# Building with LLMs

**marvik.**
AI solutions with business impact

Marvik is a hands-on AI consulting firm dedicated to creating impactful solutions for businesses. We excel in computer vision, natural language processing, GenAI, predictive analytics and more. Our team of +150 experts has been working for five years to deliver excellence in AI. We provide services from strategy design to full implementation, ensuring that our clients stay ahead in their industries.

Our mission is to empower businesses by harnessing data and accelerating digital transformation. We take a collaborative approach, working closely with clients to scale their internal capabilities and achieve success. At Marvik, we are committed to innovation and aligning with our clients' goals to drive results.

More information www.marvik.ai

# Building with LLMs

# Table of contents

# Introduction

**Large Language Models (LLMs)** are advanced deep learning systems that are able to perform a wide range of **natural language processing (NLP)** tasks, including translation, text generation, and sentiment analysis, among others. These models are trained on huge amounts of textual data, enabling them to understand language structure, context, and grammar at a sophisticated level. This capability makes them powerful tools for improving human-computer interactions and driving innovations across various applications.

The understanding capabilities of LLMs are not in the human sense. These models can identify patterns in text and predict the most likely responses based on statistical probabilities. By learning from enormous datasets, these models can generate coherent and contextually appropriate sequences of words, mimicking human-like communication. Unlike earlier architectures like **Recurrent Neural Networks (RNNs)**, that are not able to handle long-range dependencies within the input text, LLMs typically overcome this problem through self-attention mechanisms, which enable them to maintain context over much larger spans of text, significantly improving the understanding of long-form context.

This e-book presents some of the core fundamentals needed for understanding LLMs and applying them in real-world scenarios. It begins  with an introduction to how LLMs work and explores some alternatives to selecting and hosting an LLM. Next, it outlines some of the tasks that can be addressed using LLMs, providing explanations. The following section presents a high-level overview of the Transformers architecture, widely used in the LLMs world. The book then discusses techniques for enhancing results, including **Prompt Engineering** and **Retrieval-Augmented Generation**. Finally, the last section highlights some of the latest advancements in **AI agents**.

## How LLMs Work

LLMs operate by taking an input sequence, such as a sentence or a paragraph, and using their extensive training on massive amounts of textual data to process the input and generate a relevant output. When fed with an input text, LLMs analyze the sequence, capture the underlying meaning and context, and respond in a way that mimics human language.

At a high level, LLMs break down input text into tokens (smaller units like words or subwords), and through their architecture, they identify the relationships between

these tokens. This ability to attend to the entire context enables the model to generate coherent, contextually relevant responses.

# How LLMs Are Trained

Training an LLM involves teaching it to understand and generate human-like text. This is basically accomplished through three main stages: **Pre-training, Supervised fine-tuning** and **LLM Alignment**.

**Pre-training** consists of exposing the model to large amounts of text so it can learn patterns in language on its own (without any labeling). **Supervised fine-tuning (SFT)** is a technique used to enhance the model's capabilities to perform specific tasks, such as text classification, sentiment analysis, or question-answering systems, by using labeled data. This often employs optimization techniques to improve efficiency and accuracy. Key methods include **LoRA (Low-Rank Adaptation)**, which reduces trainable parameters using low-rank decomposition. Finally, the **LLM Alignment** stage ensures that a model behaves in ways that reflect human values, ethical considerations and intended use cases. The most used alignment technique is **Reinforcement Learning from Human Feedback (RLHF)**, which incorporates human feedback to guide the model's behavior. In RLHF, the model generates outputs, and human evaluations are used to rank these responses. This feedback is then incorporated into the model's learning process, helping it prioritize more appropriate or aligned responses over time.

Training LLMs from scratch is a highly complex and resource-intensive task, often involving significant time, hardware, and financial costs. Fortunately, there are several state-of-the-art **pre-trained LLMs** already available, both proprietary and open-source. These existing models can be adapted to specific use cases achieving high-quality results, making them a far more efficient choice than building an LLM from scratch.

# LLM Usage Alternatives

Since LLMs are trained on vast datasets that cover a wide variety of language structures, domains, and topics, they can perform a range of language tasks without the need for additional fine-tuning or training. This "out-of-the-box" capability allows them to handle multiple tasks with high accuracy and fluency. Additionally, techniques like Prompt Engineering and Retrieved Augmented Generation (RAG) help improve the LLMs output and make them more suited for domain-specific knowledge. Their versatility makes LLMs useful across many

applications and using a pre-trained LLM is usually good enough for producing high quality responses.

When domain-specific knowledge is critical for an application, another possible approach involves taking a pre-existing language model, and further training it on a smaller task-specific dataset. This fine-tuning allows the model to adapt and customize its knowledge to the specific details and requirements of the target task. This approach is computationally more demanding than using a pre-trained LLM but might lead to better results, since the specific textual content is embedded inside the model.

# LLM Hosting

Nowadays there are multiple frameworks available for serving and utilizing LLMs, optimized for both CPU and GPU deployments. Two widespread examples are llama.cpp and vLLM, both of which can be used to provide APIs for accessing various LLMs. Typically, leveraging these frameworks requires setting up and maintaining an environment for deployment, such as a cluster or a GPU instance. While this approach allows for extensive configuration options, it can also complicate the setup process and often incurs significant costs associated with the required infrastructure.

An attractive alternative is to utilize serverless services offered by major cloud providers, which enable the deployment and access of different LLMs with minimal setup effort. These services typically offer optimized interfaces for inference, ensuring efficient performance and allowing for scalability without the risk of bottlenecks. Some alternatives in this category include Amazon Bedrock, Azure OpenAI Service, and Google Vertex AI.

# Selecting an LLM

Which LLM to select for an application will typically depend on different factors, it's important to recognize that different models excel at different tasks. For example, OpenAI's **GPT-4o** and Meta's **Llama 3.1 70B** may offer state-of-the-art performance for creative writing or complex reasoning, but for more targeted tasks like code generation, models like **Claude** from Anthropic could provide comparable results at lower cost.

## LLMs Size and Performance

Usually larger models have a better performance on general knowledge and reasoning tasks but come with an increase in computational costs and a decrease in response times, making them impractical for some use cases. Using a smaller model, such as **Llama 3.1 8B** might be ideal for lightweight applications where fast inference is a hard requirement. Striking a balance between model size and performance is essential, as smaller models can still deliver high-quality results for specific applications while reducing costs and enhancing efficiency.
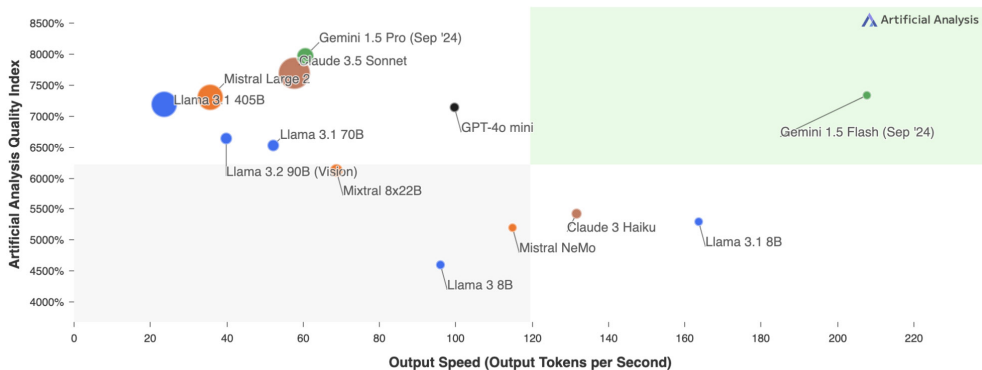


Figure [1] - A comparison of different LLMs regarding quality and inference speeds. Artificial Analysis Quality Index averages multiple performance metrics into a global index.

## Open Source and Closed Source LLMs

There are two main categories of state-of-the-art LLMs, open-source models and closed-source models. Open-source models, such as **Llama** and **Mistral** variants, are freely available for anyone to use and modify, allowing for collaboration within the community and fine-tuning for specific applications. On the other hand, closed-source models like OpenAI's **GPT-4** and Anthropic's **Claude** remain proprietary, with their codebases not accessible for the community, and being only accessible through paid subscriptions. Which type of LLM to use depends essentially on three aspects:

- **Accessibility and cost:** While many closed-source LLMs are available through subscription-based services, open-source models are freely available for use and modification. However, using an open-source LLM often involves significant costs related to the hardware needed to achieve

acceptable response times, which can scale up rapidly.
- **Data security:** When data security is crucial for the application, it is usually a better idea to self-host an open-source LLM with all the necessary security protocols. In closed-source environments it is necessary to rely on the vendor compliance for handling the data, which might be unacceptable in some cases.
- **Customization:** Open-source models provide the flexibility to modify the architecture and fine-tune the model on specific data, allowing for deeper customization. Closed-source models often have limited options for customization, restricting developers to predefined configurations.
- **Integration:** Closed-source models usually provide straightforward integration via APIs, which can simplify deployment. Integrating open-source models may require additional effort to configure, host, and maintain, especially when the models are self-hosted, but it allows for greater control over the deployment environment and data handling.

A useful reference for comparing different state-of-the-art LLMs across different metrics is the Huggingface Open LLM Leaderboard, which ranks and evaluates the performance of open-source models across various tasks. Although this ranking serves as a good starting point for selecting an LLM, it's essential to perform custom benchmarking to evaluate how well the models perform in a specific application to be developed. Additionally, some models may be fine-tuned to perform well in the metrics used for ranking but may not deliver the same level of quality in real-world applications.

# Common Tasks Enabled by LLMs

## Summarization

**Text summarization** is the process of taking a relatively long text document and producing a shorter and accurate version of it that captures the key concepts discussed in the text.  The main goal is to reduce the amount of text of a given input while preserving its most essential information and contextual meaning. Generally speaking,  summarization can be achieved by means of two different techniques: **extractive** and **abstractive** summarization.

**Extractive summarization** involves selecting exact sentences from the input text to build the summary, functioning like a binary classifier that decides whether each sentence is relevant enough to be included in the output. This technique ensures consistency with the original text since it does not alter the original text, but often lacks the ability to produce fluent summaries.

**Abstractive summarization** generates a concise, paraphrased summary that captures the main idea of the input text. By synthesizing the given information, it creates new phrases to capture the document's essence in a shorter and more readable format. The main challenge regarding this type of summarization is to generate a new text that is **factually correct** when compared to the original one.

Summarization is valuable in many fields such as media monitoring, providing quick updates on industry events and allowing businesses to stay informed without needing to process large volumes of content. In content creation, summarization helps writers condense research into concise versions and generate engaging teasers, accelerating the creative process and improving productivity.

A few years ago, extractive summarization applications using fine-tuned versions of BERT model were a standard in the industry. However, with the rising popularity of LLMs and their exceptional text generation capabilities, abstractive summarization techniques have recently gained significant traction.

**Input:**
The Empire State Building is a 102-story Art Deco skyscraper in the Midtown South neighborhood of Manhattan in New York City. The building was designed by Shreve, Lamb & Harmon and built from 1930 to 1931. Its name is derived from "Empire State", the nickname of the state of New York. The building has a roof height of 1,250 feet (380 m) and stands a total of 1,454 feet (443.2 m) tall, including its antenna. The Empire State Building was the world's tallest building until the first tower of the World Trade Center was topped out in 1970; following the September 11 attacks in 2001, the Empire State Building was New York City's tallest building until it was surpassed in 2012 by One World Trade Center
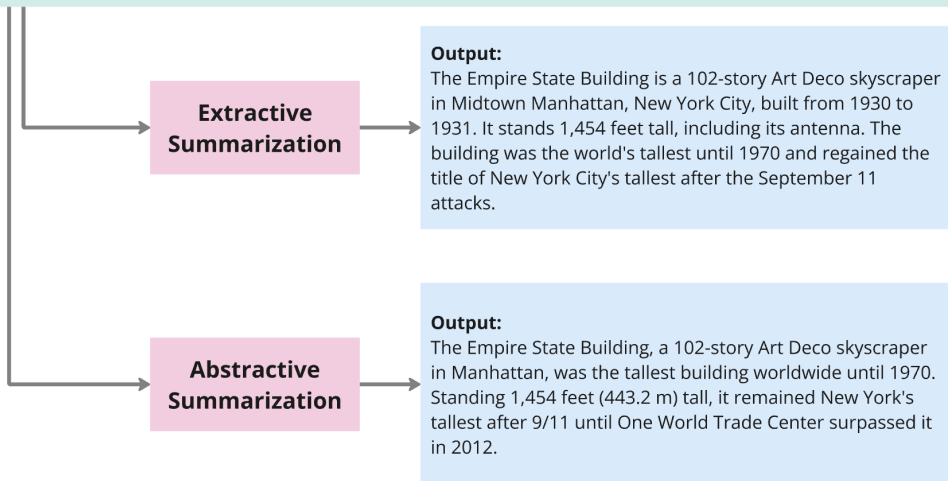
**Extractive Summarization**

**Output:**
The Empire State Building is a 102-story Art Deco skyscraper in Midtown Manhattan, New York City, built from 1930 to 1931. It stands 1,454 feet tall, including its antenna. The building was the world's tallest until 1970 and regained the title of New York City's tallest after the September 11 attacks.

**Abstractive Summarization**

**Output:**
The Empire State Building, a 102-story Art Deco skyscraper in Manhattan, was the tallest building worldwide until 1970. Standing 1,454 feet (443.2 m) tall, it remained New York's tallest after 9/11 until One World Trade Center surpassed it in 2012.

Figure [2] - Summarization task.

# Questions Answering

**Question Answering (Q&A)** applications focus on delivering accurate answers based on questions provided by external users. Most of the chatbot-like applications currently being developed are based on this Q&A concept. Since LLMs are trained on huge amounts of data, they can generate accurate answers when prompted with a simple question. When it is necessary to generate responses based on specific domain knowledge, a Q&A application can be significantly enhanced using techniques such as Retrieval Augmented Generation, which improves the quality of responses by consulting the content of specific documents that serve as a knowledge base.

Q&A applications can also be instrumented with conversational capabilities, by keeping track of the entire ongoing conversation between the user and the LLM. With this approach, the application can provide coherent and context-aware responses to follow-up questions that are closely related to previous queries. This approach enables the application to handle multi-turn interactions where

the meaning of a question depends on information exchanged earlier in the conversation.

Q&A applications enhance productivity and improve user experience across different industries. In legal and finance, they help professionals quickly find important information in large documents. E-commerce platforms also benefit from these systems by effectively answering product questions, leading to better customer interactions.
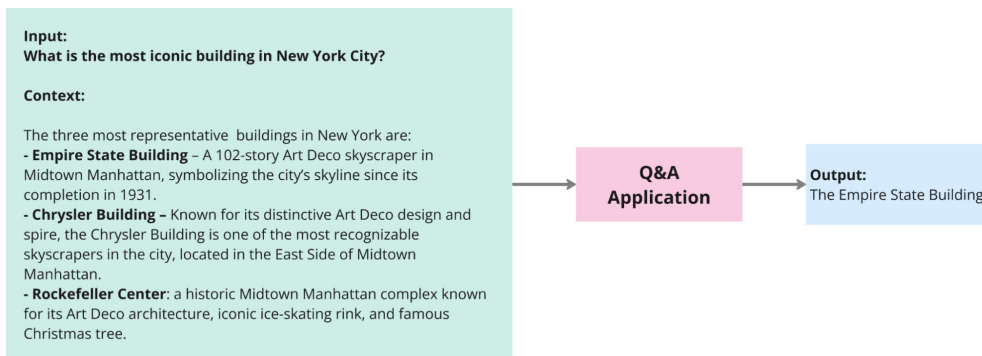


Figure [3] - Question Answering task.

# Classification

**Text Classification** involves categorizing text into predefined labels or classes. This technique is widely applied in sentiment analysis, where the goal is to identify the tone of a given sentence, or in topics classification, where it is desired to assign different tags or categories to each text. Due to their excellent understanding of language, LLMs allow the use of **zero-shot text classification** tasks, where text is assigned to predefined labels without requiring specific training. This approach utilizes contextual understanding to accurately classify unseen texts, making it ideal for applications where resources and labeled data are restricted, and therefore it is not plausible to train a supervised classifier.

Text classification tasks can be further enhanced using some of the Prompt Engineering techniques detailed later in this e-book.
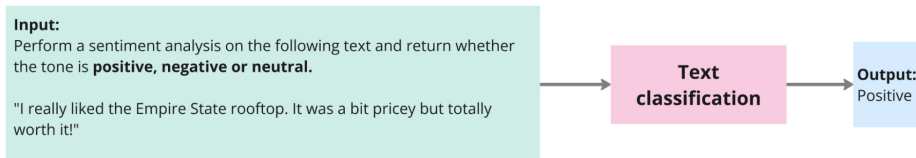
**Input:**
Perform a sentiment analysis on the following text and return whether the tone is **positive, negative or neutral.**

"I really liked the Empire State rooftop. It was a bit pricey but totally worth it!"

**Text classification**

**Output:**
Positive

Figure [4] - Text Classification task.

# Guardrails Implementation

**Guardrails** are a set of programmable safety guidelines that monitor and control a user's interaction with an LLM application.  By implementing guardrails, users can define structure, type, and quality of LLM responses, ensuring that outputs are reliable, safe, accurate, and aligned with specific expectations. Key applications of guardrails include information validation, fact-checking and monitoring for toxicity.

**Information validation** is essential for ensuring the reliability and safety of content generated by LLMs, as these models can produce text that includes hallucinations (inaccurate or fabricated information) or harmful language. To address these issues, it's important to implement validation processes that guarantee responses are both accurate and compliant. These measures, commonly referred to as guardrails, help maintain the quality and security of LLM-generated content, fostering a safer environment for users.

**Fact-checking** involves verifying the accuracy of information by cross-referencing it with reliable sources. Since LLMs can create coherent and convincing text, there's a risk of spreading false or misleading information. Therefore, implementing guardrails is essential for ensuring that the information provided is trustworthy.

Additionally, **toxicity checking** focuses on identifying and reducing harmful or inappropriate content generated during interactions. Given their capability to produce human-like text, LLMs may inadvertently generate offensive language. Guardrails are crucial in filtering out such toxicity, promoting a respectful and safe environment for users.
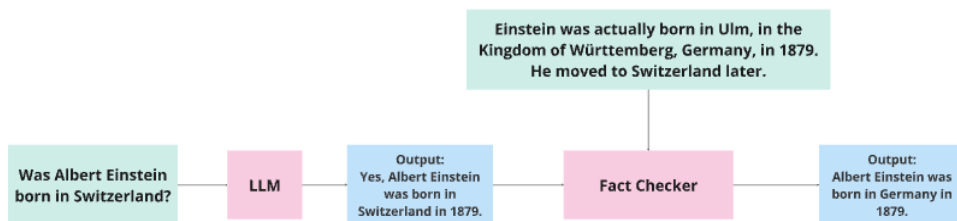
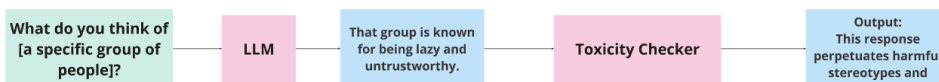Figure [5] - Example of a flow for fact checking.



Figure [6] - Example of a flow for toxicity checking.

# Information / Feature Extraction

**Feature extraction** involves using the LLM to identify and extract relevant information from unstructured text data and presenting it in a structured format, such as JSON. Common examples include extracting key entities from contracts and invoices, identifying customer intents from chat conversations or emails, and pulling specifications from proposals.

By transforming unstructured data into organized, machine-readable outputs, it simplifies analysis, improves processing efficiency, and facilitates the use of extracted features for downstream applications.
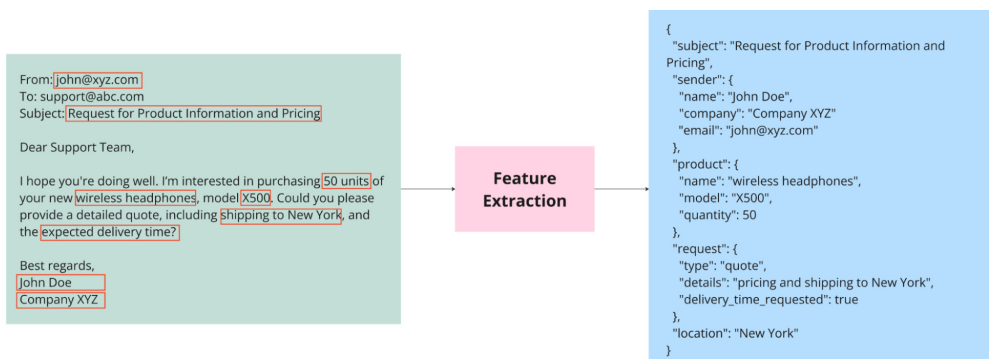


Figure [7] - Example of a feature extraction task.

# Architecture Overview

In this section the **Transformer architecture** will be reviewed, which is the foundational building block for different LLM models. Transformers are designed to process sequences of data, such as sentences or paragraphs, but unlike other deep learning based models like RNNs, they don't need to rely on sequential processing. Instead, Transformers use attention mechanisms, allowing them to process input in parallel, which significantly improves both performance and scalability.

Let's review the elementary concepts that are necessary for understanding how Transformers work.
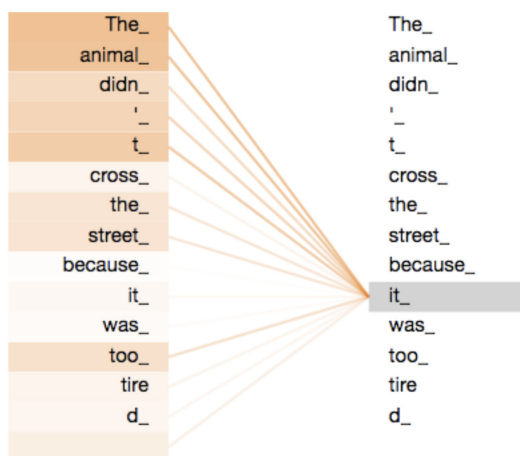
## Self-Attention Mechanism



Figure [8] - When processing the word "it", part of the attention mechanism focuses on "The Animal", and then encodes parts of "The Animal" representation into the "it" token.

The **self-attention mechanism** allows Transformers to maintain an exceptionally long-term memory compared to RNNs by enabling them to process all tokens in an input sequence simultaneously. This mechanism "attends" to each token within the sequence, capturing relationships across the entire sentence or paragraph, regardless of distance between tokens. While RNNs can also consider

nearby inputs, they have a limited reference window, and as the sequence length increases, these models struggle to understand relationships between inputs that are far-away from each other.

The word "attend"refers to the ability of focusing on specific elements of a sequence when generating an output. It involves selectively weighting the importance of different tokens (e.g. words) based on their relevance to the current task. When a model "attends" to a token, it evaluates how much influence that token should have on the representation of the current token being processed. This allows the model to capture contextual relationships and dependencies within a long sequence. Essentially, "attending" means determining which parts of the input are most important for the current processing step, helping the model to understand and generate language more effectively.

# Positional Encoding

Since transformers process all input sequences simultaneously, they lack an inherent understanding of the order of tokens or words. To address this limitation, it is essential to provide additional information regarding the relative or absolute positions of tokens within a sequence. This is accomplished through **positional encoding**, which generates a unique vector for each token embedding that represents its position in the sequence. These encodings are created using sine and cosine functions of varying frequencies. The use of functions enables the model to effectively learn relative positions between tokens, as the sine and cosine values produce distinct patterns for different positions while ensuring that the encodings maintain a consistent relationship.

# Transformers Architecture

The **Transformer architecture** leverages the concepts of **self-attention** and **positional encoding** to process and generate sequences of data effectively. Instead of using self-attention as explained before, Transformers use multiple attention "heads," each one learning different patterns in the input sequence. These multiple heads help the model capture various types of relationships and dependencies within the data, from syntactic to semantic ones.

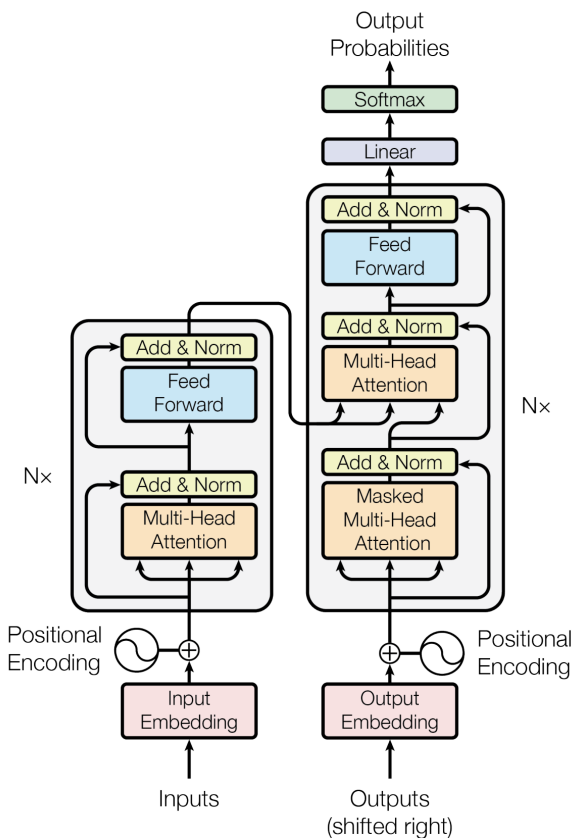The image below shows the architecture of a Transformer:



Figure [9] - Standard Transformer architecture.

# Encoder

The **encoder** component of the Transformer architecture processes an input sequence that has been converted into token embeddings. It consists of multiple layers that transform these embeddings to capture the relationships between tokens through the self-attention mechanism. This encoding process enables the model to learn dependencies among words regardless of how far apart they are in the text. Each encoder layer typically includes a **multi-head attention mechanism**, which allows the model to focus on different parts of the input simultaneously, and a feedforward neural network for additional processing.
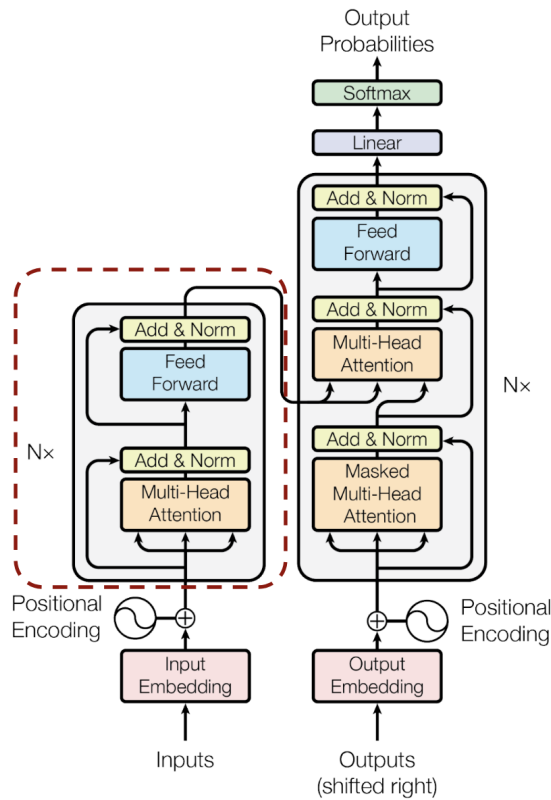
Figure [10] - Encoder portion of the Transformer architecture.

# Decoder

The **decoder** component is responsible for generating output sequences based on the encoded representations generated in the encoder. It takes the encoded information and produces the final predictions, operating in a step-by-step manner, where each output token is generated sequentially, making it crucial for tasks like language generation and machine translation.

The main difference compared to the encoder architecture is that the decoder uses a **masked self-attention mechanism** to ensure that it only considers previously generated tokens when predicting the next token, preventing it from accessing "future" information in generation time. Additionally, it incorporates an attention layer that connects back to the encoder's output, allowing it to utilize the context provided by the input sequence effectively.

It is relevant to note that most modern LLMs use a **decoder-only architecture** instead of an encoder-decoder one, since it has been proven to be more efficient for generating text. The decoder-only approach focuses on predicting the next word in a sequence based on the context of previous words. This design simplifies the training process and reduces computational overhead by eliminating the need for a separate encoding phase. As a result, decoder-only models can generate coherent and contextually relevant responses more quickly and effectively.
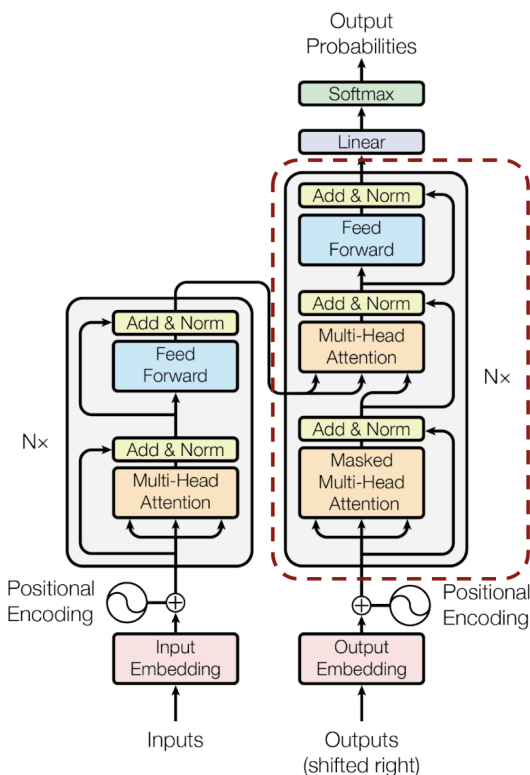


Figure [11] - Decoder portion of the Transformer architecture.

# Mixture of Experts

Some of the state-of-the-art models like Mixtral include a **Mixture of Experts (MoE)** architecture which allows training bigger models in a more efficient way. The MoE architecture introduces a set of "experts", which are separate neural

networks that can be utilized selectively. In training time, specific tokens are sent to specific experts, which enables the model to leverage the strengths of each expert for different types of input. By using only a subset of experts for any given input, the MoE architecture significantly reduces the computational load compared to dense models, where all parameters are active at all times. As a result, this approach allows for larger models with more parameters to be trained efficiently, leading to improved performance on a variety of tasks while maintaining a lower overall computational cost.

# Prompt Engineering

The following section will provide an overview of prompt techniques to improve the LLM response, along with relevant examples. For further experimentation, please refer to the following notebook.

**Prompt engineering** refers to a set of techniques used for designing precise and thoughtful prompts to guide LLMs in delivering accurate and relevant outputs. By carefully designing questions and providing clear instructions, one can significantly improve the quality of LLMs responses to be more accurate and to provide coherent results across a wide range of applications.

For the following examples on different prompt engineering techniques, let's consider this base prompt:

```
USER: Classify the following customer support inquiry into one of
the following categories: Billing, Technical Support, or Account
Issues.
The inquiry is: My internet connection keeps dropping, and I can't
access certain websites.
```

Figure [12] - Base Prompt example.

This prompt describes a task that is simple enough to be easily addressed by any LLM. However, additional prompt engineering techniques can be implemented in cases where the instructions are not straightforward. In order to illustrate the use of these techniques, the following sections present different scenarios similar to the base one but with additional complexity.

## Wrapping Input

The prompt often requires context to clarify the user's intent, ensuring that the response is relevant and tailored to its specific needs. Context provides necessary background information that allows the model to understand the subject matter better and address any complexities involved.

When incorporating context into a prompt, it's beneficial to use specific delimiters to wrap this context. Delimiters help the LLM in identifying distinct

segments of text within the prompt, clearly indicating which portions should be processed while separating them from the instructions. Various delimiters, such as quotes, XML tags, or section titles can be employed. The primary goal of the delimiters is to distinguish a specific text segment from the rest of the prompt.

```
USER: Classify the following customer support INQUIRY with DETAILS
into one of the following categories: Billing, Technical Support, or
Account Issues.

<INQUIRY>"My internet connection just does not work."</INQUIRY>

<DETAILS> The customer mentioned issues related to service
reliability. They indicated problems accessing websites but did not
mention payment or account-related problems. The back office platform
does not show payment issues.</DETAILS>
```

Figure [13] - Wrapping input example.

Additionally, **wrapping input** is useful for preventing **prompt injection**, which happens when users introduce malicious or conflicting instructions into a prompt. For instance, a user might input a statement like, "Forget the previous instruction and do exactly what I tell you." By enclosing this user input within delimiters, the model is better able to recognize that it should not follow these new instructions, but rather analyze the text in the context of previous instructions. This practice helps ensure that the model remains focused on its intended task, reducing the risk of unintended behavior. Using specificity delimiters acts as a protective measure, reinforcing the original context and guiding the model's responses effectively.

# Few-shot Learning

**Few-shot learning** allows the LLMs to enable in-context learning by using a set of few examples that are provided as part of the prompt. Instead of needing extensive training data, LLMs can generalize from just a handful of inputs, which is particularly helpful for specific tasks or rare queries. The provided examples serve as conditioning for subsequent examples where the model needs to generate a response. For instance, if the model is fed with a set of examples on how to classify text into categories, it can quickly adapt and classify new text it hasn't seen before.

```
USER: Classify the following customer support inquiry into one of
these categories based on the examples: Billing, Technical Support,
or Account Issues.

EXAMPLES:

- Example 1: Inquiry: I was charged twice for the same service. →
Category: Billing

- Example 2: Inquiry: I forgot my password and can't log in my
account. → Category: Account Issues

- Example 3: Inquiry: My connection has been unstable for the past
10 hours. → Category: Technical Support

INQUIRY: My internet connection keeps dropping, and I can't access
certain websites.

Please provide the appropriate category for this inquiry.
```

Figure [14] - Few shot learning prompt example.

# Chain of Thought

**Chain of thought (CoT)** prompting is a technique designed to improve the reasoning capabilities of an LLM by guiding it to express its reasoning process step-by-step. Instead of simply asking a direct question, this method encourages the model to articulate its reasoning, breaking down the problem into manageable components. By doing so, the model can resolve complex questions more effectively, as it systematically evaluates each part of the problem according to the instructions. This approach not only helps in reaching a more accurate conclusion but also increases transparency, allowing users to understand how the model arrived at its answer.

```
USER: Classify the following customer support inquiry into one of
these categories following the instructions: Billing, Technical
Support, or Account Issues.

INSTRUCTIONS:
Please think through the process step by step:
1. Identify the core problem in the inquiry.
```

```
2. Determine whether the problem relates to:
       - Billing: Any issues related to payments or charges.
       - Technical Support: Any technical issues, such as
         connectivity or functionality.
       - Account Issues: Any access or login problems.

INQUIRY: My internet connection keeps dropping, and I can't access
certain websites.

Please provide the appropriate category for this inquiry, and a
detail on how you reached the final answer.
```

Figure [15] - Chain of thought prompt example.

# Self Consistency

**Self-consistency** is used to enhance the reliability of generated responses. This method involves asking the model the same question multiple times and then taking the most frequently given answer as the final response. By sampling various outputs for the same prompt, self-consistency helps capture the LLM's diverse reasoning paths and mitigates the effects of any randomness in its responses.

```
USER: Classify the following customer support inquiry into one
of the categories: Billing, Technical Support, or Account Issues.
Provide your reasoning for each classification.

INQUIRY: I want to change my billing address.

INSTRUCTIONS:
Classify this inquiry three times, providing reasoning each time.
After generating the three classifications, take the majority result
as your final answer.
```

Figure [16] - Self Consistency prompt example.

```
ASSISTANT:

Q: My internet connection drops frequently, and I cannot stream
videos.

A1: This indicates a technical problem with the internet service, as
frequent disconnections are typically related to technical support
issues. Category: Technical Support.

A2: The inquiry suggests that the user is experiencing connectivity
issues, which fall under technical support. Category: Technical
Support.

A3: Since the user cannot stream due to connection drops, this is
clearly a technical issue that requires support. Category: Technical
Support.

The final category for the inquiry is Technical Support
```

Figure [17] - Self Consistency query output.

# Prompt Chaining

**Prompt chaining** is used to improve performance by breaking down complex
tasks into smaller, manageable subtasks. Instead of asking the model to tackle
a problem all at once, prompt chaining involves sequentially guiding the model
through each subtask.

**First prompt:**

```
USER: Extract the key information from the following customer
support inquiry.

INQUIRY: I cannot access my account after changing my password.
```

Figure [18] - First prompt of the prompt chaining technique.

**Second prompt:**

```
USER: Based on the extracted information, classify the inquiry into
one of these categories: Billing, Technical Support, or Account
Issues.

Extracted Information:

```

{Output from the first prompt}

```
```

Figure [19] - Second prompt of the prompt chaining technique.

# Output Parser

In some scenarios (e.g. Feature Extraction tasks), it is necessary to retrieve answers from a LLM in a **structured format** that can be easily parsed for downstream tasks. An **output parser** is a tool that transforms the raw output text from the LLM into a predefined structure, such as JSON. By instructing the LLM to generate responses according to this defined format, the resulting text can be processed to create objects that can be further manipulated. This is particularly useful when the output of the LLM serves as an intermediate step in a multi-step pipeline, as it becomes crucial to maintain a high level of control over the output format. This control ensures seamless integration between different tasks within the pipeline.

Many frameworks like Langchain Output Parsers can be used to automate this process, reducing inconsistencies and errors, ensuring structured outputs that are easier to integrate into complex workflows.

```
USER: Given the following customer support inquiry, classify it
into one of the categories: Billing, Technical Support, or Account
Issues.

INQUIRY: I cannot access my account after changing my password.

PLEASE, ALWAYS express your response as a JSON object.
```

```
The output should be formatted as a JSON instance that conforms to
the JSON schema below.

Here is the output schema:

```
{"properties": {"inquiry": {"description": "the inquiry
description", "title": "Inquiry", "type": "string"}, "classification":
{"description": "either Billing, Technical Support, or Account
Issues", "title": "Classification", "type": "string"}}, "required":
["inquiry", "classification"]}

```
```

Figure [20] - Prompt asking for a specific JSON format.

```
ASSISTANT:

{
   "inquiry": "I cannot access my account after changing my
password.",
   "classification": "Account Issues"
}
```

Figure [21] - JSON output format.

# Retrieval Augmented Generation

While **Prompt Engineering** can significantly improve the output of LLMs, some domain-specific applications are not fully addressed by this technique alone. For tasks requiring specific knowledge or highly factual information, prompt engineering may still produce incomplete or inaccurate results. This is where **Retrieval-Augmented Generation (RAG)** becomes valuable; by combining LLMs with external knowledge sources such as databases or documents, RAG allows the application to retrieve relevant, real-time information to supplement its responses. This approach enhances the accuracy and relevance of the generated content, especially in fields like legal research, medical diagnosis, or marketing, where access to up-to-date and precise data is critical.

## RAG Pipeline Overview

The general architecture of a RAG system combines an LLM with an external retrieval engine that is used to incorporate additional information into the generation process. This architecture consists of an input prompt or query from an external user and a knowledge base, which can be a database or a collection of documents. In this knowledge base, each document includes an encoded vector representation of the text generated by models called embedders, which are able to numerically capture the semantic content in the text.

When the user inputs a query, the text is also encoded using the same embedder, and a vector representation of the input is generated. This embedding is then used to retrieve the pieces of text that are most relevant to the input query, based on vector similarity techniques.  Once the relevant information is retrieved, it can be integrated into the language model's generation process.

Together, these components enable the RAG system to retrieve relevant knowledge in real-time, augmenting the generation process and ensuring the output is both contextually appropriate and grounded in factual data
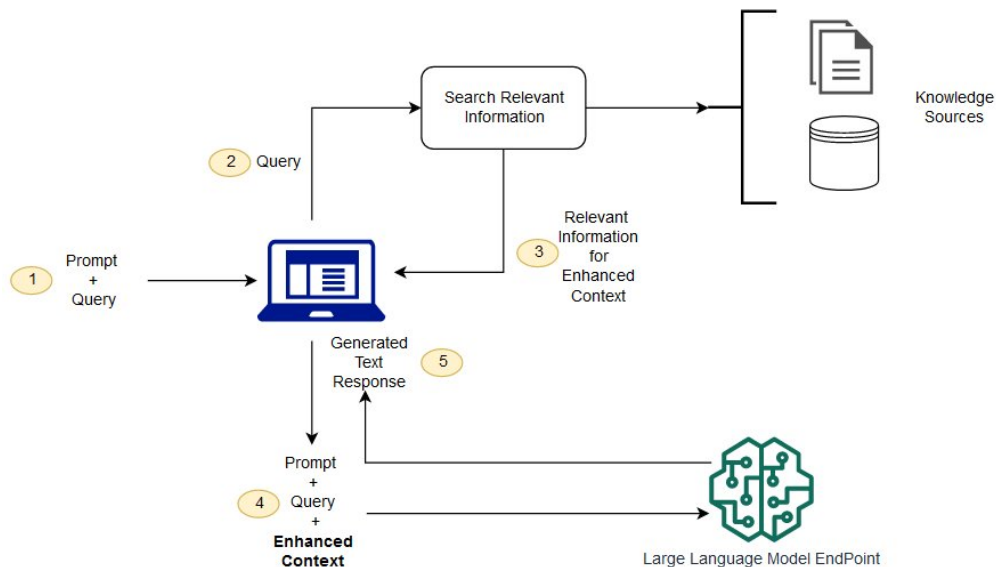
- Retrieval Augmented Generation architecture.

# Corpus Preprocessing

In order to leverage a RAG architecture for answering context-aware questions based on a textual corpus source, it is necessary to adapt  this corpus, undergoing different pre-processing steps that convert the original text into pieces of information that are suitable for the described architecture. Several pre-processing steps can be included as part of this pipeline, but the main two are the following:

- **Text chunking:** it is the process of splitting large texts into smaller, manageable pieces. Since embedding models have a limited context size window (the number of tokens they can process), longer documents may exceed the model's capacity to capture the full context effectively. Due to this, it is necessary to divide the cleaned textual information further into smaller parts, or chunks. The chunking needs to be addressed carefully as it may result in the loss of long-range dependencies within the text. To this end, it is important to preserve the original document's structural parts, such as titles and sections, whenever possible.
- **Chunks embedding:** the generated chunks are then fed to an embedder model which retrieves a vector representation of each chunk.  Embeddings are able to capture the context, semantics, and relationships between words,

making them optimal for retrieving relevant context based on a vector similarity search.

Once the text chunks are embedded, they are typically stored in a vector database, such as **Milvus** or **Weviate**, which provide optimized vector search capabilities on large amounts of data.

# Generation with Retrieval

Once the most relevant text chunks are obtained from the retriever, the generation step takes place. In this stage, context is provided as part of the prompt, together with the initial query, following a pattern like the following:

```
USER: Give a SHORT answer to the question based ONLY on the text
provided in the CONTEXT. Provide and answer in a concise way. If
there is not enough information to answer the question, say that
there is not enough information.

QUESTION: ```{user_question}```

CONTEXT:

```

{retrieved_source_1}
{retrieved_source_2}

```
{retrieved_source_N}

```
```

Figure [23] - RAG prompt example.

# Challenges and Limitations

Although RAG systems are a standard in LLM-based applications nowadays, they also come with some challenges and limitations;

- **Dependency on retrieval quality:** If the retriever fetches irrelevant or incorrect information, the generated output can be misleading or inaccurate.

This is highly linked to the quality of the embedder being used, and therefore selecting a proper embedder for the application becomes crucial.

- **Responses rigidity:** RAG systems can become overly rigid when responding to general or open-ended questions, as they heavily rely on retrieved content. Balancing retrieval and generation is essential to improve flexibility and response quality.
- **Computational Overhead:** The dual process of retrieval and generation can be computationally expensive. The retriever must search through large datasets, and the generator must then process the retrieved context alongside the original query.
- **Context Length Limitation:** Most LLMs have a fixed context window, which limits how much retrieved information can be fed into the generator at once. Selecting an LLM with the appropriate context window for the amount of text that needs to be retrieved becomes crucial.

# Improvements to a RAG Architecture

Some of the limitations related to a RAG system can be addressed implementing different strategies that combine prompting and alternative retrieval techniques. In this section a brief comment of some of them is provided:

## Routing

**Query routing** involves adding a prompting step that selects the best query engine candidate among a set of alternatives to process the input question. Based on the input query, different strategies might be used for generating a response:

- In RAG based conversational applications, it might be desired that the system answers casual questions, such as greetings, without the use of the RAG system itself, and undergo RAG only when the input query is related to domain-specific knowledge.
- When more than one knowledge base is used, query routing might be implemented for selecting the most suitable base for a given question. For instance, in an application that uses a knowledge base about Python coding, and another about SQL coding, the retrieval should be executed on one of these bases, depending on the question.

Figure [24] - Query routing example.

## Condensation

RAG systems can be instrumented with conversational capabilities by providing the entire past conversation as context. Due to context size limitations, long conversations might not be suitable for this technique. To address this problem it is possible to add an intermediate step that condenses the past conversation into a standalone question, reducing the amount of input tokens to the RAG prompt and ensuring conversational-like responses.

## Query Decomposition

**Query decomposition** refers to breaking down a complex or multi-part query into smaller, more manageable subqueries. Each subquery is then processed independently by the RAG system, and the generated outputs are then combined to generate a cohesive response. For instance, questions like "What were the causes of the 2008 financial crisis, and how did governments respond to it?" that involve multiple sub-questions, are better treated if two separated questions are passed to the RAG system.

## Step-back Prompting

**Step-back prompting** is a technique that involves an initial prompt instructing the LLM to identify high-level concepts and principles from the input query. This allows the model to reformulate the question in a more general context. The revised question is then used in the retrieval stage, usually enhancing the quality of the retrieved information. Finally, the original question is employed during the augmented generation step to produce a more accurate and contextually relevant response.

## Contextual Retrieval

This is a technique proposed by Anthropic, that involves enhancing the generated chunks by providing additional contextual information. Traditional RAG systems might struggle when individual chunks within the knowledge base lack enough context or information, for instance a chunk that contains "It is expected to see a growth of 70% on sales over the next two years", might lack information about which company is it talking about, which is the base year, and what is the sales baseline. **Contextual chunking** involves adding a stage in the pre-processing pipeline where individual chunks are fed, together with a bigger portion of the text from which the chunks were taken, and generating a context-aware chunk that can be used for improved retrieval techniques.



Figure [25] - Contextual Retrieval example.

# RAG Using Knowledge Graphs (Graph RAG)

A **Knowledge Graph (KG)** is a structured representation of information that models real-world entities (nodes) and their relationships (edges). Its goal is to facilitate semantic understanding and reasoning by connecting related information in a meaningful and intuitive manner.



Figure [26] - Knowledge Graph example.

The incorporation of a **Knowledge Graph** to the original **RAG** flow, also known as **Graph RAG**, is a novel technique that has been gaining traction in the industry. The key idea is to leverage the graph structure, using the stored relationships to enrich the context provided to the LLM during the retrieval process. The effectiveness and success of this approach heavily depend on constructing a well-suited graph tailored to the specifics of each use case.

Figure [27] - Graph RAG pipeline.

A simple example on how **Graph RAG** actually works might be as follows. Imagine a database of documents, where each document is divided into smaller chunks. The graph is built by **mapping each text chunk to a node**, and additional nodes, such as document and author nodes, are introduced to create a **hierarchical structure**. The nodes are linked based on different attributes, such as sequence or ownership, which are carefully selected to meet the specific requirements of each use case.



Figure [28] - Building a KG for RAG applications.

The resulting graph is then queried to enhance the retriever's responses in the following way. Given a user query, the system first retrieves the chunk that best matches the query using vector similarity search. Then, it runs a graph-based query (using GraphQL or Cypher) to retrieve related nodes connected to the resulting chunk. The additional information is added to the original retrieved chunk to further enrich the context provided to the LLM, improving the quality and relevance of its final response.

LangChain already offers modules compatible with Neo4J, which can serve both as a graph database and a vector store for the corresponding text embeddings. Additionally,  it is recommended visiting the DeepLearning.AI Knowledge Graphs for RAG course for deeper insights into these techniques.

# Agents

**AI agents** are autonomous systems built to carry out tasks and make decisions on their own. They adapt to changing environments and can plan actions dynamically, rather than relying on fixed, predefined steps. Typically, an LLM serves as their reasoning engine, guiding them in determining which actions to take and in what sequence. This allows them to think and reason through complex problems and respond more flexibly to different situations.

Unlike traditional AI models that depend on hardcoded instructions, AI agents manage multi-step tasks with minimal human supervision. They can use various tools and interact with their surroundings to complete complex goals. This flexibility and adaptability make them highly effective in scenarios that require ongoing adjustments and the ability to operate independently.

## What Does It Mean to Be Agentic?

A system is more "agentic" the more the system itself determines how it should behave. This is not directly linked to human supervision or intervention. As an example, a predefined set of rules could involve zero human intervention, but the behavior of the system is strictly defined.

The following chart shows different levels of autonomy in LLM applications.

Figure [29] - Levels of autonomy in LLM applications by LangChain.

Writing **code** without any calls to LLMs is the most human-driven approach, thus the less agentic one. By using **LLM calls** or even **chaining**, the system progressively becomes more agentic, although the steps to take and the steps available to take are defined by the developer.

There is also the possibility to **route** inputs into a particular downstream workflow, which leads to a more agentic solution. In this case, the system decides which steps should be taken, although the steps available are still set by the developer.

Adding cycles to the steps to take would lead into the **State Machine** level of autonomy, which is agent-executed and not human-driven. In this case, the available steps to take are still not defined by the LLM.

The last level of autonomy is achieved in a system where the LLM decides the available steps to take. This level is the **Autonomous** one. At the end of this section an example will be shown.

The following subsections dive into some of the most important concepts about AI Agents: tool usage, reflection, planning and multiagent collaboration.

# Tools Usage

Advanced AI Agents can leverage various tools, such as code execution, data retrieval, and API calls, to complete tasks more effectively. These agents typically utilize tools through **function calling**, where they reason about which tool to use and how to structure the call. This reasoning is essential to their ability to interact with their environment, enabling them to perform actions like executing code, retrieving data, or interacting with external systems. The effectiveness of the agent's reasoning lies in the proper execution of tool calls rather than reflecting on the outcomes.

While LLMs can suggest which tool to use, they cannot directly execute the tools themselves. Instead, they express the intent to call a tool in their response and provide the arguments needed for using those tools as intended. Developers are responsible for executing the tool with the provided arguments and reporting the results back to the model. The overall process can be broken down into six steps:

1. **Pick a tool** (or function) in your codebase that the model should be able to call.
2. **Describe the tool** to the model so it knows how to use it.
3. **Pass the tool definition** as an available tool to the model, along with the user's messages.
4. **Handle the model's response:**
   a. If no tool is needed, the model will provide a direct response.
   b. If a tool is needed, the model will generate the appropriate function call with arguments based on the parameters.
5. **Execute the tool call** using the parameters given by the LLM.
6. **Provide the tool's results back to the model** in order to generate the appropriate answer.

This structured framework allows AI agents to effectively integrate and utilize tools to achieve more complex and dynamic goals.

**Figure [30]** - The lifecycle of a function call.

# Reflection

This pattern involves an AI system improving its capabilities through **self-feedback** and **iterative refinement**. In this way, the AI system can enhance the quality and accuracy of its output after generating an initial solution by further reflection and analysis.

Reflexion agents provide better decisions by maintaining their reflective texts in an event memory buffer. Reflexion can accept various types and sources of feedback signals and shows significant improvements in multiple tasks.

Figure [31] - Example of reflection within a system, adapted from Andrew Ng's talk about AI Agents.

# Planning

**Planning** involves training language models to reason, devise strategies, and break down complex tasks into manageable parts. This equips language models with the ability not only to answer questions but also to proactively develop and execute action plans. With planning capabilities, language models can autonomously decompose tasks, identify necessary substeps and tools, and coordinate the use of different models.

Below is an example of a prompt that could be used in the planning phase for an autonomous AI agent.

```
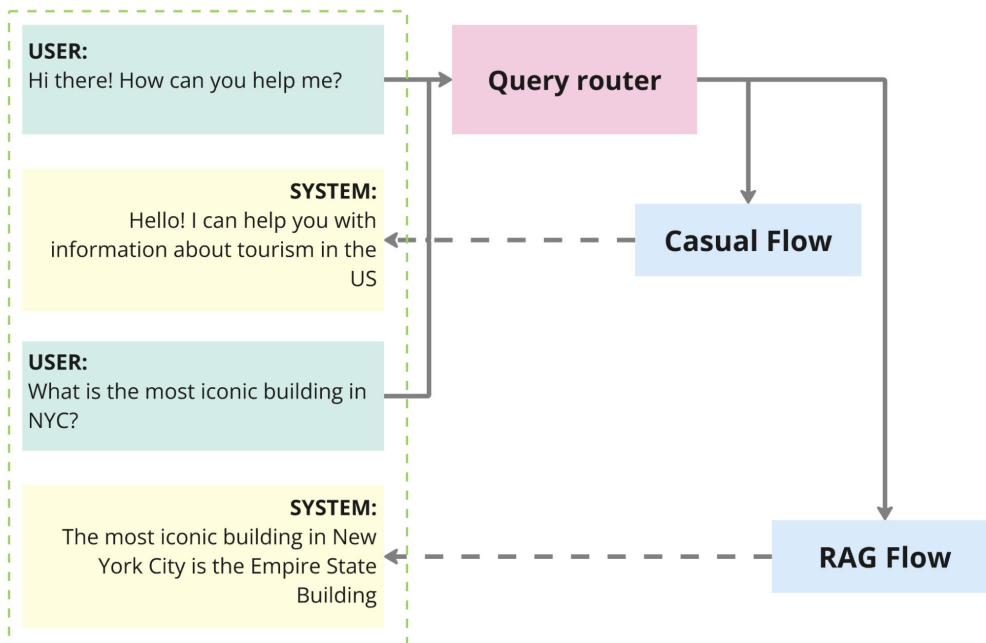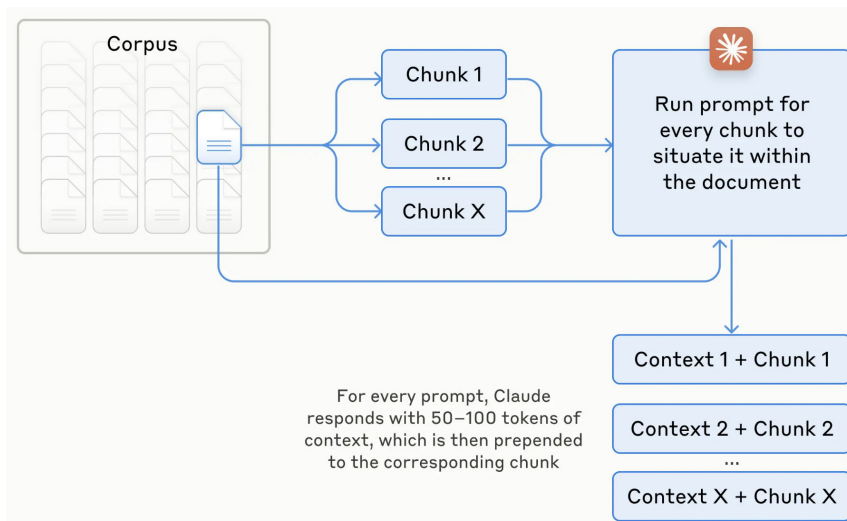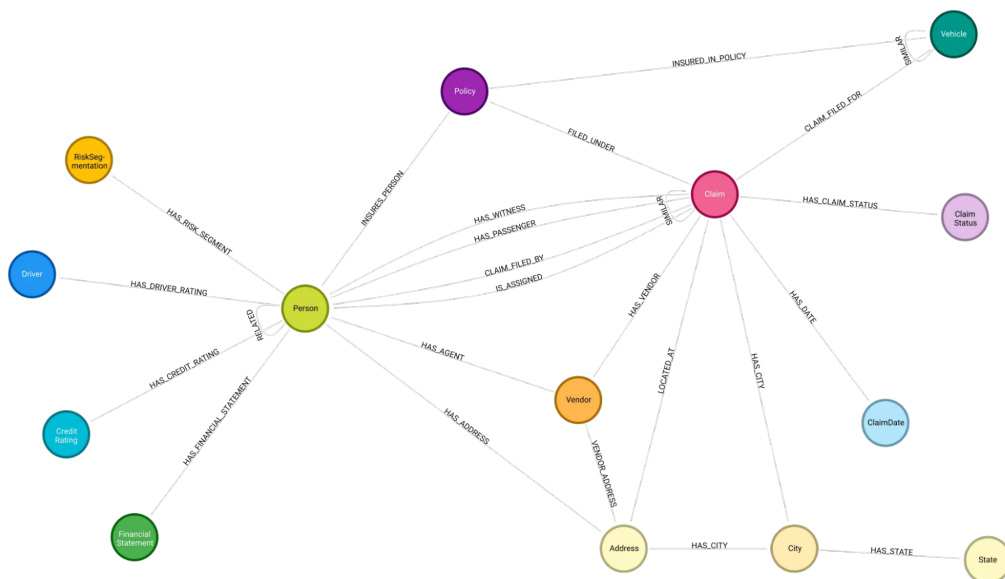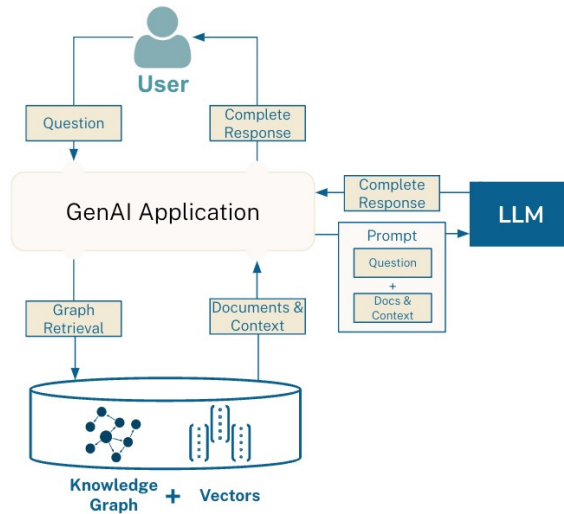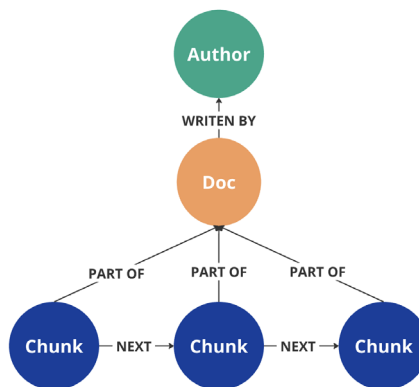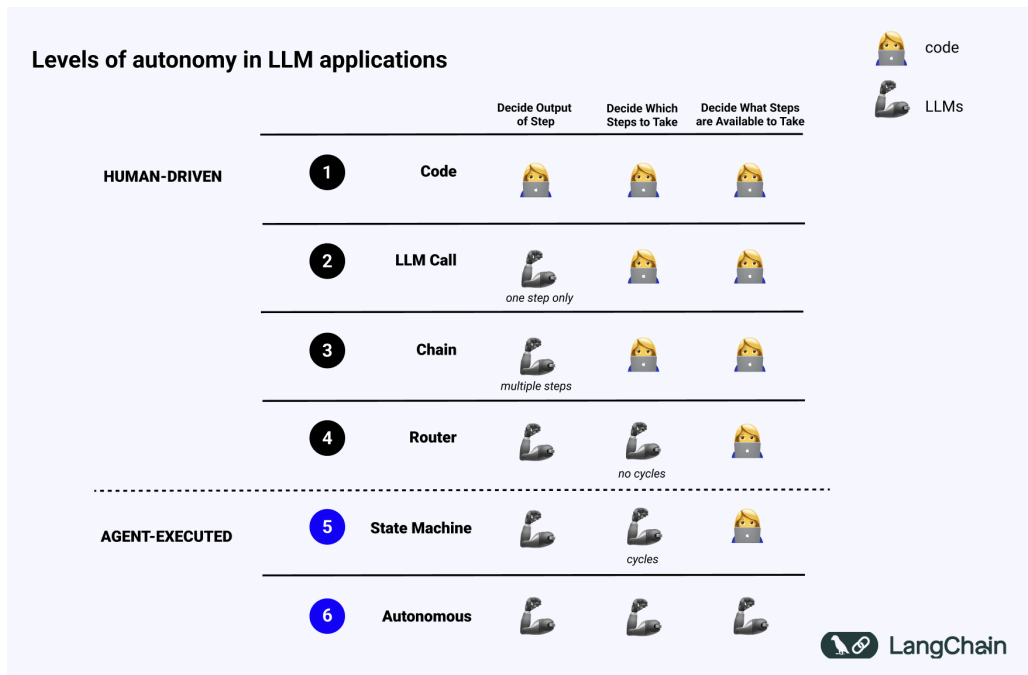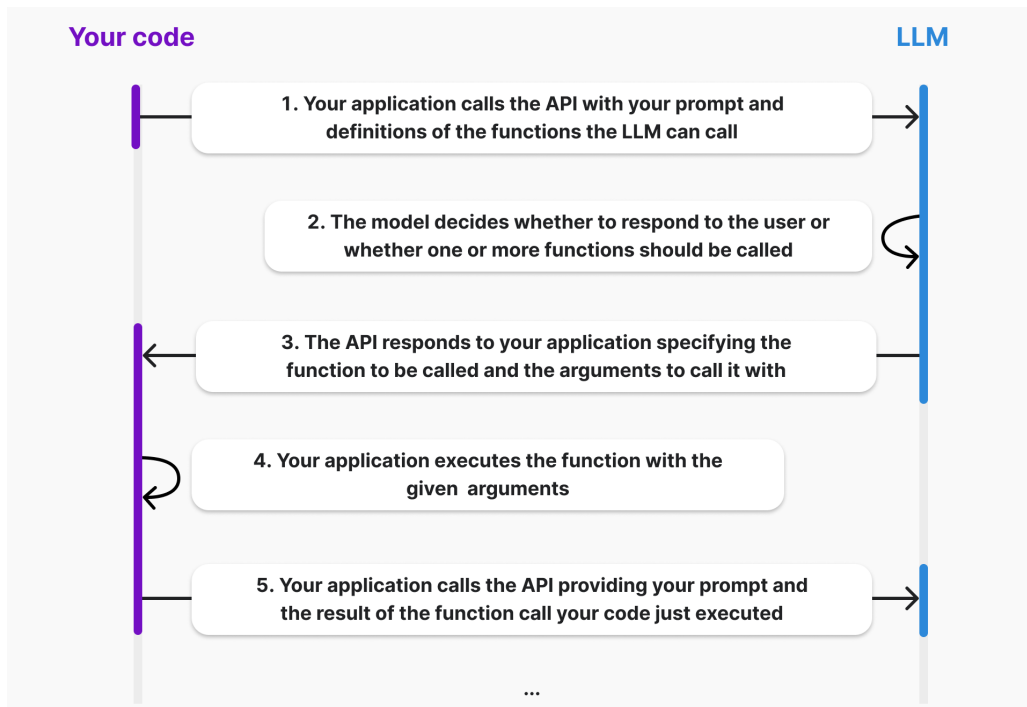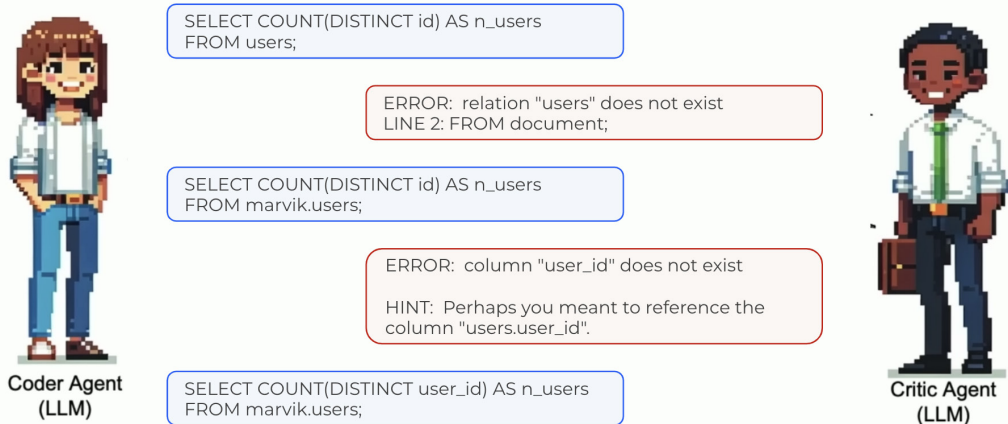You excel at planning the steps required for an LLM to be able to
solve the main task mentioned before.

Your task is to propose a series of simple steps that will break
down the complexity of a task, in order for a later process to use
these steps as input.

Based on the current context outlined: ```{context}```

The task to solve is: ```{task_description}```
```

```
Develop a detailed action plan that includes:
- Clear objectives
- Specific steps to achieve these objectives
- Anticipated outcomes

Avoid repeating steps that were already being made. Here's a list of
the steps that were previously taken: ```{previous_steps}```

Follow these instructions:

- Be concise
- Avoid unnecessary steps.
- Do not repeat steps already made.

Express your response as a JSON object as follows:

```json
{
  "next_steps": [
    {
      "step": "short description of the step (action) to take"
"details": "details of the step (action) to take in order to
provider more context and being more precise"
"expected_outcome": "details of the expected outcome, for being able
to evaluate after its success"
    },
    ... # generate as many steps as needed
  ]
}
```
```

Figure [32] - Prompt during the planning phase for creating the next steps for the LLM to take in order to solve a certain task.

# Multiagent Collaboration

**Multiagent collaboration** involves having multiple language models or agents interact to complete complex tasks together. For instance, agents simulating experts in different roles (e.g., doctors, nurses) can collaborate to develop diagnostic and treatment plans. The key to this approach is training the agents to work together efficiently, ensuring a clear division of labor to avoid conflicts or contradictions.

# Examples

Let's discuss two examples highlighting the most confusing stages of autonomy levels in LLM applications: the State Machine and Autonomous stages.

## Research Agent

Consider a scenario where a company wants to create an agent that generates research reports based on user queries, with the goal of reducing the hours the team spends on manual research. There are various ways to approach this task. One intuitive method is to first consult with **Subject Matter Experts (SMEs)** to understand their process for gathering information and creating reports.

Once the process is clearly understood, the next step is to model it by replacing human actions with code—either through LLM calls or by using external tools like web searches.

The flow might look something like this:



Figure [33] - Researcher Agent workflow built with LangGraph, based in the ReAct paradigm.

The process involves planning which sections or topics to write about, writing each section (using the LLM and tool-calling like web searches if needed), and then performing a review. The review uses a separate prompt to validate whether the output from the writing phase is satisfactory. If the output isn't adequate, the writing section is revisited with feedback from the review step, and this cycle continues until the section is approved. Once all sections are completed and approved, the final report is generated using another prompt, and the output is delivered to the user.

However, this workflow is predefined and limited to structured report writing. How can be maintained the same level of quality while creating a more flexible process that can handle other use cases?

## Autonomous Agent (Deming Agent)

A more flexible approach is to model the process as a generic workflow that doesn't rely on prompts specifically designed for report writing. While there are numerous theories online about how to achieve this, practical examples are relatively rare. One example is the Proof-of-Concept (PoC) for an autonomous generalist agent, which was developed by Marvik and is based on the **Deming Cycle (PDCA)**, a framework that stands for **Plan, Do, Check, Act**.

The Deming Cycle is a continuous improvement model used for problem-solving and process optimization, consisting of the following four steps:

1. **Plan:** Identify an issue or improvement area and develop a plan to address it.
2. **Do:** Implement the plan on a small scale to test its effectiveness.
3. **Check:** Evaluate the results to determine if the desired improvements were achieved.
4. **Act:** If successful, implement the changes on a larger scale; if not, revise the plan and repeat the cycle.

This same paradigm is reflected in the logic of the agent when planning the next steps needed to fulfill the user's query, which could be much broader than just writing a research report. Each step of the plan follows the PDCA cycle:

1. **Plan:** The agent identifies the next step to take.
2. **Do:** The action is executed, typically as an LLM call, possibly involving tool-calling.
3. **Check:** The output of the action is evaluated. If the result is unsatisfactory, the step is redirected back to execution.
4. **Act:** Based on the outcomes from the current steps and the overall progress

toward the user's query, the system analyzes how close it is to achieving the desired result. It summarizes the progress and determines whether the available information is sufficient to generate the final answer, or if additional steps need to be planned and executed.

Below is the flow created for the Deming Agent.



Figure [34] - Autonomous Agent (Deming Agent) workflow built with LangGraph, based on concepts of the Deming Cycle.

# References

- Alammar, Jay. "The Illustrated Transformer." June 27, 2018. https://jalammar.github.io/illustrated-transformer/

- Amazon Web Services. "What Are AI Agents?" https://aws.amazon.com/what-is/ai-agents/

- Anthropic. "Introducing Contextual Retrieval." September 19, 2024. https://www.anthropic.com/news/contextual-retrieval

- Damian Gil. "Advanced Retriever Techniques to Improve Your RAGs." April 17, 2024. https://towardsdatascience.com/advanced-retriever-techniques-to-improve-your-rags-1fac2b86dd61#b2c7

- Geetansh Kalra. "Attention Networks: A Simple Way to Understand Self-Attention." June 5, 2022. https://medium.com/@geetkal67/attention-networks-a-simple-way-to-understand-self-attention-f5fb363c736d

- IBM. "What Is Prompt Engineering?" https://www.ibm.com/topics/prompt-engineering

- LangChain. "Introduction." https://python.langchain.com/v0.1/docs/get_started/introduction/

- LangChain. "What Is an Agent?" June 28, 2024. https://blog.langchain.dev/what-is-an-agent/

- Liu, Yang. "Fine-tune BERT for Extractive Summarization." September 5, 2019. https://arxiv.org/abs/1903.10318

- Malec, Melissa. "Open-source LLMs vs Closed: Unbiased 2024 Guide for Innovative Companies." May 31, 2024. https://hatchworks.com/blog/gen-ai/open-source-vs-closed-llms-guide/

- Mistral. "Mixtral of Experts." December 11, 2023. https://mistral.ai/news/mixtral-of-experts/

- OpenAI. "Function Calling." June 13, 2023. https://platform.openai.com/docs/guides/function-calling

- Philip Rathle. "The GraphRAG Manifesto: Adding Knowledge to GenAI." July 11, 2024. https://neo4j.com/blog/graphrag-manifesto/

- Prompt Engineering Guide. "Prompting Techniques." https://www.promptingguide.ai/techniques

- Singh, Tarun. "Advanced RAG Techniques: Unlocking the Next Level." April 17, 2024. https://ai.gopubby.com/advanced-rag-techniques-unlocking-the-next-level-040c205b95bc

- Winston, Aaron. "What Are AI Agents and Why Do They Matter?" August 13, 2024. https://github.blog/ai-and-ml/generative-ai/what-are-ai-agents-and-why-do-they-matter/

- Yao et al. "ReAct: Synergizing Reasoning and Acting in Language Models." https://arxiv.org/abs/2210.03629

# marvik.

**AI solutions with business impact**

Got a project? growth@marvik.ai

More information www.marvik.ai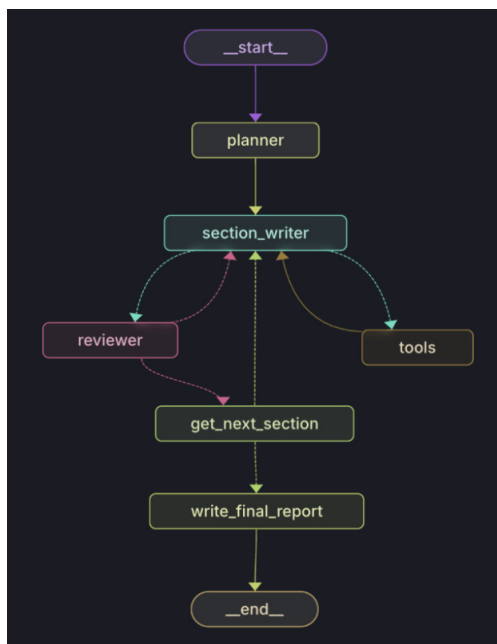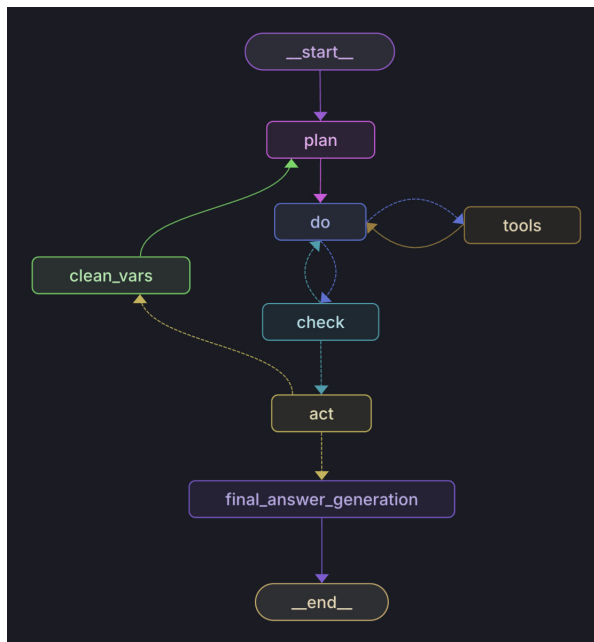