

Table of contents

[The Real Foundation of RAG Chatbots](#)[A GenAI Architecture for a Conversational Agent](#)[Step 1: Data Ingestion](#)[Step 2: Knowledge Base Creation](#)[Step 3: Document Retrieval](#)[Step 4: LLM Reasoning](#)[Step 5: User Interface](#)[How to Know if the RAG Chatbot is Delivering on its Promise?](#)[The Power of Thoughtful RAG Implementation](#)[Share](#)

Each time Generative AI models like [GPT-4](#), [Claude](#), or [Gemini](#) hallucinate or give inaccurate answers, the world questions whether these LLM-powered chatbots can truly deliver value.

Even so, building LLM chatbots is one of the most popular GenAI applications. To address some of the LLMs' limitations, companies have turned to Retrieval Augmented Generation (RAG), a framework that ensures access to accurate information by retrieving it from multiple sources.

THE REAL FOUNDATION OF RAG CHATBOTS

Most companies find RAG an easy, exciting, and one of the quickest fixes to customer and employee experience with knowledge retrieval. And that's precisely where **two of the biggest misconceptions** about building a RAG chatbot hide: assuming that it's heavily based on LLMs and that it simply works well.

Knowledge retrieval is about knowing where to find the information you need. This is a key step in developing RAG chatbots before involving a large language model. The LLM is only one part of the RAG architecture, primarily involved in reasoning and generating

Table of contents

[The Real Foundation of RAG Chatbots](#)[A GenAI Architecture for a Conversational Agent](#)[Step 1: Data Ingestion](#)[Step 2: Knowledge Base Creation](#)[Step 3: Document Retrieval](#)[Step 4: LLM Reasoning](#)[Step 5: User Interface](#)[How to Know if the RAG Chatbot is Delivering on its Promise?](#)[The Power of Thoughtful RAG Implementation](#)[Share](#)

Each time Generative AI models like [GPT-4](#), [Claude](#), or [Gemini](#) hallucinate or give inaccurate answers, the world questions whether these LLM-powered chatbots can truly deliver value.

Even so, building LLM chatbots is one of the most popular GenAI applications. To address some of the LLMs' limitations, companies have turned to Retrieval Augmented Generation (RAG), a framework that ensures access to accurate information by retrieving it from multiple sources.

THE REAL FOUNDATION OF RAG CHATBOTS

Most companies find RAG an easy, exciting, and one of the quickest fixes to customer and employee experience with knowledge retrieval. And that's precisely where **two of the biggest misconceptions** about building a RAG chatbot hide: assuming that it's heavily based on LLMs and that it simply works well.

Knowledge retrieval is about knowing where to find the information you need. This is a key step in developing RAG chatbots before involving a large language model. The LLM is only one part of the RAG architecture, primarily involved in reasoning and generating

Table of contents

[The Real Foundation of RAG Chatbots](#)[A GenAI Architecture for a Conversational Agent](#)[Step 1: Data Ingestion](#)[Step 2: Knowledge Base Creation](#)[Step 3: Document Retrieval](#)[Step 4: LLM Reasoning](#)[Step 5: User Interface](#)[How to Know if the RAG Chatbot is Delivering on its Promise?](#)[The Power of Thoughtful RAG Implementation](#)[Share](#)

Each time Generative AI models like [GPT-4](#), [Claude](#), or [Gemini](#) hallucinate or give inaccurate answers, the world questions whether these LLM-powered chatbots can truly deliver value.

Even so, building LLM chatbots is one of the most popular GenAI applications. To address some of the LLMs' limitations, companies have turned to Retrieval Augmented Generation (RAG), a framework that ensures access to accurate information by retrieving it from multiple sources.

THE REAL FOUNDATION OF RAG CHATBOTS

Most companies find RAG an easy, exciting, and one of the quickest fixes to customer and employee experience with knowledge retrieval. And that's precisely where **two of the biggest misconceptions** about building a RAG chatbot hide: assuming that it's heavily based on LLMs and that it simply works well.

Knowledge retrieval is about knowing where to find the information you need. This is a key step in developing RAG chatbots before involving a large language model. The LLM is only one part of the RAG architecture, primarily involved in reasoning and generating

Table of contents

[The Real Foundation of RAG Chatbots](#)[A GenAI Architecture for a Conversational Agent](#)[Step 1: Data Ingestion](#)[Step 2: Knowledge Base Creation](#)[Step 3: Document Retrieval](#)[Step 4: LLM Reasoning](#)[Step 5: User Interface](#)[How to Know if the RAG Chatbot is Delivering on its Promise?](#)[The Power of Thoughtful RAG Implementation](#)[Share](#)

Each time Generative AI models like [GPT-4](#), [Claude](#), or [Gemini](#) hallucinate or give inaccurate answers, the world questions whether these LLM-powered chatbots can truly deliver value.

Even so, building LLM chatbots is one of the most popular GenAI applications. To address some of the LLMs' limitations, companies have turned to Retrieval Augmented Generation (RAG), a framework that ensures access to accurate information by retrieving it from multiple sources.

THE REAL FOUNDATION OF RAG CHATBOTS

Most companies find RAG an easy, exciting, and one of the quickest fixes to customer and employee experience with knowledge retrieval. And that's precisely where **two of the biggest misconceptions** about building a RAG chatbot hide: assuming that it's heavily based on LLMs and that it simply works well.

Knowledge retrieval is about knowing where to find the information you need. This is a key step in developing RAG chatbots before involving a large language model. The LLM is only one part of the RAG architecture, primarily involved in reasoning and generating

Table of contents

[The Real Foundation of RAG Chatbots](#)[A GenAI Architecture for a Conversational Agent](#)[Step 1: Data Ingestion](#)[Step 2: Knowledge Base Creation](#)[Step 3: Document Retrieval](#)[Step 4: LLM Reasoning](#)[Step 5: User Interface](#)[How to Know if the RAG Chatbot is Delivering on its Promise?](#)[The Power of Thoughtful RAG Implementation](#)[Share](#)

Each time Generative AI models like [GPT-4](#), [Claude](#), or [Gemini](#) hallucinate or give inaccurate answers, the world questions whether these LLM-powered chatbots can truly deliver value.

Even so, building LLM chatbots is one of the most popular GenAI applications. To address some of the LLMs' limitations, companies have turned to Retrieval Augmented Generation (RAG), a framework that ensures access to accurate information by retrieving it from multiple sources.

THE REAL FOUNDATION OF RAG CHATBOTS

Most companies find RAG an easy, exciting, and one of the quickest fixes to customer and employee experience with knowledge retrieval. And that's precisely where **two of the biggest misconceptions** about building a RAG chatbot hide: assuming that it's heavily based on LLMs and that it simply works well.

Knowledge retrieval is about knowing where to find the information you need. This is a key step in developing RAG chatbots before involving a large language model. The LLM is only one part of the RAG architecture, primarily involved in reasoning and generating

Table of contents

[The Real Foundation of RAG Chatbots](#)[A GenAI Architecture for a Conversational Agent](#)[Step 1: Data Ingestion](#)[Step 2: Knowledge Base Creation](#)[Step 3: Document Retrieval](#)[Step 4: LLM Reasoning](#)[Step 5: User Interface](#)[How to Know if the RAG Chatbot is Delivering on its Promise?](#)[The Power of Thoughtful RAG Implementation](#)[Share](#)

Each time Generative AI models like [GPT-4](#), [Claude](#), or [Gemini](#) hallucinate or give inaccurate answers, the world questions whether these LLM-powered chatbots can truly deliver value.

Even so, building LLM chatbots is one of the most popular GenAI applications. To address some of the LLMs' limitations, companies have turned to Retrieval Augmented Generation (RAG), a framework that ensures access to accurate information by retrieving it from multiple sources.

THE REAL FOUNDATION OF RAG CHATBOTS

Most companies find RAG an easy, exciting, and one of the quickest fixes to customer and employee experience with knowledge retrieval. And that's precisely where **two of the biggest misconceptions** about building a RAG chatbot hide: assuming that it's heavily based on LLMs and that it simply works well.

Knowledge retrieval is about knowing where to find the information you need. This is a key step in developing RAG chatbots before involving a large language model. The LLM is only one part of the RAG architecture, primarily involved in reasoning and generating

Table of contents

[The Real Foundation of RAG Chatbots](#)[A GenAI Architecture for a Conversational Agent](#)[Step 1: Data Ingestion](#)[Step 2: Knowledge Base Creation](#)[Step 3: Document Retrieval](#)[Step 4: LLM Reasoning](#)[Step 5: User Interface](#)[How to Know if the RAG Chatbot is Delivering on its Promise?](#)[The Power of Thoughtful RAG Implementation](#)[Share](#)

Each time Generative AI models like [GPT-4](#), [Claude](#), or [Gemini](#) hallucinate or give inaccurate answers, the world questions whether these LLM-powered chatbots can truly deliver value.

Even so, building LLM chatbots is one of the most popular GenAI applications. To address some of the LLMs' limitations, companies have turned to Retrieval Augmented Generation (RAG), a framework that ensures access to accurate information by retrieving it from multiple sources.

THE REAL FOUNDATION OF RAG CHATBOTS

Most companies find RAG an easy, exciting, and one of the quickest fixes to customer and employee experience with knowledge retrieval. And that's precisely where **two of the biggest misconceptions** about building a RAG chatbot hide: assuming that it's heavily based on LLMs and that it simply works well.

Knowledge retrieval is about knowing where to find the information you need. This is a key step in developing RAG chatbots before involving a large language model. The LLM is only one part of the RAG architecture, primarily involved in reasoning and generating

Table of contents

[The Real Foundation of RAG Chatbots](#)[A GenAI Architecture for a Conversational Agent](#)[Step 1: Data Ingestion](#)[Step 2: Knowledge Base Creation](#)[Step 3: Document Retrieval](#)[Step 4: LLM Reasoning](#)[Step 5: User Interface](#)[How to Know if the RAG Chatbot is Delivering on its Promise?](#)[The Power of Thoughtful RAG Implementation](#)[Share](#)

Each time Generative AI models like [GPT-4](#), [Claude](#), or [Gemini](#) hallucinate or give inaccurate answers, the world questions whether these LLM-powered chatbots can truly deliver value.

Even so, building LLM chatbots is one of the most popular GenAI applications. To address some of the LLMs' limitations, companies have turned to Retrieval Augmented Generation (RAG), a framework that ensures access to accurate information by retrieving it from multiple sources.

THE REAL FOUNDATION OF RAG CHATBOTS

Most companies find RAG an easy, exciting, and one of the quickest fixes to customer and employee experience with knowledge retrieval. And that's precisely where **two of the biggest misconceptions** about building a RAG chatbot hide: assuming that it's heavily based on LLMs and that it simply works well.

Knowledge retrieval is about knowing where to find the information you need. This is a key step in developing RAG chatbots before involving a large language model. The LLM is only one part of the RAG architecture, primarily involved in reasoning and generating

Table of contents

[The Real Foundation of RAG Chatbots](#)[A GenAI Architecture for a Conversational Agent](#)[Step 1: Data Ingestion](#)[Step 2: Knowledge Base Creation](#)[Step 3: Document Retrieval](#)[Step 4: LLM Reasoning](#)[Step 5: User Interface](#)[How to Know if the RAG Chatbot is Delivering on its Promise?](#)[The Power of Thoughtful RAG Implementation](#)[Share](#)

Each time Generative AI models like [GPT-4](#), [Claude](#), or [Gemini](#) hallucinate or give inaccurate answers, the world questions whether these LLM-powered chatbots can truly deliver value.

Even so, building LLM chatbots is one of the most popular GenAI applications. To address some of the LLMs' limitations, companies have turned to Retrieval Augmented Generation (RAG), a framework that ensures access to accurate information by retrieving it from multiple sources.

THE REAL FOUNDATION OF RAG CHATBOTS

Most companies find RAG an easy, exciting, and one of the quickest fixes to customer and employee experience with knowledge retrieval. And that's precisely where **two of the biggest misconceptions** about building a RAG chatbot hide: assuming that it's heavily based on LLMs and that it simply works well.

Knowledge retrieval is about knowing where to find the information you need. This is a key step in developing RAG chatbots before involving a large language model. The LLM is only one part of the RAG architecture, primarily involved in reasoning and generating

Table of contents

[The Real Foundation of RAG Chatbots](#)[A GenAI Architecture for a Conversational Agent](#)[Step 1: Data Ingestion](#)[Step 2: Knowledge Base Creation](#)[Step 3: Document Retrieval](#)[Step 4: LLM Reasoning](#)[Step 5: User Interface](#)[How to Know if the RAG Chatbot is Delivering on its Promise?](#)[The Power of Thoughtful RAG Implementation](#)[Share](#)

Each time Generative AI models like [GPT-4](#), [Claude](#), or [Gemini](#) hallucinate or give inaccurate answers, the world questions whether these LLM-powered chatbots can truly deliver value.

Even so, building LLM chatbots is one of the most popular GenAI applications. To address some of the LLMs' limitations, companies have turned to Retrieval Augmented Generation (RAG), a framework that ensures access to accurate information by retrieving it from multiple sources.

THE REAL FOUNDATION OF RAG CHATBOTS

Most companies find RAG an easy, exciting, and one of the quickest fixes to customer and employee experience with knowledge retrieval. And that's precisely where **two of the biggest misconceptions** about building a RAG chatbot hide: assuming that it's heavily based on LLMs and that it simply works well.

Knowledge retrieval is about knowing where to find the information you need. This is a key step in developing RAG chatbots before involving a large language model. The LLM is only one part of the RAG architecture, primarily involved in reasoning and generating

Table of contents

[The Real Foundation of RAG Chatbots](#)[A GenAI Architecture for a Conversational Agent](#)[Step 1: Data Ingestion](#)[Step 2: Knowledge Base Creation](#)[Step 3: Document Retrieval](#)[Step 4: LLM Reasoning](#)[Step 5: User Interface](#)[How to Know if the RAG Chatbot is Delivering on its Promise?](#)[The Power of Thoughtful RAG Implementation](#)[Share](#)

Each time Generative AI models like [GPT-4](#), [Claude](#), or [Gemini](#) hallucinate or give inaccurate answers, the world questions whether these LLM-powered chatbots can truly deliver value.

Even so, building LLM chatbots is one of the most popular GenAI applications. To address some of the LLMs' limitations, companies have turned to Retrieval Augmented Generation (RAG), a framework that ensures access to accurate information by retrieving it from multiple sources.

THE REAL FOUNDATION OF RAG CHATBOTS

Most companies find RAG an easy, exciting, and one of the quickest fixes to customer and employee experience with knowledge retrieval. And that's precisely where **two of the biggest misconceptions** about building a RAG chatbot hide: assuming that it's heavily based on LLMs and that it simply works well.

Knowledge retrieval is about knowing where to find the information you need. This is a key step in developing RAG chatbots before involving a large language model. The LLM is only one part of the RAG architecture, primarily involved in reasoning and generating

Table of contents

[The Real Foundation of RAG Chatbots](#)[A GenAI Architecture for a Conversational Agent](#)[Step 1: Data Ingestion](#)[Step 2: Knowledge Base Creation](#)[Step 3: Document Retrieval](#)[Step 4: LLM Reasoning](#)[Step 5: User Interface](#)[How to Know if the RAG Chatbot is Delivering on its Promise?](#)[The Power of Thoughtful RAG Implementation](#)[Share](#)

Each time Generative AI models like [GPT-4](#), [Claude](#), or [Gemini](#) hallucinate or give inaccurate answers, the world questions whether these LLM-powered chatbots can truly deliver value.

Even so, building LLM chatbots is one of the most popular GenAI applications. To address some of the LLMs' limitations, companies have turned to Retrieval Augmented Generation (RAG), a framework that ensures access to accurate information by retrieving it from multiple sources.

THE REAL FOUNDATION OF RAG CHATBOTS

Most companies find RAG an easy, exciting, and one of the quickest fixes to customer and employee experience with knowledge retrieval. And that's precisely where **two of the biggest misconceptions** about building a RAG chatbot hide: assuming that it's heavily based on LLMs and that it simply works well.

Knowledge retrieval is about knowing where to find the information you need. This is a key step in developing RAG chatbots before involving a large language model. The LLM is only one part of the RAG architecture, primarily involved in reasoning and generating

Table of contents

[The Real Foundation of RAG Chatbots](#)[A GenAI Architecture for a Conversational Agent](#)[Step 1: Data Ingestion](#)[Step 2: Knowledge Base Creation](#)[Step 3: Document Retrieval](#)[Step 4: LLM Reasoning](#)[Step 5: User Interface](#)[How to Know if the RAG Chatbot is Delivering on its Promise?](#)[The Power of Thoughtful RAG Implementation](#)[Share](#)

Each time Generative AI models like [GPT-4](#), [Claude](#), or [Gemini](#) hallucinate or give inaccurate answers, the world questions whether these LLM-powered chatbots can truly deliver value.

Even so, building LLM chatbots is one of the most popular GenAI applications. To address some of the LLMs' limitations, companies have turned to Retrieval Augmented Generation (RAG), a framework that ensures access to accurate information by retrieving it from multiple sources.

THE REAL FOUNDATION OF RAG CHATBOTS

Most companies find RAG an easy, exciting, and one of the quickest fixes to customer and employee experience with knowledge retrieval. And that's precisely where **two of the biggest misconceptions** about building a RAG chatbot hide: assuming that it's heavily based on LLMs and that it simply works well.

Knowledge retrieval is about knowing where to find the information you need. This is a key step in developing RAG chatbots before involving a large language model. The LLM is only one part of the RAG architecture, primarily involved in reasoning and generating

Table of contents

[The Real Foundation of RAG Chatbots](#)[A GenAI Architecture for a Conversational Agent](#)[Step 1: Data Ingestion](#)[Step 2: Knowledge Base Creation](#)[Step 3: Document Retrieval](#)[Step 4: LLM Reasoning](#)[Step 5: User Interface](#)[How to Know if the RAG Chatbot is Delivering on its Promise?](#)[The Power of Thoughtful RAG Implementation](#)[Share](#)

Each time Generative AI models like [GPT-4](#), [Claude](#), or [Gemini](#) hallucinate or give inaccurate answers, the world questions whether these LLM-powered chatbots can truly deliver value.

Even so, building LLM chatbots is one of the most popular GenAI applications. To address some of the LLMs' limitations, companies have turned to Retrieval Augmented Generation (RAG), a framework that ensures access to accurate information by retrieving it from multiple sources.

THE REAL FOUNDATION OF RAG CHATBOTS

Most companies find RAG an easy, exciting, and one of the quickest fixes to customer and employee experience with knowledge retrieval. And that's precisely where **two of the biggest misconceptions** about building a RAG chatbot hide: assuming that it's heavily based on LLMs and that it simply works well.

Knowledge retrieval is about knowing where to find the information you need. This is a key step in developing RAG chatbots before involving a large language model. The LLM is only one part of the RAG architecture, primarily involved in reasoning and generating

Table of contents

[The Real Foundation of RAG Chatbots](#)[A GenAI Architecture for a Conversational Agent](#)[Step 1: Data Ingestion](#)[Step 2: Knowledge Base Creation](#)[Step 3: Document Retrieval](#)[Step 4: LLM Reasoning](#)[Step 5: User Interface](#)[How to Know if the RAG Chatbot is Delivering on its Promise?](#)[The Power of Thoughtful RAG Implementation](#)[Share](#)

Each time Generative AI models like [GPT-4](#), [Claude](#), or [Gemini](#) hallucinate or give inaccurate answers, the world questions whether these LLM-powered chatbots can truly deliver value.

Even so, building LLM chatbots is one of the most popular GenAI applications. To address some of the LLMs' limitations, companies have turned to Retrieval Augmented Generation (RAG), a framework that ensures access to accurate information by retrieving it from multiple sources.

THE REAL FOUNDATION OF RAG CHATBOTS

Most companies find RAG an easy, exciting, and one of the quickest fixes to customer and employee experience with knowledge retrieval. And that's precisely where **two of the biggest misconceptions** about building a RAG chatbot hide: assuming that it's heavily based on LLMs and that it simply works well.

Knowledge retrieval is about knowing where to find the information you need. This is a key step in developing RAG chatbots before involving a large language model. The LLM is only one part of the RAG architecture, primarily involved in reasoning and generating

Table of contents

[The Real Foundation of RAG Chatbots](#)[A GenAI Architecture for a Conversational Agent](#)[Step 1: Data Ingestion](#)[Step 2: Knowledge Base Creation](#)[Step 3: Document Retrieval](#)[Step 4: LLM Reasoning](#)[Step 5: User Interface](#)[How to Know if the RAG Chatbot is Delivering on its Promise?](#)[The Power of Thoughtful RAG Implementation](#)[Share](#)

Each time Generative AI models like [GPT-4](#), [Claude](#), or [Gemini](#) hallucinate or give inaccurate answers, the world questions whether these LLM-powered chatbots can truly deliver value.

Even so, building LLM chatbots is one of the most popular GenAI applications. To address some of the LLMs' limitations, companies have turned to Retrieval Augmented Generation (RAG), a framework that ensures access to accurate information by retrieving it from multiple sources.

THE REAL FOUNDATION OF RAG CHATBOTS

Most companies find RAG an easy, exciting, and one of the quickest fixes to customer and employee experience with knowledge retrieval. And that's precisely where **two of the biggest misconceptions** about building a RAG chatbot hide: assuming that it's heavily based on LLMs and that it simply works well.

Knowledge retrieval is about knowing where to find the information you need. This is a key step in developing RAG chatbots before involving a large language model. The LLM is only one part of the RAG architecture, primarily involved in reasoning and generating

Table of contents

[The Real Foundation of RAG Chatbots](#)[A GenAI Architecture for a Conversational Agent](#)[Step 1: Data Ingestion](#)[Step 2: Knowledge Base Creation](#)[Step 3: Document Retrieval](#)[Step 4: LLM Reasoning](#)[Step 5: User Interface](#)[How to Know if the RAG Chatbot is Delivering on its Promise?](#)[The Power of Thoughtful RAG Implementation](#)[Share](#)

Each time Generative AI models like [GPT-4](#), [Claude](#), or [Gemini](#) hallucinate or give inaccurate answers, the world questions whether these LLM-powered chatbots can truly deliver value.

Even so, building LLM chatbots is one of the most popular GenAI applications. To address some of the LLMs' limitations, companies have turned to Retrieval Augmented Generation (RAG), a framework that ensures access to accurate information by retrieving it from multiple sources.

THE REAL FOUNDATION OF RAG CHATBOTS

Most companies find RAG an easy, exciting, and one of the quickest fixes to customer and employee experience with knowledge retrieval. And that's precisely where **two of the biggest misconceptions** about building a RAG chatbot hide: assuming that it's heavily based on LLMs and that it simply works well.

Knowledge retrieval is about knowing where to find the information you need. This is a key step in developing RAG chatbots before involving a large language model. The LLM is only one part of the RAG architecture, primarily involved in reasoning and generating

Table of contents

[The Real Foundation of RAG Chatbots](#)[A GenAI Architecture for a Conversational Agent](#)[Step 1: Data Ingestion](#)[Step 2: Knowledge Base Creation](#)[Step 3: Document Retrieval](#)[Step 4: LLM Reasoning](#)[Step 5: User Interface](#)[How to Know if the RAG Chatbot is Delivering on its Promise?](#)[The Power of Thoughtful RAG Implementation](#)[Share](#)

Each time Generative AI models like [GPT-4](#), [Claude](#), or [Gemini](#) hallucinate or give inaccurate answers, the world questions whether these LLM-powered chatbots can truly deliver value.

Even so, building LLM chatbots is one of the most popular GenAI applications. To address some of the LLMs' limitations, companies have turned to Retrieval Augmented Generation (RAG), a framework that ensures access to accurate information by retrieving it from multiple sources.

THE REAL FOUNDATION OF RAG CHATBOTS

Most companies find RAG an easy, exciting, and one of the quickest fixes to customer and employee experience with knowledge retrieval. And that's precisely where **two of the biggest misconceptions** about building a RAG chatbot hide: assuming that it's heavily based on LLMs and that it simply works well.

Knowledge retrieval is about knowing where to find the information you need. This is a key step in developing RAG chatbots before involving a large language model. The LLM is only one part of the RAG architecture, primarily involved in reasoning and generating

Table of contents

[The Real Foundation of RAG Chatbots](#)[A GenAI Architecture for a Conversational Agent](#)[Step 1: Data Ingestion](#)[Step 2: Knowledge Base Creation](#)[Step 3: Document Retrieval](#)[Step 4: LLM Reasoning](#)[Step 5: User Interface](#)[How to Know if the RAG Chatbot is Delivering on its Promise?](#)[The Power of Thoughtful RAG Implementation](#)[Share](#)

Each time Generative AI models like [GPT-4](#), [Claude](#), or [Gemini](#) hallucinate or give inaccurate answers, the world questions whether these LLM-powered chatbots can truly deliver value.

Even so, building LLM chatbots is one of the most popular GenAI applications. To address some of the LLMs' limitations, companies have turned to Retrieval Augmented Generation (RAG), a framework that ensures access to accurate information by retrieving it from multiple sources.

THE REAL FOUNDATION OF RAG CHATBOTS

Most companies find RAG an easy, exciting, and one of the quickest fixes to customer and employee experience with knowledge retrieval. And that's precisely where **two of the biggest misconceptions** about building a RAG chatbot hide: assuming that it's heavily based on LLMs and that it simply works well.

Knowledge retrieval is about knowing where to find the information you need. This is a key step in developing RAG chatbots before involving a large language model. The LLM is only one part of the RAG architecture, primarily involved in reasoning and generating

