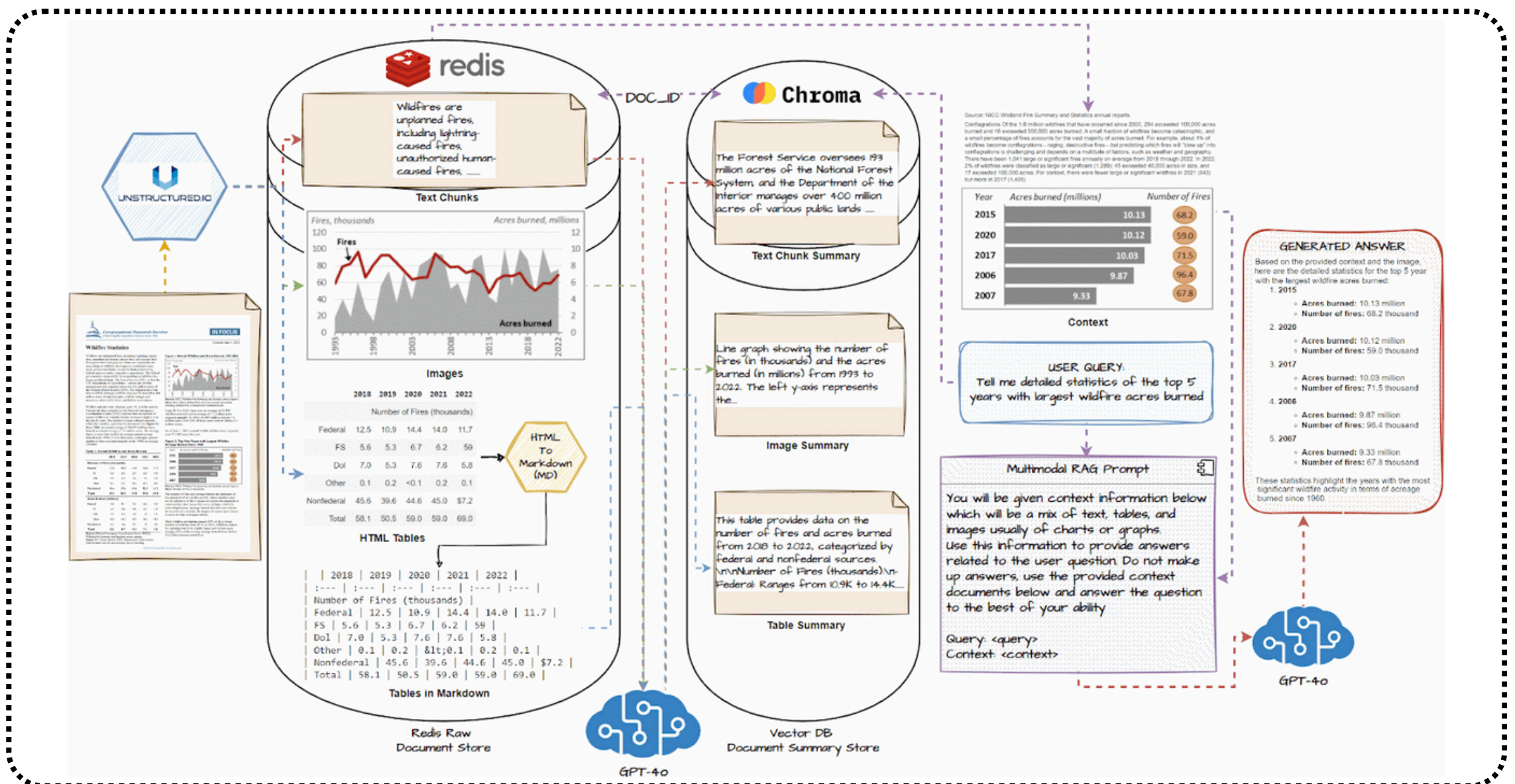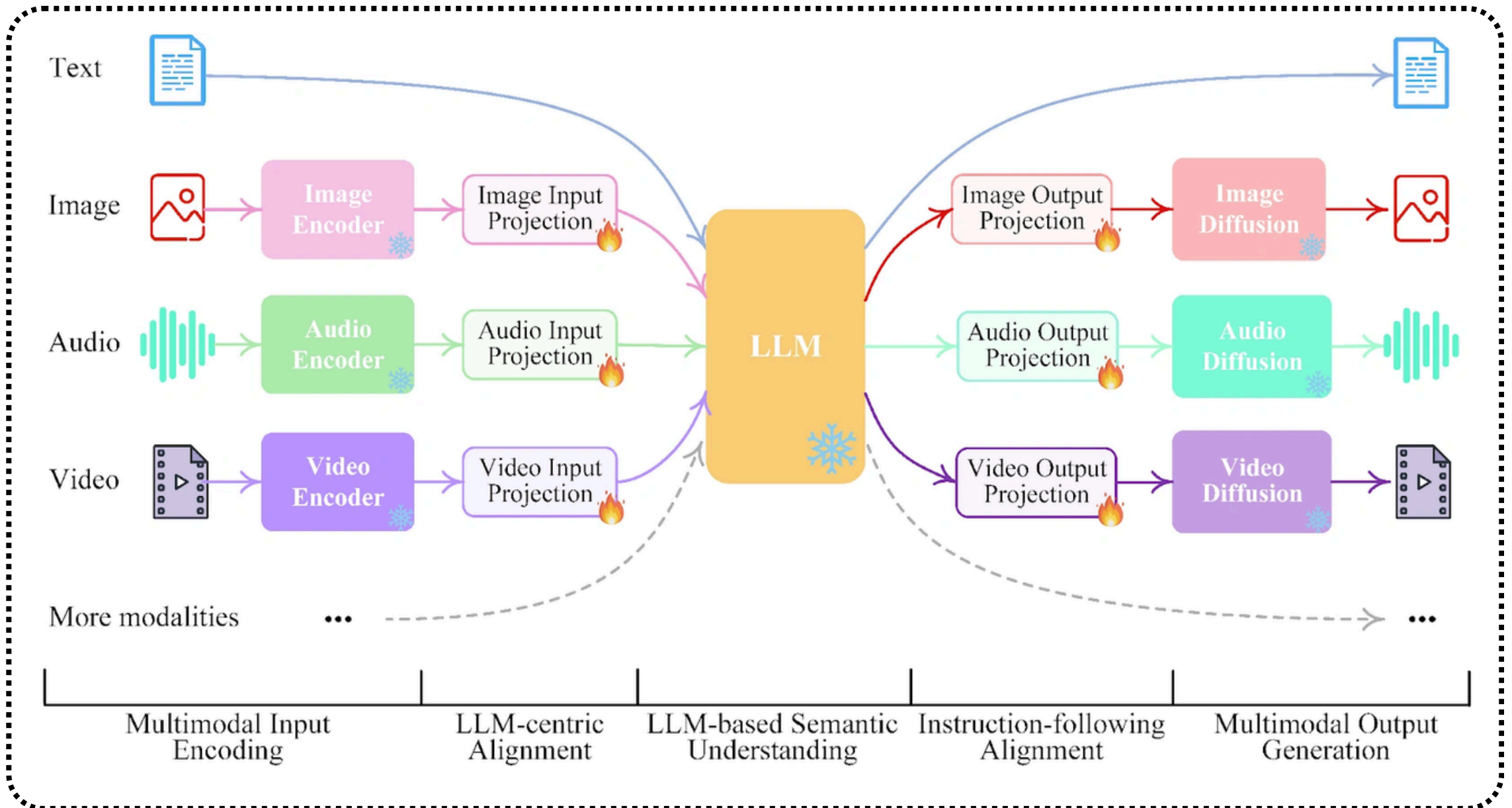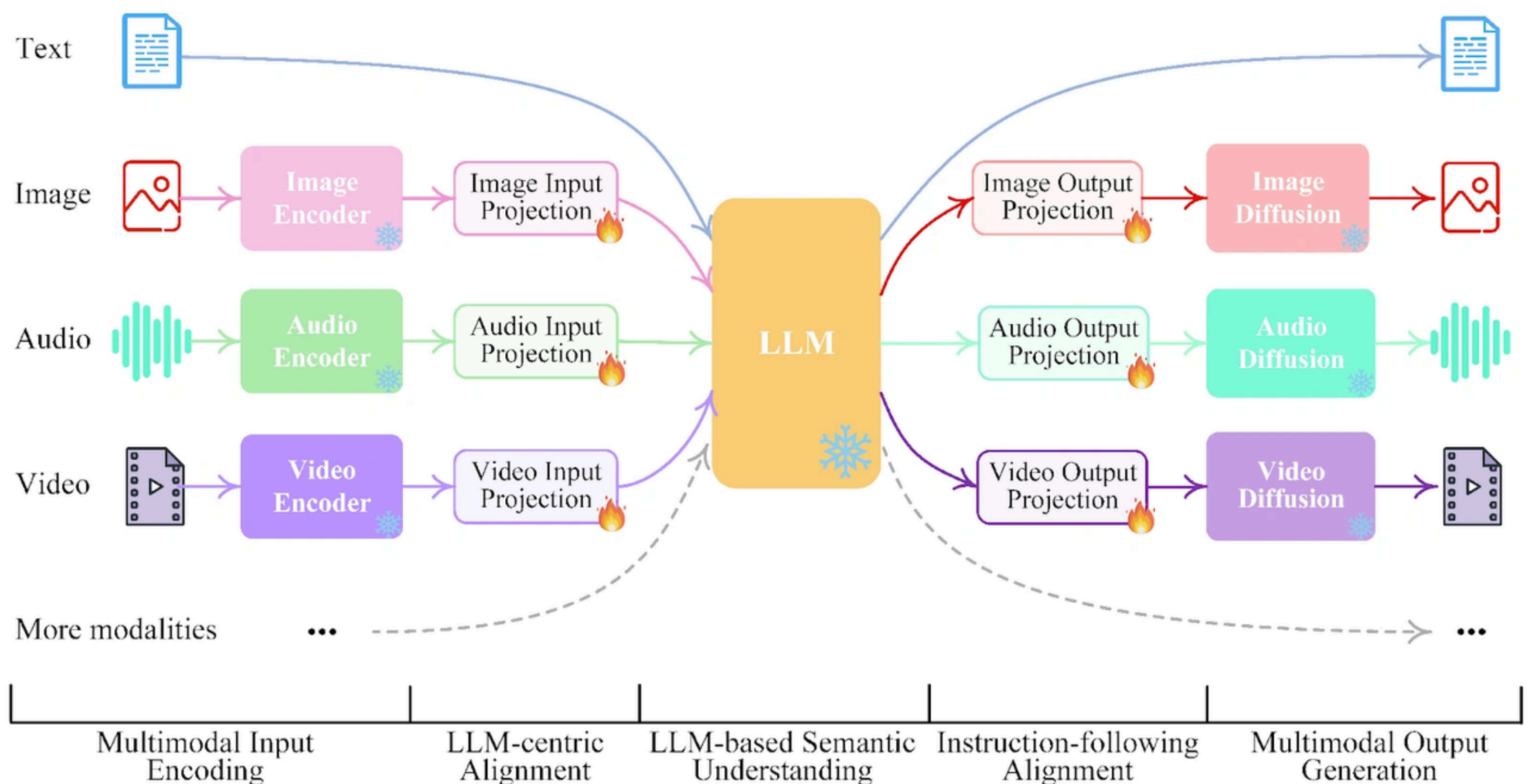# A Guide to Multimodal RAG System

# What is a Multimodal Large Language Model?

- Multimodal Large Language Models (LLMs) are essentially transformer-based LLMs that have been pre-trained and fine-tuned on multimodal data to analyze and understand various data formats, including text, images, tables, audio, and video.

- A true multimodal model ideally should be able not just to understand mixed data formats but also generate the same as shown in the following workflow illustration of NExT-GPT, published as a paper, **NExT-GPT: Any-to-Any Multimodal Large Language Model**

From the paper on NExT-GPT, any true multimodal model would typically have the following stages:

- Multimodal Encoding Stage. Leveraging existing well-established models to encode inputs of various modalities.

- LLM Understanding and Reasoning Stage. An LLM is used as the core agent of NExT-GPT. Technically, they employ the **Vicuna LLM** which takes as input the representations from different modalities and carries out semantic understanding and reasoning over the inputs. It outputs 1) the textual responses directly and 2) signal tokens of each modality that serve as instructions to dictate the decoding layers whether to generate multimodal content and what content to produce if yes.

- Multimodal Generation Stage. Receiving the multimodal signals with specific instructions from LLM (if any), the Transformer-based output projection layers map the signal token representations into the ones that are understandable to following multimodal decoders.
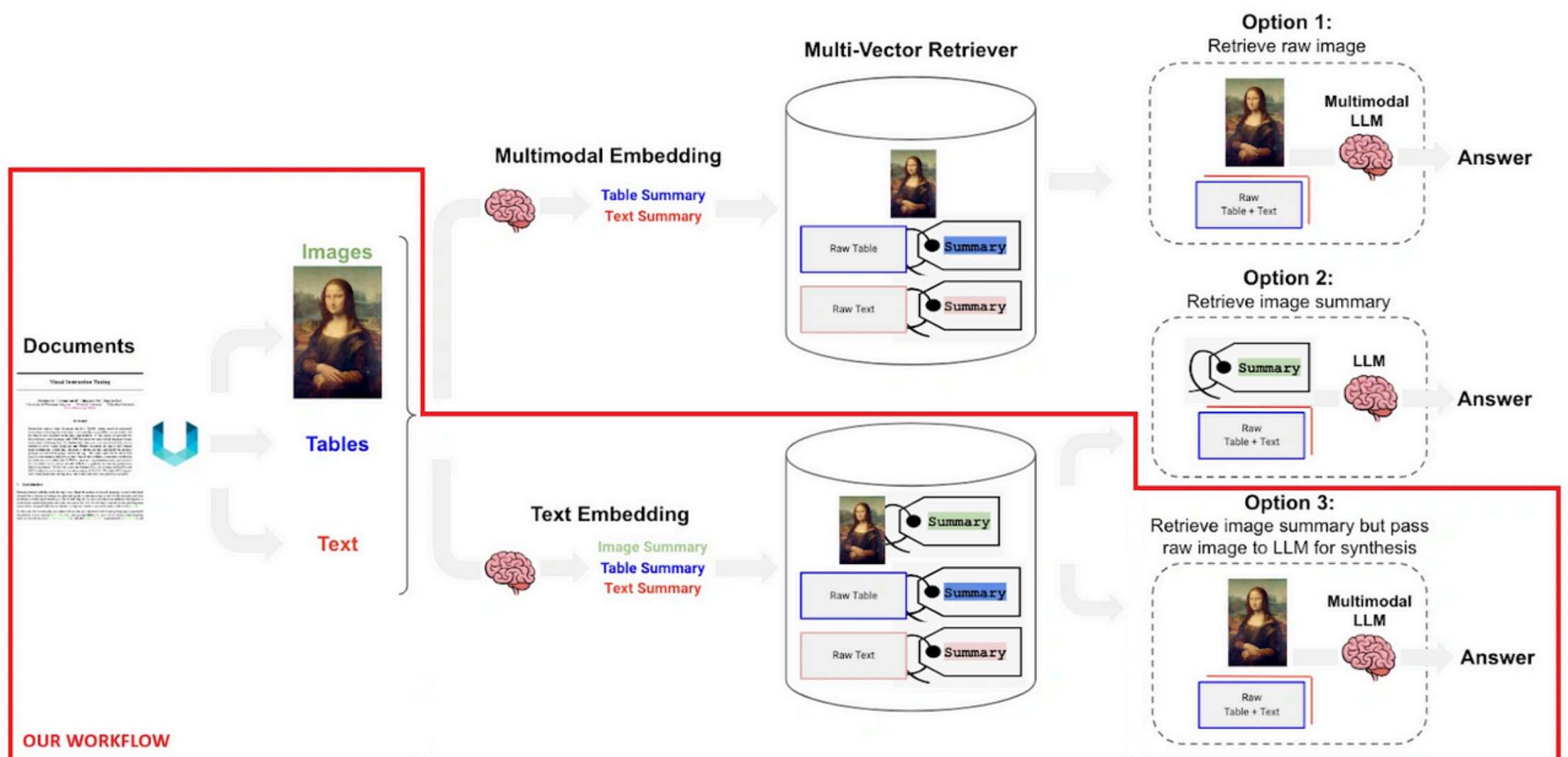
However, most current Multimodal LLMs available for practical use are one-sided, which means they can understand mixed data formats but only generate text responses. The most popular commercial multimodal models are as follows:

- **GPT-4V & GPT-4o (OpenAI)**: GPT-4o can understand text, images, audio, and video, although audio and video analysis are still not open to the public.

- **Gemini (Google)**: A multimodal LLM from Google with true multimodal capabilities where it can understand text, audio, video, and images.

- **Claude (Anthropic)**: A highly capable commercial LLM that includes multimodal capabilities in its latest versions, such as handling text and image inputs.

For our Multimodal RAG System, we will leverage **GPT-4o**, one of the most powerful multimodal models currently available.

# Multimodal RAG System Workflow

- In this section, we will explore potential ways to build the architecture and workflow of a multimodal RAG system. The following figure illustrates potential approaches in detail and highlights the one we will use in this guide.

# End-to-End Workflow

Multimodal RAG Systems can be implemented in various ways, the above figure illustrates three possible workflows as recommended in the <u>LangChain blog</u>, this include:

- Option 1: Use multimodal embeddings (such as <u>CLIP</u>) to embed images and text together. Retrieve either using similarity search, but simply link to images in a docstore. Pass raw images and text chunks to a multimodal LLM for synthesis.

- Option 2: Use a multimodal LLM (such as <u>GPT-4o</u>, <u>GPT4-V</u>, <u>LLaVA</u>) to produce text summaries from images. Embed and retrieve text summaries using a text embedding model. Again, reference raw text chunks or tables from a docstore for answer synthesis by a regular LLM; in this case, we exclude images from the docstore.

- Option 3: Use a multimodal LLM (such as GPT-4o, GPT4-V, LLaVA) to produce text, table and image summaries (text chunk summaries are optional). Embed and retrieved text, table, and image summaries with reference to the raw elements, as we did above in option 1. Again, raw images, tables, and text chunks will be passed to a multimodal LLM for answer synthesis. This option is sensible if we don't want to use multimodal embeddings, which don't work well when working with images that are more charts and visuals. However, we can also use multimodal embedding models here to embed images and summary descriptions together if necessary.

There are limitations in Option 1 as we cannot use images, which are charts and visuals, which is often the case with a lot of documents. The reason is that multimodal embedding models can't often encode granular information like numbers in these visual images and compress them into meaningful embedding. Option 2 is severely limited because we do not end up using images at all in this system even if it might contain valuable information and it is not truly a multimodal RAG system.

# For more information, visit this <u>article</u>



Advanced    Best of Tech    Generative AI    Guide    Large Language Models

## A Comprehensive Guide to Building Multimodal RAG Systems

Build a Multimodal RAG system to handle text, images & tables, overcoming traditional RAG limita

*Dipanjan (DJ) Sarkar*    08 Jan, 2025