# DocuQuery

## Introduction

DocuQuery is an intelligent chatbot designed to enhance the learning experience by meticulously analyzing the content of uploaded PDF files. It assists users in delving deeply into specific concepts, providing detailed explanations and insights as requested. This tool streamlines the acquisition of knowledge by presenting complex information in an accessible and user-centric manner.

## Purpose

The current system of physical document navigation poses considerable challenges, particularly when it comes to the efficient retrieval of information. Traditional methods require manual scanning of pages, a process that is both time-consuming and cumbersome, especially when attempting to cross-reference or revisit specific concepts. **DocuQuery** revolutionizes this approach by offering an advanced, user-friendly platform that simplifies the search and learning process. With **DocuQuery,** users gain the ability to instantly access targeted information within any provided document, thereby enhancing their educational and informational experience. This system not only optimizes time management but also promotes an enriched, on-demand learning environment.

**Target Audience:** This chatbot would appeal to anyone looking for an optimized document navigation experience that saves time and makes the retrieval of information as convenient as possible.

**Platform Used:** Flowise

Prompt patterns:
1. **What is Linear Regression? (if source data mining textbook)**
2. **How to make a nut raisin cake ? (cook book recipe source)**

## Basic Foundations for Flowise

## What is LLM?

LLMs, or Large Language Models, are advanced artificial intelligence programs designed to understand, generate, and sometimes translate natural language text. They are trained on vast datasets of human language, which allows them to predict and produce text that is coherent and contextually relevant. LLMs use deep learning techniques, particularly neural networks, to capture the nuances of language, including grammar, style, and even some elements of common sense reasoning and knowledge.

These models are used in a variety of applications, such as chatbots, translation services, content creation tools, and more. They are capable of performing tasks like answering questions, composing essays, summarizing documents, and engaging in conversation, making them valuable across many domains including customer service, education, and content generation.

# What is Flowise?

Flowise is an open-source, low-code tool that allows developers to build customized applications harnessing the capabilities of Large Language Models (LLMs). It provides a user-friendly drag-and-drop interface, enabling the creation of applications with complex LLM flows and AI agents efficiently and with minimal coding. The platform supports local hosting of LLMs and offers seamless deployment across various cloud platforms like AWS, Azure, and GCP.

It comes with a range of features, including ready-to-use app templates, conversational agents that can maintain context, and the ability to integrate and extend functionality through APIs, SDKs, and embedded chat widgets. Flowise also boasts community support and documentation to assist developers in building and deploying their apps.

Flowise aims to democratize the use of LLMs, making it accessible to developers who may not have deep expertise in machine learning or AI, while still offering powerful features for those who do. With Flowise, you can build a chatbot, connect with multiple workflows and AI agents, and even create project management tasks integrated with tools like Slack and Notion

Website: [Flowise - Low code LLM Apps Builder](#)
GitRepo: [GitHub](#)
Flowise Documentation: [Welcome to Flowise | FlowiseAI](#)

# Why Flowise?

Flowise stands out as a tool designed for building custom Large Language Model (LLM) applications, and here are several reasons why it might be a particularly good choice:

1.Low-Code Environment: Flowise offers a low-code environment that enables users to construct complex LLM applications without the need for extensive programming knowledge. This can be particularly beneficial for users with limited technical backgrounds or for projects that need to be prototyped and iterated quickly (Flowise - Low code LLM Apps Builder) (Y Combinator).

2. Drag-and-Drop Interface: The platform's drag-and-drop interface simplifies the process of creating workflows. This visual approach can help you understand the architecture of your applications better and make it easier to experiment with different configurations (Flowise - Low code LLM Apps Builder) (CODE Online).

3. Integration Capabilities: Flowise allows for seamless integration with various APIs, enabling the incorporation of LLM functionalities into existing front-end applications. This can enhance the capabilities of your applications by adding conversational chatbots, language translation tools, and other language processing functionalities (Homepage).
4. Extensibility: With its open-source nature, Flowise is built to be extensible. It supports various libraries and frameworks like LangChain, LlamaIndex, and HuggingFace, among others,

allowing developers to customize their integrations and leverage a broad range of functionalities (Y Combinator).
5. Community and Support: Being a tool that's trending and backed by a growing community, Flowise has the advantage of collective knowledge sharing and support. It's a dynamic platform that benefits from the contributions and feedback of its user base (Flowise - Low code LLM Apps Builder).

6. Deployment Flexibility: The platform supports deployment on a variety of cloud platforms, which gives you the flexibility to host your applications in a way that best suits your needs (Flowise - Low code LLM Apps Builder)

7. No-Cost and Open-Source Advantage: Flowise is open-source and available for both personal and commercial use without any cost. This can significantly lower the barriers to entry for creating sophisticated LLM applications (Homepage)
Given these features, Flowise can be an attractive option for anyone looking to harness the power of LLMs to build applications that are both sophisticated and user-friendly, without the need for deep coding expertise or significant resources. Whether you're a developer looking to streamline your workflow or a non-technical user aiming to leverage LLMs, Flowise provides a platform that is accessible, adaptable, and supportive of creative LLM application development.


**How to build DocuQuery chatbot to chat with your own data with Flowise and Pinecone Vector database?**
**Prequisites:**
1. Open AI API key
2. Pinecone API key and Pinecone Index name
   - To create a pine cone index first create account then select index give a name for the index and the dimension i.e 1536

**Procedure Steps to build a rag chatbot to chat with your own data**
Task-1:
Open command prompt or Terminal
1. Install flowise
   npm install -g flowise
2. Start Flowise
   npx  flowise start
3. Open http://localhost:3000
Task-2:
1. Now click on + button on the top right corner of the page to create a new chatflow and give it a name, say project-1
2. Drag and Drop the following from the left side menu by clicking on + button on the top left corner of the page:
   - Add Pdf file from Document Loader section under LangChain
   - Add OpenAI from chat models section
   - Add RecursiveCharacterText Splitter from text splitter
   - Add pinecone upsert vector DB from vector stores
   - Add open ai embeddings from embeddings sections

- Add conversational retrieval QA chain from LLM Chains section
- Add in memory cache to the chat open ai input cache

Task-3:

Now provide API keys of open ai model and pinecone database. Also provide name of the model for embedding model of what is being used.

Task-4:

Connect each tile to another as following
- Connect PDF file to recursive character text splitter to split the document into chunks
- Now connect the output of PDF to the inputs of pinecone vector store
- Now connect open ai embeddings to inputs of pinecone vector store.
- Connect the output of the chat open ai and the output of pinecone to the input of the conversational retrieval Q&A chain input

**Task-5:**

Now save the chatflow

**TESTING**

- **Objective:** To evaluate the performance of DocuQuery and ensure its functionality meets the intended requirements.
- **Test Cases:**
  i. **Document Upload Test:**
     - **Test Case:** Upload a Textbook of a subject or a cookbook recipe PDF.
     - **Expected Result:** DocuQuery explains the topic of requested.
     - **Actual Result:** Key topics are successfully extracted and presented to the user.
  ii. **Question-Answering Test:**
     - **Test Case:** Ask DocuQuery about the explanation of a topic.
     - **Expected Result:** DocuQuery accurately identifies the specific topic and provides relevant details.
     - **Actual Result:** DocuQuery retrieves the correct topic information from the uploaded documents.
  iii. **Error Handling Test:**
     - **Test Case:** Input a vague query with ambiguous terms.
     - **Expected Result:** DocuQuery responds with a helpful error message, guiding the user to provide clearer input.
     - **Actual Result:** DocuQuery recognizes the ambiguity and prompts the user Hmm, I dont know.

**Encountered Challenges and Solutions**

**Document Extraction Inconsistencies**

- **Challenge**: There were irregularities in the results from document extraction.
  o **Solution**: Enhancement of document parsing algorithms has been implemented, resulting in increased accuracy of data extraction.

**Response Time Lags**

- **Challenge**: Users experienced occasional lags in receiving responses.
  o **Solution**: Backend processes have been optimized and server configurations adjusted, leading to quicker response times.

## Complex Query Misinterpretations

- **Challenge**: The system sometimes misinterpreted queries of a complex nature.
  - **Solution**: The training datasets have been expanded and the natural language processing models have been refined to better comprehend a wider array of query structures.

## Quality and Documentation for DocuQuery
## Quality Assessment
## User Experience

- **Report**: Users of DocuQuery find the platform's interface to be highly intuitive and user-friendly.
- **Impact**: The prompt and accurate responses provided by DocuQuery significantly enrich the document browsing and learning experience for users.

## Accuracy

- **Consistency**: DocuQuery maintains a high degree of precision in extracting and interpreting information from documents.
- **Outcome**: The system delivers relevant answers to user inquiries with a minimized rate of error, thereby reducing misinformation.

## Error Handling

- **Management**: DocuQuery efficiently manages errors, offering users clear and actionable feedback.
- **Guidance**: Upon receiving ambiguous queries, DocuQuery prompts users to rephrase their questions, ensuring a smooth and continuous user experience.

**Some basic quiz questions to have basic understanding of the concepts we have implemented:**

**What is a vector database commonly used for in the context of machine learning?**
A) Storing relational data
B) Optimizing search queries
C) Storing and searching embeddings
D) Managing user permissions
Correct Answer: C
Explanation: Vector databases are specifically designed to store and manage embeddings, which are high-dimensional vectors used in machine learning to represent data.

**Which tool is best suited for automating the creation of large language models?**
A) LangSmith
B) LangFlow
C) Lang Chain
D) Flowise
Correct Answer: A
Explanation: LangSmith is geared towards automating the creation of large language models, providing tools and workflows to streamline this process.

**What is the primary purpose of LangFlow?**
A) Database management
B) Workflow automation
C) Language model training

D) Embeddings optimization
Correct Answer: B
Explanation: LangFlow is primarily used for workflow automation, helping users manage and automate their data processing workflows efficiently.
Which of the following is NOT a characteristic of large language models?
A) Ability to generate human-like text
B) Use of embeddings
C) Low resource requirements
D) Trained on diverse internet text
Correct Answer: C
Explanation: Large language models typically require significant computational resources for training and operation, contrary to option C.
Flowise is designed to improve which of the following aspects?
A) Text generation speed
B) Model accuracy
C) Workflow efficiency
D) Data storage
Correct Answer: C
Explanation: Flowise is designed to enhance workflow efficiency, enabling smoother and more effective management of machine learning workflows.
What is a primary use of embeddings in machine learning?
A) Improving the security of data
B) Reducing storage needs for data
C) Facilitating accurate and fast searches
D) Generating automated email responses
Correct Answer: C
Explanation: Embeddings are used to transform data into a high-dimensional space, making it easier to perform fast and accurate searches.
Lang Chain primarily focuses on which of the following?
A) Chain of command in organizational structures
B) Connecting different language processing tools into a cohesive workflow
C) Encryption and security solutions
D) Human resources management
Correct Answer: B
Explanation: Lang Chain helps in integrating different language processing tools and technologies into a cohesive workflow.
Which of the following best describes Langsmith?
A) A tool for improving physical manufacturing processes
B) A language model training facilitator
C) A type of vector database
D) A cybersecurity service
Correct Answer: B
Explanation: Langsmith assists in the facilitation and automation of training language models.
In the context of AI, what is the role of a vector database like Pinecone or Milvus?
A) Managing network infrastructure

**B) Storing discrete mathematical vectors**
**C) Serving and searching embeddings for similarity searches**
**D) Optimizing financial models**
**Correct Answer: C**
**Explanation: Pinecone and Milvus are vector databases designed to store embeddings and perform similarity searches efficiently.**
**What advantage do embeddings provide in natural language processing tasks?**
**A) They speed up the network**
**B) They enhance the physical security of data centers**
**C) They represent words in a way that captures semantic meaning**
**D) They decrease the accuracy of models**
**Correct Answer: C**
**Explanation: Embeddings are beneficial in natural language processing as they can represent words or phrases in vector form that encapsulates semantic meanings, aiding in various tasks like sentiment analysis, translation, and more.**