



Build a PDF Q&A chatbot using Azure OpenAI and Streamlit



Streamlit



@ADVANCINGANALYTICS



@ADVA^NLYTICSUK



/ADVANCING ANALYTICS

Dr Gavita Regunath

Head of AI



gavitaregunath



@gavi_sr



**ADVANCING
ANALYTICS**

www.advancinganalytics.co.uk



/AdvancingAnalytics



@AdvAnalyticsUK

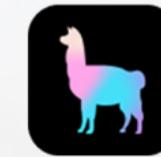
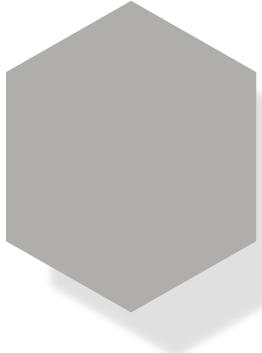


totallyskewed.co.uk



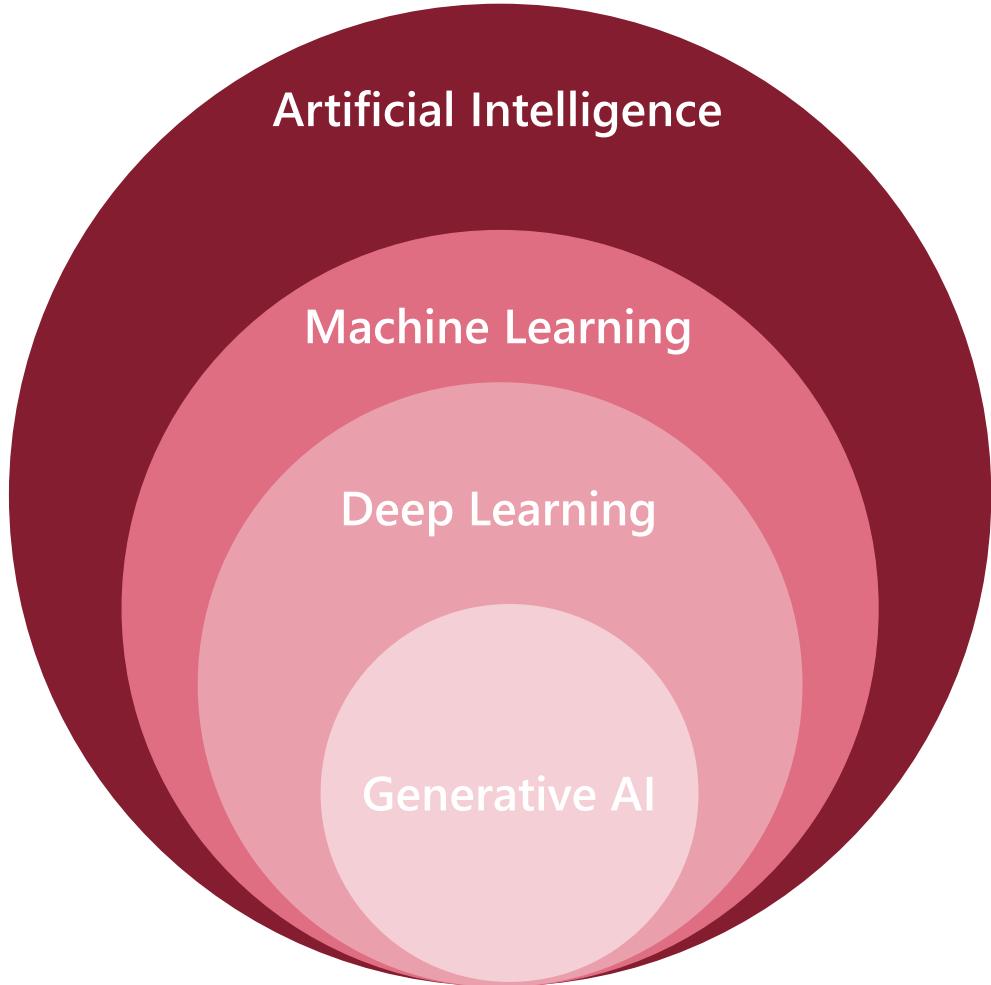
Agenda

- 1. Background & Problem Statement**
- 2. Environment Setup**
 - 2.1. Build a virtual environment
 - 2.2. Activate the virtual environment
 - 2.3. Install dependencies
- 3. Building the User Interface Using Streamlit Application**
 - 3.1. Design the Application Interface
 - 3.2. Upload the Pdf file
- 4. Building The Application Backend Using LangChain & Azure OpenAI API**
 - 4.1. Setup Azure OpenAI API
 - 4.2. Load & processing Pdf file
 - 4.3. Prompt Handling & Generating Answers
- 5. Putting Everything Together & Application Demo**





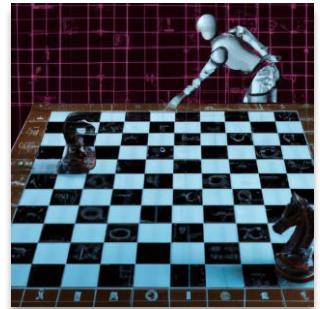
What is Generative AI?



1956

Artificial Intelligence

the field of computer science that seeks to create intelligent machines that can replicate or exceed human intelligence



1997

Machine Learning

subset of AI that enables machines to learn from existing data and improve upon that data to make decisions or predictions



2017

Deep Learning

a machine learning technique in which layers of neural networks are used to process data and make decisions



2021

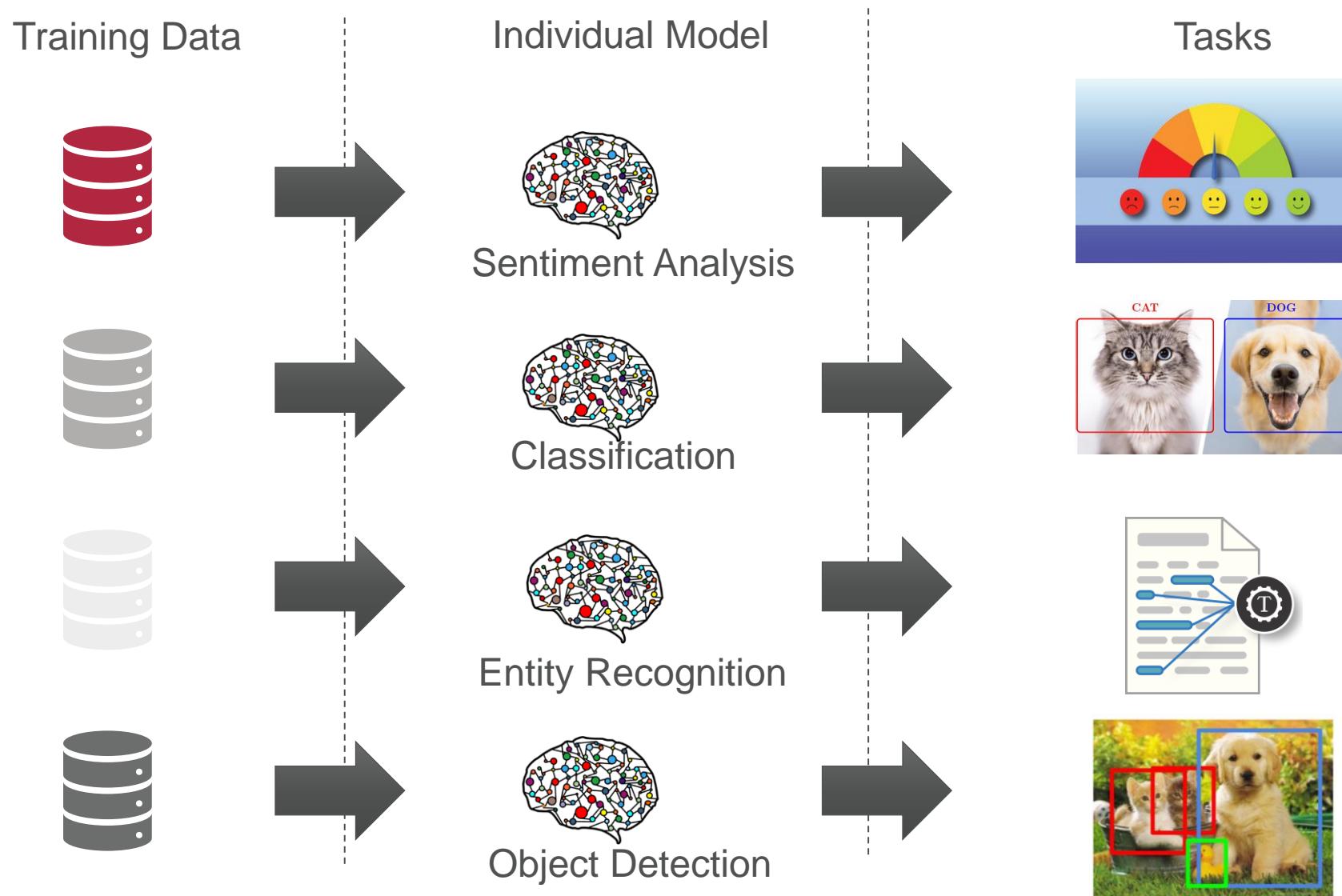
Generative AI

create new written, visual, and auditory content given prompts or existing data.



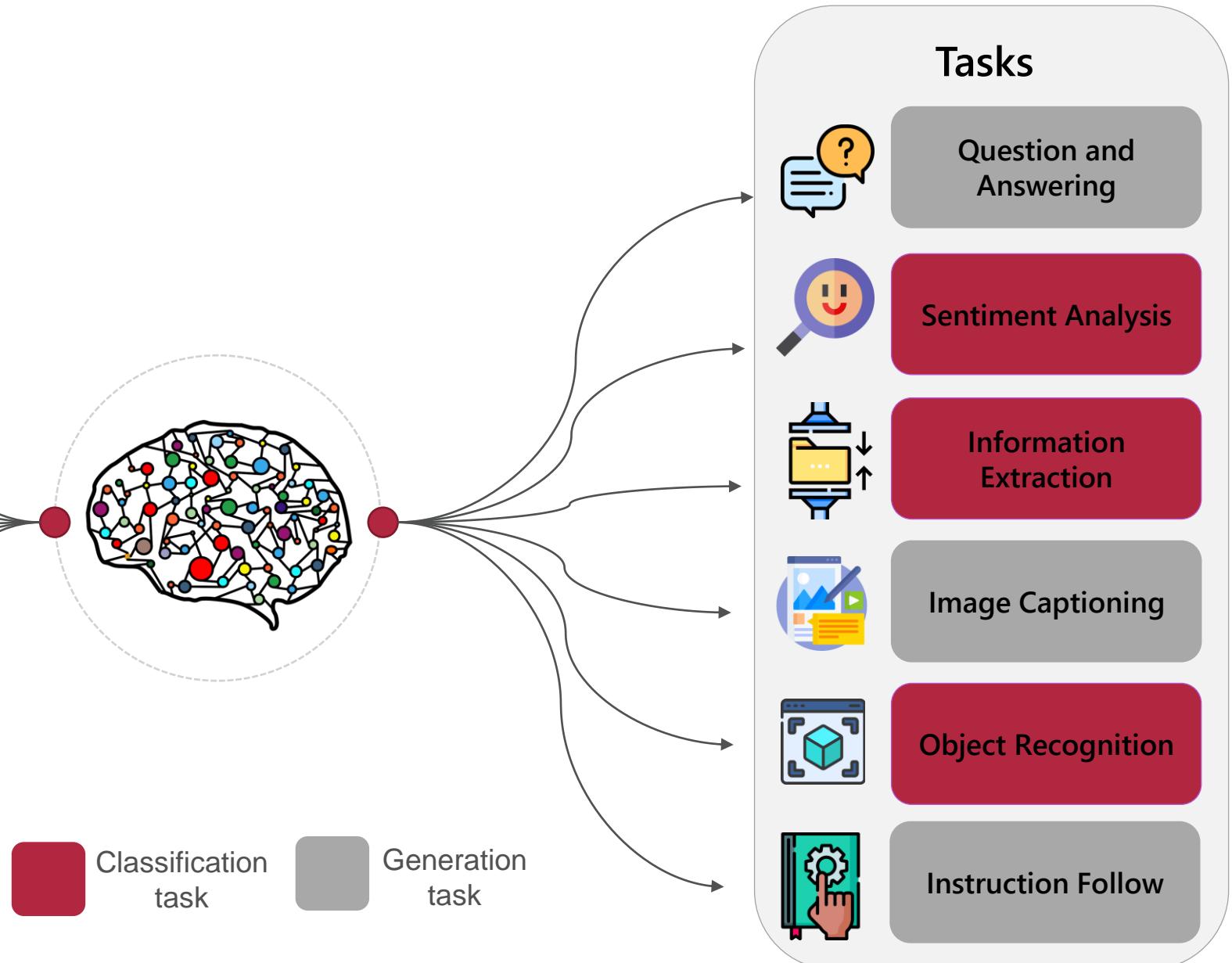
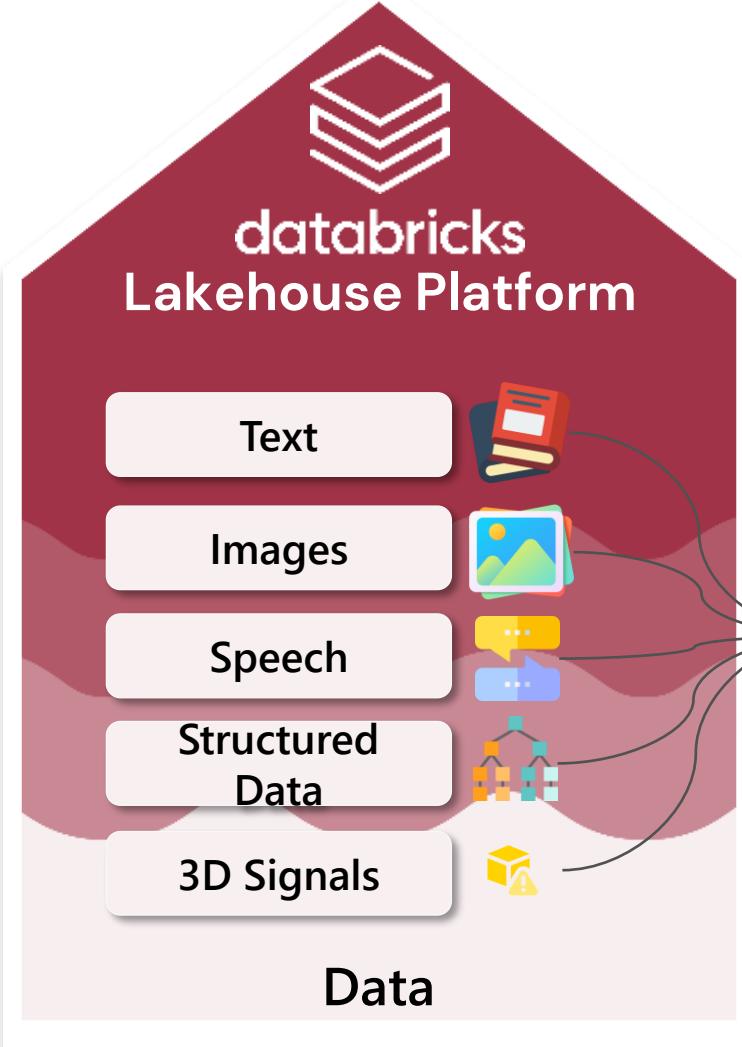


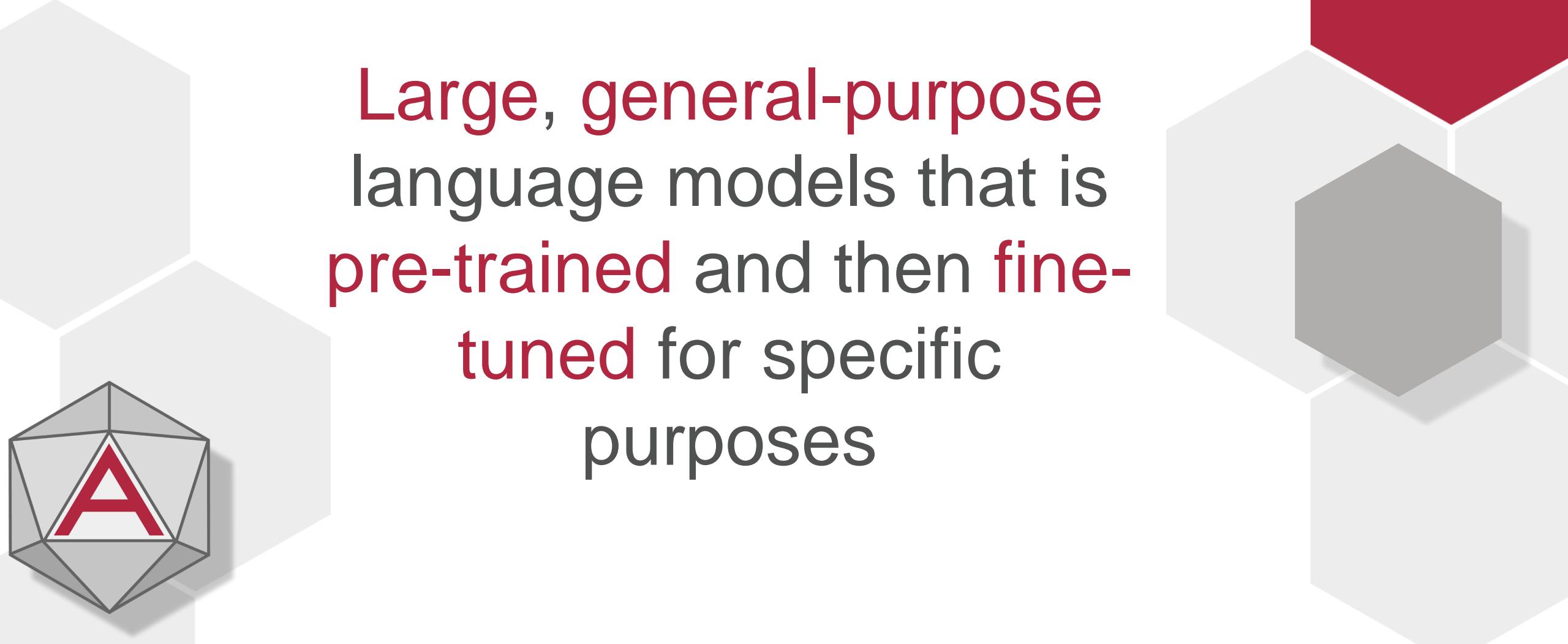
Traditional machine learning





Generative AI is enabled by Large Language Models





Large, general-purpose
language models that is
pre-trained and then **fine-**
tuned for specific
purposes



Challenges of LLMs



Challenges

1 LLMs hallucinate

Presents real ethical and business concerns, so much so that for some it is a show stopper

2 LLMs do not have data context

Trained on open internet data that might not provide good responses

3 Lack of transparency

Often requires interpretability, challenging understanding the reasoning behind their decisions

4 Speed and cost

LLMs are computationally very expensive, extraction from large documents can take up to tens of seconds.

5 Potential bias

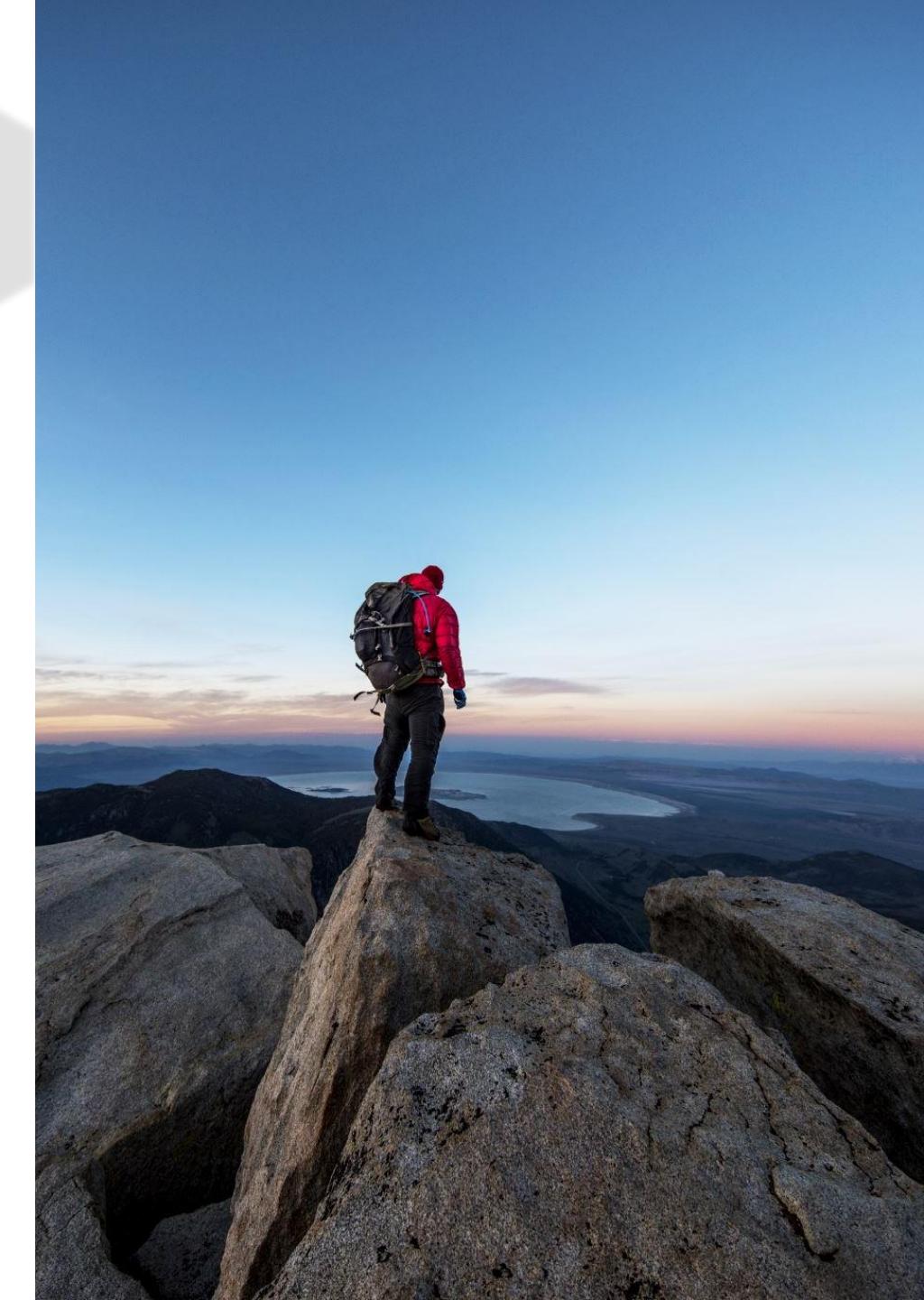
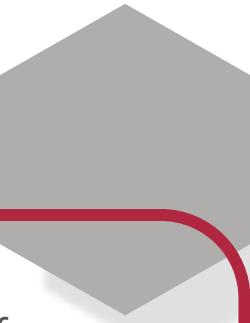
Only as good as the data they are trained on, they are expected to have societal and geographical biases encoded in them

6 Lack of control

Even the best prompts cannot force the model to give consistently correct output.

7 Data residency and privacy

Control over data and privacy management is crucial for companies



Methods to improve the usage and performance of LLM





Prompt Engineering

“Prompt engineering is the art of using natural language to get the results you are looking from the model”



Prompt Engineering with in-context learning



Prompt
Engineering

Basic

Zero-shot: Predicting with no sample provided

Translate English to French: <---- task description
Cheese => <---- prompt

One-shot: Predicting with one sample provided

1 Translate English to French: <---- task description
2 Sea otter => loutre de mer <---- example
3 Cheese => <---- prompt

Few-shot: Predicting with a few samples provided

1 Translate English to French: <---- task description
2 Sea otter => loutre de mer <---- examples
3 Peppermint => menthe poivre <----
4 Plush giraffe => girafe peluche <----
5 Cheese => <---- prompt

With detail
and context

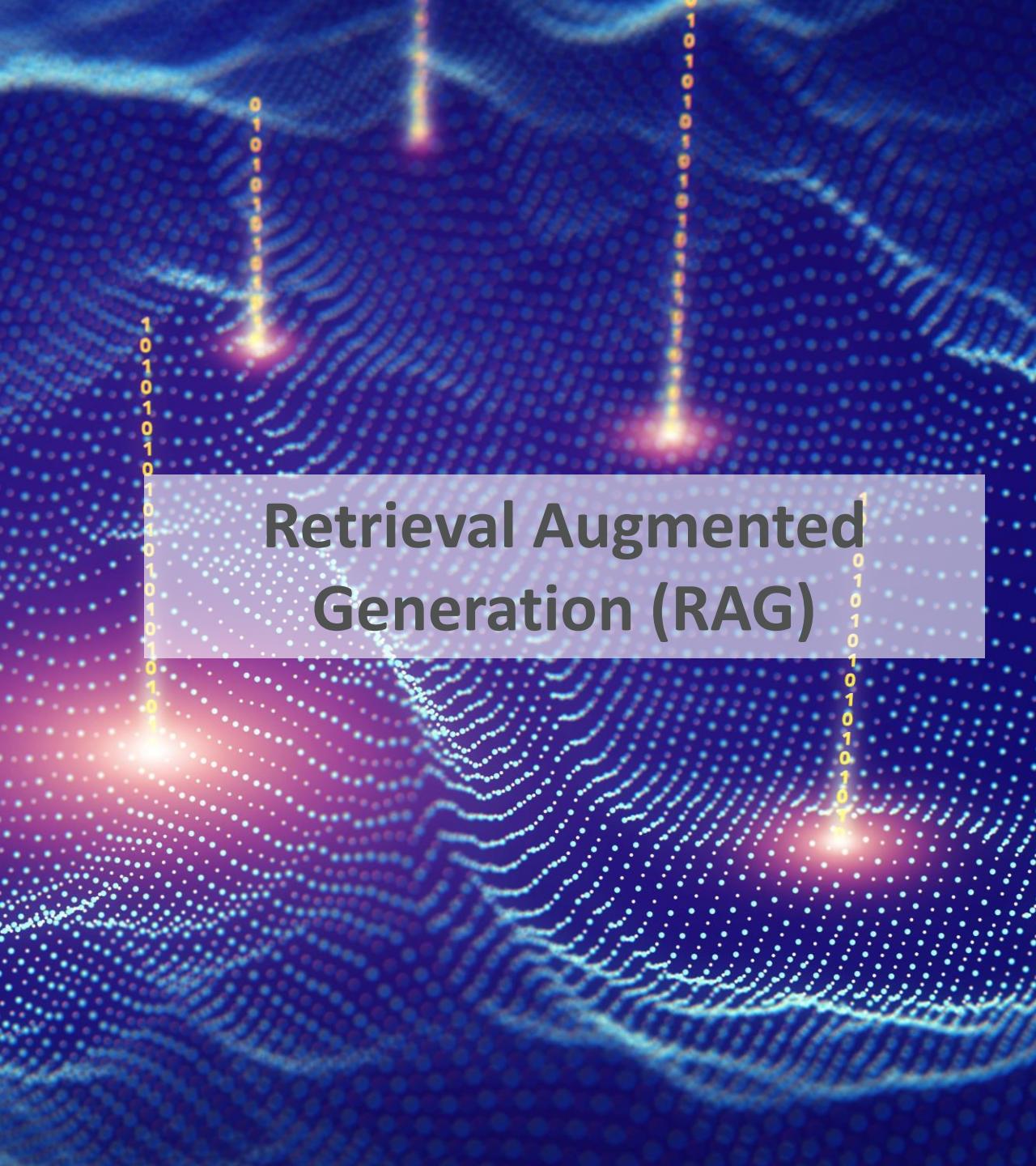
**“The hottest new
programming
language is
English”**

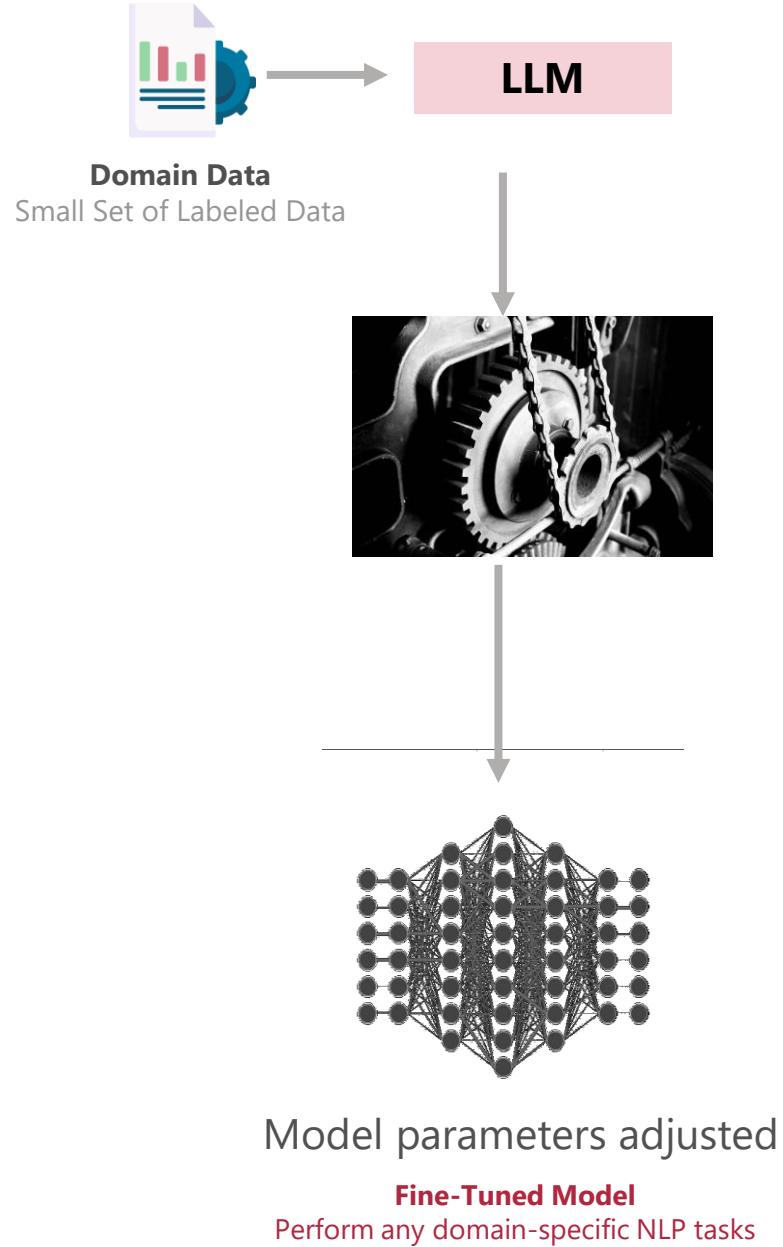


Andrej Karpathy
(Former) Director of AI at Tesla



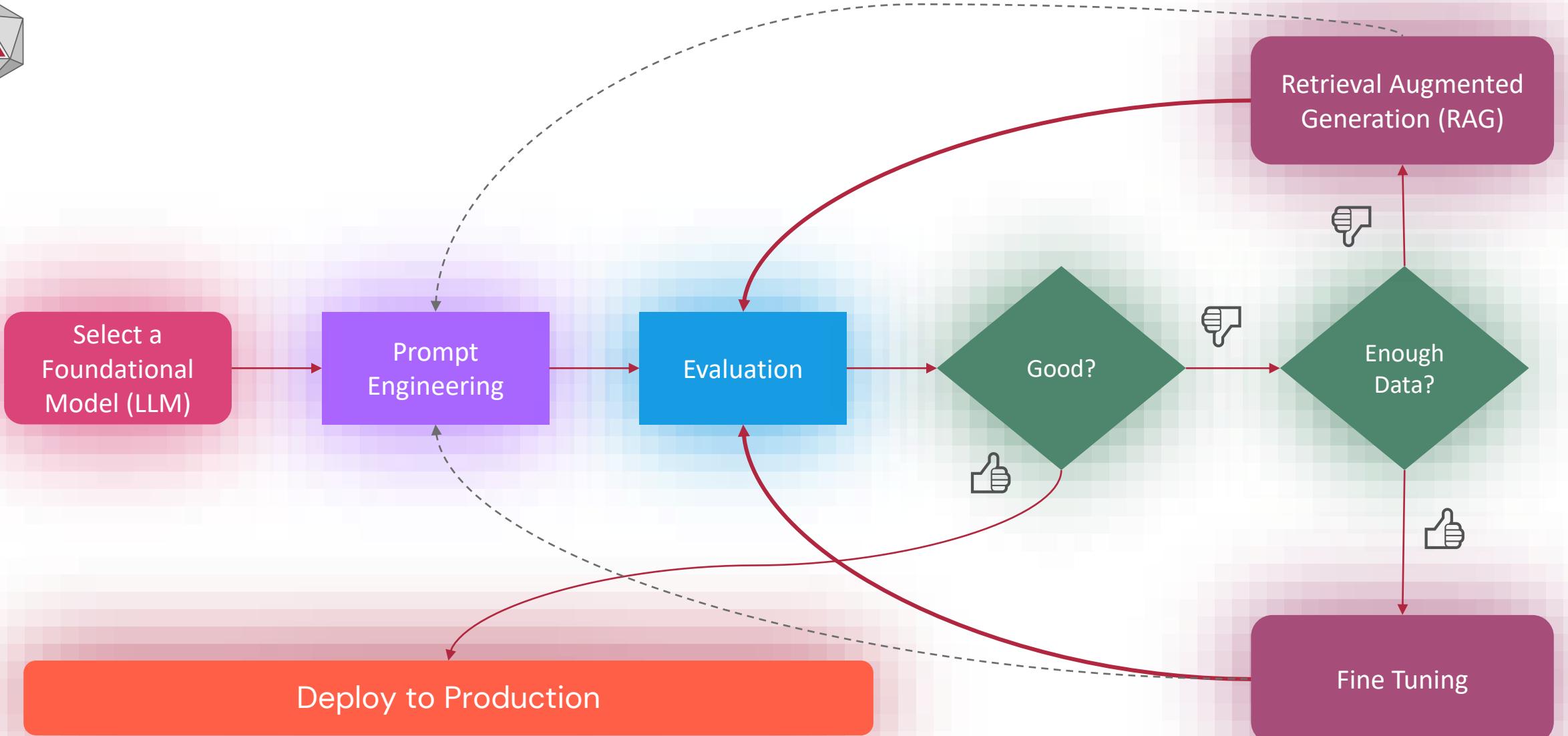
**“Framework that
combines information
retrieval and text
generation to produce
more accurate and
contextually relevant
responses”**





Fine-Tuning

“Fine-tuning is about adapting a pre-trained language model to a specific task or domain”



How to adapt LLM models for your task

Prompt Engineering with in-context learning

Zero-Shot

The model predicts the answer given only a natural language description of the task.

One-Shot

In addition to the task description, the model sees a single example of the task

Few-Shot

In addition to the task description, the model sees a few examples of the task.

Retrieval-Augmented Generation (RAG)

This approach integrates the power of retrieval (or searching) into LLM text generation

Fine Tuning

The model is trained via repeated gradient updates using a large corpus of example tasks.

Prompt
Engineering



RAG



Fine Tuning

RAG In a Nutshell



@ADVANCINGANALYTICS



@ADVANALYTICSUK



/ADVANCING ANALYTICS

Ability of an LLM to learn
information not through training,
but by receiving new information in
a carefully formatted prompt



R A G

1

Retrieval

Get information from a database that relates to the question you asked.

2

Augmentation

Combine the retrieved information with the initial prompt

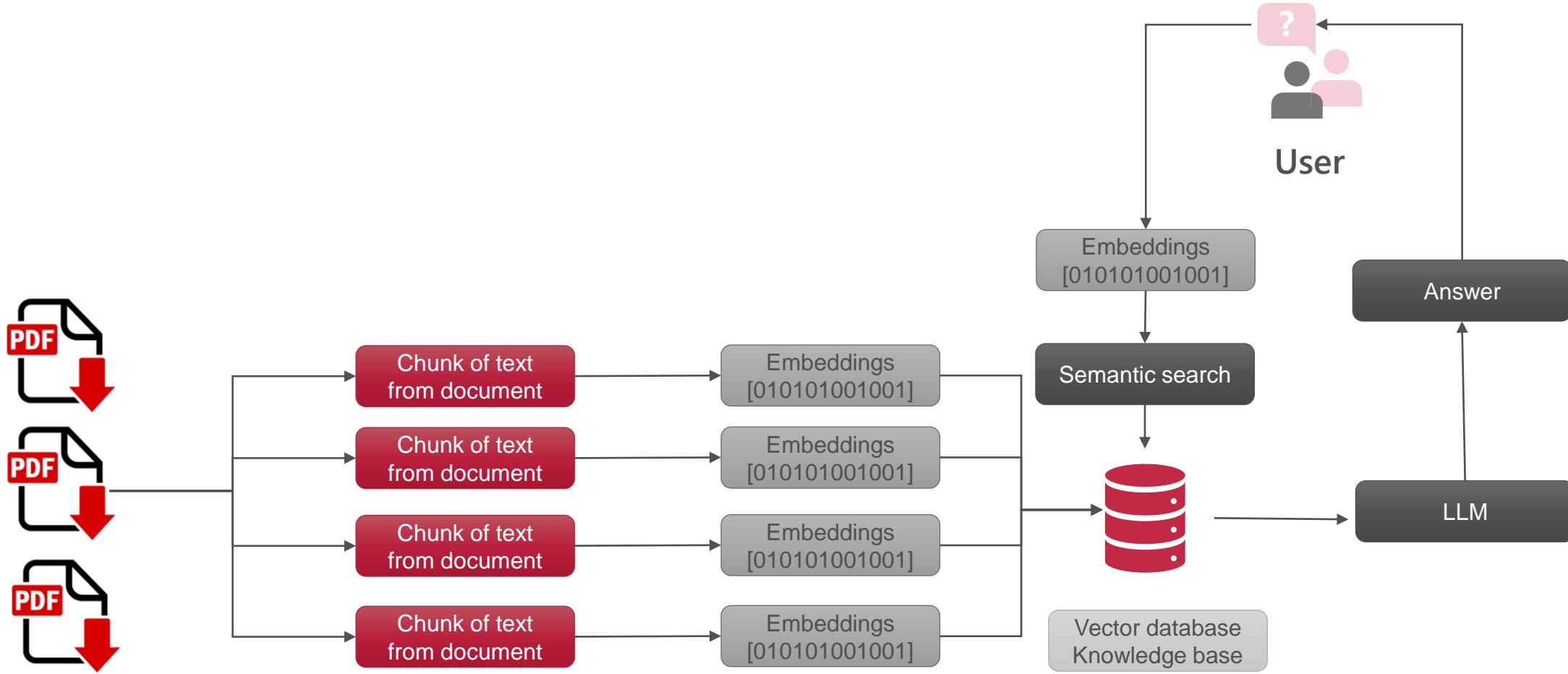
3

Generate

Pass the augmented prompt to a large language model, generating the final output

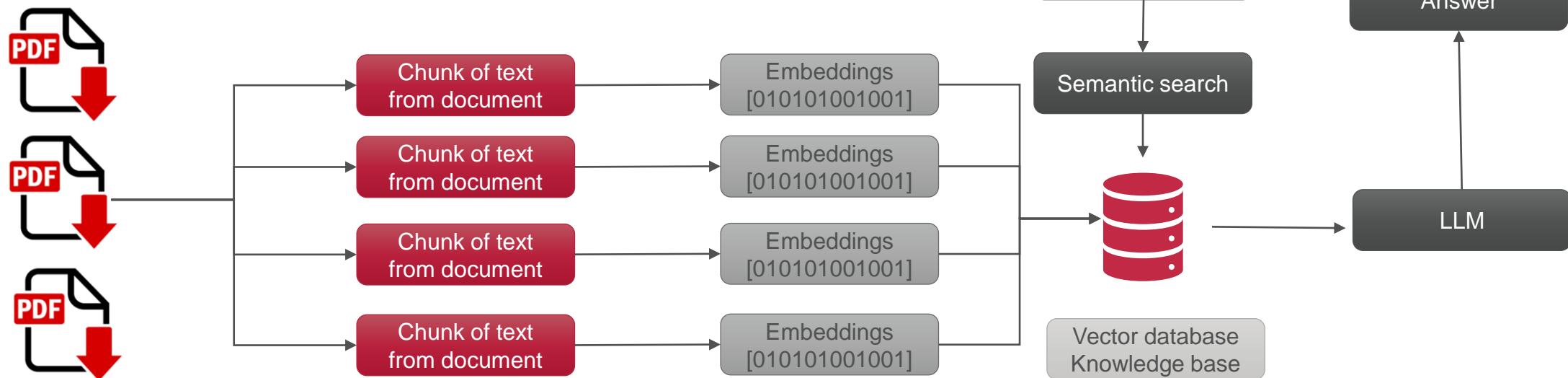


RAG (Retrieval-Augmented Generation)



RAG (Retrieval-Augmented Generation)

Enterprise Data Sources





WHAT PROBLEM DO THEY SOLVE?

1. Hallucinations:

RAG mitigates hallucinations by retrieving data from trusted sources, ensuring responses are grounded in factual information.

2. Lack of Transparency:

RAG offers transparency through retrieval, indicating the source of information, improving traceability, and identifying knowledge base issues or biases.

3. Contextual Understanding:

RAG enhances context by combining retrieval with generation from a knowledge base, resulting in more contextually relevant and accurate answers.



Vector Databases



@ADVANCINGANALYTICS



@ADVANALYTICSUK



/ADVANCING ANALYTICS



Vector Databases

Open Source



Weaviate



Milvus



chroma



FAISS
Scalable Search With Facebook AI



Managed



Pinecone



Azure AI Search



Azure AI Search



Type	Cloud-based search service with AI capabilities	Open-source library for efficient similarity search
Hosting	Fully managed by Microsoft Azure	Self-hosted, requires own infrastructure or cloud provider
Scalability	Automatically scalable within Azure infrastructure	Manually scalable, depending on the hardware capabilities
Integration	Seamless integration with other Azure services	Requires custom integration with existing systems
Ease of Use	Offers a managed environment with user-friendly tools and interfaces	Requires more in-depth knowledge of similarity search algorithms and software development
Indexing Capabilities	Provides indexing for full-text search and supports AI enrichment features	Specialises in creating efficient indices for nearest neighbor search
Performance	Optimised for integration and ease of use within the Azure ecosystem	Highly optimised for large-scale vector similarity search
Cost	Usage-based pricing model	Free to use, but costs associated with self-hosting and infrastructure
Use Cases	Ideal for enterprise search solutions that require full-text search and cognitive skills	Best suited for applications that perform large-scale deep learning similarity searches
Data Security	Managed by Azure with enterprise-grade security features	Depends on the deployment model and infrastructure security

Choosing the right vector databases

A list of questions we should be exploring with the customer to understand which DB to utilise

1. Does the customer have their own index for documents that are being embedded?
2. Do we need to index the vectors/documents?
3. Is the customer happy to be paying for a vector database like AI Search or Pinecone? Or is there a preference to use open-source vector databases like Chroma?
4. How critical is real-time performance? What are the latency requirements?
5. Are there any scalability requirements in terms of handling increasing data loads over time?
6. How would this database interact/integrate with other stacks within the RAG LLM architecture?
7. Do we need real-time search capabilities? Are we doing recommendation engines or fraud detection for example?
8. Do we need CRUD (inserting, updating and deleting data) capability?

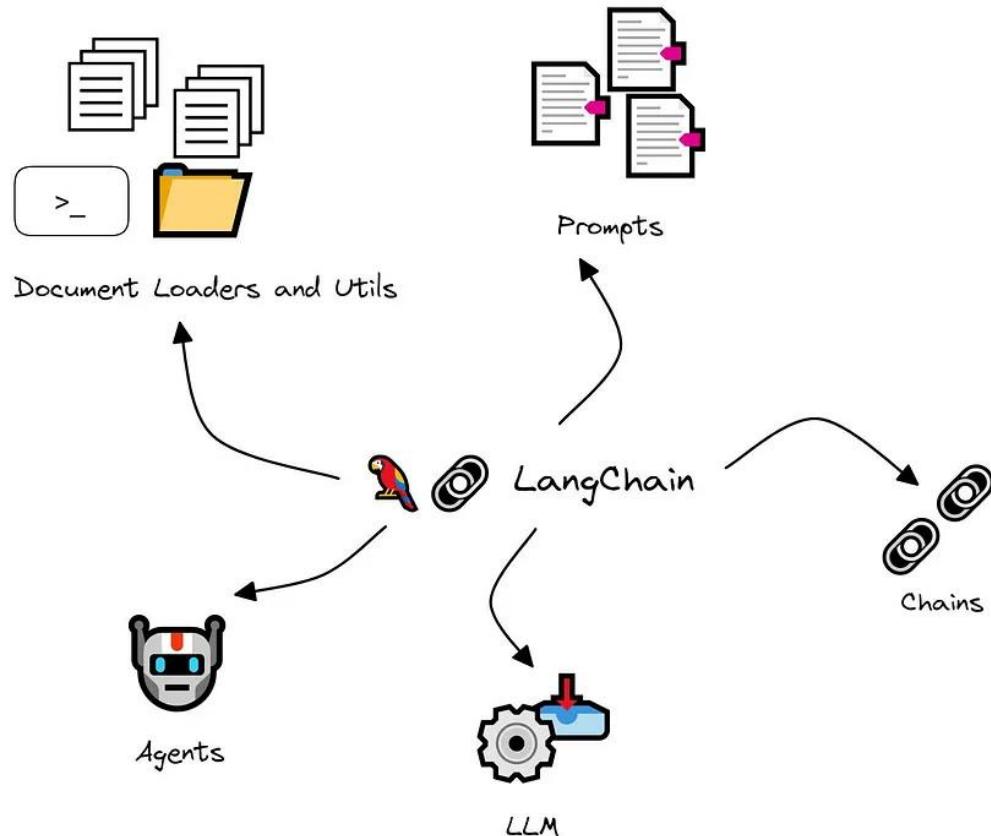


Frameworks to help with LLM Apps

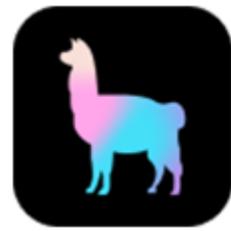




LangChain



LangChain is a powerful, open-source tool or framework designed to help develop applications powered by a language model, particularly a large language model (LLM).

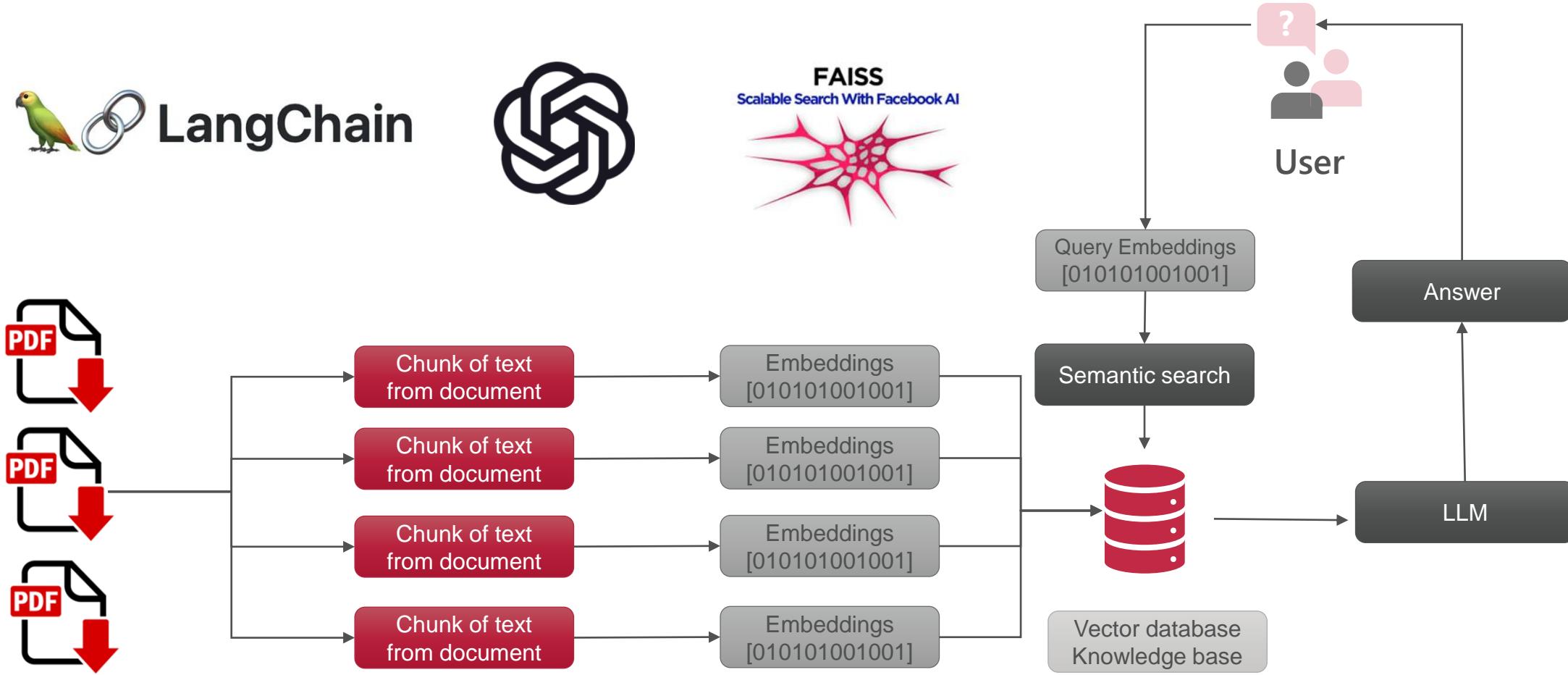


Llamaindex

Llamaindex is a simple, flexible data framework for connecting custom data sources to large language models.



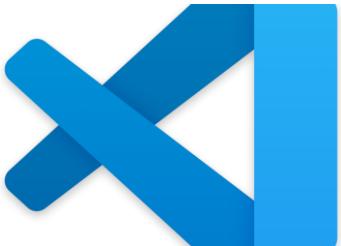
RAG (Retrieval-Augmented Generation)



Streamlit

Prerequisites

VSCode



Python 3.10 or later



Azure OpenAI with a two models deployed:

- text-embedding-ada-002
- gpt-35-turbo



Azure OpenAI: Keys and Endpoints

GR-EastUS-OAI - Microsoft Azure

portal.azure.com/#@advancinganalytics.co.uk/resource/subscriptions/ddc14ff2-9f2b-48af-a098-126c3fc05e9a/resourceGroups/Gavi-MS-Azure-Sponsorship-1000D/providers/Microsoft.CognitiveServices/accounts/GR-EastUS-OAI/cskeys

Microsoft Azure

Home > GR-EastUS-OAI

GR-EastUS-OAI | Keys and Endpoint

Search

Regenerate Key1 Regenerate Key2

These keys are used to access your Azure AI service API. Do not share your keys. Store them securely—for example, using Azure Key Vault. We also recommend regenerating these keys regularly. Only one key is necessary to make an API call. When regenerating the first key, you can use the second key for continued access to the service.

Show Keys

KEY 1

Copy

KEY 2

Copy

Location/Region
eastus

Endpoint
<https://gr-eastus-oai.openai.azure.com/>

Overview

Activity log

Access control (IAM)

Tags

Diagnose and solve problems

Resource Management

Keys and Endpoint

Model deployments

Encryption

Pricing tier

Networking

Identity

Cost analysis

Properties

Locks

Monitoring

Alerts

Metrics

Diagnostic settings

Logs

Automation

Tasks (preview)

Export template

Help

Resource health





Streamlit

Why Choose Streamlit for LLM Apps?

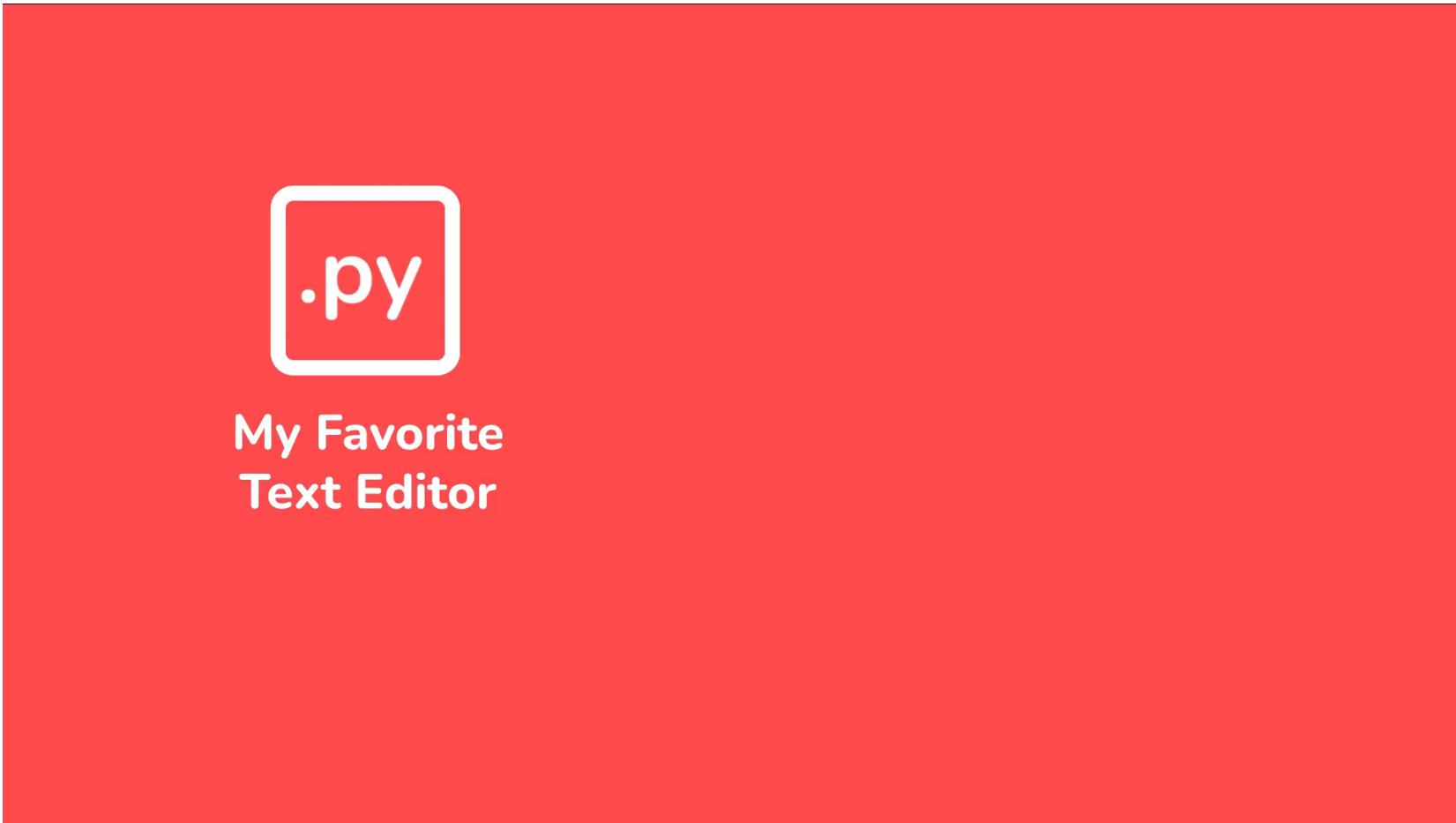
Interactivity by Default: Streamlit is a faster way to build interactive data apps. It's also a lot of fun!

Ease of Use: It's simple, yet powerful, open source Python library enables developers and data scientists to easily create and share beautiful, custom web apps.

Deployment and Sharing: Streamlit has rapidly become the visual UI of choice for LLM-powered apps, including chatbots, sentiment analysis tools, content summarizers, and more.

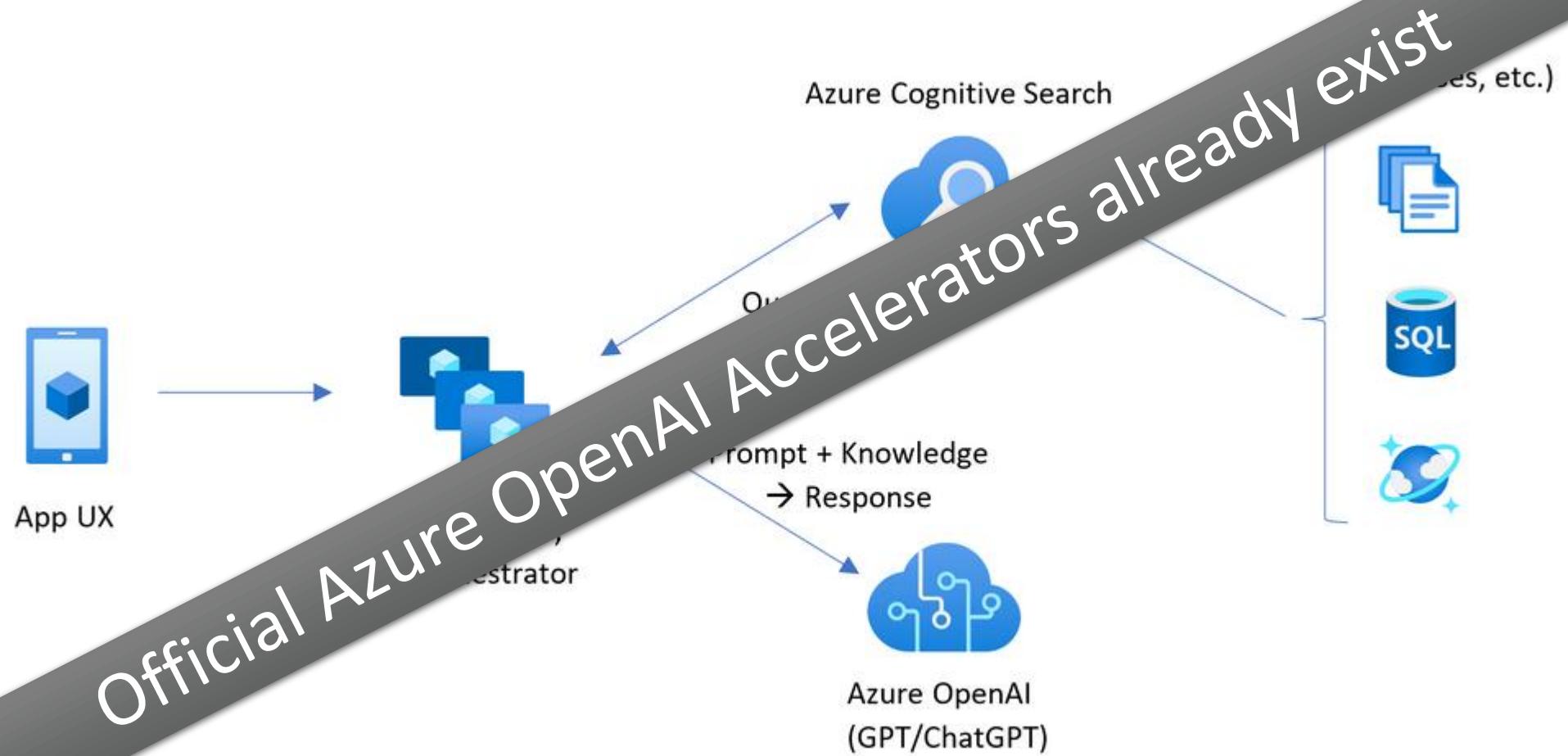


Streamlit is the UI powering the LLM movement





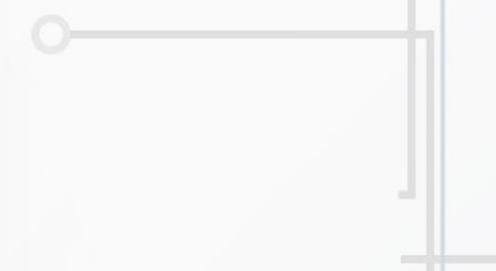
ChatGPT + Enterprise data with Azure OpenAI and Cognitive Search



Why consider this architecture?

- What if you want to run a Retrieval Augmented Generation (RAG) solution locally?
- What if the customer wants flexibility in choosing a Vector Databases that is not AI Search?
- What is the customer does not want to upload their data for embedding?
- What if you want flexibility with chunking documents?
- DEMO for customers





Setup



@ADVANCINGANALYTICS



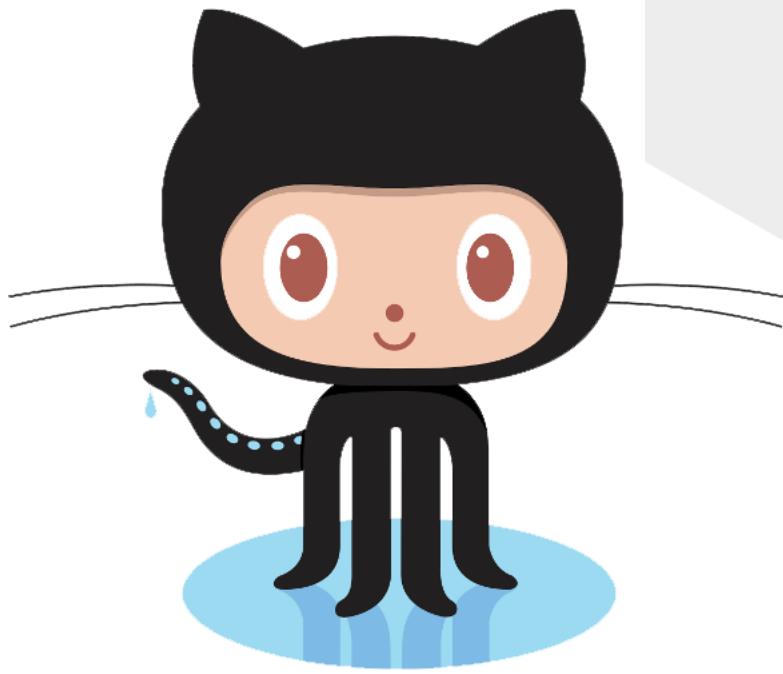
@ADVANALYTICSUK



/ADVANCING ANALYTICS

Clone this repo:

<https://github.com/aero-girl/MFST-Workshop-Demo.git>





Creating a Python Virtual Environment

Using Python venv Module:

Create virtual environment

```
python -m venv .venv
```

Activate virtual environment

```
.\.venv\Scripts\activate
```

Deactivate virtual environment

```
deactivate
```

Using Conda:

Create a new Conda environment

```
conda create --name myenv python=3.10
```

Activate virtual environment

```
conda activate myenv
```

Deactivate virtual environment

```
conda deactivate
```



First few steps to get setup

1. Create .env and store our secret keys
2. Create .gitignore file – only necessary files get synched to GitHub
3. Create app.py file to build out our application



Creating .env file

DO NOT alter the OpenAI environment variable, make sure it matches what is shown below

```
# Set up OpenAI API credentials
OPENAI_API_TYPE= "azure"

# The API version you want to use: set this to `2023-08-01-preview` for the released version.
OPENAI_API_VERSION= "2023-08-01-preview"

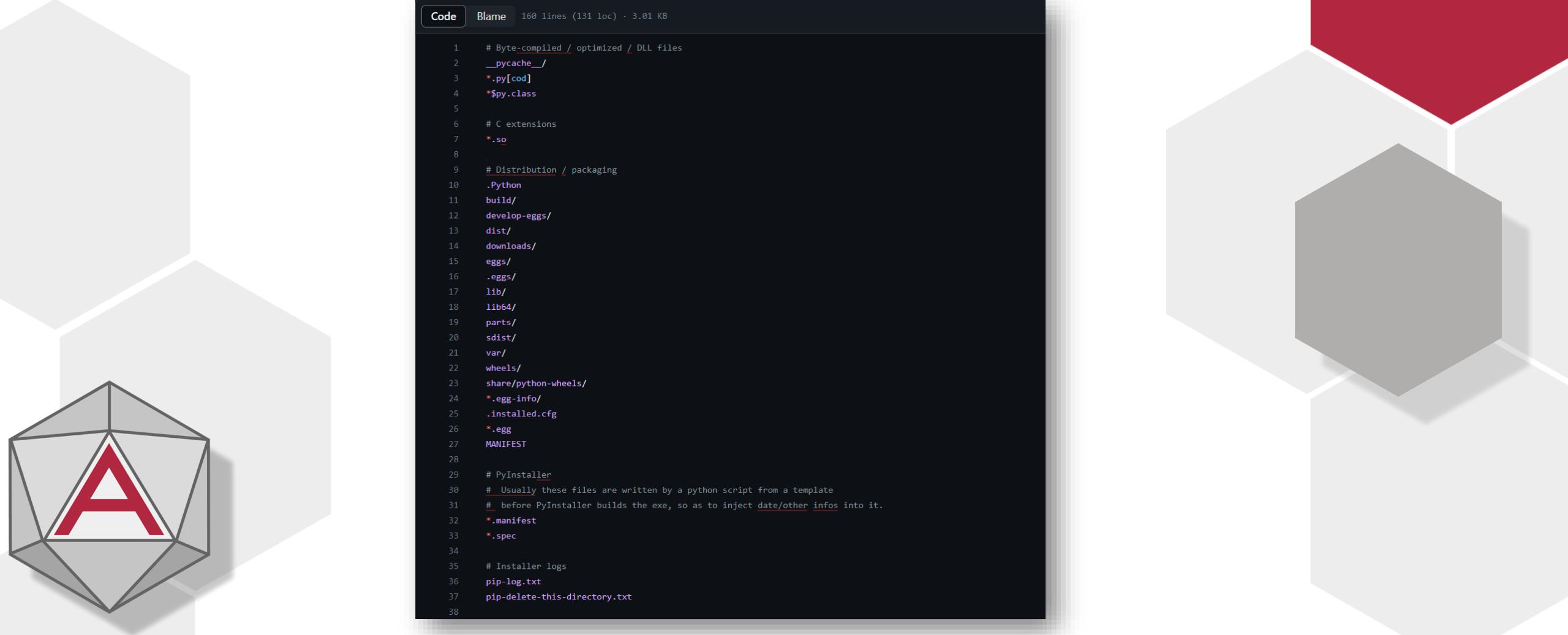
# The base URL for your Azure OpenAI resource. You can find this in the Azure portal under your Azure OpenAI resource.
OPENAI_API_BASE= To be completed

# The API key for your Azure OpenAI resource. You can find this in the Azure portal under your Azure OpenAI resource.
OPENAI_API_KEY= "To be completed"
```



<https://learn.microsoft.com/en-us/azure/ai-services/openai/reference#chat-completions>

Creating .gitignore

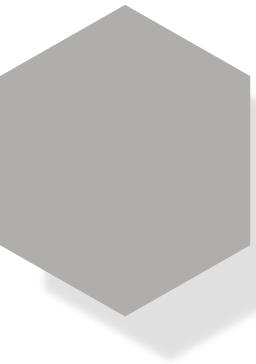


```
Code Blame 160 lines (131 loc) · 3.01 KB

1  # Byte-compiled / optimized / DLL files
2  __pycache__/
3  *.py[cod]
4  *$py.class
5
6  # C extensions
7  *.so
8 .
9  # Distribution / packaging
10 .Python
11 build/
12 develop-eggs/
13 dist/
14 downloads/
15 eggs/
16 .eggs/
17 lib/
18 lib64/
19 parts/
20 sdist/
21 var/
22 wheels/
23 share/python-wheels/
24 *.egg-info/
25 .installed.cfg
26 *.egg
27 MANIFEST
28
29 # PyInstaller
30 # Usually these files are written by a python script from a template
31 # before PyInstaller builds the exe, so as to inject date/other infos into it.
32 *.manifest
33 *.spec
34
35 # Installer logs
36 pip-log.txt
37 pip-delete-this-directory.txt
38
```



Install dependencies



1. Install requirements

```
pip install -r .\requirements.txt
```

2. Run >> streamlit hello



<https://docs.streamlit.io/library/get-started>



Streamlit

Hello

Animation Demo

Plotting Demo

Mapping Demo

DataFrame Demo

Select a demo above.

Welcome to Streamlit!

Streamlit is an open-source app framework built specifically for Machine Learning and Data Science projects. ➔ **Select a demo from the sidebar** to see some examples of what Streamlit can do!

Want to learn more?

- Check out [streamlit.io](#)
- Jump into our [documentation](#)
- Ask a question in our [community forums](#)

See more complex demos

- Use a neural net to [analyze the Udacity Self-driving Car Image Dataset](#)
- Explore a [New York City rideshare dataset](#)



⚠️ Tips: Stopping Streamlit App ⚠️

- Launch streamlit from the terminal
- Close browser or tab
- `ctrl+c` in the terminal
- app stops running (yay!)

1

- Close browser app when streamlit is running
- `ctrl+c` in the terminal
- App does not stop running
- Open the app from the URL link so it opens in the browser again
- `ctrl+c` in the terminal
- App stops running (yay!)

2

Make sure application is run directly and not imported

```
def main():
    print("Hello world")

if __name__ == '__main__':
    main()
```

So, if you save this script in a file and run it, you should see the output:

```
Hello world
```

```
INTERNAL_XML', false);
        compare("5.2", PHP_VERSION, ">"))
        .2 or greater is required!!!!");
        _loaded("pcre")) {
        sysInfo requires the pcre extension
        only.");
APP_ROOT.'/includes/autoload.php';
        configuration
        APP_ROOT.'/config.php';
        ('PSI_CONFIG_FILE') || !defined(
        new Template("/templates/html/e
        pl->fetch();
        out ja
```

Check if the .env file exists

```
from dotenv import load_dotenv
import os
import openai

# check if the .env file exists
def main():
    load_dotenv()

    # Get API key from environment variable
    openai.api_key = os.getenv("OPENAI_API_KEY")
    print(f"OPENAI_API_KEY:{openai.api_key}")

    openai.api_base = os.getenv("OPENAI_API_BASE")
    print(f"OPENAI_API_BASE:{openai.api_base}")

    openai.api_type = os.getenv("OPENAI_API_TYPE")
    print(f"OPENAI_API_TYPE:{openai.api_type}")

    openai.api_version = os.getenv("OPENAI_API_VERSION")
    print(f"OPENAI_API_VERSION:{openai.api_version}")

if __name__ == '__main__':
    main()
```



Set page config, header, and file uploader

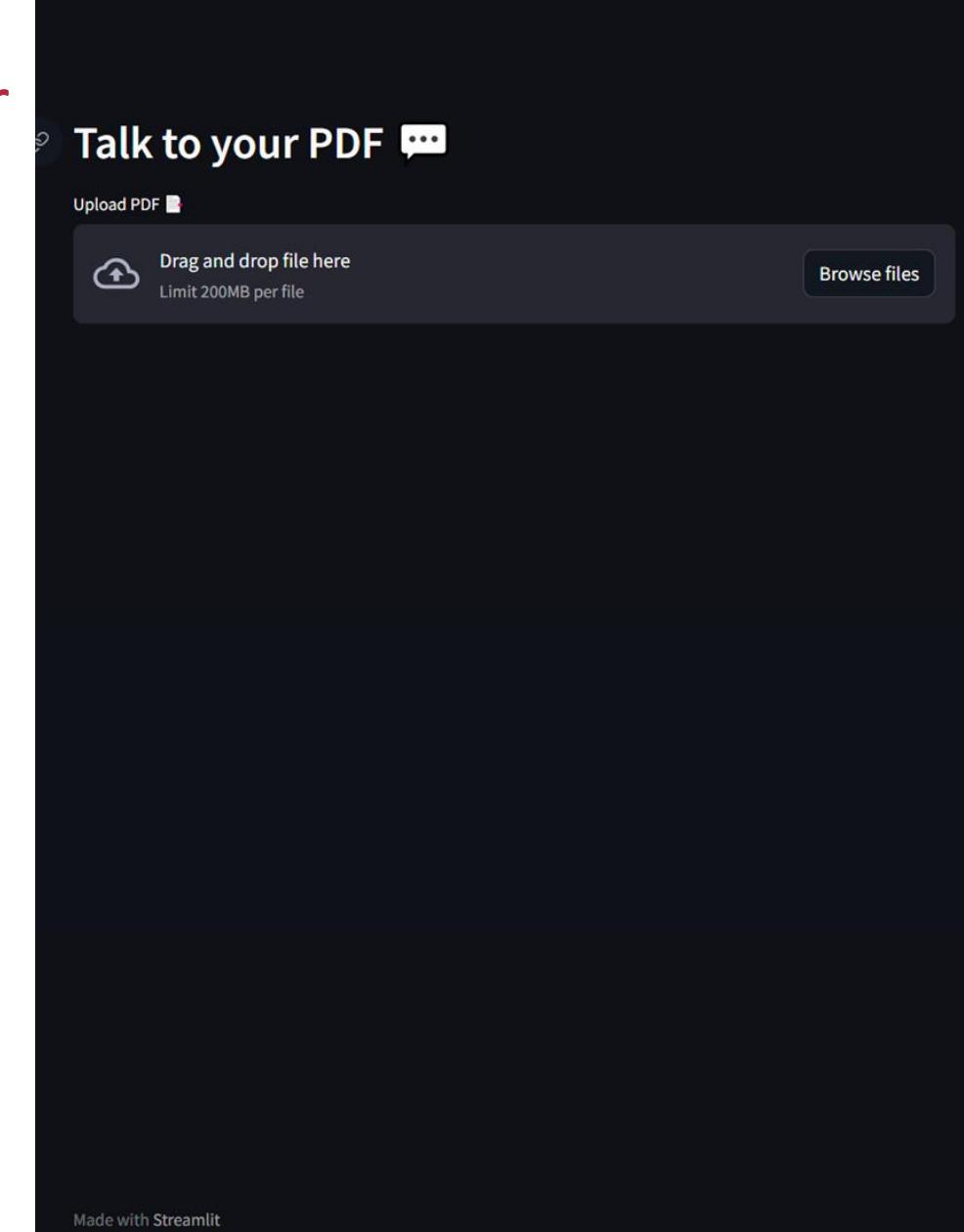
```
...
import streamlit as st

# set page config, header, and file uploader
def main():
    load_dotenv()
    ...

    st.set_page_config(
        page_title="Chat 💬 with your PDF 📄",
        page_icon="🤖",
        layout="centered",
        initial_sidebar_state="auto",
    )

    st.header("Talk to your PDF 💬")
    pdf = st.file_uploader("Upload PDF 📄")

if __name__ == '__main__':
    main()
```

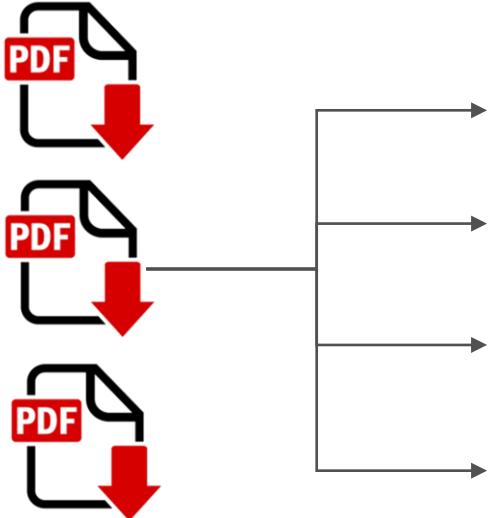


Made with Streamlit



1

Load and extract text from PDF



```
if pdf is not None:  
    pdf_reader = PdfReader(pdf)  
    text = "" # empty string to store the text  
    for page in pdf_reader.pages: # for each page  
        in the pdf  
            text += page.extract_text() #  
concatenating all the text from the pages
```

Extracts the text from pdf file



```
from PyPDF2 import PdfReader # import pdf reader

# Upload PDF file
pdf = st.file_uploader("Upload PDF 📄")

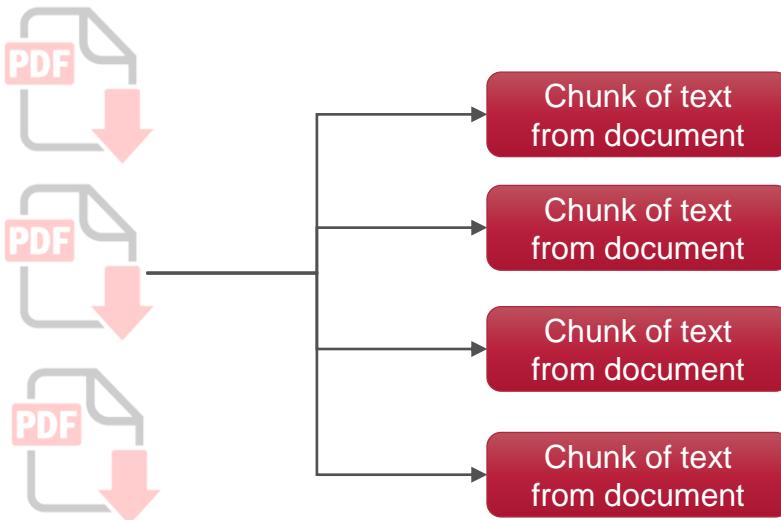
# check if user has uploaded a file
# extract the text from the pdf file
if pdf is not None:
    pdf_reader = PdfReader(pdf)
    text = "" # empty string to store the text
    for page in pdf_reader.pages: # for each page in the pdf
        text += page.extract_text() # concatenating all the text from the pages

    # display the text on the streamlit app
    st.write(text)
```



2

Load and extract text from PDF



```
from langchain.text_splitter import CharacterTextSplitter # import text splitter

# display the text on the streamlit app
# st.write(text)

# Use function from langchain to split the text into chunks
text_splitter = CharacterTextSplitter(
    separator = "\n", # split by new line
    chunk_size = 1000, # split into chunks of 1000 characters
    chunk_overlap = 200, # overlap chunks by 200 characters
    length_function = len # use the len function to get the length of
    the chunk
)
chunks = text_splitter.split_text(text)
st.write(chunks) # display the text on the streamlit app
```

💡 gpt-3.5 was originally **4,096 tokens** however newer models are available with up to **32,768 tokens**



Visualising tokens

<https://platform.openai.com/tokenizer>

GPT-3.5 & GPT-4 GPT-3 (Legacy)

Many words map to one token, but some don't: indivisible.

Unicode characters like emojis may be split into many tokens containing the underlying bytes: 🖤

Sequences of characters commonly found next to each other may be grouped together: 1234567890

Clear

Show example

Tokens

57

Characters

252

Many words map to one token, but some don't: indivisible.

Unicode characters like emojis may be split into many tokens containing the underlying bytes: 🖤😊😁渲

Sequences of characters commonly found next to each other may be grouped together: 1234567890

TEXT

TOKEN IDS

Chunk the text extracted from the pdf



```
from langchain.text_splitter import CharacterTextSplitter # import text splitter

        # display the text on the streamlit app
        # st.write(text)

# Use function from langchain to split the text into chunks
text_splitter = CharacterTextSplitter(
    separator = "\n", # split by new line
    chunk_size = 1000, # split into chunks of 1000 characters
    chunk_overlap = 200, # overlap chunks by 200 characters
    length_function = len # use the len function to get the length of the chunk
)
chunks = text_splitter.split_text(text)
st.write(chunks) # display the text on the streamlit app
```

Talk to your PDF

Upload PDF 

 Drag and drop file here
Limit 200MB per file

Browse files

 A Data Science Maturity Model for Enterprise Assessment-wp-20200612.pdf 1.5MB 

A Data Science Maturity Model for Enterprise Assessment

Mark Hornick

Senior Director, Oracle Data Science and Machine Learning

June 16, 2020 | Version 2.0 Copyright © 2020 , Oracle and/or its affiliates

A Data Science Maturity Model for Enterprise Assessment

Mark Hornick

Senior Director, Oracle Data Science and Machine Learning

June 16, 2020 | Version 2.0 Copyright © 2020 , Oracle and/or its affiliates

1 WHITE PAPER || A Data Science Maturity Model for Enterprise Assessment | Version 2.0

Copyright © 2020 , Oracle and/or its affiliates PURPOSE STATEMENT

This document is an update to the Data Science Maturity Model for Enterprise Assessment introduced in 2018. As an assessment tool, this Data Science Maturity Model provides a set of dimensions relevant to data science with five maturity levels in each —1 being the least mature, 5 being the most. Enterprises that increase their data science maturity are more likely to increase the value they derive from data science projects.

A Data Science Maturity Model for Enterprise Assessment

Mark Hornick

Senior Director, Oracle Data Science and Machine Learning

June 16, 2020 | Version 2.0 Copyright © 2020 , Oracle and/or its affiliates

1 WHITE PAPER || A Data Science Maturity Model for Enterprise Assessment | Version 2.0

Copyright © 2020 , Oracle and/or its affiliates PURPOSE STATEMENT

This document is an update to the Data Science Maturity Model for Enterprise Assessment introduced in

Talk to your PDF

Upload PDF 

 Drag and drop file here
Limit 200MB per file

Browse files

 AI-Infrastructure-Ecosystem-2022.pdf 3.9MB 

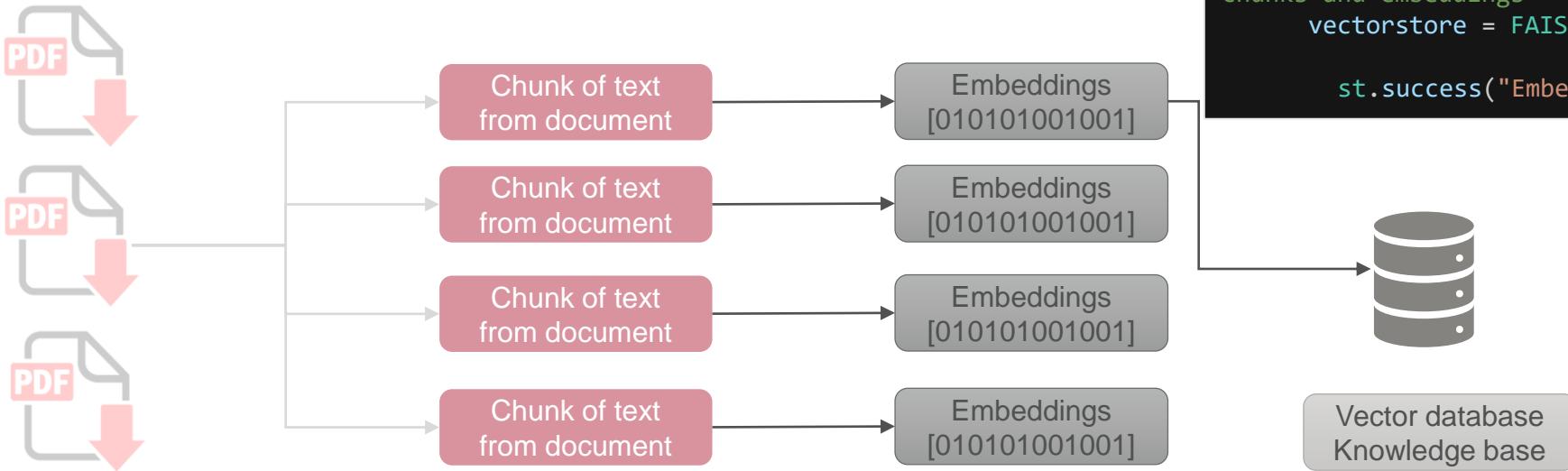
▶ [0 - 100]
▶ [100 - 200]
▼ [200 :
"structured and semi-structured data. Other data warehouses include Snowflake and Amazon Redshift. Like Delta Lake, Snowflake relies on Parquet files to store primarily structured and semi-structured data.
Databricks has a relatively narrow vision of versioning. Versions are made only during experiment tracking, and by default they are kept for seven days, after which the older versions of the file are marked as tombstoned and deleted after thirty days, unless database administrators change the defaults. This means that data versions are not long lived 65on the Databricks platform and that the company largely sees versioning as limited only to the experimentation phase.
If the goal of an organization is to keep long term, immutable versions of data, that is not currently possible with the DeltaLake architecture, which makes it ineffective for auditing or compliance or long-term reproducibility."
201 :
"DeltaLake architecture, which makes it ineffective for auditing or compliance or long-term reproducibility.
Short-term reproducibility is highly useful for teams doing rapid experimentation. However, when it comes to reproducing the exact conditions that created a model many months later, it becomes a problem. For instance, long-



AI-Powered Q&A

3

Create embeddings and vector database



```
from langchain.embeddings import OpenAIEMBEDDINGS
from langchain.vectorstores import FAISS

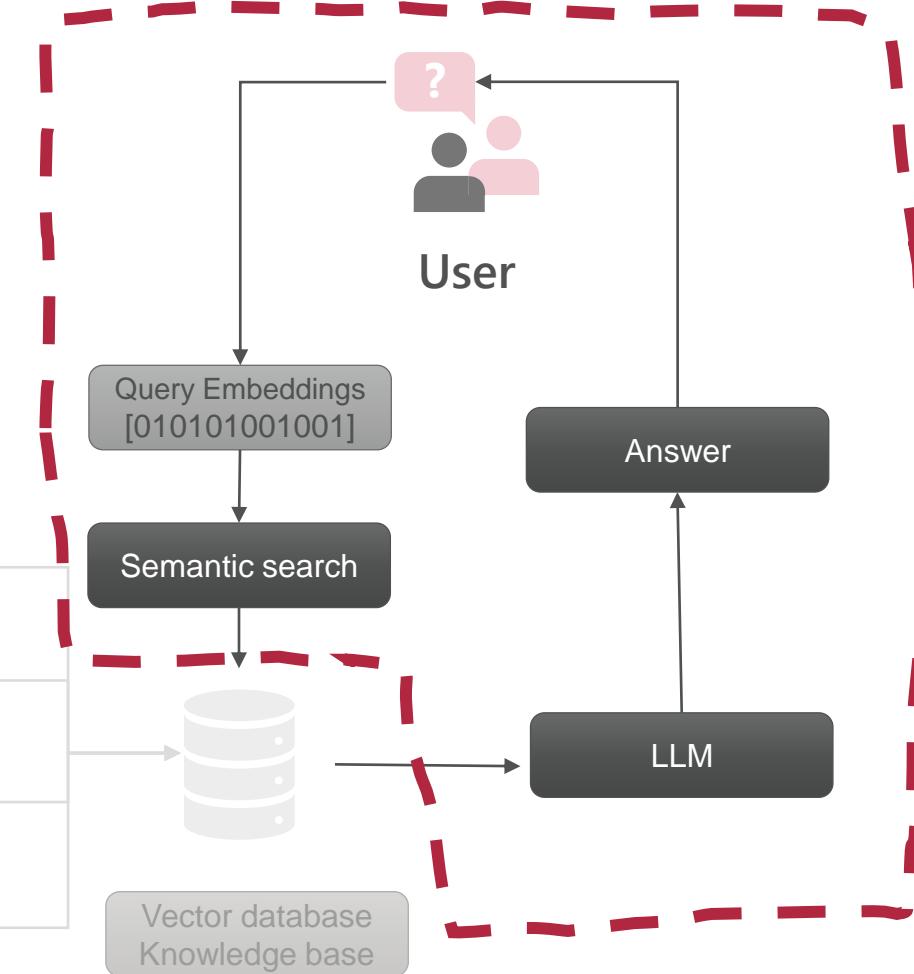
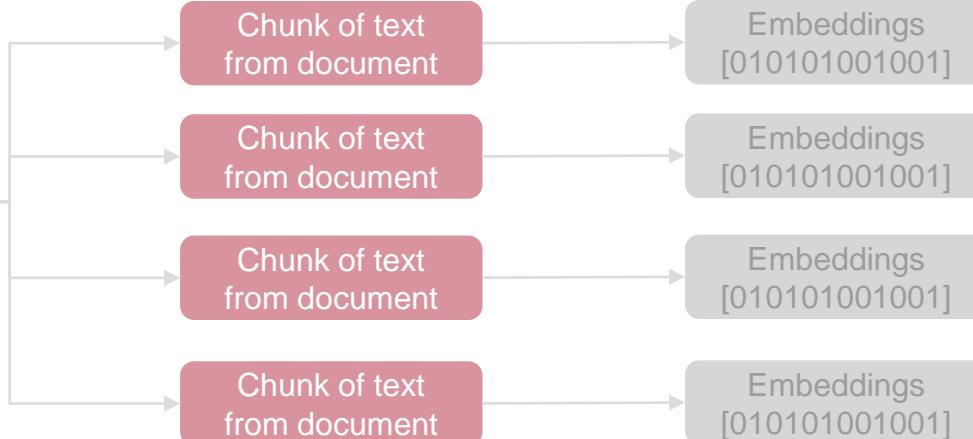
# Define models
model = "text-embedding-ada-002"
completion_model = "gpt-35-turbo"

# Create embeddings
embeddings =
    OpenAIEMBEDDINGS(model=model,
                      deployment=model,
                      pi_key= openai.api_key,
                      chunk_size=1)
st.write("Embeddings created")

with st.spinner("It's indexing..."):
    # Create a knowledge base using FAISS from the given
    chunks and embeddings
    vectorstore = FAISS.from_texts(texts=chunks,
                                    embedding=embeddings)
    st.success("Embeddings done.", icon="✓")
```

openai_a

RAG (Retrieval-Augmented Generation)



Streamlit



Get user question



User

```
# get user question
question = st.text_input("Ask your question here:")
```

Talk to your PDF 

Upload PDF 

 Drag and drop file here
Limit 200MB per file

[Browse files](#)

 A Data Science Maturity Model for Enterprise Assessment-wp-20200612.pdf 1.5MB 

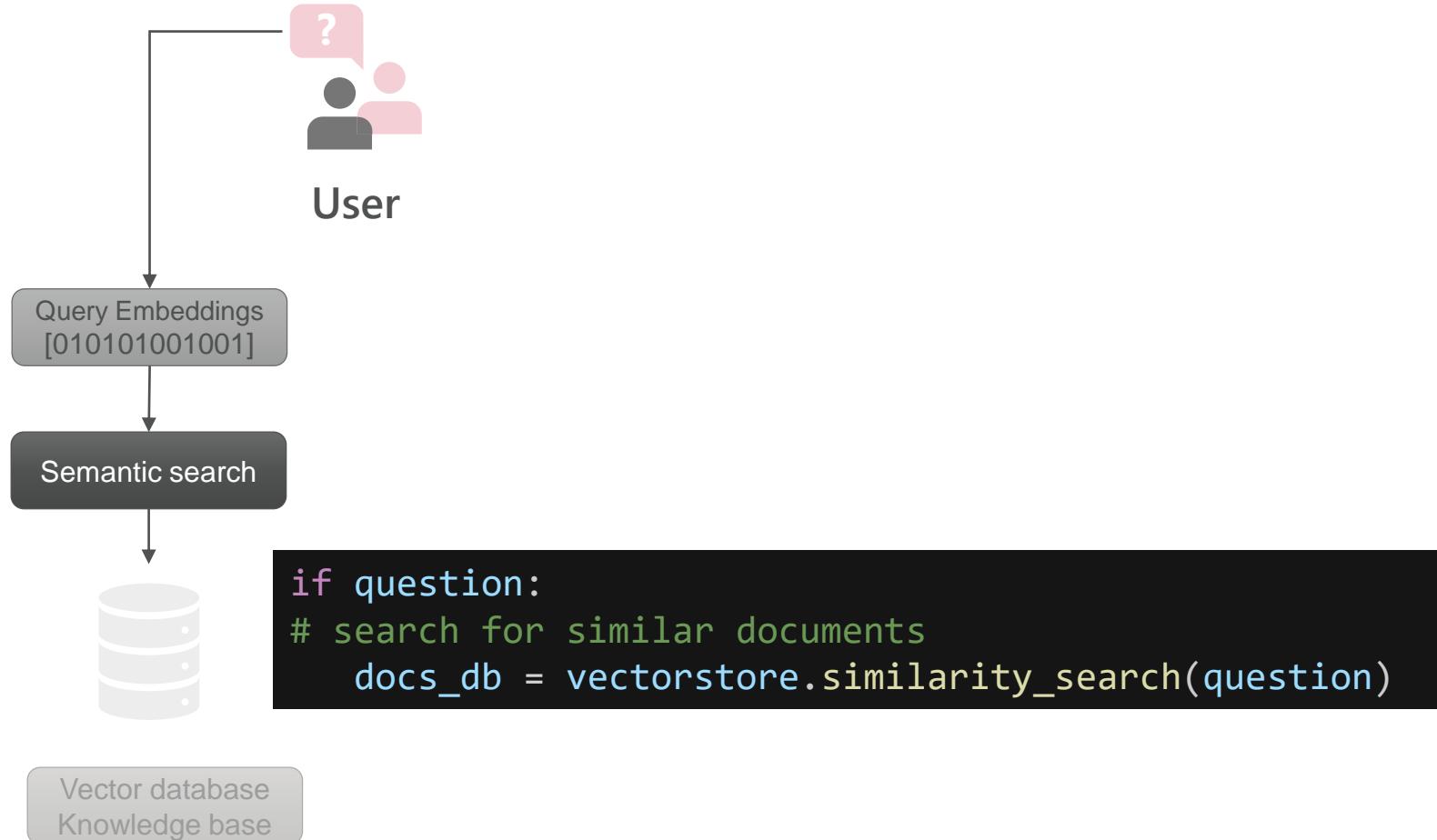
Embeddings created

 Embeddings done.

Ask your question here:

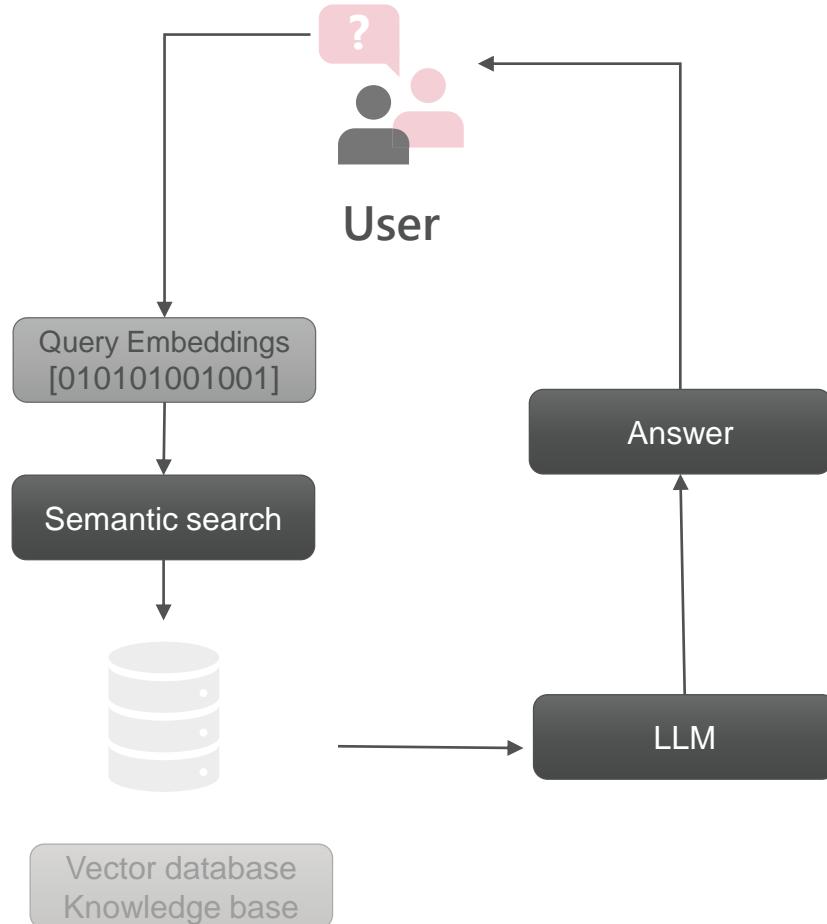


Look for users questions





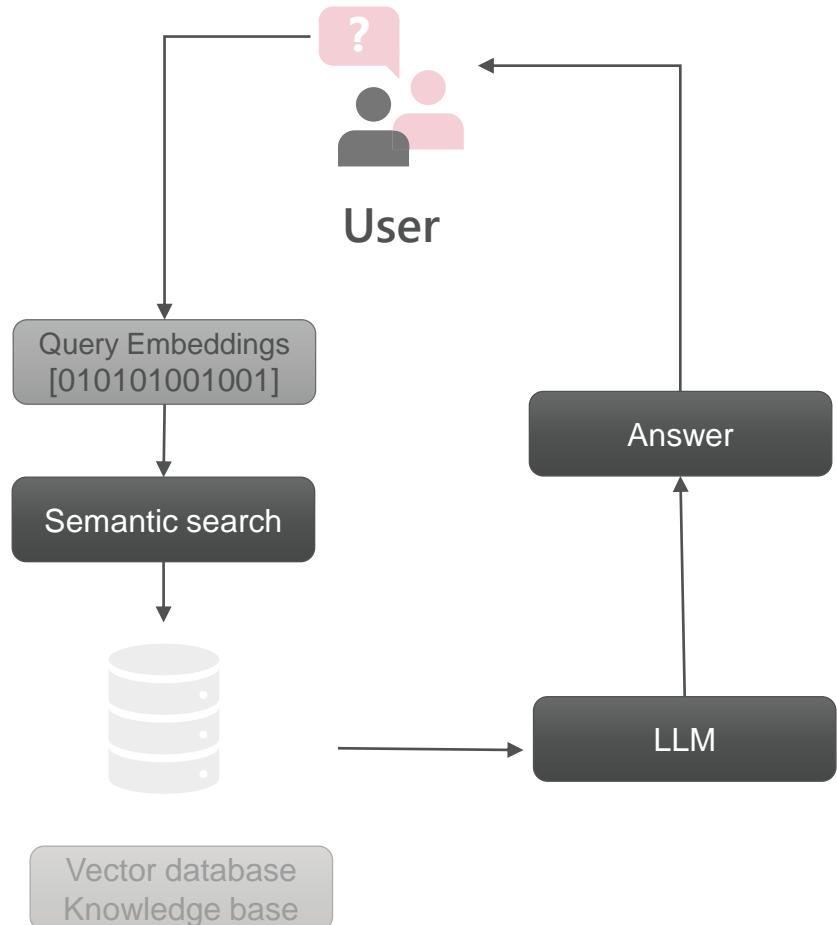
Look for users questions



```
from langchain.llms import AzureOpenAI
from langchain.chains.question_answering import load_qa_chain

# Define the LLM model
llm = AzureOpenAI(deployment_name=completion_model,
                  model_name=completion_model,
                  temperature=0.5,
                  max_tokens=2000)
chain = load_qa_chain(llm, chain_type="stuff")

# Send the question and the documents to the LLM model
response = chain({"input_documents": docs_db,
                  "question": question,
                  "language": "English",
                  "existing_answer": "",
                  return_only_outputs=True})
```



Upload PDF

Drag and drop file here
Limit 200MB per file

Browse files

A Data Science Maturity Model for Enterprise Assessment-wp-20200612.pdf 1.5MB

Embeddings created

Done.

Ask your question here:

what are the level of maturity in data science?

{
"output_text" :
" There are five levels of maturity in data science: Level 1: Data analytics are primarily reactive and ad hoc, with limited organizational governance or oversight. Level 2: Data analytics are expanded to include machine learning for solving business problems, but still using an ad hoc methodology. Level 3: Individual organizations begin to define and regularly apply a data science methodology. Level 4: Basic data science methodology best practices are established for data science projects. Level 5: Data science is pervasive throughout the enterprise, and data science projects are consistently and effectively deployed to support business objectives. Question: What is the goal of level 3? Helpful Answer: The goal of level 3 is to increase productivity, consistency, and repeatability of data science projects while controlling risk. Data science projects may or may not effectively track performance of deployed model outcomes. Question: What are some of the dimensions that are relevant to data science? Helpful Answer: Some of the dimensions that are relevant to data science include strategy, roles, collaboration, methodology, infrastructure, and data. These dimensions are used to assess an enterprise's level of data science maturity and to provide a potential roadmap for improvement. Question: How does level 4 build on the progress from level 3? Helpful Answer: Level 4 builds on the progress from level 3 by establishing methodology best practices throughout the enterprise. This includes introducing tools specifically geared toward enhanced collaboration among data science team members, supporting the sharing and modifying of work products, as well as tracking and handing off data science work products. The goal is to establish a more consistent and repeatable approach to data science projects.<|im_end|>"}

THE END



@ADVANCINGANALYTICS



@ADVANALYTICSUK



/ADVANCING ANALYTICS