statement	6
1.3 Report Structure	6
2 Prerequisite: ANN, Transformer, LLM	7
2.1 Artificial Neural Network	7
2.1.1 Overall Architecture of Neural Networks	7
2.1.2 Mathematical Model of an Artificial Neuron	7
2.1.3 Activation Functions	8
2.1.4 Single-layer and Multi-layer Perceptrons	8
2.1.5 Forward Propagation	9
2.1.6 Loss Function	10
2.1.7 Backpropagation and Weight Update	11
2.2 Transformer architecture	12
2.2.1 Introduction and Background	12
2.2.2 Self-Attention Mechanism	13
2.2.3 Multi-Head Attention (MHA)	14
2.2.4 Feed-Forward Network and Add-Norm	14
2.2.5 Embeddings and Positional Encoding	14
2.2.6 Operational Flow	15
2.2.7 Comparative Analysis	15
2.3 Large language model	15
2.3.1 Evolution and Architecture	16
2.3.2 Pretraining and Fine-Tuning	16
2.3.3 LLMs in Chatbot Systems	17
3 Agent, URAG, LangChain	19
3.1 AI Agent	19
3.1.1 Introduction to AI Agents	19
3.1.2 Core Components of an AI Agent	19
3.1.3 Levels of Agent Agency	19
3.1.4 Thought–Action–Observation Cycle	20
3.1.5 Practical Illustration of the TAO Cycle	20
3.1.6 Observation and Response Generation	21
3.1.7 Internal Reasoning: Chain-of-Thought and ReAct	21
3.1.8 Actions: Enabling Interaction with the Environment	22
3.1.9 Observation and Adaptation	22
3.2 Unified Hybrid RAG (URAG)	22
3.2.1 Challenges in Conventional RAG	22



3.2.2	Evolution of Enhancement Strategies	22
3.2.3	The URAG Framework: Tiered Architecture	23
3.2.4	URAG Preparation Workflow	24
3.3	Introduction to LangChain	25
4	Retrieval-Augmented Generation for Large Language Models: A Survey	27
4.1	Overview of RAG	27
4.1.1	Naive RAG	27
4.1.2	Advanced RAG	27
4.1.3	Modular RAG	28
4.2	Retrieval	28
4.2.1	Retrieval Sources and Granularity	28
4.2.2	Indexing Optimization	29
4.2.3	Query Optimization	29
4.2.4	Embeddings and Adapters	29
4.3	Generation	29
4.3.1	Overview of the Generation Stage	29
4.3.2	Context Curation for Improved Generation	30
4.3.3	LLM Fine-tuning for Generation	31
4.4	Augmentation	32
4.4.1	Overview of the Augmentation Process	32
4.4.2	Iterative Retrieval	33
4.4.3	Recursive Retrieval	34
4.4.4	Adaptive Retrieval	34
4.4.5	Summary	35
4.5	Tasks and Evaluation Framework	35
4.5.1	Downstream Tasks and Benchmarks	35
4.5.2	Evolution of Evaluation Targets	36
4.5.3	Evaluation Aspects: The "RAG Triad" and Key Abilities	36
4.6	Discussion	36
5	Reranking, RAG-Reasoning, RAG-RL	38
5.1	Reranking	38
5.1.1	Limitations of Vector-Based Retrieval	38
5.1.2	Reranking as a Post-Retrieval Refinement Step	38
5.1.3	Bi-Encoder versus Cross-Encoder for Reranking	39
5.1.4	Summary	40
5.2	Towards Agentic RAG with Deep Reasoning: A Survey of RAG Reasoning Systems in LLMs	40
5.3	RAG-RL: Advancing Retrieval Augmented Generation via RL and Curriculum Learning	43
5.3.1	The Problem of Imperfect Retrieval	43
5.3.2	Foundational Concepts: RL and CL	43
5.3.3	The Synergy of RL and CL	44
5.3.4	Methodology	44
5.3.5	Training Schedules and Results	45



6 Implementation: LangChain, Text-to-speech, Speech-to-text	46
6.1 Building an AI Agent Chatbot with LangChain	46
6.1.1 RAG Agent Architecture	46
6.1.2 Agent Tools	46
6.1.3 Agent Configuration and Memory	47
6.2 Text-to-Speech	47
6.2.1 Text-to-Speech Pipeline	47
6.2.2 gTTS: Google Text-to-Speech	47
6.3 Speech-to-text	48
6.3.1 SpeechRecognition with Google Web Speech API	48
6.3.2 OpenAI Whisper	49
6.3.3 Experimental Comparison and Evaluation	49
6.3.4 Conclusion for Prototype Implementation	50
6.4 Python Flask app for the chatbot	50
6.4.1 Backend Core	50
6.4.2 Conversational Agent & RAG Pipeline	50
7 Evaluation	52
8 Conclusion	60
References	61



List of Figures

1	Basic architecture of a feedforward neural network	7
2	Mathematical structure of an artificial neuron	8
3	Common activation functions: Sigmoid, Tanh, and ReLU	8
4	Comparison between single-layer and multi-layer perceptrons	9
5	Forward propagation through hidden layers	10
6	Illustration of prediction error and loss minimization	11
7	Training process: forward propagation, backpropagation, and weight update	11
8	Transformer architecture	13
9	MultiHead Attention	14
10	Transformer architectures	16
11	RLHF for ChatGPT	17
12	Retrieval augmented generation	17
13	High-level AI agent workflow: user request, reasoning, and tool-based action	19
14	The Thought–Action–Observation (TAO) cycle in an AI agent	20
15	Reasoning and action steps using an external weather API	21
16	Observation feedback and final response generation	21
17	URAG idea	23
18	URAG workflow	24
19	Preparation of URAG	25
20	RAG - Indexing	26
21	RAG - Retrieval and Generation	26
22	RAG overview	27
23	Types of RAG	28
24	Generation stage in Naive RAG, where retrieved context is directly passed to a frozen LLM via a prompt	30
25	Advanced RAG with post-retrieval context curation before generation	31
26	Generation stage using a fine-tuned LLM in an Advanced RAG pipeline	32
27	Overview of the iterative augmentation process in RAG	33
28	Iterative retrieval with repeated retrieval-generation cycles	33
29	Recursive retrieval with query transformation and decomposition	34
30	Adaptive retrieval where the LLM actively controls retrieval decisions	35
31	Approximate nearest neighbor search may return vectors that are close in embedding space but not truly semantically relevant	38
32	Reranking as an intermediate step between retrieval and generation in Advanced RAG	39
33	Comparison between bi-encoder retrieval and cross-encoder reranking architectures	39
34	Model performance under the distractor and ideal retrieval evaluation settings across different curriculum construction settings	45
35	AI Agent	46
36	General Text-to-Speech pipeline from text input to synthesized speech output	47
37	Example of using gTTS in Python to generate an MP3 file from text	48

List of Tables

1	Member list & workload	5
2	Comparison: LSTM/GRU vs. Transformer	15
4	Comparison: SpeechRecognition API vs. Whisper Model	49



Member list & Workload

No.	Fullname	Student ID	Problems	% done
1	Nguyễn Thiện Minh	2312097	- ANN, Agent - Reranking - Text to speech	100%
2	Huỳnh Đức Nhân	2312420	- Transformer, URAG - RAG-RL - Speech to text	100%
3	Phạm Trần Minh Trí	2313622	- LLM, LangChain - RAG-Reasoning - Evaluation	100%

Table 1: Member list & workload

Source code for the RAG chat app

The source code can be accessed here: https://github.com/phamtranminhtri/music_rag



1 Introduction

1.1 Background and Motivation

In recent years, the field of Artificial Intelligence has witnessed a paradigm shift with the emergence of Large Language Models (LLMs) such as GPT-4, Gemini, and Llama. These models, built upon the Transformer architecture, have demonstrated remarkable capabilities in natural language understanding, generation, and reasoning. They can draft emails, write code, and even engage in complex conversations that mimic human interaction.

Despite their impressive performance, LLMs face inherent limitations. One significant challenge is "hallucination," where models generate plausible but factually incorrect information. Additionally, purely parametric models are limited by their training data cut-off dates and lack access to private or real-time proprietary data. For specialized domains, such as a dedicated music platform, generic LLMs may not provide specific, deep, or curated knowledge required to satisfy user queries accurately.

1.2 Problem Statement

To address these limitations, Retrieval-Augmented Generation (RAG) has emerged as a powerful framework. By grounding the generation process in external, verifiable knowledge sources, RAG significantly reduces hallucinations and allows models to access up-to-date information without the need for frequent retraining.

Furthermore, user interaction with AI systems is evolving beyond text-based interfaces. The integration of Speech-to-Text (STT) and Text-to-Speech (TTS) technologies enables more natural, hands-free, and accessible communication. Building a comprehensive AI system that combines advanced retrieval mechanisms with multi-modal interaction capabilities represents a significant engineering challenge that requires integrating various state-of-the-art components.

1.3 Report Structure

The remainder of this report is organized as follows:

- **Section 2** provides the theoretical prerequisites, covering Artificial Neural Networks, Transformer architectures, and Large Language Models.
- **Section 3** introduces the concepts of AI Agents, Universal RAG (URAG), and the LangChain framework.
- **Section 4** presents a survey of Retrieval-Augmented Generation, detailing the standard "Retrieve-Read" pipeline.
- **Section 5** discusses advanced RAG techniques, such as Reranking, to improve retrieval quality and reduce noise.
- **Section 6** details the implementation of our music chatbot, utilizing LangChain, Wikipedia data loaders, and voice technologies.
- **Section 7** presents the evaluation methodology and results, using metrics like Cosine Similarity and LLM-as-a-Judge.
- **Section 8** concludes the report and discusses potential future improvements.

2 Prerequisite: ANN, Transformer, LLM

2.1 Artificial Neural Network

Artificial Neural Networks (ANNs) are computational models inspired by the biological neural system. An ANN consists of multiple artificial neurons interconnected through weighted connections, enabling the model to learn complex mappings between input data and output predictions. Neural networks are widely used in classification, regression, and pattern recognition tasks due to their strong representation capability.

2.1.1 Overall Architecture of Neural Networks

A typical feedforward neural network is composed of three main components: an input layer, one or more hidden layers, and an output layer. Information flows from the input layer through hidden layers to produce the final output.

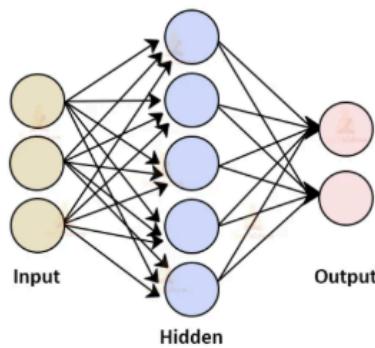


Figure 1: Basic architecture of a feedforward neural network

Each neuron in one layer is usually connected to all neurons in the next layer through weighted edges, allowing the network to model complex relationships.

2.1.2 Mathematical Model of an Artificial Neuron

An artificial neuron performs a weighted summation of its inputs followed by a nonlinear transformation. Given n input features x_i , the neuron computes an intermediate value z as:

$$z = \sum_{i=1}^n w_i x_i + b \quad (1)$$

where w_i denotes the weight associated with input x_i , and b is the bias term. The output of the neuron is obtained by applying an activation function $f(\cdot)$:

$$y = f(z) \quad (2)$$

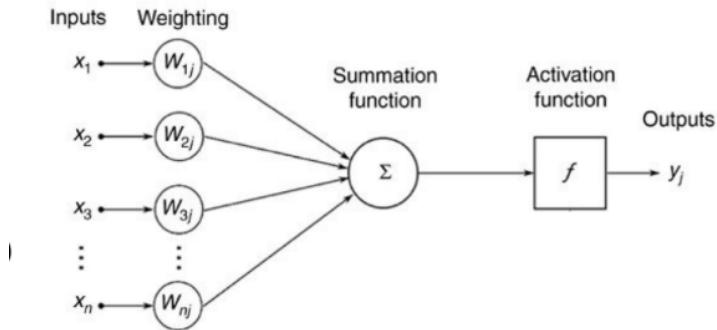


Figure 2: Mathematical structure of an artificial neuron

2.1.3 Activation Functions

Activation functions introduce non-linearity into the neural network. Without them, the network would behave as a linear model and fail to capture complex patterns in data. Common activation functions include:

- **Sigmoid:** maps input values into the range $(0, 1)$, often used in binary classification.
- **Tanh:** outputs values in the range $(-1, 1)$, providing zero-centered activations.
- **ReLU (Rectified Linear Unit):** defined as $f(z) = \max(0, z)$, widely used due to its simplicity and effectiveness.

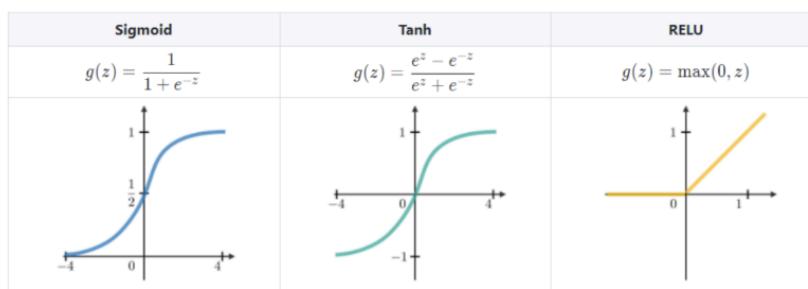


Figure 3: Common activation functions: Sigmoid, Tanh, and ReLU

2.1.4 Single-layer and Multi-layer Perceptrons

A **single-layer perceptron** contains no hidden layers and can only solve linearly separable problems. In contrast, a **multi-layer perceptron (MLP)** consists of one or more hidden layers, enabling the network to learn complex nonlinear relationships.

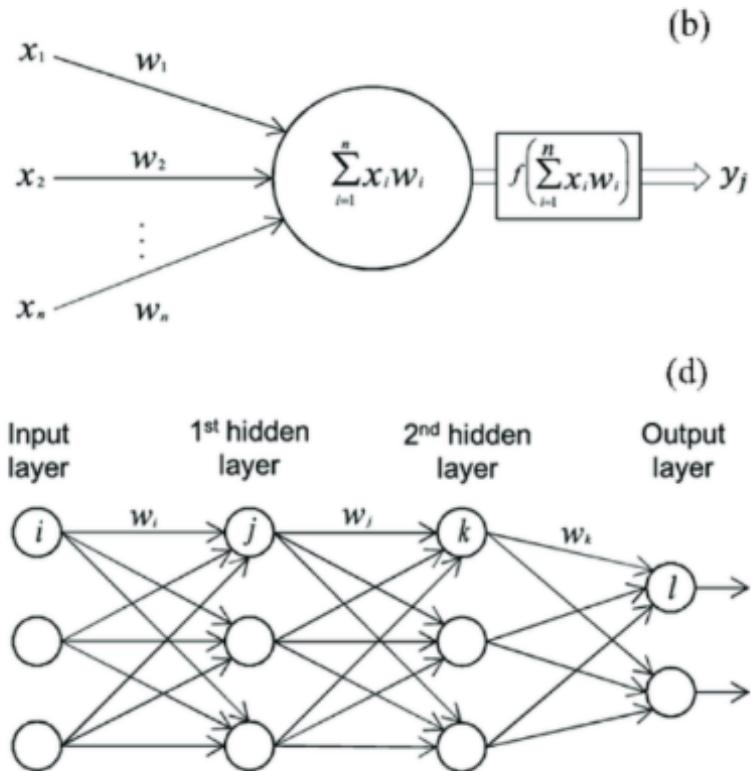


Figure 4: Comparison between single-layer and multi-layer perceptrons

2.1.5 Forward Propagation

Forward propagation is the process in which input data is passed through the network layer by layer to compute the output prediction. At each layer, weighted sums and activation functions are applied to transform the data.

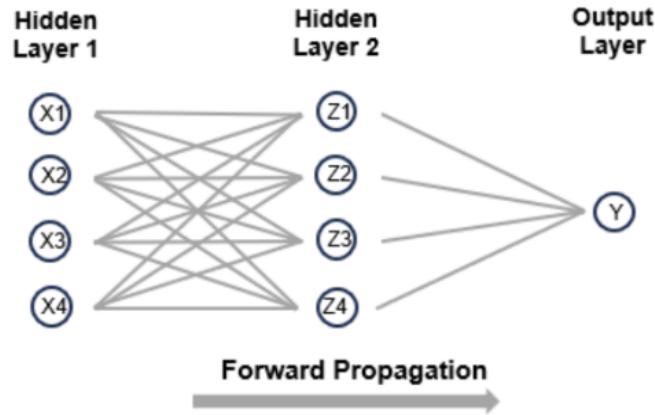


Figure 5: Forward propagation through hidden layers

2.1.6 Loss Function

The loss function measures the discrepancy between the predicted output \hat{y} and the true target value y . It acts as a guidance signal for the learning process, where a smaller loss indicates better model performance. Common loss functions include:

- Mean Squared Error (MSE) for regression tasks:

$$L = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (3)$$

- Cross-Entropy Loss for classification tasks:

$$L = - \sum_{i=1}^n y_i \log(\hat{y}_i) \quad (4)$$

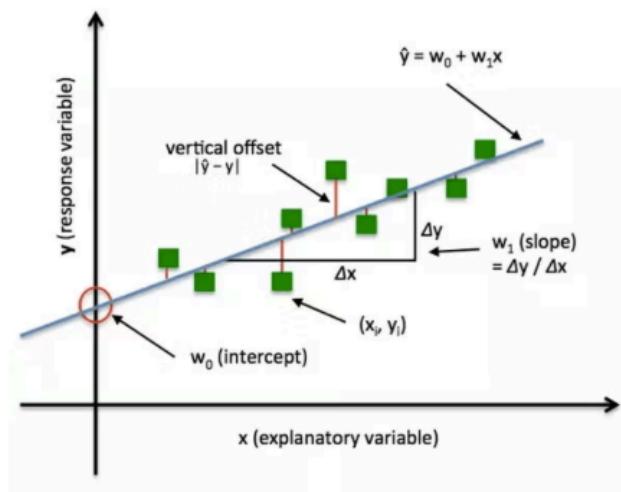


Figure 6: Illustration of prediction error and loss minimization

2.1.7 Backpropagation and Weight Update

Backpropagation is an algorithm used to compute the gradient of the loss function with respect to each weight in the network. Based on these gradients, weights are updated iteratively using optimization methods such as gradient descent.

The training process of a neural network consists of three repeated steps:

1. Forward propagation
2. Backpropagation
3. Weight update

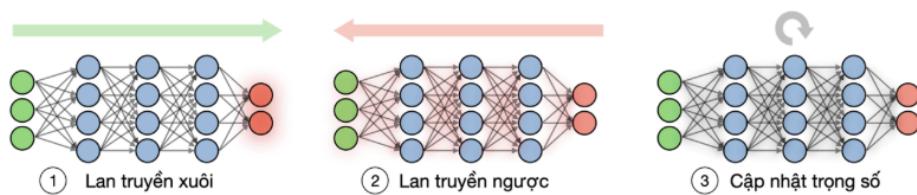


Figure 7: Training process: forward propagation, backpropagation, and weight update

2.2 Transformer architecture

2.2.1 Introduction and Background

Prior to the advent of the Transformer, sequence modeling and transduction problems, such as machine translation and language modeling, were dominated by Recurrent Neural Networks (RNNs), Long Short-Term Memory (LSTM), and Gated Recurrent Units (GRUs). However, these architectures faced several inherent limitations:

- **Sequential Computation:** They require $O(n)$ steps to process a sequence of length n , preventing parallelization during training.
- **Vanishing and Exploding Gradients:** Despite improvements in LSTMs, these models still struggle with gradient stability in very long sequences.
- **Long-range Dependencies:** Difficulty in capturing relationships between distant tokens in a sentence.
- **Efficiency Decay:** Performance significantly drops as the sequence length increases.

In 2015, the Attention mechanism was introduced as an enhancement for Seq2seq models. By 2017, the landmark paper "*Attention is All You Need*" proposed the Transformer architecture, which entirely replaced recurrence with attention. The Transformer follows an **Encoder-Decoder** structure:

- **Encoder:** Processes the input sequence and generates continuous representations.
- **Decoder:** Generates the output sequence token by token, attending to the encoder's output.
- Both components consist of multiple identical layers stacked to increase the model's capacity and representational power.

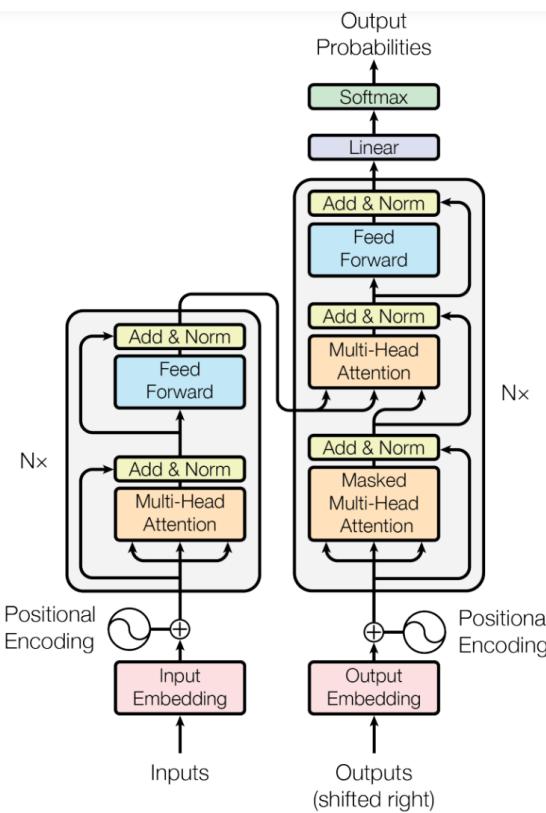


Figure 8: Transformer architecture

2.2.2 Self-Attention Mechanism

The core of the Transformer is the **Scaled Dot-Product Attention**. Unlike standard attention, Self-Attention relates different positions of a single sequence to compute a representation of the same sequence.

For an input embedding X , we compute three vectors: **Query (Q)**, **Key (K)**, and **Value (V)** using learned weight matrices W^Q, W^K, W^V :

$$Q = XW^Q, \quad K = XW^K, \quad V = XW^V \quad (5)$$

The attention scores are calculated as follows:

$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V \quad (6)$$

where d_k is the dimension of the keys.

Advantages:

- **Parallelization:** Unlike RNNs, all tokens are processed simultaneously.
- **Global Context:** It captures both local and global dependencies regardless of their distance in the sequence.
- **Training Speed:** When combined with masking, the training process becomes highly efficient.

2.2.3 Multi-Head Attention (MHA)

Instead of performing a single attention function, Multi-Head Attention allows the model to jointly attend to information from different representation subspaces at different positions.

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O \quad (7)$$

where each $\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$.

Purpose:

- Enables the model to learn multiple types of relationships (e.g., syntactic vs. semantic) within a single block.
- Further enhances parallel computation capabilities.

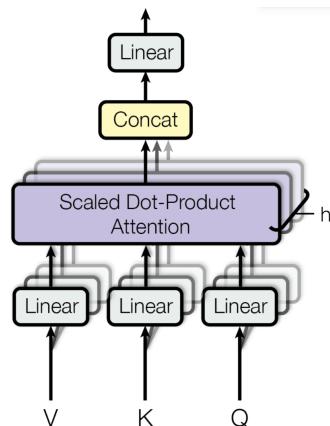


Figure 9: MultiHead Attention

2.2.4 Feed-Forward Network and Add-Norm

Each layer in our encoder and decoder contains a fully connected **Position-wise Feed-Forward Network (FFN)**:

$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2 \quad (8)$$

This consists of two linear transformations with a ReLU activation in between, projecting the data to a higher dimension and then back to stabilize representations.

To ensure deep networks can be trained effectively, each sub-layer (Attention and FFN) is wrapped in an **Add-Norm** layer:

- **Residual Connection (Add):** Helps avoid the vanishing gradient problem and retains early information.
- **Layer Normalization (Norm):** Normalizes the inputs to stabilize the learning process.

2.2.5 Embeddings and Positional Encoding

Since the Transformer contains no recurrence or convolution, it has no inherent sense of the relative or absolute position of tokens.

- **Token Embedding:** Converts text into fixed-size vectors.
- **Positional Encoding:** Unique vectors are added to the input embeddings to provide information about the order of tokens.
- **Shift Right:** In the decoder, the output is shifted to the right to ensure that the prediction for a token only depends on the previously generated tokens (preventing "looking into the future").

2.2.6 Operational Flow

Training Phase:

1. Source and target sequences are embedded and augmented with positional encodings.
2. The Encoder processes the source sequence through multiple layers. The final output serves as the *Key* and *Value* for the Decoder's *Cross-Attention* mechanism.
3. The Decoder receives the masked target input, allowing it to compute all positions in parallel via *Masked Multi-Head Attention*.
4. The final output passes through a Linear layer and a Softmax function to predict tokens, calculating loss for backpropagation.

Inference Phase: The process is similar but **autoregressive**. The target sequence starts empty (or with a start token). Each predicted token is appended to the input of the decoder for the next step until an end-of-sequence token is generated.

2.2.7 Comparative Analysis

The following table summarizes the differences between traditional recurrent models and the Transformer:

Table 2: Comparison: LSTM/GRU vs. Transformer

Feature	LSTM / GRU	Transformer
Training	Slow (Sequential computation)	Fast (Highly parallelizable)
Inference	Similar speed to Transformer	Efficient
Long-range Dependency	Poor (Memory fades over time)	Excellent (Global attention)
Short/Medium Sequences	Performs very well	Performs exceptionally
Very Long Sequences	Significant performance drop	Superior (but requires memory optimization)

2.3 Large language model

Large Language Models (LLMs) are advanced neural network models designed to understand, generate, and manipulate natural language at scale. They are typically built upon the Transformer architecture and trained on massive text corpora, enabling them to perform a wide range of Natural Language Processing (NLP) tasks such as text generation, question answering, summarization, and dialogue systems.

2.3.1 Evolution and Architecture

The development of LLMs has progressed from traditional statistical language models (e.g., n-gram models) to recurrent architectures such as Recurrent Neural Networks (RNNs) and Long Short-Term Memory (LSTM) networks, and ultimately to Transformer-based models. Transformers address the limitations of sequential computation and long-term dependency learning through the self-attention mechanism, making them highly scalable and effective for large datasets.

Based on their architectural design, LLMs can be categorized into three main types:

- **Decoder-only models (causal language models):** These models, such as GPT and LLaMA, predict the next token in a sequence and are primarily used for text generation tasks.
- **Encoder-only models (bidirectional language models):** Models like BERT, RoBERTa, and DistilBERT focus on learning contextual representations of text and are widely used for text understanding tasks such as classification and semantic similarity.
- **Encoder–Decoder models (sequence-to-sequence):** Examples include T5 and BART, which encode the input sequence and then decode it into an output sequence, making them suitable for machine translation, summarization, and question answering.

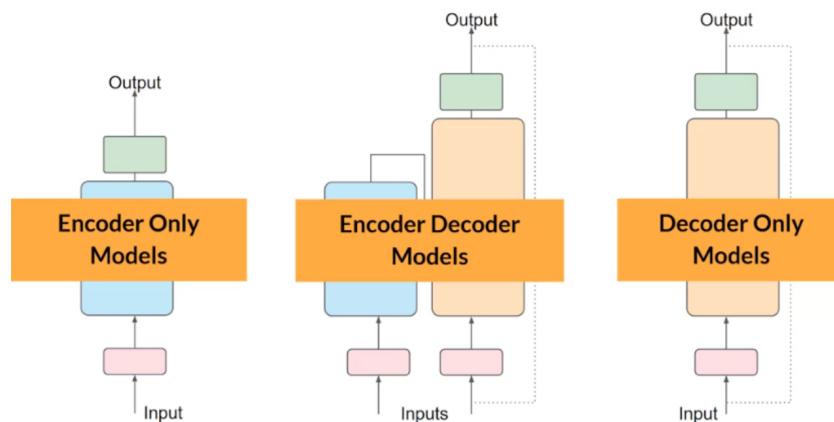


Figure 10: Transformer architectures

2.3.2 Pretraining and Fine-Tuning

LLMs are first trained during a pretraining phase using large-scale datasets such as Common Crawl, Wikipedia, and book corpora. Common pretraining objectives include autoregressive language modeling, masked language modeling, and denoising tasks. These objectives allow the model to learn grammar, semantics, and world knowledge from raw text.

After pretraining, LLMs are adapted to specific tasks through fine-tuning. This process may include Supervised Fine-Tuning (SFT), where human-labeled prompt-response pairs are used, and Reinforcement Learning from Human Feedback (RLHF), where human preferences guide the optimization of model outputs. Parameter-Efficient Fine-Tuning (PEFT) techniques such as LoRA and QLoRA are often employed to reduce computational cost while maintaining performance.

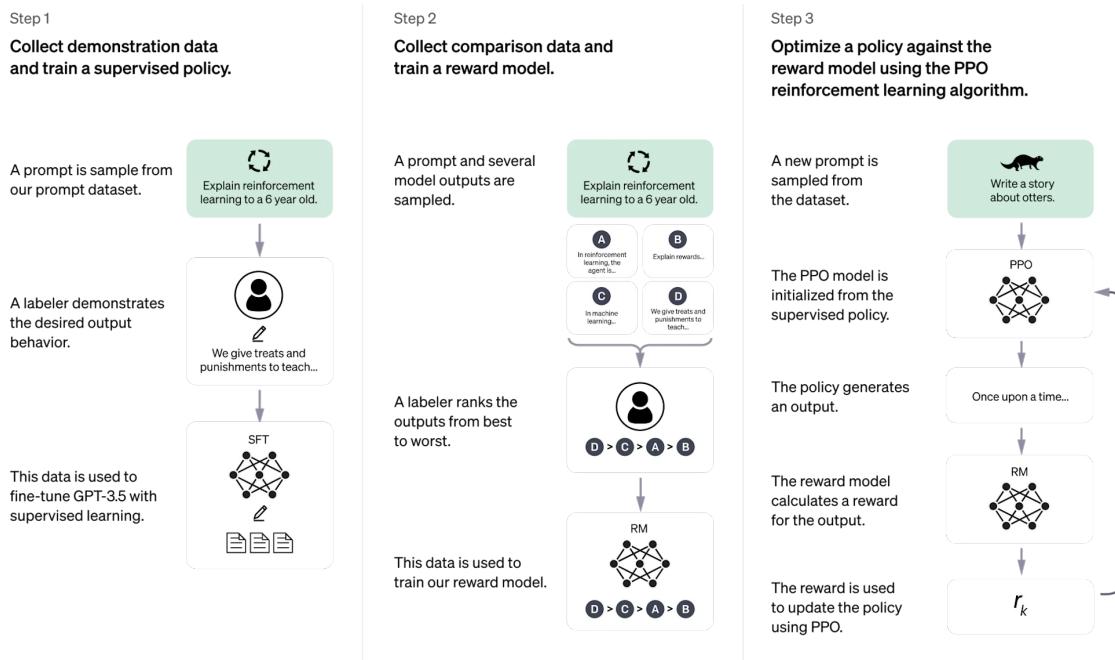


Figure 11: RLHF for ChatGPT

2.3.3 LLMs in Chatbot Systems

LLMs play a central role in modern chatbot systems. To improve factual accuracy and domain specificity, Retrieval-Augmented Generation (RAG) is commonly used. In a RAG-based system, relevant documents are embedded into a vector space and stored in a vector database. When a user submits a query, the system retrieves the most relevant documents and provides them as additional context to the LLM before generating a response. This approach significantly enhances the reliability and explainability of chatbot answers.

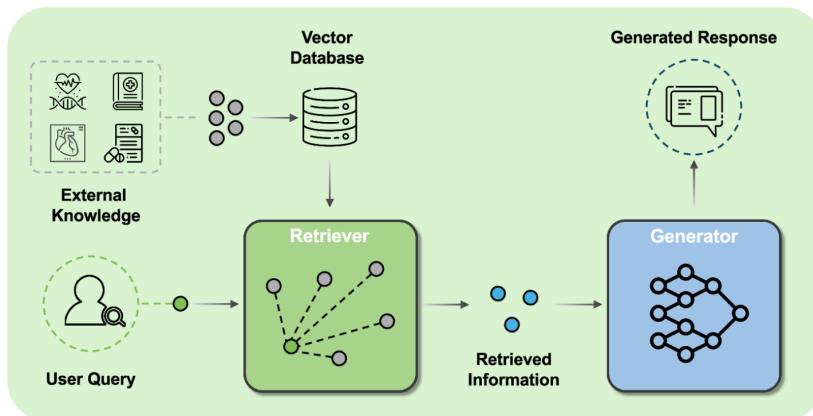


Figure 12: Retrieval augmented generation

In addition, prompt engineering and system instructions are used to control the behavior,



style, and scope of the chatbot. For domain-specific applications such as a music information website, these techniques ensure that responses remain relevant, concise, and aligned with pre-defined constraints.

3 Agent, URAG, LangChain

3.1 AI Agent

An AI agent is a system that leverages an artificial intelligence model to interact with its environment in order to achieve a user-defined objective. Unlike traditional AI systems that produce static outputs, an AI agent integrates reasoning, planning, and action execution to complete tasks in a dynamic and autonomous manner.

3.1.1 Introduction to AI Agents

An AI agent operates by receiving a user request, reasoning about the task, selecting appropriate tools, and executing actions to fulfill the objective. This process allows agents to solve multi-step problems that require interaction with external systems and continuous adaptation based on feedback.

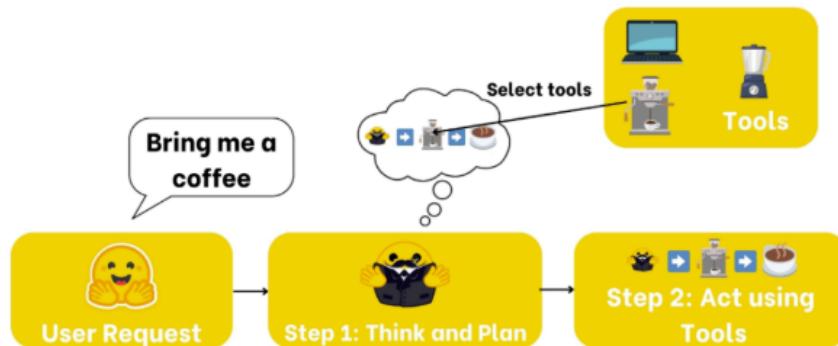


Figure 13: High-level AI agent workflow: user request, reasoning, and tool-based action

3.1.2 Core Components of an AI Agent

AI agents consist of two main components:

- **The Brain (AI Model):** This component is typically a Large Language Model (LLM) that performs reasoning, understands context, plans future steps, and decides which actions or tools should be used.
- **The Body (Capabilities and Tools):** This component enables the agent to act within its environment. It includes external tools and APIs such as search engines, code execution environments, and databases, as well as memory and observation mechanisms.

3.1.3 Levels of Agent Agency

AI agents can be categorized according to their level of autonomy, commonly referred to as their *agency level*. As the agency level increases, the agent gains greater control over program execution and decision-making.

3.1.4 Thought–Action–Observation Cycle

The behavior of an AI agent follows a continuous loop known as the *Thought–Action–Observation* cycle. In this cycle, the agent first reasons about the problem, then performs an action, and finally observes the result of that action. The observation is used to update the agent's internal state and guide subsequent reasoning steps.

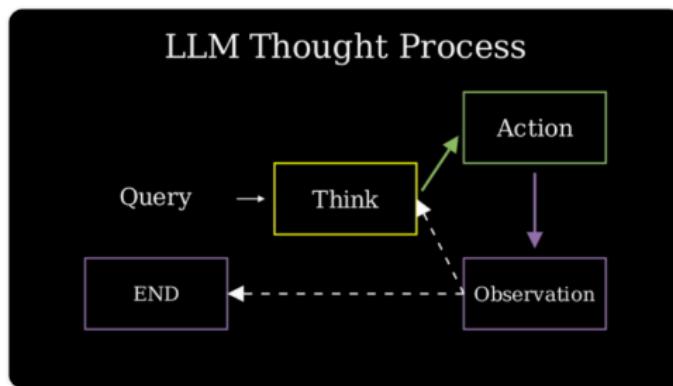


Figure 14: The Thought–Action–Observation (TAO) cycle in an AI agent

This iterative cycle enables agents to adapt their behavior dynamically based on environmental feedback.

3.1.5 Practical Illustration of the TAO Cycle

A common example of the Thought–Action–Observation cycle is a task that requires external information, such as retrieving real-time weather data. The agent reasons about the user's request, calls an appropriate external API, observes the returned data, and then produces a response.



Figure 15: Reasoning and action steps using an external weather API

3.1.6 Observation and Response Generation

After executing an action, the agent observes the outcome and integrates the feedback into its internal context. Based on this updated information, the agent performs additional reasoning to generate the final response for the user.



Figure 16: Observation feedback and final response generation

3.1.7 Internal Reasoning: Chain-of-Thought and ReAct

Reasoning is a critical capability of AI agents. Two widely used prompting approaches are:

- **Chain-of-Thought (CoT):** A technique that encourages the model to reason step-by-step before generating a final answer.
- **ReAct (Reasoning and Acting):** A technique that interleaves reasoning steps with actions, allowing the agent to think, act using tools, observe results, and continue reasoning until the task is completed.

3.1.8 Actions: Enabling Interaction with the Environment

Actions allow an AI agent to engage with its environment and accomplish tasks. Common categories of actions include information gathering, tool usage, environment interaction, and communication.

3.1.9 Observation and Adaptation

After executing an action, the agent observes the outcome and integrates the feedback into its internal context. This enables the agent to update its memory, refine its strategy, and improve performance in future interactions.

3.2 Unified Hybrid RAG (URAG)

3.2.1 Challenges in Conventional RAG

Despite the success of Retrieval-Augmented Generation (RAG) in grounding LLM responses, several critical challenges remain:

- **Dangerous Hallucinations:** LLMs may still generate plausible-looking but factually incorrect information, which is particularly risky in specialized domains.
- **Accuracy Limitations in Lightweight LLMs:** Smaller, cost-effective models often struggle to maintain high accuracy compared to their larger counterparts.
- **Noise in Retrieval Results:** Irrelevant or noisy documents retrieved from the vector database can distract the model, leading to suboptimal generation.
- **Decoupling of Retrieval and Generation:** A significant disconnect often exists between the retrieval phase and the final generation phase, preventing seamless information flow.

3.2.2 Evolution of Enhancement Strategies

Various strategies have been proposed to mitigate these issues, though they often come with trade-offs:

- **Retriever Enhancements:** Incorporating structured data, such as Knowledge Graphs, to improve context retrieval.
- **Reflective/Corrective RAG:** Implementing self-correction loops to detect and fix hallucinations by enriching inputs.
- **Speculative RAG:** Utilizing parallel processing to generate multiple draft responses simultaneously.

Drawbacks: While effective, these methods significantly increase **system complexity**, lead to **longer latency** (processing times), and require **higher computational resources**.

3.2.3 The URAG Framework: Tiered Architecture

The core concept of URAG is inspired by **Human Advisor Systems**. Just as a human advisor first consults a Frequently Asked Questions (FAQ) list before searching through thick manuals, URAG operates on a tiered retrieval strategy:

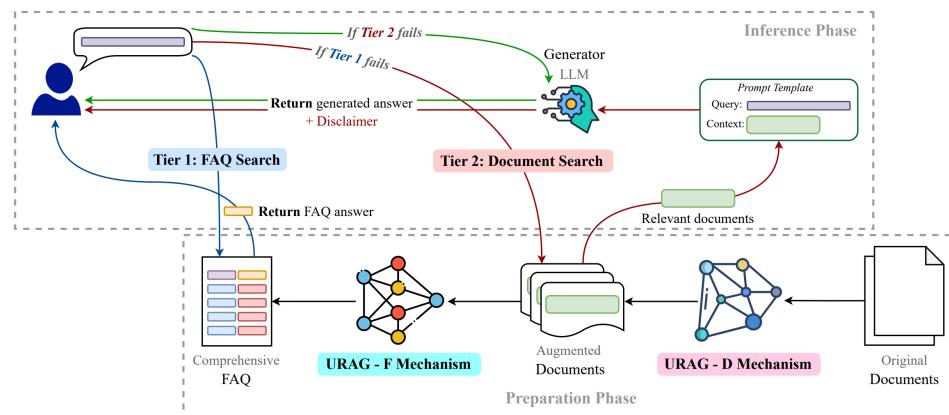


Figure 17: URAG idea

1. **Tier 1: Enriched FAQ Database:** This tier utilizes a database that integrates data from various text corpora with existing FAQs. Automated via the **URAG-F** mechanism, it ensures that common and critical queries receive direct, high-precision responses.
2. **Tier 2: Augmented Document Retrieval:** If Tier 1 fails to find a match, the system searches a document corpus enhanced by **URAG-D**. This mimics traditional RAG but uses optimized prompt templates to guide the LLM.
3. **Fallback Mechanism:** In cases where neither tier retrieves sufficient information, the system generates a response directly from the LLM's internal knowledge, accompanied by a **disclaimer** to manage user expectations regarding potential inaccuracies.

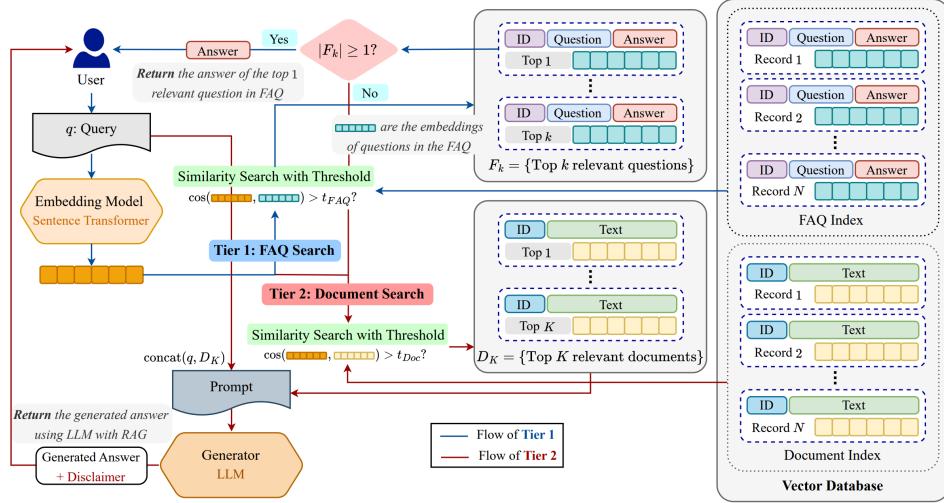


Figure 18: URAG workflow

3.2.4 URAG Preparation Workflow

The URAG architecture unifies two preparatory mechanisms to optimize performance. A key component used throughout is **Chain-of-Thought (CoT) Prompting**, which enhances the reasoning capabilities of the LLM to ensure reliable output.

URAG-D: Document Database Augmentation Unlike "Naive RAG" which uses raw document chunks, URAG-D improves retrieval through:

- **Strategic Segmentation:** Dividing documents into coherent, logical chunks.
- **Contextual Rewriting:** Each chunk is rewritten for consistency and contextual relevance. The process extracts the general context of the original document to maintain logical coherence across chunks.
- **Summarization:** A brief summary sentence is prepended to each rewritten chunk to facilitate better matching during retrieval.

URAG-F: FAQ Database Enrichment The URAG-F mechanism leverages the augmented corpus from URAG-D and the initial FAQ set to:

- **Q-A Pair Generation:** Extract and generate new, high-quality Question-Answer pairs from the processed documents.
- **Linguistic Diversity:** Paraphrase these pairs into multiple variations to ensure the system can handle diverse user phrasing and linguistic styles.

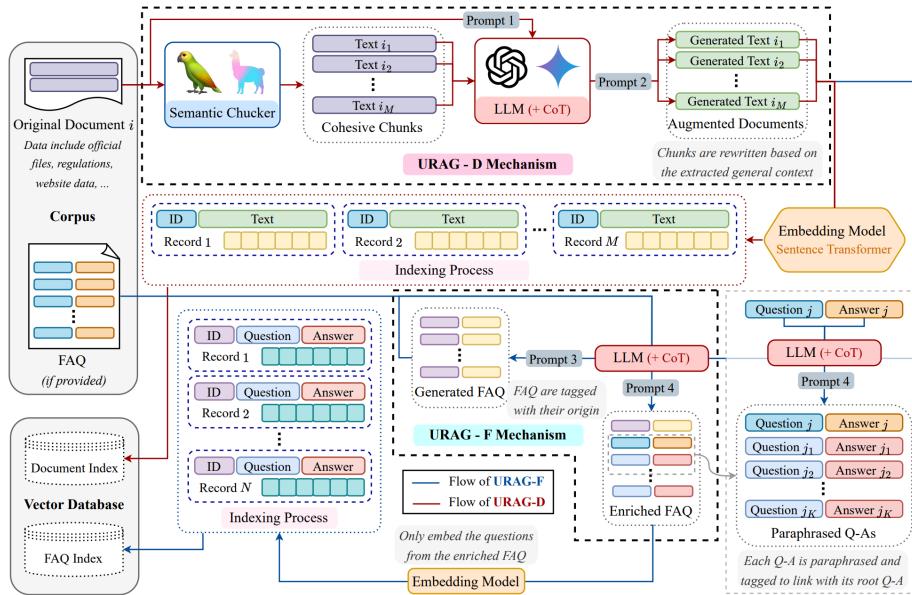


Figure 19: Preparation of URAG

3.3 Introduction to LangChain

LangChain is a comprehensive framework designed to support the development, deployment, and monitoring of applications powered by Large Language Models (LLMs). It provides modular components and abstractions that simplify the construction of complex LLM-based systems such as Retrieval-Augmented Generation (RAG) pipelines and agentic workflows.

LangChain supports the full lifecycle of an LLM application. During the development phase, developers can build applications using LangChain's core components, including prompt templates, chains, retrievers, vector stores, and integrations with third-party tools and APIs. For more advanced agentic behaviors, LangChain introduces LangGraph, which enables the definition of multi-step agents with explicit states, nodes, and control flows. This graph-based design allows agents to reason, invoke tools, and iterate over multiple steps in a structured and controllable manner.

In the production phase, LangChain is complemented by LangSmith, a platform that provides observability, debugging, and evaluation capabilities. LangSmith allows developers to inspect prompt execution, track intermediate steps, monitor latency and costs, and systematically evaluate model outputs. These features are particularly important for RAG systems, where both retrieval quality and generation accuracy must be continuously assessed.

For deployment, LangChain offers the LangGraph Platform, which facilitates the deployment and scaling of agentic workflows. This platform-oriented approach makes LangChain suitable not only for experimentation but also for real-world applications.

Within a RAG architecture, LangChain structures the workflow into two main stages. The first stage is indexing, an offline process that involves loading documents, splitting them into chunks, generating embeddings using an embedding model, and storing these embeddings in a vector database. The second stage is retrieval and generation, which occurs at runtime: relevant document chunks are retrieved from the vector store based on the user query, combined with the query into a structured prompt, and passed to an LLM to generate a final response.

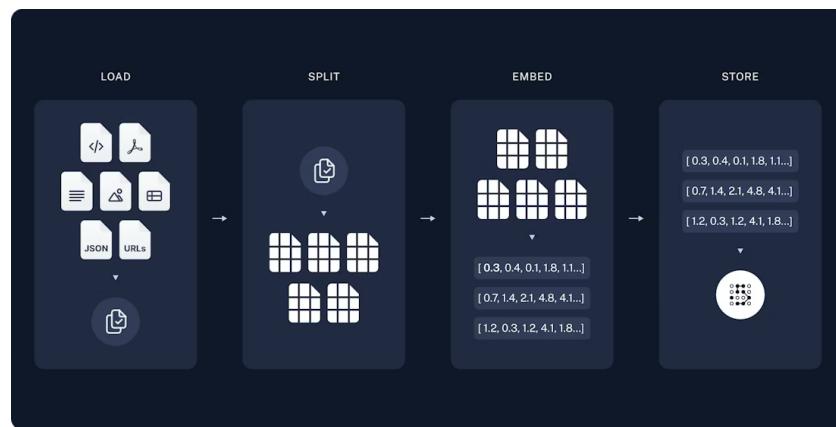


Figure 20: RAG - Indexing

LangChain provides a flexible and extensible foundation for building RAG-based and agent-driven applications, enabling developers to integrate LLMs, retrieval systems, and external tools into a unified and production-ready framework.

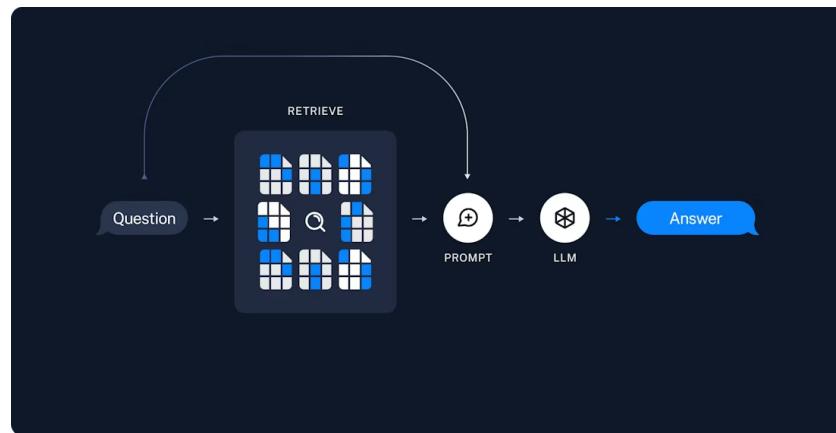


Figure 21: RAG - Retrieval and Generation

4 Retrieval-Augmented Generation for Large Language Models: A Survey

Link to the paper: <https://arxiv.org/abs/2312.10997> [1]

4.1 Overview of RAG

Retrieval-Augmented Generation (RAG) is a paradigm that enhances large language models (LLMs) by incorporating external knowledge retrieved at inference time. Instead of relying solely on parametric knowledge stored in model weights, RAG dynamically retrieves relevant documents and uses them as additional context for response generation.

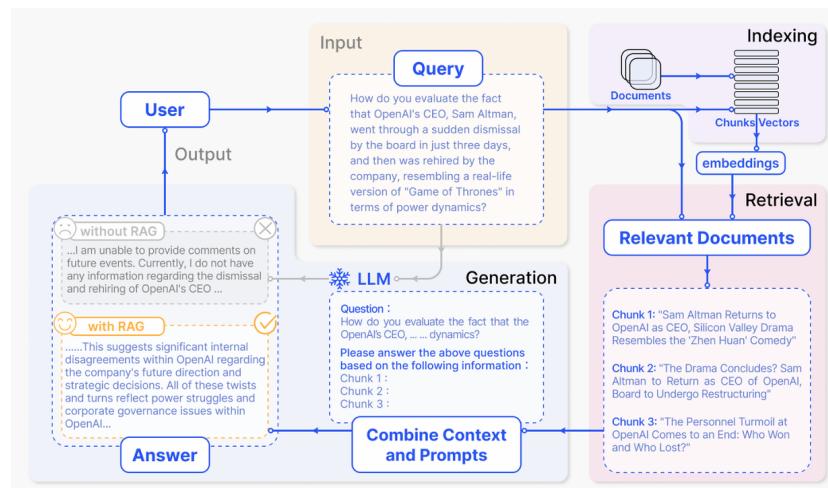


Figure 22: RAG overview

4.1.1 Naive RAG

The Naive RAG framework follows a simple *Retrieve–Read* pipeline consisting of three main stages: indexing, retrieval, and generation. During indexing, raw data sources such as PDFs, HTML pages, or Word documents are converted into plain text, segmented into smaller chunks, and encoded into vector representations stored in a vector database. In the retrieval stage, the user query is embedded and compared with stored vectors using similarity metrics to obtain the top- K most relevant chunks. Finally, in the generation stage, the LLM produces an answer based on the user query and the retrieved context, optionally incorporating conversation history.

Despite its simplicity, Naive RAG suffers from several limitations. Retrieval may lack precision and recall, resulting in irrelevant or missing information. During generation, the model may hallucinate content not supported by the retrieved context or produce outputs with bias or irrelevance. Moreover, effectively integrating retrieved information across different tasks remains challenging, often leading to redundant, incoherent, or overly extractive responses.

4.1.2 Advanced RAG

Advanced RAG aims to improve retrieval quality and context utilization through enhanced pre-retrieval and post-retrieval techniques. Pre-retrieval optimization focuses on improving indexing

structures and query formulation, including data granularity control, metadata alignment, mixed retrieval strategies, and query rewriting or expansion. Post-retrieval optimization emphasizes effective context integration, such as re-ranking retrieved documents to prioritize relevance and compressing context to reduce noise and prompt length.

4.1.3 Modular RAG

Modular RAG extends beyond fixed retrieval-generation pipelines by introducing specialized, interchangeable modules. These include search modules for heterogeneous data sources, RAG-Fusion for multi-query expansion and re-ranking, and memory modules that maintain a continuously updated retrieval memory pool. Additional components such as routing, prediction, and task adapters allow RAG systems to dynamically select retrieval pathways and adapt to downstream tasks.

This modular design enables flexible retrieval patterns, including Rewrite–Retrieve–Read, Generate–Read, and hybrid retrieval strategies that combine keyword-based, semantic, and vector-based search. As a result, Modular RAG exhibits strong adaptability and scalability across diverse applications.

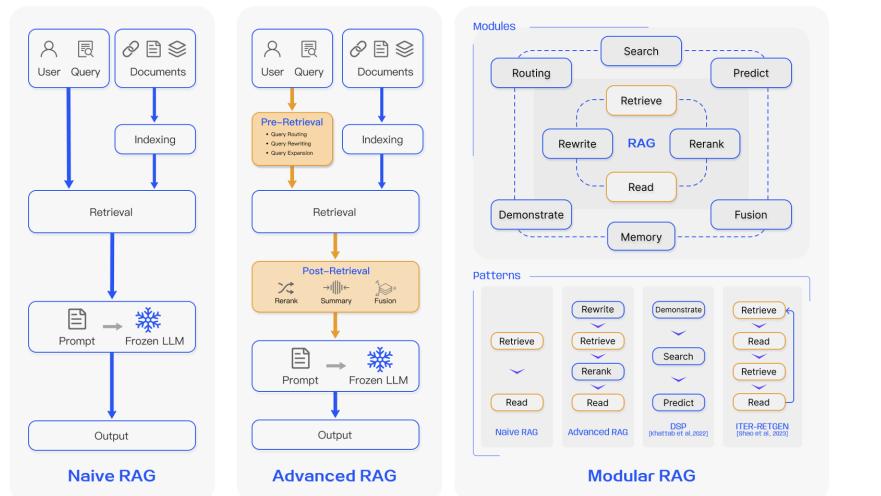


Figure 23: Types of RAG

4.2 Retrieval

Retrieval is a core component of RAG, responsible for identifying and supplying relevant external knowledge to the generation model. The effectiveness of a RAG system heavily depends on retrieval source selection, indexing strategies, query optimization, and embedding quality.

4.2.1 Retrieval Sources and Granularity

Retrieval sources can be categorized into unstructured, semi-structured, and structured data. Unstructured text data, such as Wikipedia articles or domain-specific documents, is the most common source. Semi-structured data, including PDFs with tables, presents challenges due to structural complexity, while structured sources like knowledge graphs offer precise and verified information at the cost of higher construction and maintenance effort.



Retrieval granularity ranges from tokens and sentences to chunks and full documents. Coarse-grained retrieval provides richer context but may introduce redundancy and noise, whereas fine-grained retrieval improves precision but risks losing essential semantic information.

4.2.2 Indexing Optimization

Indexing optimization techniques aim to balance context richness and efficiency. Chunking strategies play a critical role, where large chunks capture broader context but increase noise and computational cost, while small chunks reduce noise but may lack sufficient information. The *Small-to-Big* approach mitigates this trade-off by retrieving smaller units and expanding context hierarchically.

Metadata attachments, such as page numbers or timestamps, enable filtered retrieval and scoped search. Structural indexing methods, including hierarchical document structures and knowledge graph indices, further enhance retrieval speed and relevance. Techniques like Reverse HyDE leverage LLMs to generate potential questions that each chunk can answer, improving retrievability.

4.2.3 Query Optimization

Query optimization improves retrieval effectiveness by refining or expanding user queries. Query expansion and multi-query techniques enrich the query with additional context, while sub-query decomposition breaks complex questions into simpler ones. Query transformation methods include rewriting queries, generating hypothetical answers (HyDE), and step-back prompting to retrieve higher-level contextual information.

Query routing mechanisms further enhance retrieval by directing queries to appropriate data sources or pipelines using metadata-based or semantic routing strategies.

4.2.4 Embeddings and Adapters

Modern RAG systems often employ hybrid retrieval that combines sparse retrievers, such as BM25 for keyword matching, with dense retrievers based on neural embeddings for semantic understanding. Embedding models can be fine-tuned for domain-specific tasks, with LM-supervised retrievers aligning retrieval objectives with generation outcomes using LLM feedback.

When fine-tuning is impractical, adapter-based methods provide lightweight alternatives. These include prompt retrievers, bridging modules that transform retrieved content into LLM-friendly formats, and plug-in knowledge generators that replace or augment traditional retrievers in white-box settings.

4.3 Generation

Generation is the final stage in a Retrieval-Augmented Generation (RAG) pipeline, where the Large Language Model (LLM) produces the final response based on the user query and the augmented context retrieved from external knowledge sources. This stage directly determines the quality, coherence, and usefulness of the system's output.

4.3.1 Overview of the Generation Stage

In the basic RAG setting, also referred to as *Naive RAG*, the generation stage receives the user query together with the retrieved documents as input. These elements are combined into a single prompt, which is then passed to a frozen LLM to generate the final answer.

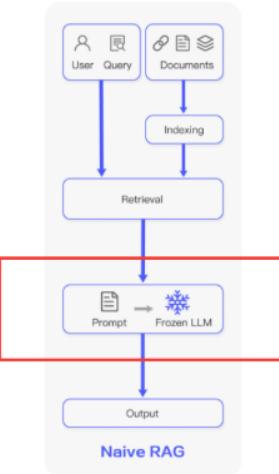


Figure 24: Generation stage in Naive RAG, where retrieved context is directly passed to a frozen LLM via a prompt

The generated response is expected to be fluent, context-grounded, and aligned with the retrieved evidence. However, this approach may suffer from noisy or redundant context, which can degrade answer quality.

4.3.2 Context Curation for Improved Generation

To address the limitations of Naive RAG, advanced RAG systems introduce a *context curation* step before generation. The goal of context curation is to improve relevance and reduce noise in the retrieved context supplied to the LLM.

Key components of context curation include:

- **Reranking:** reorders retrieved documents to prioritize the most relevant ones.
- **Context Selection or Compression:** filters, summarizes, or shortens retrieved content to avoid overly long prompts.

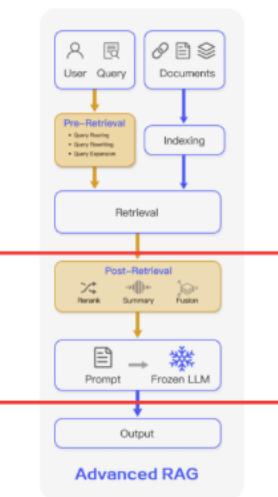


Figure 25: Advanced RAG with post-retrieval context curation before generation

By providing a more concise and focused augmented prompt, context curation enables the LLM to generate higher-quality and more accurate responses.

4.3.3 LLM Fine-tuning for Generation

Beyond improving the quality of the input context, generation performance can be further enhanced by fine-tuning the LLM itself. The objective of LLM fine-tuning is to adapt the model to domain-specific knowledge and desired response styles.

The key benefits of LLM fine-tuning include:

- **Domain Adaptation:** improves understanding and reasoning in specialized domains.
- **Output Alignment:** aligns tone, structure, and formatting with predefined guidelines.
- **Reduced Hallucination:** encourages stronger reliance on retrieved evidence.

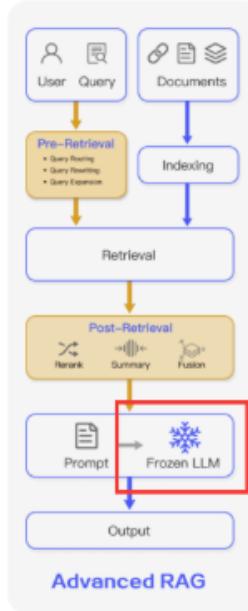


Figure 26: Generation stage using a fine-tuned LLM in an Advanced RAG pipeline

As a result, fine-tuned generation produces responses that are more accurate, coherent, and context-grounded, making it suitable for domain-specific and high-stakes applications.

4.4 Augmentation

Augmentation is a key process in Retrieval-Augmented Generation (RAG) systems that enhances answer quality by iteratively refining retrieval and generation. Instead of performing retrieval only once, the system can retrieve additional information, reformulate queries, or adjust context dynamically during the answering process.

4.4.1 Overview of the Augmentation Process

In the augmentation process, retrieval and generation are tightly coupled in an iterative loop. After each generation step, the system evaluates whether the current information is sufficient or if additional retrieval is required before producing the final response.

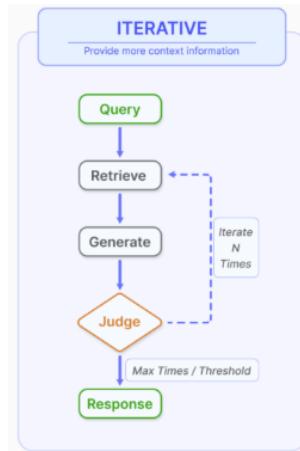


Figure 27: Overview of the iterative augmentation process in RAG

The main purposes of augmentation are:

- Handling complex or multi-step reasoning questions.
- Producing more accurate and context-grounded answers.

4.4.2 Iterative Retrieval

Iterative retrieval repeatedly alternates between retrieving new context and generating partial answers. Each generation step provides new signals that guide the next retrieval, gradually enriching the available context.

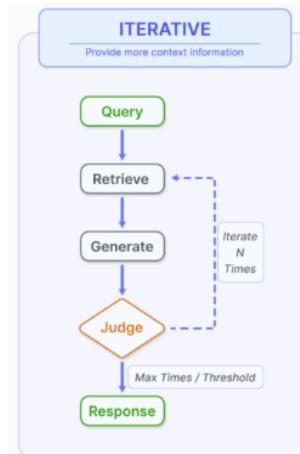


Figure 28: Iterative retrieval with repeated retrieval-generation cycles

The key ideas behind iterative retrieval include:

- Retrieval is performed again based on what has been generated so far.

- Context is progressively enriched to support step-by-step reasoning.

Pros and Cons:

- **Advantages:** more comprehensive and targeted information.
- **Disadvantages:** risk of semantic drift or accumulation of irrelevant context.

4.4.3 Recursive Retrieval

Recursive retrieval refines the query step-by-step using feedback from previous retrieval results. Each iteration clarifies what information is still missing, leading to progressively more relevant context.

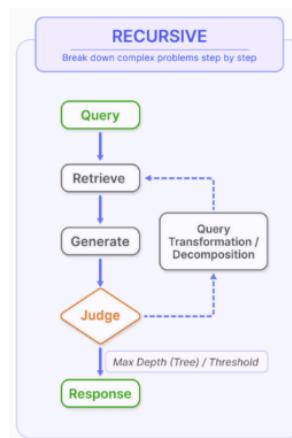


Figure 29: Recursive retrieval with query transformation and decomposition

Key characteristics of recursive retrieval include:

- Construction of intermediate reasoning structures such as chain-of-thought or clarification trees.
- Continuous reformulation of the query to target more specific information.

Pros and Cons:

- **Advantages:** improved accuracy and relevance over multiple iterations.
- **Disadvantages:** increased latency due to multiple refinement steps.

4.4.4 Adaptive Retrieval

Adaptive retrieval allows the LLM to dynamically decide when and what to retrieve during generation. Instead of enforcing retrieval at fixed stages, the model monitors its own confidence and triggers retrieval only when additional context is required.

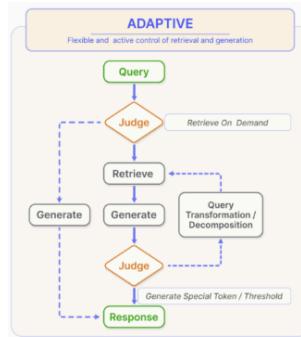


Figure 30: Adaptive retrieval where the LLM actively controls retrieval decisions

The key ideas of adaptive retrieval are:

- Retrieval is performed on demand rather than at predefined stages.
- The LLM behaves like an agent, evaluating intermediate outputs and deciding whether to retrieve more context or proceed to the final answer.

4.4.5 Summary

In summary, augmentation enhances the RAG pipeline by enabling iterative and adaptive interaction between retrieval and generation. Iterative, recursive, and adaptive retrieval strategies provide increasing levels of flexibility and reasoning capability, allowing RAG systems to better handle complex queries and produce more accurate, context-aware responses.

4.5 Tasks and Evaluation Framework

4.5.1 Downstream Tasks and Benchmarks

The primary objective of RAG systems is to support various Natural Language Processing (NLP) tasks by providing external context. While Question Answering (QA) remains the most prominent application, the scope of RAG is rapidly expanding.

- **Question Answering (QA):** This is the core application, encompassing several complexities:
 - *Single-hop QA*: Requires retrieving information from a single document (e.g., **Natural Questions (NQ)**, **TriviaQA**, **SQuAD**).
 - *Multi-hop QA*: Requires synthesizing information across multiple documents to answer complex queries (e.g., **HotpotQA**, **2WikiMultiHopQA**, **MuSiQue**).
 - *Domain-specific QA*: Tailored for specialized fields like medicine or research (e.g., **Qasper**, **COVID-QA**, **MMCU Medical**).
- **Expanded Tasks:** Beyond QA, RAG is utilized in **Information Extraction (IE)**, **Dialogue Generation**, and **Code Search**, where external documentation is vital for accuracy.



4.5.2 Evolution of Evaluation Targets

Traditional evaluation relied on lexical overlap metrics, but modern RAG research shifts toward component-level analysis.

1. **Traditional Metrics:** Metrics such as **EM (Exact Match)**, **F1 score**, and **ROUGE** were standard. However, they fail to evaluate the internal reasoning of the model or the quality of the retrieved context.
2. **Modern RAG Metrics:** Evaluation is now bifurcated into two primary targets:
 - **Retrieval Quality:** Measured by **Hit Rate**, **Mean Reciprocal Rank (MRR)**, and **Normalized Discounted Cumulative Gain (NDCG)**.
 - **Generation Quality:** Focused on **Faithfulness** (factuality), **Relevance** (answering the intent), and **Non-harmfulness**.

4.5.3 Evaluation Aspects: The "RAG Triad" and Key Abilities

Current research emphasizes three quality scores (often referred to as the RAG Triad) and four essential robustness abilities.

The Three Quality Scores:

- **Context Relevance:** Evaluates whether the retrieved snippets are precise and specific to the query, minimizing "noise" or extraneous content.
- **Answer Faithfulness:** Ensures the generated output is derived strictly from the retrieved context without "hallucinating" external or contradictory information.
- **Answer Relevance:** Measures how well the response directly addresses the user's question.

Four Essential Abilities:

- **Information Integration:** The capacity to synthesize insights from multiple disparate documents.
- **Noise Robustness:** The ability to filter out irrelevant or misleading information within the retrieved set.
- **Counterfactual Robustness:** The ability to recognize and disregard known inaccuracies or "fake news" injected into documents.
- **Negative Rejection:** The system's ability to decline answering when the provided documents do not contain the necessary knowledge.

4.6 Discussion

- **RAG vs. Long-Context LLMs:** With modern LLMs supporting contexts over 200,000 tokens, the necessity of RAG is often debated. However, RAG remains superior due to:
 - **Efficiency:** Processing massive contexts per request is computationally expensive and slow. RAG's chunk-based retrieval is significantly faster and more cost-effective.
 - **Observability:** RAG provides clear citations and references, whereas Long-Context generation operates as a "black box," making verification difficult.



- **The Noise Paradox:** Recent studies have identified a "Robustness Paradox": in certain scenarios, including a small amount of irrelevant (noisy) documents can unexpectedly improve model accuracy by over 30%, suggesting that some noise may act as a regularizer during the reasoning process.
- **Hybrid Approaches:** The most effective current strategy is a hybrid model: using **Fine-tuning** to adapt the model's behavior and style, while using **RAG** to provide updated, dynamic knowledge.
- **Scaling Laws:** While LLMs follow predictable scaling laws, it is unclear if these apply to RAG. Some research suggests an "**Inverse Scaling Law**", where smaller, specialized models may outperform larger general-purpose models when equipped with high-quality RAG pipelines.

5 Reranking, RAG-Reasoning, RAG-RL

5.1 Reranking

5.1.1 Limitations of Vector-Based Retrieval

In large-scale Retrieval-Augmented Generation (RAG) systems, retrieval is typically performed using dense vector embeddings and approximate nearest neighbor (ANN) search. While this approach is highly efficient and scalable, it introduces fundamental limitations.

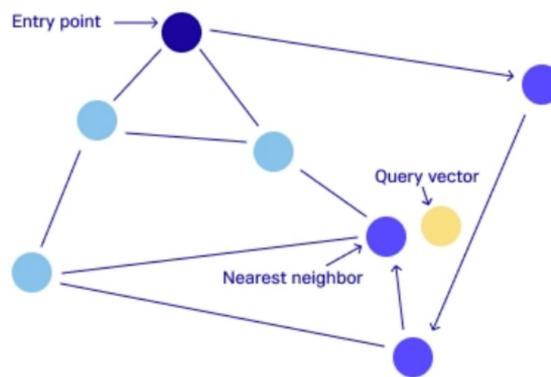


Figure 31: Approximate nearest neighbor search may return vectors that are close in embedding space but not truly semantically relevant

Dense retrieval commonly relies on a bi-encoder architecture, where queries and documents are encoded independently into vector representations. ANN search then finds the nearest vectors in the embedding space. However, vector similarity does not necessarily correspond to true semantic relevance. As a result, the retrieved top- k documents may lack sufficient reasoning, logical structure, or supporting evidence required by the LLM. Noisy or weakly relevant context can increase confusion and hallucination during generation.

5.1.2 Reranking as a Post-Retrieval Refinement Step

To mitigate the limitations of vector-based retrieval, reranking is introduced as a post-retrieval refinement step in the RAG pipeline. Instead of directly passing the retrieved documents to the generator, reranking evaluates and reorders them based on more accurate relevance estimation.

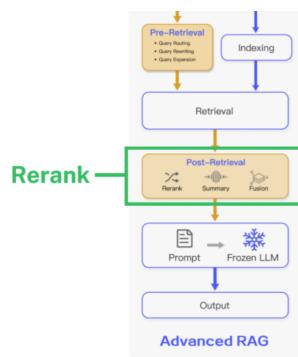


Figure 32: Reranking as an intermediate step between retrieval and generation in Advanced RAG

In this stage, each retrieved document is assessed together with the query, and the documents are reordered to prioritize those that are most relevant. By selecting a smaller set of high-quality passages, reranking reduces noise in the prompt and improves factual accuracy, coherence, and reliability of the final LLM-generated response.

5.1.3 Bi-Encoder versus Cross-Encoder for Reranking

A key architectural distinction in RAG systems lies between bi-encoders used for retrieval and cross-encoders used for reranking.

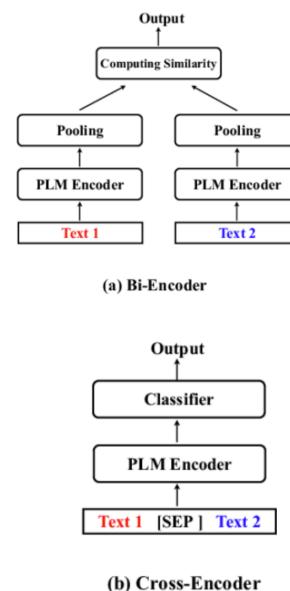


Figure 33: Comparison between bi-encoder retrieval and cross-encoder reranking architectures

Bi-encoders encode the query and documents separately and compute relevance using vector similarity. This design enables fast retrieval over large document collections but captures only coarse semantic similarity. In contrast, cross-encoders jointly encode the query and document



within a single model, allowing full cross-attention between them. This results in significantly more accurate relevance estimation.

Due to their higher computational cost, cross-encoders are typically applied only to the top- k candidates returned by the bi-encoder. This two-stage design effectively balances scalability and accuracy in modern RAG systems.

5.1.4 Summary

In summary, reranking plays a critical role in bridging the gap between efficient retrieval and accurate generation. By refining retrieved results using more precise relevance modeling, reranking improves context quality, reduces hallucination, and significantly enhances the overall performance of Retrieval-Augmented Generation systems.

5.2 Towards Agentic RAG with Deep Reasoning: A Survey of RAG Reasoning Systems in LLMs

Link to paper: <https://arxiv.org/abs/2507.09477> [5]

Retrieval-Augmented Generation (RAG) has emerged as a powerful paradigm for enhancing large language models (LLMs) by grounding generation in external knowledge sources. However, traditional RAG pipelines primarily focus on surface-level semantic retrieval and often struggle with multi-hop reasoning, noisy contexts, and complex decision-making tasks. The survey by Li et al. proposes a comprehensive framework that systematically analyzes how reasoning capabilities can be deeply integrated into RAG systems, moving toward more agentic and intelligent architectures.

The survey categorizes existing approaches into three major paradigms. The first is Reasoning-Enhanced RAG, where reasoning is explicitly incorporated to improve retrieval, integration, and generation. At the retrieval stage, techniques such as reasoning-aware query reformulation, retrieval planning, and retriever model enhancement aim to obtain evidence that is more relevant to downstream reasoning tasks. During integration, retrieved documents are assessed, filtered, and fused using reasoning-driven relevance assessment and information synthesis mechanisms. At the generation stage, context-aware and grounded generation methods ensure that the model's outputs remain faithful to retrieved evidence and follow coherent reasoning paths.

The second paradigm is RAG-Enhanced Reasoning, which treats retrieval as a tool to directly support the reasoning process of LLMs. In this setting, models retrieve external knowledge such as structured knowledge bases, web content, tools, or prior experiences to assist in complex reasoning tasks including mathematical problem solving, fact verification, and planning. In-context retrieval of examples and historical interactions further enables LLMs to adapt reasoning strategies dynamically based on retrieved demonstrations or memories.

The third paradigm, Synergized RAG-Reasoning, represents the most advanced integration, where retrieval and reasoning are tightly interwoven within an agentic workflow. These systems interleave reasoning steps with retrieval actions using chain-based, tree-based, or graph-based reasoning workflows. Moreover, agent orchestration techniques, including single-agent and multi-agent systems, allow LLMs to autonomously plan, retrieve, reason, and verify information. Such agentic RAG systems exhibit improved robustness, interpretability, and adaptability in complex tasks.

Category	Method summary	Related papers
----------	----------------	----------------



Reasoning-Aware Query Reformulation (§3.1.1)	Reformulates the original query to better retrieve reasoning-relevant context. This includes query decomposition (breaking complex queries into simpler ones) , reformulation (recasting ambiguous queries) , and expansion (enriching the query via CoT).	e.g., Collab-RAG (Xu et al., 2025b), DynQR (Anonymous, 2025), DeepRetrieval (Jiang et al., 2025)
Retrieval Strategy and Planning (§3.1.2)	Covers global retrieval guidance. This involves advance planning to generate a retrieval blueprint before execution or adaptive retrieval methods that predict whether and how to retrieve based on query complexity.	e.g., PAR-RAG (Zhang et al., 2025d), LPKG (Wang et al., 2024b), FIND (Jia et al., 2025)
Retrieval Model Enhancement (§3.1.3)	Enhances retrievers with reasoning. This is done by leveraging structured knowledge (like KGs with GNNs or symbolic rules) or integrating explicit reasoning (like CoT) with the query.	e.g., GNN-RAG (Mavromatis & Karypis, 2024), RuleRAG (Chen et al., 2024c)
Relevance Assessment & Filtering (§3.2.1)	Uses deeper reasoning to assess the relevance of retrieved fragments. This can involve using "assessor experts" to select faithful evidence or models to filter non-entailing passages.	e.g., SEER (Zhao et al., 2024c), M-RAG-R (Yoran et al., 2024)
Information Synthesis & Fusion (§3.2.2)	Fuses relevant snippets into a coherent evidence set after they are identified. Methods include aggregating sub-question answers or building a reasoning graph to evaluate and aggregate knowledge.	e.g., BeamAggR (Chu et al., 2024), DualRAG (Cheng et al., 2025), CRP-RAG (Xu et al., 2024)
Context-Aware Generation (§3.3.1)	Ensures outputs remain relevant and reduces noise. This includes selective-context utilization (pruning or re-weighting content) and reasoning path generation (building explicit logical chains).	e.g., Open-RAG (Islam et al., 2024), RARE (Wang et al., 2025d), Self-Reasoning (Xia et al., 2025b)
Grounded Generation Control (§3.3.2)	Introduces verification mechanisms to anchor outputs to retrieved evidence. This is done via fact verification , citation generation , and faithful reasoning (ensuring steps adhere to evidence).	e.g., RARR (Gao et al., 2023a), TRACE (Fang et al., 2024), AlignRAG (Wei et al., 2025b)
Knowledge Base (§4.1.1)	Retrieves from KBs storing arithmetic, commonsense, or logical knowledge. This can include formal lemmas for math , legal precedents , or code snippets.	e.g., Premise-Retrieval (Tao et al., 2025), ReaRAG (Lee et al., 2025), CBR-RAG (Wiratunga et al., 2024)



Web Retrieval (\$4.1.2)	Accesses dynamic online content like web pages, news, or social media. It is used for fact-checking by verifying claims step-by-step or for QA by iteratively refining reasoning.	e.g., ALR ² (Li et al., 2024d), RARE (Tran et al., 2024), Open-RAG (Islam et al., 2024)
Tool Using (\$4.1.3)	Leverages external resources like calculators, libraries, or APIs to enhance reasoning interactively. This improves numerical accuracy and computational robustness.	e.g., TATU (Li et al., 2024g), TRICE (Qiao et al., 2024), Re-Invoke (Chen et al., 2024a)
Prior Experience (\$4.2.1)	Retrieves past interactions or successful strategies stored in a model's internal memory. This includes leveraging past decisions for planning or recalling conversational histories for adaptive reasoning.	e.g., RAP (Kagaya et al., 2024), JARVIS-1 (Wang et al., 2024f), EM-LLM (Fountas et al., 2024)
Example or Training Data (\$4.2.2)	Retrieves external examples from demonstrations or training data. This provides relevant exemplars to guide the model in emulating specific reasoning patterns.	e.g., MoD (Wang et al., 2024c), RE4 (Li et al., 2024c), UPRISE (Cheng et al., 2023)
Chain-based (\$5.1.1)	Interleaves retrieval operations between the linear "step-by-step" reasoning of a Chain-of-Thought (CoT) to avoid error propagation. Methods can also add verification or filtering steps.	e.g., IRCOT (Trivedi et al., 2023), Rat (Wang et al., 2024g), CoV-RAG (He et al., 2024a), RAFT (Zhang et al., 2024a)
Tree-based (\$5.1.2)	Explores multiple reasoning pathways. Tree-of-Thought (ToT) methods build a deterministic reasoning tree. Monte Carlo Tree Search (MCTS) methods use probabilistic tree search to dynamically prioritize exploration.	ToT: e.g., RATT (Zhang et al., 2025a), Tree of Clarifications (Kim et al., 2023) MCTS: e.g., AirRAG (Feng et al., 2025), MCTS-RAG (Hu et al., 2025)
Graph-based (\$5.1.3)	Walk-on-Graph uses graph learning techniques (like GNNs) to retrieve and reason over graph-structured data. Think-on-Graph integrates graph structures into the LLM's reasoning loop, letting the LLM decide which node to explore next.	Walk-on-Graph: e.g., QA-GNN (Yasunaga et al., 2021) Think-on-Graph: e.g., ToG (Sun et al., 2024b), Graph-CoT (Jin et al., 2024)
Single-Agent (\$5.2.1)	A single agent interweaves retrieval into its reasoning loop. This is achieved via Prompting (e.g., ReAct), Supervised Fine-Tuning (SFT), or Reinforcement Learning (RL).	Prompting: e.g., ReAct (Yao et al., 2023b) SFT: e.g., Toolformer (Schick et al., 2023) RL: e.g., Search-R1 (Jin et al., 2025)



Multi-Agent (\$5.2.2)	Uses multiple agents for collaboration. Decentralized systems use specialized agents that work together. Centralized systems use a hierarchical (e.g., manager-worker) pattern for task decomposition.	Decentralized: e.g., M-RAG (Wang et al., 2024) Centralized: e.g., HM-RAG (Liu et al., 2025), Chain of Agents (Zhang et al., 2024c)
--------------------------	--	---

The survey highlights a clear evolution from static retrieval pipelines toward dynamic, agent-based RAG systems with deep reasoning capabilities. It identifies key challenges such as efficiency, evaluation, and controllability, while outlining future research directions that aim to unify reasoning, retrieval, and agent learning into a coherent framework for next-generation LLM systems.

5.3 RAG-RL: Advancing Retrieval Augmented Generation via RL and Curriculum Learning

Link to paper: <https://arxiv.org/abs/2503.12759> [3]

5.3.1 The Problem of Imperfect Retrieval

Modern RAG systems are highly dependent on the precision of the retriever. However, in practical scenarios, retrievers often return a mix of relevant documents and irrelevant "distractors" (noise). This leads to two major issues:

- **Generator Hallucination:** The LLM (Generator) may struggle to distinguish between gold context and noise, leading to incorrect or fabricated outputs.
- **Bottleneck of Retrieval Optimization:** Achieving 100% retriever precision is computationally prohibitive and technically challenging.

Proposed Solution: Instead of solely focusing on fixing the retriever, the **RAG-RL** approach aims to train a **Robust Generator**. This empowers the LLM to autonomously identify relevant contexts and ignore noise by combining **Reinforcement Learning (RL)** for reasoning and **Curriculum Learning (CL)** for noise robustness.

5.3.2 Foundational Concepts: RL and CL

Reinforcement Learning (RL) RL is a paradigm where an **Agent** interacts with an **Environment** to maximize a cumulative **Reward**.

- **Agent:** The decision-maker (in this case, the LLM).
- **State (s):** The current situation (Query + Retrieved Documents).
- **Action (a):** The generated response and citations.
- **Reward (r):** Feedback indicating the quality of the action.

The RL feedback loop enables the agent to update its policy π by balancing *exploration* (trying new generation patterns) and *exploitation* (using known high-reward patterns).

Curriculum Learning (CL) Inspired by human education, CL involves training a model by starting with simple concepts and gradually increasing difficulty.

- **Synthetic Difficulty:** Samples are scored based on complexity (e.g., number of distractors).
- **Scheduler:** Manages the transition from *Easy Samples* (low noise) to *Hard Samples* (high noise) to prevent the model from being overwhelmed early in training.

5.3.3 The Synergy of RL and CL

Combining these two paradigms addresses the following:

- **Solving the "Cold Start" Problem:** In noisy environments, an untrained model fails to generate correct citations, leading to zero rewards and stagnant learning. CL provides "easy" (clean) samples first to establish a positive reward baseline.
- **Two-Phase Robustness:** The model learns **how** to cite (formatting/syntax) in Phase 1, and **what** to cite (relevance/filtering) in Phase 2.

5.3.4 Methodology

Agent and Environment In the RAG-RL framework, the **Generator (LLM)** acts as the agent. The **Environment** consists of the retrieved documents (Gold context + Distractors). The **Action** is the final generated answer bundled with specific citations.

Optimization via GRPO The framework utilizes **Group Relative Policy Optimization (GRPO)**. For each query, the LLM generates a group of G candidate answers. The model prioritizes answers with higher relative rewards within the group, effectively refining its policy without requiring a separate value function.

Reward Design To quantify the quality of the generation, the total reward R_{total} is defined as a multi-objective function:

$$R_{total} = R_{answer} + R_{citation} + R_{format} \quad (9)$$

Where the individual components are defined as:

- **Answer Reward:** Uses an indicator function $\mathbb{1}$ to check if the generated answer o_{ans} matches the ground truth G_{ans} .

$$R_{answer} = \gamma_{ans} \cdot \mathbb{1}(o_{ans} = G_{ans}) \quad (10)$$

- **Citation Reward:** Measures the accuracy of the sources cited.

$$R_{citation} = \gamma_{cor} \cdot \text{Recall}(o_{cit}, G_{cit}) - \gamma_{inc} \cdot c_{inc} \quad (11)$$

where c_{inc} is the count of incorrect citations.

- **Formatting Reward:**

$$R_{format} = \begin{cases} \gamma_{fmt} & \text{if format is correct} \\ -p & \text{otherwise} \end{cases} \quad (12)$$

Constants are set as: $\gamma_{ans} = 5$, $\gamma_{cor} = 5$, $\gamma_{inc} = 2$.

5.3.5 Training Schedules and Results

The study evaluates three curriculum strategies to manage the number of distractors (K):

1. **Max Curriculum (Baseline):** Static training on the hardest samples (K distractors) throughout.
2. **Linear Curriculum:** Gradually increases the number of distractors from 1 to K linearly over training steps.
3. **Min-Max Curriculum:** A step-function approach where the first half of training uses only 1 distractor, and the second half jumps to K distractors.

Experimental Results: Using **Qwen2.5-7B-Instruct** as the base model and testing across **HotpotQA**, **MuSiQue**, and **2WikiMultiHopQA** datasets, the results indicate that **Min-Max** and **Linear Curriculum** schedules consistently and significantly outperform the Max Curriculum baseline, demonstrating the importance of gradual noise introduction in training robust RAG agents.

Eval. Setting	Curriculum	HotpotQA			MuSiQue			2Wiki		
		Answer F1	Citation F1	Joint F1	Answer F1	Citation F1	Joint F1	Answer F1	Citation F1	Joint F1
Distractor	Baseline	60.65	36.47	45.55	25.88	25.35	25.61	48.71	40.18	44.03
	Max	66.04	73.93	69.76	40.91	53.07	46.20	67.99	77.08	72.25
	Linear	68.71	78.54	73.24	44.68	50.76	51.14	68.23	83.23	75.40
	Min-Max	68.87	81.44	74.72	47.18	64.48	54.49	70.77	76.37	73.46
	Base-Answerable	66.19	72.81	69.34	41.40	54.28	46.97	68.56	76.36	72.25
	Base-Unanswerable	65.25	71.13	68.06	38.84	52.15	44.52	66.59	72.04	69.21
Ideal Retrieval	Baseline	67.90	63.26	65.50	41.16	58.16	48.21	70.29	54.49	61.39
	Max	74.25	86.26	79.81	54.64	68.84	60.92	71.82	78.53	75.02
	Linear	75.67	89.34	81.94	61.10	73.90	66.89	74.31	88.82	80.92
	Min-Max	76.18	93.13	83.81	65.06	81.51	72.37	75.06	80.37	77.63
	Base-Answerable	74.76	83.69	78.98	57.53	67.75	62.22	72.61	79.98	76.11
	Base-Unanswerable	75.23	82.99	78.92	54.53	68.37	60.67	71.04	74.30	72.65

Figure 34: Model performance under the distractor and ideal retrieval evaluation settings across different curriculum construction settings

6 Implementation: LangChain, Text-to-speech, Speech-to-text

The full implementation code can be accessed at https://github.com/phamtranminhtri/music_rag

6.1 Building an AI Agent Chatbot with LangChain

We implement AI agent chatbot using the LangChain framework, designed to support a music-related website through Retrieval-Augmented Generation (RAG) combined with speech-based interaction. LangChain enables the construction of agentic systems by integrating large language models (LLMs) with external tools, memory, and retrieval mechanisms, allowing the chatbot to reason, decide actions, and iteratively solve user queries.

An **AI agent** in LangChain is defined as a system that combines a language model with a set of tools, enabling it to select and invoke appropriate tools based on the task context. Tools act as functional interfaces that extend the model's capabilities beyond pure text generation, such as searching documents or querying the web.

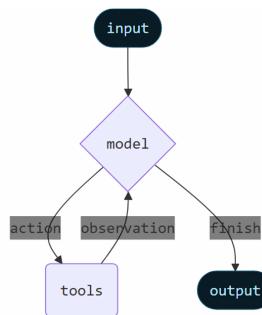


Figure 35: AI Agent

6.1.1 RAG Agent Architecture

The chatbot is built around a RAG agent architecture consisting of three main stages: indexing, retrieval, and response generation. During the indexing phase, nearly 100 Wikipedia articles related to Music are loaded using LangChain's `WikipediaLoader`. These documents are embedded using the `gemini-embedding-001` model and stored persistently in a `Chroma` vector database, enabling efficient semantic similarity search.

6.1.2 Agent Tools

The agent is equipped with multiple tools to enhance its reasoning and information access:

- **Context Retrieval Tool:** Retrieves the top- k (with $k = 10$) most relevant documents from the Chroma vector store based on semantic similarity.
- **Web Search Tool:** Uses DuckDuckGo to fetch real-time search results, including URLs, titles, and snippets, enabling access to up-to-date information.
- **Web Loader Tool:** Employs LangChain's `WebBaseLoader` to extract and process raw HTML content from selected web pages.

These tools are wrapped and exposed to the agent, allowing it to dynamically decide whether to rely on internal knowledge, retrieved documents, or external web sources.

6.1.3 Agent Configuration and Memory

The agent is initialized with a system prompt that defines its role and behavior. The conversational backbone uses the `gemini-2.5-flash` chat model for response generation. To maintain contextual coherence, short-term memory is incorporated, enabling the agent to remember and reference previous turns within a single conversation thread.

6.2 Text-to-Speech

Text-to-Speech (TTS) is a technology that converts written text into synthesized speech, enabling AI systems to produce audible responses for users. In conversational AI and Retrieval-Augmented Generation (RAG) systems, TTS plays an important role in enhancing user interaction by delivering information in a natural and accessible audio format.

6.2.1 Text-to-Speech Pipeline

A typical Text-to-Speech system follows a multi-stage processing pipeline that transforms raw text into an audio signal. The main stages of this pipeline are illustrated in Figure 36.

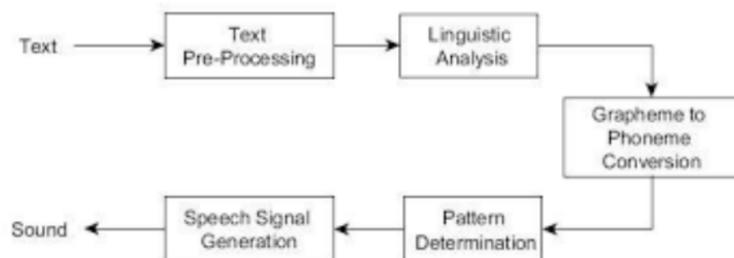


Figure 36: General Text-to-Speech pipeline from text input to synthesized speech output

The process begins with **text pre-processing**, where the input text is normalized by handling numbers, abbreviations, and formatting. Next, **linguistic analysis** examines grammatical structure and syntactic features. The **grapheme-to-phoneme conversion** stage maps characters to phonemes to determine correct pronunciation. Based on this phonetic representation, **pattern determination** defines prosodic features such as rhythm, pitch, and stress. Finally, **speech signal generation** synthesizes the waveform that produces the audible speech output.

6.2.2 gTTS: Google Text-to-Speech

gTTS (Google Text-to-Speech) is a lightweight Python library that provides an easy-to-use interface for converting text into speech. It leverages Google Translate's text-to-speech service to generate spoken audio from textual input and outputs the result as an MP3 file.

The simplicity of gTTS makes it suitable for rapid prototyping and integration into AI applications where basic speech output is required. Figure 37 shows a minimal example of using gTTS to synthesize Vietnamese speech from text.

```
1 from gtts import gTTS
2
3 text = "Xin chào, đây là gTTS"
4 tts = gTTS(text, lang='vi')
5 tts.save("output.mp3")
```



Figure 37: Example of using gTTS in Python to generate an MP3 file from text

In this example, the input text is passed to the gTTS library along with the target language code. The generated speech is then saved as an audio file, which can be played back by the system. While gTTS is easy to use and freely accessible, it provides limited control over voice characteristics and requires an internet connection.

Overall, Text-to-Speech systems such as gTTS complement language models by enabling multimodal interaction, allowing AI-driven applications to communicate information not only through text but also through natural-sounding speech.

6.3 Speech-to-text

Automatic Speech Recognition (ASR), commonly known as Speech-to-Text (STT), is the process of converting spoken language into written text. In this study, two distinct approaches were evaluated: the **SpeechRecognition** library (utilizing Google Web Speech API) and OpenAI's **Whisper** model.

6.3.1 SpeechRecognition with Google Web Speech API

The **SpeechRecognition** library serves as a flexible wrapper for several high-level STT APIs. For this implementation, the **Google Web Speech API** was utilized.

- **Mechanism:** The process follows a client-server architecture. The client records or loads an audio file (typically in .wav format), which is then transmitted to Google's cloud servers. The server processes the audio using proprietary deep learning models and returns the transcribed text in a JSON format.
- **Technical Characteristics:** Since the heavy computation is handled server-side, the local machine requires minimal processing power. It supports a vast array of languages, including high-accuracy recognition for Vietnamese.
- **Strengths:** Extremely lightweight for the host device, high inference speed, and minimal setup requirements.

6.3.2 OpenAI Whisper

Whisper is a state-of-the-art, multi-purpose speech recognition model trained on 680,000 hours of multilingual and multitask supervised data. Unlike API-based solutions, Whisper is a **Transformer-based encoder-decoder** model designed for local deployment.

- **Architecture:**

- **Input:** The raw audio is first converted into a **Log-Mel Spectrogram**, which represents the frequency content of the audio over time.
- **Encoder:** A series of convolutional layers and Transformer blocks encode the audio features into a latent representation.
- **Decoder:** An autoregressive Transformer decoder predicts the next text token based on the previously generated tokens and the encoder's audio features.
- **Strengths:** Excellent understanding of context, automatic punctuation, and the ability to function entirely offline.

6.3.3 Experimental Comparison and Evaluation

Through practical implementation in a demonstration application, both technologies were compared based on performance, resource consumption, and output quality.

Table 4: Comparison: SpeechRecognition API vs. Whisper Model

Feature	SpeechRecognition (Google)	OpenAI Whisper
Inference Speed	Extremely Fast (Cloud-dependent)	Slower (Hardware-dependent)
Hardware Req.	Minimal (Low CPU/RAM)	High (Requires GPU/Strong CPU)
Connectivity	Requires Internet	Offline capable
Formatting	Basic (No punctuation/casing)	Advanced (Punctuation/Casing)
Vietnamese	Excellent recognition	Good, but prone to hallucinations
Customization	None (Closed API)	Support for Fine-tuning

Pros and Cons Analysis:

- **SpeechRecognition:**

- *Pros:* Ideal for low-end hardware; exceptional Vietnamese recognition with zero local training.
- *Cons:* Strictly requires an internet connection; limited by request quotas; lacks advanced text formatting in the free tier.

- **Whisper:**

- *Pros:* Provides human-like transcriptions with proper capitalization and grammar; provides privacy through local processing.
- *Cons:* Computationally "heavy"; may suffer from *hallucinations* (generating text when there is only background noise).

6.3.4 Conclusion for Prototype Implementation

Based on the experimental results, **SpeechRecognition (Google Web Speech API)** was selected for the final demonstration. This decision was driven by its **ease of integration, high inference speed on low-resource hardware**, and **superior Vietnamese recognition** without the need for complex local environment configurations. While Whisper offers higher contextual quality, its hardware demands outweigh its benefits for this specific lightweight application demo.

6.4 Python Flask app for the chatbot

The application is built as a web-based conversational agent specifically designed for music-related queries. The architecture leverages the Flask microframework for the backend, orchestrating interaction between the user interface, external APIs, and a Retrieval-Augmented Generation (RAG) pipeline. The system components are detailed below:

6.4.1 Backend Core

The backend is implemented in Python using Flask. It exposes several RESTful endpoints to handle user interactions:

- **/** (GET): Serves the main web interface and initializes a unique session identifier for conversation history tracking.
- **/chat** (POST): The primary endpoint for processing text-based user queries. It invokes the LangChain agent to generate responses.
- **/stt** (POST): Handles Speech-to-Text conversion. It accepts audio files (WebM format), converts them to WAV using *pydub*, and utilizes the *SpeechRecognition* library with the Google Speech Recognition API to transcribe user speech.
- **/tts** (POST): Handles Text-to-Speech conversion. It detects the language of the agent's response using *langdetect* and synthesizes audio using the Google Text-to-Speech (*gTTS*) library.

6.4.2 Conversational Agent & RAG Pipeline

The core intelligence is powered by **LangChain** and the **Google Gemini 2.5 Flash** model. The agent is configured with a persistent memory mechanism (**InMemorySaver**) to maintain context within a session. The system employs a decision-making strategy defined in its system prompt to select between three distinct tools:

1. **Context Retrieval (RAG)**: For queries regarding classical music or established knowledge, the agent utilizes a local **Chroma** vector database. Queries are embedded using **GoogleGenerativeAIEmbeddings** (model: *gemini-embedding-001*) to perform similarity searches against a pre-indexed collection of music-related Wikipedia documents.
2. **Web Search**: For current events or information not present in the local database, the agent employs *DuckDuckGoSearchResults* to perform live web searches.
3. **Web Loader**: When specific external web pages are identified as satisfied sources, this tool utilizes **WebBaseLoader** to scrape and ingest the full textual content of the URL for synthesis.



The agent logic prioritizes the local RAG tool for efficiency and reliability, falling back to web search tools only when necessary information is absent or requires real-time data.



7 Evaluation

We need to evaluate how effective our RAG system is. We will be rating the RAG feature, not the web search tool here. As discussed in the previous section, the RAG vector store is taken from 76 Wikipedia articles about Music. We will also measure the time it takes to retrieve the relevance context from the database.

We use 2 types of evaluation score:

1. **Cosine Similarity (0-1):** Semantic similarity between query embeddings and document embeddings
2. **LLM-as-a-Judge (1-5 scale):**
 - 5 = Highly relevant, directly answers query
 - 4 = Relevant, provides useful information
 - 3 = Somewhat relevant, tangentially related
 - 2 = Minimally relevant
 - 1 = Not relevant

We use 10 text queries: Basic factual queries (1-2), Temporal queries (3-5), Event-based queries (6-7), Complex analytical queries (8-10)

1. What are the most common musical instruments?
2. Explain the concept of harmony in music
3. How has music evolved over the centuries?
4. What are the characteristics of music in the Renaissance period?
5. When did electronic music emerge?
6. What impact did the invention of recording technology have on music?
7. How did jazz influence modern music genres?
8. Compare Western and Eastern musical traditions
9. What is the relationship between rhythm and cultural identity in music?
10. How do musical scales affect emotional perception?

For each of the 10 text queries, we first retrieve the top 5 relevant documents from the vector store. Then we use the 2 methods above to evaluate the relevant score of the retrieved docs.

1. For method 1, we calculate the embedding of the query and of each document, then we calculate the cosine similarity between the embedding of query with the embedding of each doc.
2. For method 2, we feed a prompt to the LLM to ask it to give a score based on how relevant the document is to the query. The model used is **gemini-2.5-flash**.



```
1 for i, doc in enumerate(retrieved_docs):
2     prompt = f"""Rate the relevance of the following document to the query on
3         a scale of 1-5.
4
5 Query: {query}
6 Document Content: {doc.page_content[:500]}...
7
8 Scoring scale:
9 5 - Highly relevant, directly answers the query
10 4 - Relevant, provides useful information
11 3 - Somewhat relevant, tangentially related
12 2 - Minimally relevant, barely related
13 1 - Not relevant, unrelated to query
14
15 Respond with ONLY the numeric score (1-5)."""
16
17     try:
18         response = model.invoke(prompt)
19         score = int(response.content.strip())
20         scores.append(min(max(score, 1), 5)) # Ensure score is between 1-5
21     except:
22         scores.append(3) # Default to neutral if parsing fails
23
```

Listing 1: LLM-as-a-judge

Here is the result:

```
=====
2 Query 1: What are the most common musical instruments?
3 =====
4
5 Metrics:
6     Retrieval time: 1.732s
7     Avg similarity: 0.697
8     Avg relevance: 2.60/5
9
10 Top 3 Retrieved Documents:
11
12     Document 1:
13         Similarity: 0.715
14         Relevance: 3/5
15         Source: https://en.wikipedia.org/wiki/Musical_instrument
16         Content: as the only system that applies to any culture and, more importantly,
17             provides the only possible classification for each instrument. The most
18             common c...
19
20     Document 2:
21         Similarity: 0.703
22         Relevance: 4/5
23         Source: https://en.wikipedia.org/wiki/Percussion_instrument
24         Content: === By prevalence in common knowledge ===
25 It is difficult to define what is common knowledge but there are instruments
26             percussionists and composers us...
27
28     Document 3:
29         Similarity: 0.697
30         Relevance: 2/5
31         Source: https://en.wikipedia.org/wiki/Musical_instrument
32         Content: There are many different methods of classifying musical instruments.
33             Various methods examine aspects such as the physical properties of the
```



```
    instrument...  
30  
31 =====  
32 Query 2: Explain the concept of harmony in music  
33 =====  
34  
35 Metrics:  
36     Retrieval time: 0.483s  
37     Avg similarity: 0.755  
38     Avg relevance: 5.00/5  
39  
40 Top 3 Retrieved Documents:  
41  
42 Document 1:  
43     Similarity: 0.773  
44     Relevance: 5/5  
45     Source: https://en.wikipedia.org/wiki/Music  
46     Content: Harmony refers to the "vertical" sounds of pitches in music, which  
means pitches that are played or sung together at the same time creates a  
chord. Us...  
47  
48 Document 2:  
49     Similarity: 0.764  
50     Relevance: 5/5  
51     Source: https://en.wikipedia.org/wiki/Music\_theory  
52     Content: In music, harmony is the use of simultaneous pitches (tones, notes),  
or chords. The study of harmony involves chords and their construction and  
chord ...  
53  
54 Document 3:  
55     Similarity: 0.751  
56     Relevance: 5/5  
57     Source: https://en.wikipedia.org/wiki/Music  
58     Content: === Harmony ===.  
59  
60 =====  
61 Query 3: How has music evolved over the centuries?  
62 =====  
63  
64 Metrics:  
65     Retrieval time: 0.465s  
66     Avg similarity: 0.720  
67     Avg relevance: 4.20/5  
68  
69 Top 3 Retrieved Documents:  
70  
71 Document 1:  
72     Similarity: 0.731  
73     Relevance: 5/5  
74     Source: https://en.wikipedia.org/wiki/Popular\_music  
75     Content: There are multiple possible explanations for many of these changes.  
One reason for the brevity of songs in the past was the physical capability of  
rec...  
76  
77 Document 2:  
78     Similarity: 0.730  
79     Relevance: 5/5  
80     Source: https://en.wikipedia.org/wiki/Classical\_music  
81     Content: By the 20th century, stylistic unification gradually dissipated while  
the prominence of popular music greatly increased. Many composers actively  
avoid...
```



```
83 Document 3:  
84     Similarity: 0.718  
85     Relevance: 3/5  
86     Source: https://en.wikipedia.org/wiki/Popular\_music  
87     Content: In addition to many changes in specific sounds and technologies used,  
     there has been a shift in the content and key elements of popular music since  
     th...  
88 =====  
89 Query 4: What are the characteristics of music in the Renaissance period?  
90 =====  
91  
92 Metrics:  
93     Retrieval time: 0.479s  
94     Avg similarity: 0.724  
95     Avg relevance: 4.40/5  
96  
97 Top 3 Retrieved Documents:  
98  
99 Document 1:  
100     Similarity: 0.744  
101     Relevance: 5/5  
102     Source: https://en.wikipedia.org/wiki/Music  
103     Content: Renaissance music (c.1400 to 1600) was more focused on secular themes  
     , such as courtly love. Around 1450, the printing press was invented, which  
     made...  
104  
105 Document 2:  
106     Similarity: 0.734  
107     Relevance: 5/5  
108     Source: https://en.wikipedia.org/wiki/Classical\_music  
109     Content: === Renaissance ===  
110  
111 The musical Renaissance era lasted from 1400 to 1600. It was characterized by  
     greater use of instrumentation, multiple interwea...  
112  
113 Document 3:  
114     Similarity: 0.722  
115     Relevance: 4/5  
116     Source: https://en.wikipedia.org/wiki/Baroque\_music  
117     Content: tritone, perceived as an unstable interval, to create dissonance (it  
     was used in the dominant seventh chord and the diminished chord). An interest  
     in ...  
118  
119 =====  
120 Query 5: When did electronic music emerge?  
121 =====  
122  
123 Metrics:  
124     Retrieval time: 0.571s  
125     Avg similarity: 0.728  
126     Avg relevance: 4.60/5  
127  
128 Top 3 Retrieved Documents:  
129  
130 Document 1:  
131     Similarity: 0.741  
132     Relevance: 5/5  
133     Source: https://en.wikipedia.org/wiki/Electronic\_music  
134     Content: The first electronic musical devices were developed at the end of the  
     19th century. During the 1920s and 1930s, some electronic instruments were  
     intro...  
135
```



```
136
137 Document 2:
138   Similarity: 0.740
139   Relevance: 5/5
140   Source: https://en.wikipedia.org/wiki/Electronic\_music
141   Content: During the 1960s, digital computer music was pioneered, innovation in
142     live electronics took place, and Japanese electronic musical instruments
143     began t...
144
145 Document 3:
146   Similarity: 0.724
147   Relevance: 5/5
148   Source: https://en.wikipedia.org/wiki/Electronic\_music
149   Content: === United States ===
150 In the United States, electronic music was being created as early as 1939, when
151   John Cage published Imaginary Landscape, No. 1, ...
152 =====
153
154 Metrics:
155   Retrieval time: 0.472s
156   Avg similarity: 0.708
157   Avg relevance: 5.00/5
158
159 Top 3 Retrieved Documents:
160
161 Document 1:
162   Similarity: 0.715
163   Relevance: 5/5
164   Source: https://en.wikipedia.org/wiki/Music
165   Content: distributed. The introduction of the multitrack recording system had
166     a major influence on rock music, because it could do more than record a band's
167     pe...
168
169 Document 2:
170   Similarity: 0.712
171   Relevance: 5/5
172   Source: https://en.wikipedia.org/wiki/Popular\_music
173   Content: In the 1950s and 1960s, the new invention of television began to play
174     an increasingly important role in disseminating new popular music. Variety
175     shows...
176
177 Document 3:
178   Similarity: 0.709
179   Relevance: 5/5
180   Source: https://en.wikipedia.org/wiki/Music
181   Content: In the 19th century, a key way new compositions became known to the
182     public was by the sales of sheet music, which middle class amateur music
183     lovers wo...
184 =====
185
186 Query 6: What impact did the invention of recording technology have on music?
187 =====
188 Metrics:
189   Retrieval time: 0.539s
190   Avg similarity: 0.705
191   Avg relevance: 3.60/5
192
193 Top 3 Retrieved Documents:
```



```
189
190 Document 1:
191 Similarity: 0.719
192 Relevance: 5/5
193 Source: https://en.wikipedia.org/wiki/Rock\_music
194 Content: fusion began to take its audience, but acts like Steely Dan, Frank Zappa and Joni Mitchell recorded significant jazz-influenced albums in this period, ...
195
196 Document 2:
197 Similarity: 0.712
198 Relevance: 3/5
199 Source: https://en.wikipedia.org/wiki/Music
200 Content: Jazz evolved and became an important genre of music over the course of the 20th century, and during the second half, rock music did the same. Jazz is ...
201
202 Document 3:
203 Similarity: 0.702
204 Relevance: 5/5
205 Source: https://en.wikipedia.org/wiki/Rock\_music
206 Content: British acts to emerge in the same period from the blues scene, to make use of the tonal and improvisational aspects of jazz, included Nucleus and the...
207 =====
208 Query 8: Compare Western and Eastern musical traditions
209 =====
210 =====
211
212 Metrics:
213 Retrieval time: 0.499s
214 Avg similarity: 0.720
215 Avg relevance: 4.40/5
216
217 Top 3 Retrieved Documents:
218
219 Document 1:
220 Similarity: 0.734
221 Relevance: 4/5
222 Source: https://en.wikipedia.org/wiki/Music
223 Content: In the West, much of the history of music that is taught deals with the Western civilization's art music, known as classical music. The history of mus...
224
225 Document 2:
226 Similarity: 0.720
227 Relevance: 4/5
228 Source: https://en.wikipedia.org/wiki/Classical\_music
229 Content: Classical music generally refers to the art music of the Western world, considered to be distinct from Western folk music or popular music traditions....
230
231 Document 3:
232 Similarity: 0.715
233 Relevance: 4/5
234 Source: https://en.wikipedia.org/wiki/Music
235 Content: === Asian cultures ===
236
237 Asian music covers a swath of music cultures surveyed in the articles on Arabia, Central Asia, East Asia, South Asia, and Sout...
238 =====
```



```
240 Query 9: What is the relationship between rhythm and cultural identity in music?  
241 =====  
242  
243 Metrics:  
244   Retrieval time: 0.477s  
245   Avg similarity: 0.696  
246   Avg relevance: 3.20/5  
247  
248 Top 3 Retrieved Documents:  
249  
250   Document 1:  
251     Similarity: 0.709  
252     Relevance: 5/5  
253     Source: https://en.wikipedia.org/wiki/Tempo  
254     Content: song (although this would be less likely with an experienced  
255       bandleader). Differences in tempo and its interpretation can differ between  
256       cultures, as ...  
257  
258   Document 2:  
259     Similarity: 0.698  
260     Relevance: 3/5  
261     Source: https://en.wikipedia.org/wiki/Tempo  
262     Content: This context-dependent perception of tempo and rhythm is explained by  
263       the principle of correlative perception, according to which data are  
264       perceived i...  
265  
266   Document 3:  
267     Similarity: 0.695  
268     Relevance: 2/5  
269     Source: https://en.wikipedia.org/wiki/Music  
270     Content: === Rhythm ===...  
271  
272 =====  
273 Query 10: How do musical scales affect emotional perception?  
274 =====  
275  
276 Metrics:  
277   Retrieval time: 0.451s  
278   Avg similarity: 0.705  
279   Avg relevance: 3.00/5  
280  
281 Top 3 Retrieved Documents:  
282  
283   Document 1:  
284     Similarity: 0.709  
285     Relevance: 5/5  
286     Source: https://en.wikipedia.org/wiki/Music\_theory  
287     Content: The interrelationship of the keys most commonly used in Western tonal  
288       music is conveniently shown by the circle of fifths. Unique key signatures  
289       are a...  
290  
291   Document 2:  
292     Similarity: 0.708  
293     Relevance: 4/5  
294     Source: https://en.wikipedia.org/wiki/Scale\_\(music\)  
295     Content: Tetratonic (4 notes), tritonic (3 notes), and ditonic (2 notes):  
296       generally limited to prehistoric ("primitive") music  
297     Scales may also be described by ...  
298  
299   Document 3:  
300     Similarity: 0.704  
301     Relevance: 2/5
```



```
295     Source: https://en.wikipedia.org/wiki/Soundtrack
296     Content: plot anticipations, and moral judgement of the characters.
297     Furthermore, eyetracking and pupillometry studies found that film music is
298     able to influenc...
299 =====
```

Listing 2: RAG score result

Overall, we can see that the time it take to retrieve the context is acceptable (less than 1 second for most cases). The average similarity score is around 0.7-0.75; and the average LLM judge rating varied more, around 2.6 to 5.0. LLM judge score can be pretty inconsistance and biased, or even inaccurate, so we use extra manual inspection to check the relevance of the retrieved context to the query; but generally, the LLM can be quite good at deciding which infomation is important to answer the query.

In conclusion, we can say that with our vector database of 76 Wiki articles is quite sufficient for our RAG system, and RAG did a good job at retrieving the neccesary context for the queries.



8 Conclusion

In this project, we have successfully designed and implemented an AI Agent chatbot specializing in the music domain, leveraging the power of LLM and RAG. By integrating LangChain, a vector database constructed from 76 Wikipedia articles, and speech interfaces (Text-to-Speech and Speech-to-Text), we created an interactive system capable of understanding and answering complex user queries with contextually relevant information.

Our theoretical exploration covered the fundamental building blocks of modern AI, including Artificial Neural Networks, Transformer architectures, and the evolution of LLMs. We delved into the specifics of RAG, examining advanced techniques such as reranking and reasoning to overcome the limitations of standard vector-based retrieval.

The evaluation of our system demonstrated the effectiveness of the RAG approach. The system achieved sub-second retrieval times and maintained consistent relevance scores, with cosine similarity averaging between 0.7 and 0.75. The "LLM-as-a-Judge" evaluation, complemented by manual inspection, confirmed that the retrieved context significantly enhanced the quality and accuracy of the generated responses compared to a standalone model.

Looking ahead, there are several avenues for future improvement. Expanding the knowledge base beyond the initial set of Wikipedia articles would broaden the agent's expertise. Implementing more sophisticated reranking algorithms and exploring reinforcement learning techniques could further refine retrieval accuracy. Additionally, optimizing the latency of the speech interfaces would contribute to a more seamless real-time user experience. Overall, this project validates the potential of domain-specific RAG agents in bridging the gap between general-purpose LLMs and specialized information needs.

References

- [1] Yunfan Gao et al. *Retrieval-Augmented Generation for Large Language Models: A Survey*. Preprint. 2023. DOI: [10.48550/arXiv.2312.10997](https://doi.org/10.48550/arXiv.2312.10997). arXiv: [2312.10997 \[cs.CL\]](https://arxiv.org/abs/2312.10997). URL: <https://arxiv.org/abs/2312.10997>.
- [2] Google AI for Developers. *Google AI Studio*. Accessed: 2026-01-08. Google. 2025. URL: <https://ai.google.dev/aistudio>.
- [3] Jerry Huang et al. *RAG-RL: Advancing Retrieval-Augmented Generation via RL and Curriculum Learning*. Preprint. 2025. DOI: [10.48550/arXiv.2503.12759](https://doi.org/10.48550/arXiv.2503.12759). arXiv: [2503.12759 \[cs.CL\]](https://arxiv.org/abs/2503.12759). URL: <https://arxiv.org/abs/2503.12759>.
- [4] LangChain Documentation. *Build a RAG agent with LangChain*. Accessed: 2026-01-05. LangChain. 2025. URL: <https://docs.langchain.com/oss/python/langchain/rag>.
- [5] Yangning Li et al. *Towards Agentic RAG with Deep Reasoning: A Survey of RAG-Reasoning Systems in LLMs*. Preprint. 2025. DOI: [10.48550/arXiv.2507.09477](https://doi.org/10.48550/arXiv.2507.09477). arXiv: [2507.09477 \[cs.CL\]](https://arxiv.org/abs/2507.09477). URL: <https://arxiv.org/abs/2507.09477>.
- [6] OpenAI. *Introducing ChatGPT*. Accessed: 2026-01-08. OpenAI. 2022. URL: <https://openai.com/index/chatgpt/>.