# Chat with documents – RAG Assistant

Alejandra Caggiano

We have seen that GeneXus Enterprise AI allows creating different types of Artificial Intelligence assistants. In particular, we have already succeeded in creating an interactive conversation.

Now we want to be able to define an assistant that allows us to chat with documents. For that, we are going to work with RAG Assistants.

# Retrieval Augmented Generation (RAG)

Retrieval Augmented Generation (RAG) is an approach that combines information retrieval from unstructured data and text generation to improve performance in tasks such as question answering.

# Retrieval Augmented Generation (RAG)

➢ Data Ingestion

➢ Retrieval

➢ Generation

➢ End user interaction

This process consists of the following four phases:

• The first phase is **data ingestion**: It involves loading various types of documents, in different formats and from multiple sources.

• It is followed by the **Retrieval phase:** In this stage, the data retrieval process is started, taking advantage of the previously loaded and organized information. A selective search is performed on a set of documents, identifying related information and efficiently reducing the search space. This approach ensures that focus is placed on the most relevant and meaningful information.

• The next phase is **Generation**: It focuses on generating relevant and contextually consistent answers. In this process, the system uses the RAG Assistant configuration to know which model to access and with which parameters. This assistant incorporates the necessary elements to define the search strategy, and to achieve consistency and relevance in the generated content.

• The last phase is **Interaction with the end user**: GeneXus Enterprise AI facilitates smooth and efficient communication between end users and RAG Assistants, completing the cycle and providing answers to queries in an efficient manner.

## Retrieval Augmented Generation (RAG)

Example:

Chat with GeneXus Training documents

Good. As we said at the beginning, our goal now is to create an assistant that allows us to chat with a set of documents, and we will do it with GeneXus Training documents.
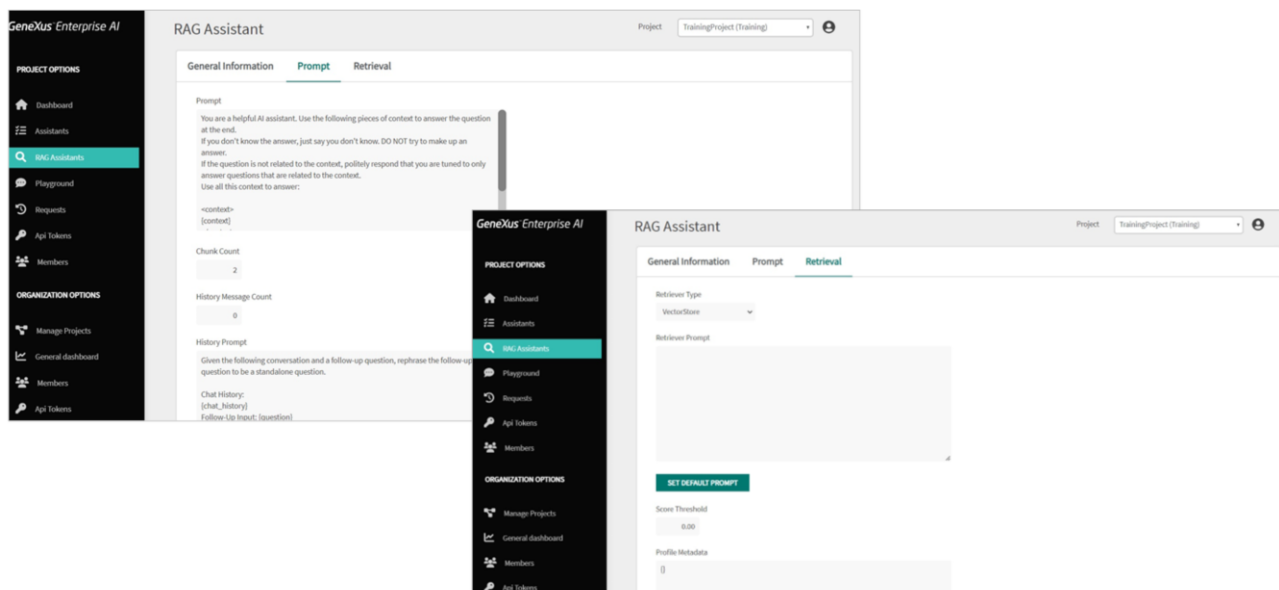
# RAG Assistants



So we enter the platform, select the project we are going to work on, and in the menu we choose RAG Assistant. By default, when opening this section we see a RAG Assistant named Default, which can be customized, or we can create new ones. We click on Create new, name it ChatWithGXTraining, and enter a brief description. We Confirm. From the Update option, we can customize the assistant definition as needed. The settings are organized in the following tabs:

# RAG Assistants



General Information, Prompt, containing instructions to guide the assistant on how to approach and answer questions. These instructions establish clear guidelines for the assistant to provide relevant and useful answers based on the context provided.
This option indicates the number of fragments that are retrieved to expand the context.

Then this Message History option sets the number of historical messages that are taken into account in the conversation. It is useful for tracking the interaction history and understanding the context gathered in the conversation.

If this value is set to 4, it means that the last 4 interactions will be considered. The minimum value it can take is 0, which indicates that the conversation history is of no interest.   When the value is greater than 0, it is used together with the message indicated in the following History Prompt option.

# RAG Assistants



Next, we see the options for configuring the model used by the assistant to generate the answer. This includes the service provider, model name, temperature, maximum token limit, and other parameters that affect how answers are generated.

Lastly, the Retrieval tab specifies how the information is retrieved. We leave the default values.
OK. Now let's load the documents. We click on Add Documents. The Add Files button allows loading files of various formats:. .txt, .pdf, .docx, .pptx, .xlsx, .odt, .odp, .ods, .xlsx, .epub, .json, .jsonl, and .csv.

In our example, we are going to load a small set of PDF and Word documents corresponding to GeneXus Training material. We click on Add Files.

# RAG Assistants



Once the documents have been added, we click on Start upload. To view the uploaded documents, each with its details, we click on View Documents.

# RAG Assistants



By selecting the name of the file we see all the associated information, and we can view and download the file from the URL. Good. We have already created our RAG Assistant and uploaded the corresponding files. We are now ready to test it.

We will do it next from the Playground option of the menu.

GeneXus™
by Globant

training.genexus.com