

WEBEX and Generative AI

Generative AI and Webex integration - Omer Ilyas

Omer Ilyas

Omer Ilyas - Technical Marketing Engineer - oilyas@cisco.com

Table of contents

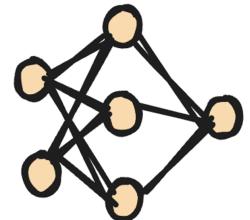
1. Lab	3
1.1 Overview - Understanding AI and Its Integration with Webex	3
1.2 Pre-Requisites and Setup	10
1.3 AI/ML Revolution Unveiled	32
1.4 Tokenization	36
1.5 Task 4: Embeddings and Vector Database	54
1.6 Task 5: Context Windows and Retrieval Augmented Generation (RAG)	78
1.7 Multimodal RAG	96
1.8 GenAI- Frameworks	0
1.9 Understanding Fine-Tuning - Large Language Models	0
1.10 Fine Tuning - Deep Dive into Quantization , LoRA and SFT	0
1.11 Task 8a - Configuring and Fine Tuning - Using llama2	0
1.12 Task 8b - Configuring and Fine Tuning - Using llama3	0
1.13 Task 8c - Configuring and Fine Tuning - Using ChatGPT	0
1.14 Language Model Merging	0
1.15 Logging Out and Ending the Lab Session	0
1.16 Conclusion	0

1. Lab

1.1 Overview - Understanding AI and Its Integration with Webex

Technology and AI

Technology and AI are influencing every aspect of our lives, making the world around us more complex by the day



Artificial Intelligence (AI) is transforming the way we work, enabling innovative solutions and enhancing productivity. Cisco is committed to innovating responsibly, with responsible AI being non-negotiable. Our approach is grounded in the principles of transparency, fairness, accountability, reliability, security, and privacy. Cisco's Responsible AI Framework, aligned with the AI Risk Management Framework, ensures that all AI initiatives undergo rigorous assessments, particularly around privacy.

Remember data has value, Cisco's AI strategy emphasizes building ethical and trustworthy systems that mitigate risks while enhancing innovation across products and services.

**IN THE AGE OF INFORMATION,
I DID NOT KNOW
IS NOT AN EXCUSE**

Please refer to the below info for more detailed insights about Cisco Responsible AI:

- Responsible AI is built on a foundation of privacy
- Cisco's Responsible Artificial Intelligence Principles
- Cisco's Responsible Approach to Governing Artificial Intelligence
- Cisco's Responsible Artificial Intelligence Framework

Responsible AI Framework

Cisco's goal is to provide clarity and consistency in informing users when AI is employed in our technologies

Cisco Trust Portal

trustportal.cisco.com



Your Data Belongs to You



Privacy, Security, and Human Rights by Design, Not by Chance



Trust Is Earned, Transparency Is Paramount

1.1.1 The Evolution of AI at Cisco

At Cisco, innovation is woven into the fabric of everything we do. We've been at the forefront of AI development long before it became a buzzword. Our journey spans decades, from early developments in audio and video intelligence, like echo cancellation, media resilience, to recent breakthroughs like noise removal, face recognition, and immersive experiences, just to name a few. We continue to push boundaries, enhancing AI across multiple domains to deliver immersive and intelligent experiences. Innovation in AI isn't just a trend for us—it's a longstanding commitment to excellence and progress.



What we want everyone to understand is that AI is the core fabric that powers our platforms, enabling reimagining of work, workspaces, and customer experiences. Whether through the Webex Suite, advanced devices, or contact center solutions, AI drives the seamless, intelligent, and immersive experiences that define modern collaboration. Cisco's AI-powered Webex platform is designed to transform the way we work, ensuring that every interaction is smarter, more efficient, and more personalized.



REIMAGINE WORK WITH WEBEX SUITE



REIMAGINE WORKSPACES WITH DEVICES



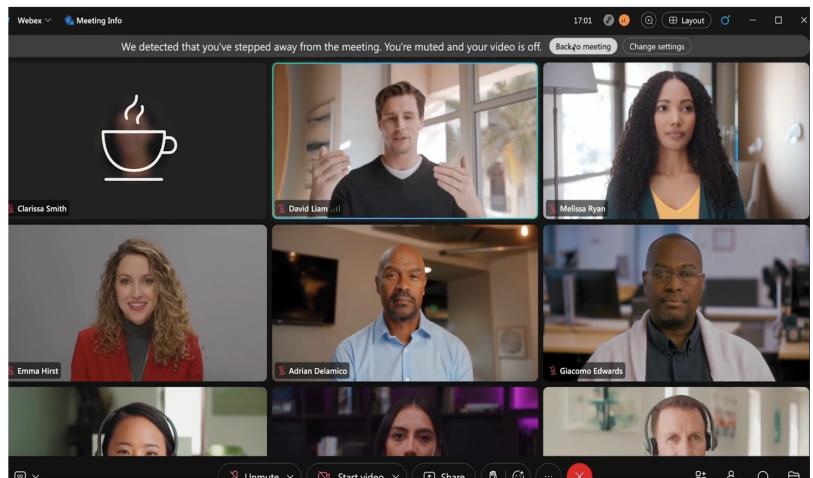
REIMAGINE CUSTOMER EXPERIENCE WITH CONTACT CENTRE

AI-powered Webex Platform

In recent years, we've expanded our platform with several new AI capabilities. Here are a few highlights:

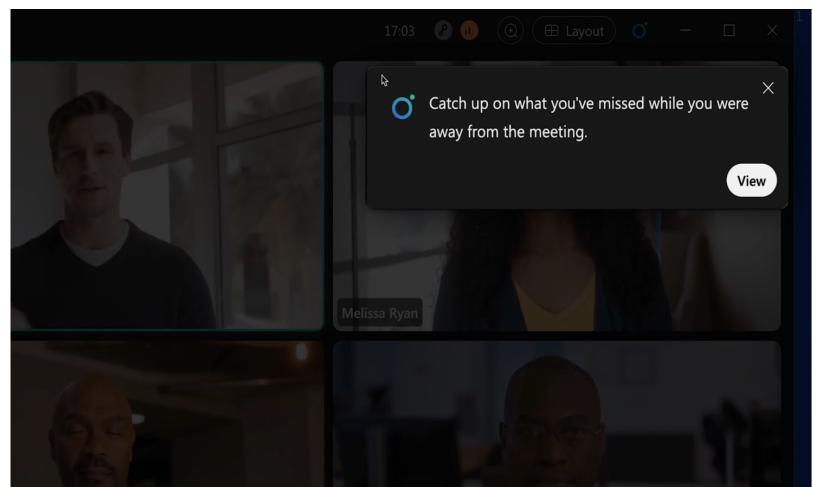
Catch me up in a meeting

Late to the meeting?
Stepped away from the meeting? You can quickly catch up on what you missed.



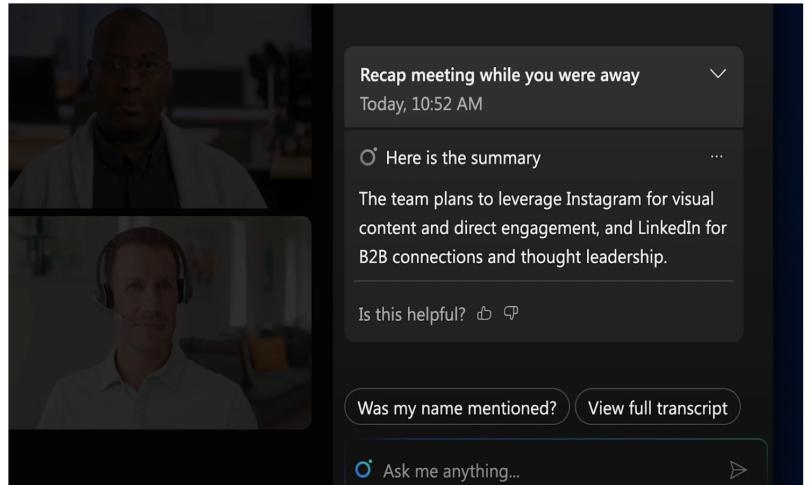
Catch me up in a meeting

Late to the meeting?
Stepped away from the meeting? You can quickly catch up on what you missed.



Catch me up in a meeting

Late to the meeting?
Stepped away from the meeting? You can quickly catch up on what you missed.



Translate Messages with Ease



Cinematic Meetings

Keep everyone engaged and show the best view of people in the room from different angles through adaptive AI directed framing

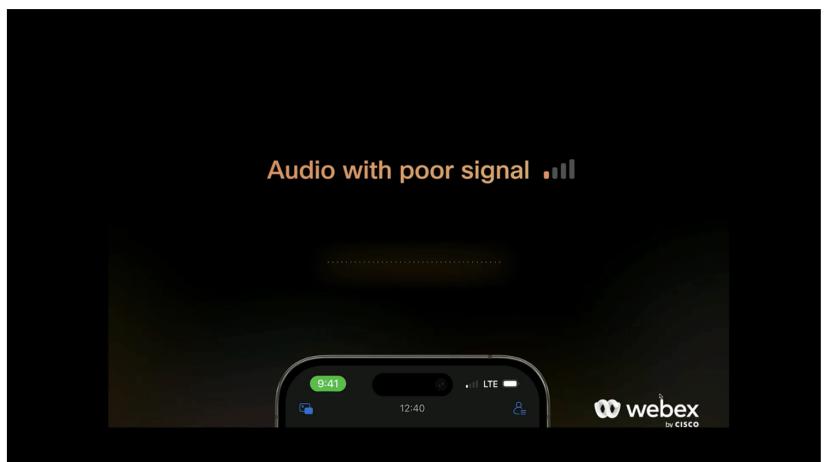


AI Codec Improve audio quality in poor network conditions

Webex AI Codec

Be heard and seen clearly

New AI codec and Super resolution transform audio and video even in the poorest network connection



Webex AI Codec

Compatible
Devices +
Webex

Min
Version

App (44.8)

+ Phase1

(App to
App only)

Macbook

Windows

iPhone 11+

Android

iPad 13+

Be heard and seen clearly

New AI codec and Super resolution transform audio and video even in the poorest network connection



1.1.2 Summary of newly integrated AI features within our platform

Reimagine work and workspaces

Optimize for voice	Meeting Summaries
Smart Chapters	Space Summaries
Background Noise Removal	Rewrite/Translate messages
AI Codec	Multi Stream
HD Voice	Smart Relighting
Be right back	People Recognition in Rooms
Super Resolutions	<u>Slido</u> Topic Summaries
	Video Highlights

Reimagine Customer Experience

New Agent Desktop
New Supervisor Desktop
Topic Analysis
Intelligent Summary
Suggested Responses
Coaching Highlights
Agent Burnout Detection

REMINDER: Our product team is hosting AI sessions. Please review the agenda and join the ones that interest you.

1.1.3 Lab Guide Overview

In this lab guide you will go through the fundamental concepts of AI, various techniques, and how AI can be integrated with Webex to create efficient workflows.

1.1.4 Upon completion of this lab you will be able to

- Understand the basics of AI and its applications
- Learn about embedding techniques
- Vector databases
- Gain insights into Generative AI models
- Familiarize yourself with different AI frameworks
- Integrate AI with Webex to create seamless workflows - Call Recording
- Develop hands-on skills through practical exercises
- Understand Fine-tuning and Quantization
- Deploy Fine Tuning techniques
- Retrieval Augmented Generation (RAG) and Multimodal RAG

1.1.5 Prerequisites

- Basic understanding of AI concepts is helpful but not required.
- We will be using Google Colab for this lab. However, if you have your own machine with an Nvidia chipset and prefer to run the lab locally, feel free to do so.

1.1.6 Disclaimer

The lab design and configuration examples provided are for reference/learning purposes only. This is a sample deployment, and not all recommended features are used or enabled optimally. For design-related questions, please contact your representative at Cisco or a TME team.

1.1.7 Lab Overview - Generative AI and Webex integration

- Lab Setup - Google Colab and Hugging Face
- AI/ML Revolution Unveiled: Task2.md
- Tokenization:
- Embeddings and Vector Database:
- Context Windows and Retrieval Augmented Generation (RAG)
- Multimodal RAG
- Understanding Fine-Tuning for Large Language Models
- Fine Tuning - Deep Dive into Quantization , LoRA and SFT
- Configuring and Fine Tuning - Using llama2
- Configuring and Fine Tuning - Using llama3
- Configuring and Fine Tuning - Using ChatGPT
- Model Merging (Theory Only)

Let's get started! Click on Task 1 - Google Collab- Accessing Google Collab and creating account.

1.2 Pre-Requisites and Setup

1.2.1 Google Collab- Accessing Google Collab and creating account

Google Colab is a free, cloud-based platform that provides a convenient environment for running notebooks. If you want to create a machine learning model but don't have a computer that can handle the workload, Google Colab is the platform for you. In our lab, we will be using Google Colab to test and run our code. However, if you have your own Python environment and prefer to run the code on your local machine, please feel free to do so.

Here are some reasons why using Google Colab can be beneficial for this lab:

- Free Access to GPUs and TPUs: Google Colab offers free access to powerful GPUs and TPUs, which can significantly accelerate the training and fine-tuning of machine learning models.
- No Setup Required: With Colab, there is no need to set up your local environment. Everything runs in the cloud, which saves time and avoids configuration issues.
- Easy Collaboration: Colab notebooks can be easily shared and collaborated on with team members, making it an ideal tool for collaborative projects.
- Integration with Google Drive: Colab integrates seamlessly with Google Drive, allowing you to save and manage your work conveniently.
- Pre-installed Libraries: Many popular machine learning libraries, including TensorFlow and PyTorch, come pre-installed in Colab, making it easy to start working on your projects immediately.

1.2.2 Getting Started With Google Colab

To start working with Google Collaboratory Notebook you first need to log in to your Google account, then go to this link [Google Colab](#)

- Create a new Jupyter Notebook



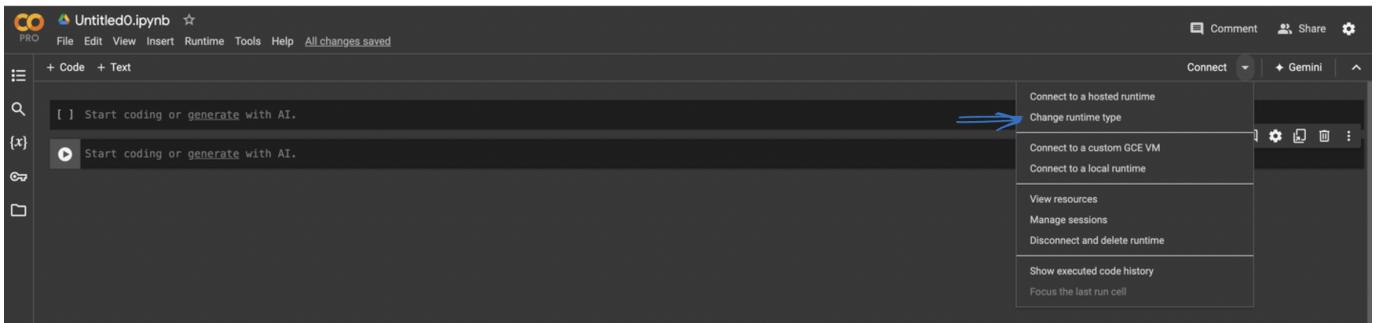
- On creating a new notebook, it will create a Jupyter notebook with Untitled0.ipynb and save it to your google drive in a folder named Colab Notebooks. Now as it is essentially a Jupyter Notebook, all commands of Jupyter Notebooks will work here.



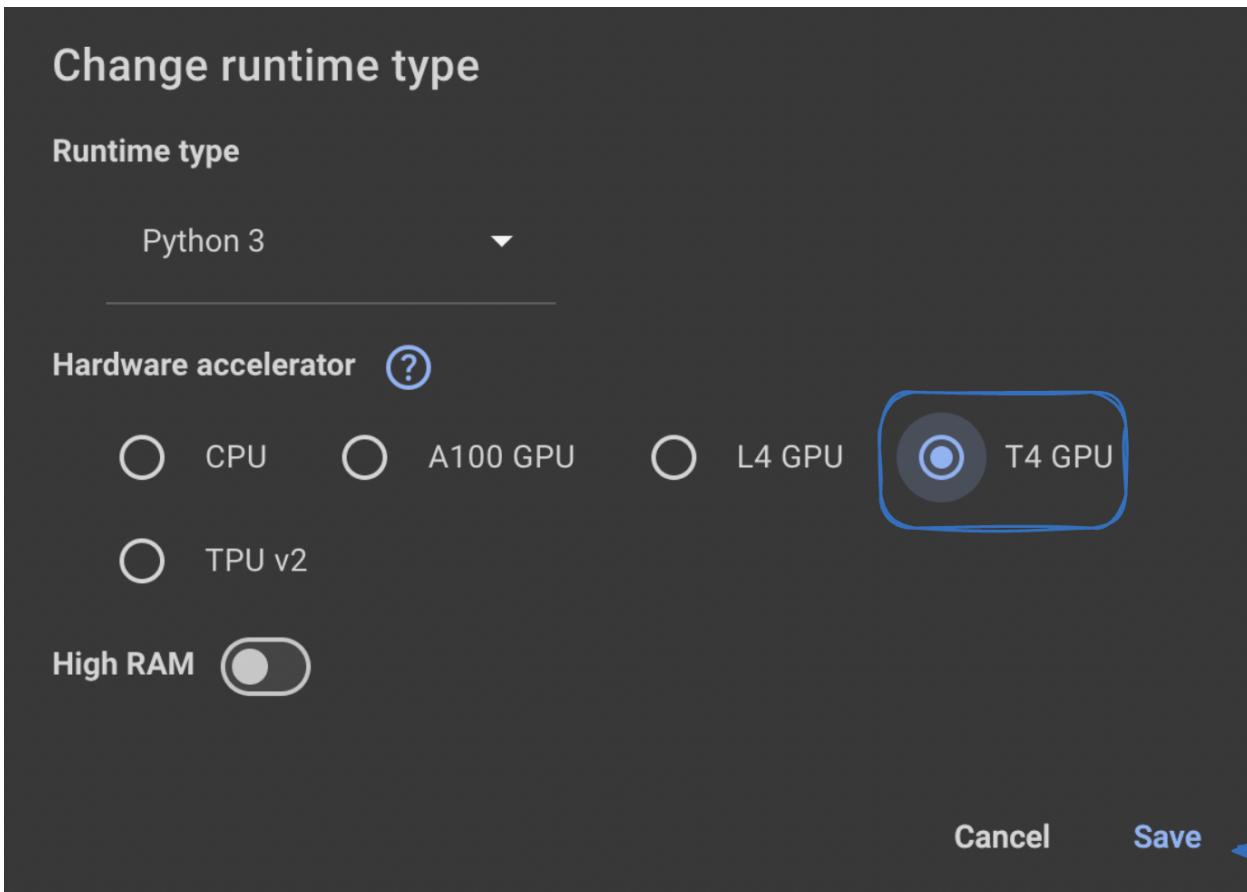
- There might be times when we need to fine-tune models or perform specific tasks that require changing the runtime environment in Colab. Google Colab offers different runtime environments that can be selected based on your requirements:
- Python Versions: You can select between different versions of Python (e.g., Python 2 or Python 3) depending on the compatibility of the code and libraries. We will be using Python3 for our lab.
- Hardware Accelerators: Colab provides access to hardware accelerators, which can be particularly useful for intensive computations. You can choose between:

```
None: No hardware acceleration, suitable for basic tasks.  
GPU: Accelerate your computations with a Graphics Processing Unit.  
TPU: Use a Tensor Processing Unit for even faster performance, especially beneficial for deep learning tasks.
```

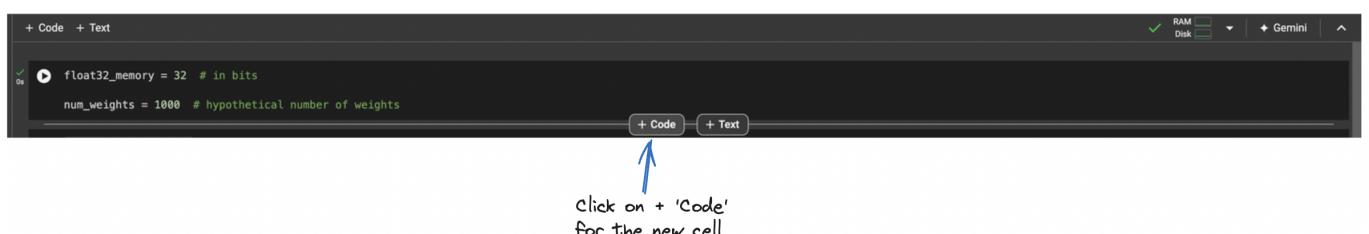
- Change Runtime Environment: Click the “Runtime” dropdown menu at the top of the Colab interface.



- Select “Change runtime type”: This will open a dialog box where you can configure the runtime environment.
- Select Python Version: Choose Python 3 from the “Runtime type” dropdown menu.
- Select Hardware Accelerator: From the “Hardware accelerator” dropdown menu, choose GPU, or TPU based on your needs.



- Save Settings: Click “Save” to apply the changes.
- New Cell: Whenever you want to copy the code in Google Colab and run it, be sure to click on + Code to add a new code cell.



- Execute Code: Click the play button to the left of the code, or use the keyboard shortcut "Command/Ctrl+Enter" while the cell is selected.



1.2.3 Notes: On GPU and TPU Access:

While Google Colab offers free access to GPUs and TPUs, there are limitations. For more consistent access to high-performance GPUs and TPUs, you might need to subscribe to Colab Pro or Colab Pro+ accounts. These paid plans provide priority access to better hardware, longer runtimes, and more memory.

1.2.4 Using Huggingface Hub to share our Datasets

In this lab, we will be utilizing the Hugging Face Hub to load our custom datasets. Hugging Face provides an extensive repository of datasets that can be easily integrated into your machine learning workflows. For the purposes of this lab, we will demonstrate how to access/upload and use our custom datasets effectively.

However, when fine-tuning models in your own work environment, especially if you are using private data, there are important considerations to keep in mind:

- Private Datastores: If you are working with proprietary or sensitive data, it is crucial to use your organization's secure datastores. Ensure that all data handling complies with your organization's data privacy policies and regulations.
- Hugging Face Datasets: If you prefer to use Hugging Face for dataset storage and management, make sure to mark your datasets as private. This setting ensures that your data cannot be accessed by anyone outside your organization, maintaining the confidentiality and integrity of your information. Please refer to Huggingface documentation for more info.

Few more Consideration

- Upload Dataset: When uploading your dataset to Hugging Face, choose the appropriate privacy settings. You can set your dataset to private during the upload process.
- Check Permissions: Regularly review and manage the permissions of your datasets to ensure they remain private and secure.
- Collaborator Access: If you need to share the dataset with specific team members, use the Hugging Face interface to grant access to trusted collaborators only.

By following these guidelines, you can ensure that your data remains secure while leveraging the powerful tools and resources provided by Hugging Face. This approach not only enhances your workflow efficiency but also upholds the best practices in data security and privacy.

1.2.5 Accessing Hugging Face Hub and creating account

Hugging Face can be accessed by browsing to [huggingface](https://huggingface.co)

The screenshot shows the Hugging Face Hub homepage. At the top, there's a search bar with placeholder text "Search models, datasets, users...". Below the search bar, there are navigation links for "Models", "Datasets", "Spaces", "Posts", "Docs", "Solutions", "Pricing", "Log In", and "Sign Up". A banner at the top left says "NEW! AI Tools are now available in HuggingChat". The main content area features a large yellow emoji of a smiling face. Below it, the text "The AI community building the future." is displayed in a large, bold, white font. To the right, there's a sidebar with categories like "Tasks", "Libraries", "Datasets", "Languages", "Licenses", and "Other". Under "Tasks", there are sections for "Multimodal", "Computer Vision", "Natural Language Processing", "Audio", "Tabular", and "Reinforcement Learning". On the right side, there's a list of models with their names, descriptions, and statistics. Some examples include "meta-llama/Llama-2-70b", "stabilityai/stable-diffusion-xl-base-0.9", and "tiiuae/falcon-40b-instruct".

Signing up

- Browse to Hugging Face home page and click on Sign up. Follow the instructions as per below images

The screenshot shows the "Join Hugging Face" sign-up form. It has a yellow emoji at the top. The form asks for an "Email Address" (with placeholder "omer.lyas@boldbetz.com") and a "Password" (with placeholder "*****"). Below the password field are three validation rules: "Must contain at least 8 characters", "Must contain uppercase, lowercase letters, and numbers", and "If less than 12 characters, must contain a special character". At the bottom of the form are "Next" and "Already have an account? Log in" buttons. Blue arrows point from the text "Follow the instructions as per below images" to the "Email Address" and "Password" fields.

Complete your profile

One last step to join the community

Username

Full name

Avatar (optional)

Twitter username (optional)

Upload file

Twitter account

GitHub username (optional)

LinkedIn profile (optional)

GitHub username

LinkedIn profile

Homepage (optional)

Homepage

AI & ML interests (optional)

AI & ML interests

I have read and agree with the [Terms of Service](#) and the [Code of Conduct](#)

Create Account

- Please check your email address for a confirmation link and click to verify your account

Your email address has been verified successfully.

Join an organization

Hugging Face is way more fun with friends and colleagues!

Being part of an organization lets you:

- Show your university, lab or company work to the Hugging Face community
- Collaborate on private datasets, models and spaces
- Subscribe to a paid plan and use Hugging Face hosted training and inference services

Search for an organization

or Create a new organization

- Organization Creation (Optional): While you can upload datasets and fine-tune models directly on Hugging Face without creating an organization, you have the option to create an organization on Hugging Face. This can be particularly useful for team collaboration, as it allows you to upload all your datasets and models in one centralized location.

Optional Step

New Organization

Complete your organization profile

Organization Username: WebexOne

Organization type: Classroom

Homepage (optional): Homepage

Logo (optional): Upload a file

Organization Full name: WebexOne

GitHub username (optional): GitHub alias

Twitter username (optional): Twitter account

AI & ML interests (optional): AI & ML interests

Create Organization

- At this stage, you will see no models or datasets created under your account.

The screenshot shows the Hugging Face Hub organization page for 'WebexOne'. At the top, there's a search bar and a navigation bar with links for Models, Datasets, Spaces, Posts, Docs, Solutions, Pricing, and a profile icon. Below the header, it says 'Classroom' and 'Upgrade to Enterprise'. There's a section for 'AI & ML interests' with a note 'None defined yet.' and a 'Create a Card' button. Under 'Team members', there's one member represented by a purple profile picture. The main content area is titled 'No organization card' with a sub-note 'Customize this page and let people know more about your organization by creating an organization card!' and a 'Create a Card' button. Below this, there are sections for 'Models' (None yet) and 'Datasets' (None yet), each with a blue arrow pointing to the text 'No models or Datasets created yet'.

Models ← No models or Datasets created yet

Datasets ← No models or Datasets created yet

The screenshot shows the user profile menu on the right side of the header. It includes a profile picture (purple circle), the name 'Webex', and a 'Profile' link. Below that are sections for 'Notifications' (Inbox (0)), 'New Model', 'New Dataset', 'New Space', and 'New Collection'. There's also a 'Create organization' button, a 'Settings' link, and a 'Sign Out' link. A blue arrow points from the top of the page to the profile picture in the menu.

Hugging Face API Keys

Create an API Key: As we will be uploading our datasets to the Hugging Face Hub, we need to create an API key for our account. This API key will be used to authenticate and interact with the Hugging Face services programmatically.

you can browse to [huggingface API Key](#)

or

Click on your profile picture > Settings > Access Tokens



The screenshot shows the 'Access Tokens' section of the Hugging Face Hub settings. It includes a 'User Access Tokens' section with a note about token authentication and a 'Create new token' button. The 'Access Tokens' link in the sidebar is circled with a blue circle labeled '1'. A blue arrow points from the 'Create new token' button to another circled '2'.

Under the "Access Tokens" section, click on "Create new token." You will see options to select the token type and provide a token name. For example, you might name your token "Webexone" and select the appropriate permissions.

- fine-grained: tokens with this role can be used to provide fine-grained access to specific resources, such as a specific model or models in a specific organization. This type of token is useful in production environments, as you can use your own token without sharing access to all your resources.
- read: tokens with this role can only be used to provide read access to repositories you could read. That includes public and private repositories that you, or an organization you're a member of, own. Use this role if you only need to read content from the Hugging Face Hub (e.g. when downloading private models or doing inference).
- write: tokens with this role additionally grant write access to the repositories you have write access to. Use this token if you need to create or push content to a repository (e.g., when training a model or modifying a model card).

As we have a lab environment we will be using the "Write" permission. This token will have read and write access to all your resources and can make calls to inference API on your behalf, as shown in the image below.

Create new Access Token

Token type

Fine-grained Read Write ← 1

! This cannot be changed after token creation.

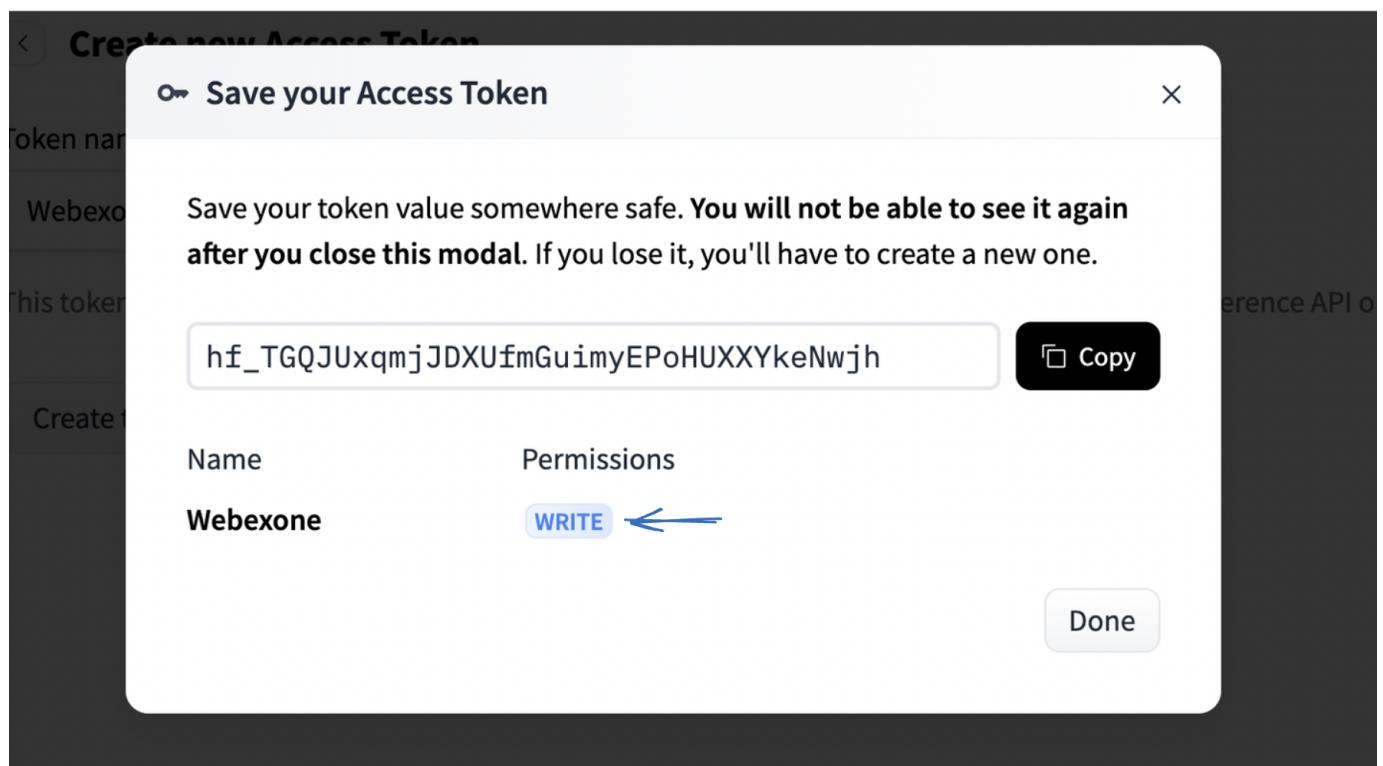
Token name

Webexone

This token has read and write access to all your and your orgs resources and can make calls to inference API on your behalf.

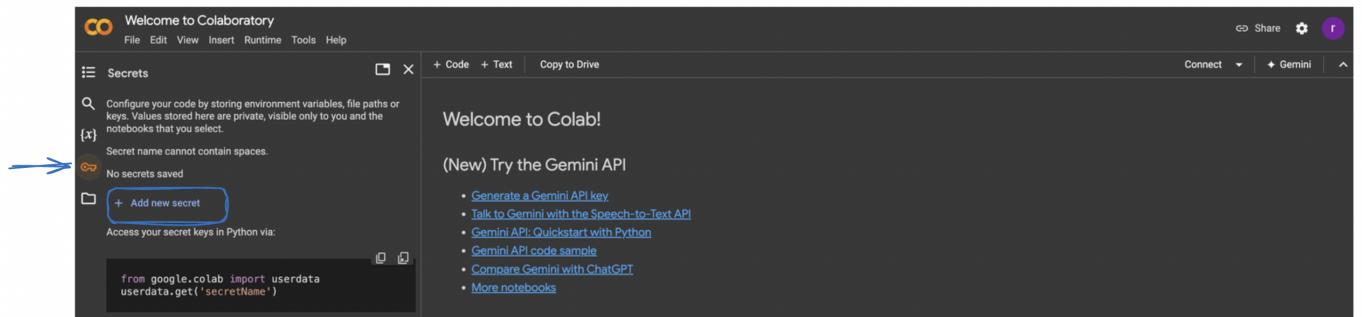
Create token ← 2

Save and Secure the Token: Once the token is generated, save it securely. This token will be required for accessing and managing your datasets via the API. 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



1.2.6 Accessing Hugging Face API in Google Colab

- Open the Google Colab notebook and navigate to the new “Secrets” section in the sidebar.



- Click on “Add a new secret.” Enter the name example: HF_TOKEN and value of the secret. Note: The name is permanent once set.
- The list of secrets is global across all your notebooks.
- Use the “Notebook access” toggle to grant or revoke access to a secret for each notebook.

The screenshot shows the 'Secrets' sidebar in Google Colab. It includes a search icon, a note about secret names, and a table for managing secrets. The table has columns for 'Notebook access', 'Name', 'Value', and 'Actions'. A secret named 'HF_TOKEN' is listed with its value partially visible as 'hf_TGQJUxq...' and a series of icons for edit, copy, and delete. Handwritten annotations include arrows pointing to the '+ Add new secret' button and the 'Value' column, with text below them: 'Give it a meaningful name' and 'The secret from Huggingface you created earlier'.

Notebook access	Name	Value	Actions
<input checked="" type="checkbox"/>	HF_TOKEN	hf_TGQJUxq...	

Give it a
meaningful
name

The secret from Huggingface
you created earlier

Optional Steps below

Incorporating Secrets into Your Code - We will use it later in our lab

- To use a secret in your notebook, use the following code snippet

```
from google.colab import userdata  
my_secret_key = userdata.get('HF_TOKEN')
```

- Replace with your secret's name.

Using Secrets as Environment Variables - Optional Step , we will use it later in our lab

- For Python modules requiring API keys as environment variables, use the below snippet:

```
# Import Colab Secrets userdata module  
  
from google.colab import userdata  
import os  
  
# Set other API keys similarly  
os.environ["HF_TOKEN"] = userdata.get('HF_TOKEN')
```

1.2.7 Getting Started with Langchain

Let's explore how to create a LangChain account and obtain the API key, which we will use later in our lab.

To start building applications with LangChain, you'll need an API key. This key allows your applications to securely connect with LangChain's services, ensuring proper authentication and usage tracking.

Lets start by visiting the official LangChain website and create an account. You'll need to enter some basic information about yourself or your organization. In this example, I will be using my Google account to sign up.



LangSmith

Data Region  US ▾

Create an account

 Continue with Github

 Continue with Google

 Continue with Discord

OR

[Sign up with email](#)

Already have an account? [Log in](#)

By continuing, you agree to our [Terms of Service](#) and [Privacy Policy](#).
Data Security is important to us. To learn more about our policies and certifications, visit trust.langchain.com for more information.

To create an API key head to the Settings page. Then click Create API Key.



Setup Langchain envoirnment - Sample Code

- In the code snippet below, we'll configure our environment to enable tracing of AI calls and set up the API key for interacting with LangChain's services, as we will be using these in the upcoming section.

```
1 !pip install langchain-community
2 export LANGCHAIN_TRACING_V2=true
3 export LANGCHAIN_API_KEY=<your-api-key>
```

1.2.8 Optional Step - Running Ollama locally

Have you ever used GPT and amazed at its ability to understand and respond to your queries? But did you know that you can also harness the power of open-source models using Ollama? With Ollama, you can download and interact with these models directly, getting responses that are tailored to your needs.

In this lab, I'll show you how to install Ollama locally on your machine and start using it to generate responses. You'll learn how to download and load open-source models, and then use Ollama's intuitive interface to interact with them. More info for Ollama can be found [here](#)

Note: The steps below are provided for informational purposes. If you are using the demo laptop in this lab, feel free to follow these instructions. However, if you are using your personal or work machine, please ensure that you have the necessary privileges and authorization from your organization's administrator to install Ollama locally.

Prerequisites for Installing Ollama Locally

- Ensure that Docker is installed and running on your machine. It is available for various operating systems, including macOS, Windows, and Linux. You can download it from the official Docker website and follow the installation instructions for your specific OS.

Installing Ollama Locally

Note: Since we're using a Mac, I'll demonstrate how you can install it on macOS.

- First, you need to install Ollama. You can download it here and clicking on the download button. Follow the instructions.
- Once installed, Ollama functions as a command-line application, allowing you to interact with it directly through the terminal. To get started, open the terminal and enter the following command:

```
SHELL
```

```
ollama
```

It will output the below commands

Usage:

```
ollama [flags]  
ollama [command]
```

Available Commands:

serve	Start ollama
create	Create a model from a Modelfile
show	Show information for a model
run	Run a model
pull	Pull a model from a registry
push	Push a model to a registry
list	List models
cp	Copy a model
rm	Remove a model
help	Help about any command

Flags:

-h, --help	help for ollama
-v, --version	Show version information

- Ollama supports a list of models available here. Below are some example models that can be downloaded:

Model	Parameters	Size	Download
Llama 3.1	8B	4.7GB	<code>ollama run llama3.1</code>
Llama 3.1	70B	40GB	<code>ollama run llama3.1:70b</code>
Llama 3.1	405B	231GB	<code>ollama run llama3.1:405b</code>
Phi 3 Mini	3.8B	2.3GB	<code>ollama run phi3</code>
Phi 3 Medium	14B	7.9GB	<code>ollama run phi3:medium</code>
Gemma 2	2B	1.6GB	<code>ollama run gemma2:2b</code>
Gemma 2	9B	5.5GB	<code>ollama run gemma2</code>
Gemma 2	27B	16GB	<code>ollama run gemma2:27b</code>
Mistral	7B	4.1GB	<code>ollama run mistral</code>
Moondream 2	1.4B	829MB	<code>ollama run moondream</code>
Neural Chat	7B	4.1GB	<code>ollama run neural-chat</code>
Starling	7B	4.1GB	<code>ollama run starling-lm</code>
Code Llama	7B	3.8GB	<code>ollama run codellama</code>
Llama 2 Uncensored	7B	3.8GB	<code>ollama run llama2-uncensored</code>
LLaVA	7B	4.5GB	<code>ollama run llava</code>
Solar	10.7B	6.1GB	<code>ollama run solar</code>

- To download a model, for example, gemma:2b, which is a lighter model, we can simply type:

SHELL

```
ollama pull gemma:2b
```

```
Available Commands:
  serve      Start ollama
  create     Create a model from a Modelfile
  show       Show information for a model
  run        Run a model
  pull       Pull a model from a registry
  push       Push a model to a registry
  list       List models
  ps         List running models
  cp         Copy a model
  rm         Remove a model
  help      Help about any command
```

Flags:

```
-h, --help      help for ollama
-v, --version   Show version information
```

Use "ollama [command] --help" for more information about a command.

```
[malek@Malek-iMacs-iMac ~ % ollama run
Error: requires at least 1 arg(s), only received 0
[malek@Malek-iMacs-iMac ~ % ollama pull gemma:2b
pulling manifest
pulling c1864a5eb193... 89% ━━━━━━ | 1.5 GB/1.7 GB 20 MB/s 9s
```

- After pulling the model we can interact with it in the terminal by typing:

```
SHELL
ollama run gemma:2b

[malek@Malek-iMac ~ % ollama pull gemma:2b
pulling manifest
pulling c1864a5eb193... 100% ━━━━━━ 1.7 GB
pulling 097a36493f71... 100% ━━━━━━ 8.4 KB
pulling 109037bec39c... 100% ━━━━ 136 B
pulling 22a838ceb7fb... 100% ━━━━ 84 B
pulling 887433bb9a90... 100% ━━━━ 483 B
verifying sha256 digest
writing manifest
success
[malek@Malek-iMac ~ % ollama run gemma:2b
```

- We can now ask our question directly in the terminal



```

malek@Malek-iMacs-iMac ~ % ollama run gemma:2b
>>> Can you tell me about Cisco
Sure, here's a summary about Cisco:

**Cisco Systems** is a global leader in networking and digital solutions. They develop, manufacture, and sell networking equipment, software, and services to various industries, including service providers, governments, and enterprises.

**Key facts about Cisco:**

* Founded in 1984
* Headquartered in San Jose, California
* Has over 100,000 employees worldwide
* Offers a wide range of networking products and services
* Is known for innovation, quality, and customer support

**Some of the most popular Cisco products and services include:**

* **Networking equipment:** Switches, routers, firewalls, VPNs, and more
* **Software:** Cisco IOS, Cisco Catalyst, and other network management tools
* **Cloud services:** Cisco Cloud and Cisco Virtual Network
* **Security solutions:** Cisco SecureNet, Cisco TrustRadius, and other security products

**Here are some of the benefits of using Cisco products and services:**

* **Increased network efficiency and performance**
* **Enhanced security and compliance**
* **Improved collaboration and communication**
* **Reduced IT costs**

**Some potential drawbacks of using Cisco products and services:**

* **High price point for some products and services**
* **Complex technology can be difficult to manage**
* **Customer service can be expensive or difficult to reach**

Overall, Cisco is a leading provider of networking solutions that can meet the needs of a variety of businesses and organizations.

I hope this information is helpful. Please let me know if you have any other questions.

>>> End a message (/? for help)

```

- To remove model you can type

SHELL

```
ollama rm gemma:2b
```

- To see the models installed, just enter

SHELL

```
ollama list
```

Web Interface

Once we've installed Ollama, we can interact with it directly through the terminal, as demonstrated earlier. But what if you prefer a web interface? There are several open-source tools available, but I'll show you how to use Open WebUI, which was previously known as Ollama WebUI.

Open WebUI is a versatile, feature-rich, and user-friendly web interface that you can host yourself and use entirely offline. It offers a ChatGPT-style interface, allowing you to interact with language models running on locally (Ollama). This tool is especially useful for those who want to run language models locally or in a self-hosted environment, ensuring both data privacy and control.

- After installing Ollama, simply run the following Docker command to set up the interface

SHELL

```
docker run -d -p 3000:8080 --add-host=host.docker.internal:host-gateway -v open-webui:/app/backend/data --name open-webui --restart always ghcr.io/open-webui/open-webui:main
```

```

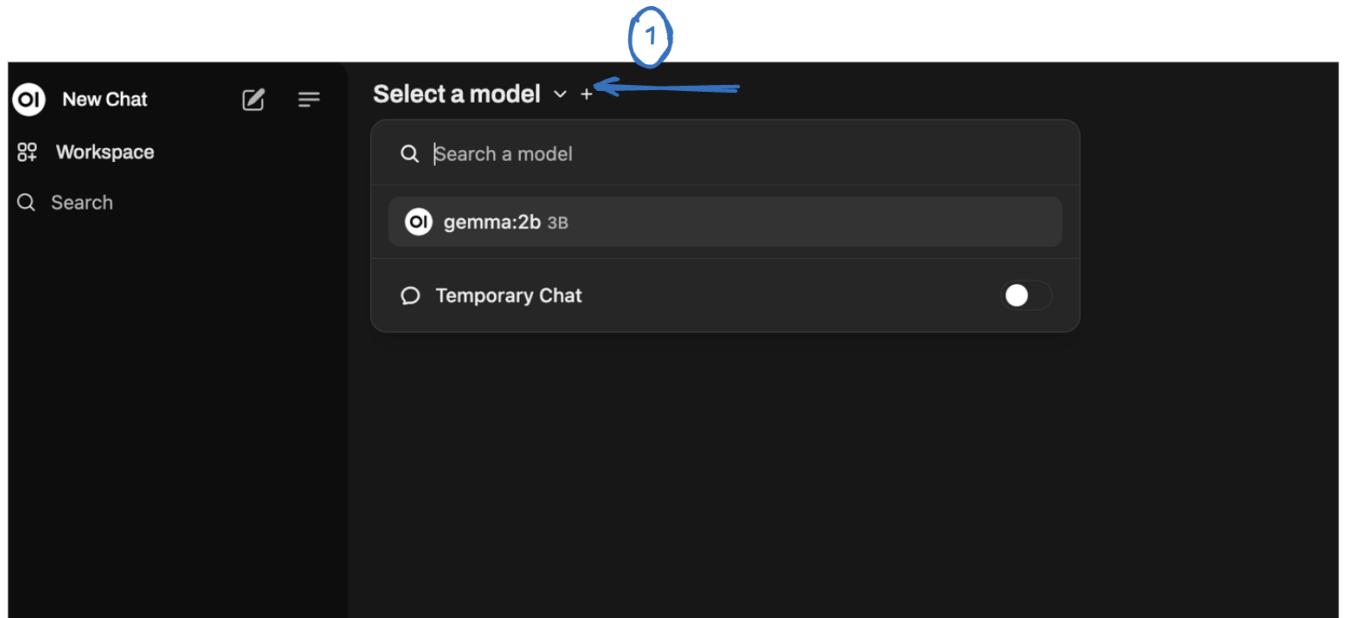
malek@Malek-iMacs-iMac ~ % docker run -d -p 3000:8080 --add-host=host.docker.internal:host-gateway -v open-webui:/app/backend/data --name open-webui --restart always ghcr.io/open-webui/open-webui:main
Unable to find image 'ghcr.io/open-webui/open-webui:main' locally
main: Pulling from open-webui/open-webui
e4fff0779a6d: Pull complete
d97016d0706d: Pull complete
53db1713e5d9: Pull complete
a8cd795d9cc0: Pull complete
de3ba92de392: Pull complete
6f4d87c224b0: Pull complete
4f4fb700ef54: Pull complete
dd92a6022ddb: Pull complete
bbbfded48a772: Pull complete
a825beebdb5b: Downloading [=====] 256.8MB/294.4MB
fd0f6bc6022b: Download complete

```

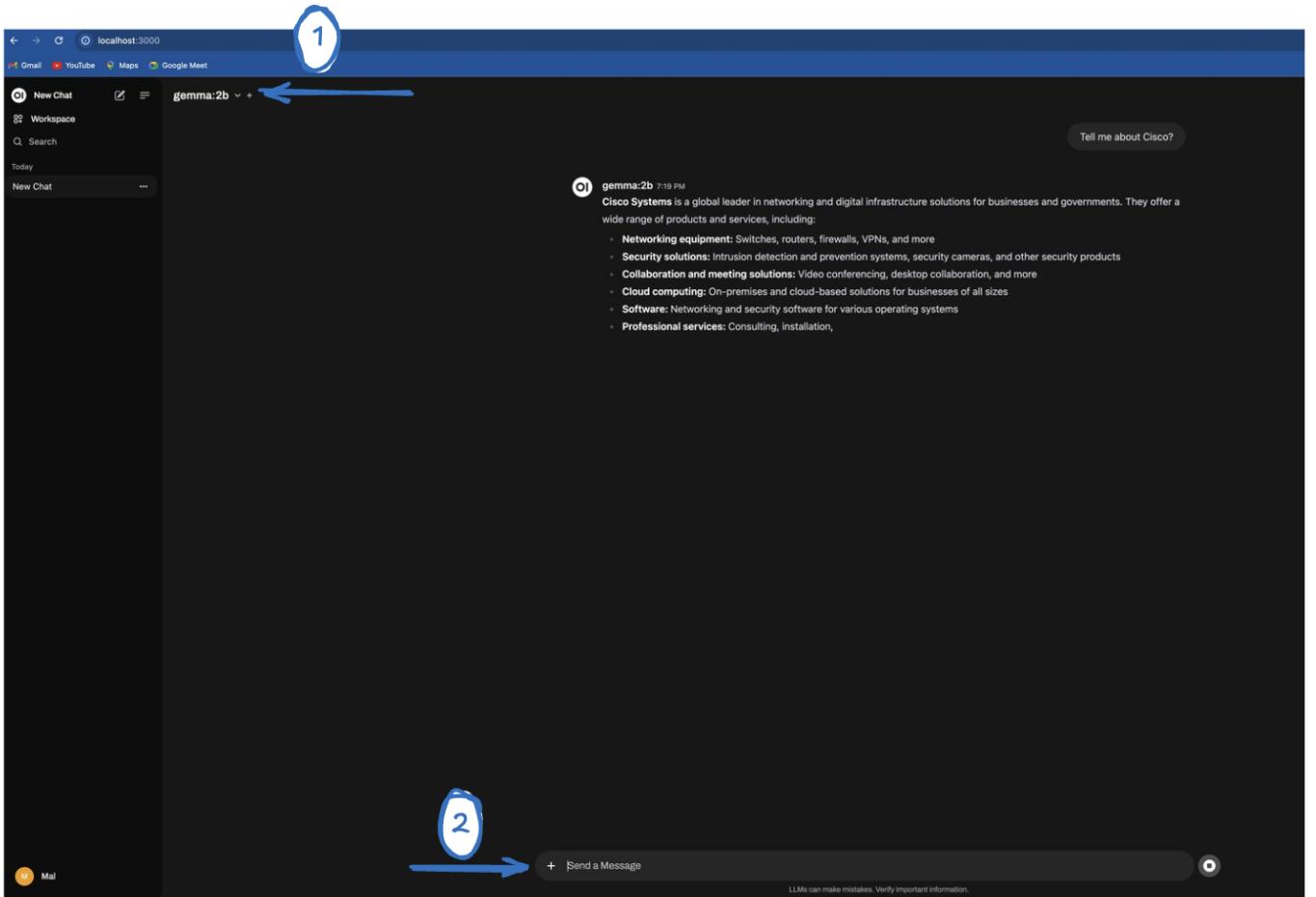
- Once installed, you can access Open WebUI at <http://localhost:3000>.

Note: Before accessing the web interface, you will be prompted to create an account. Please follow the instructions.

- In the dropdown menu, you'll see the model you previously installed via the terminal (gemma:2b).



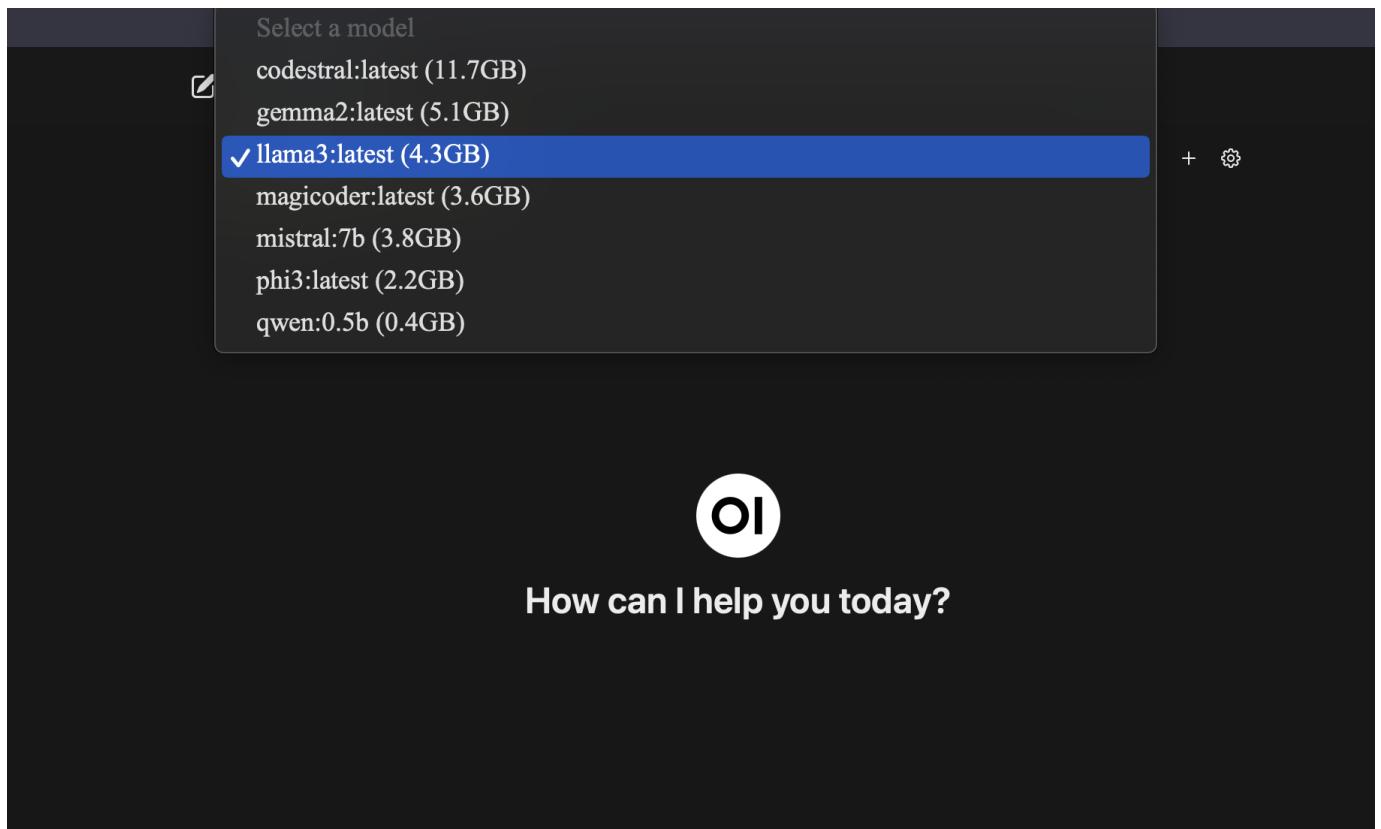
- You can now interact with the model via WebUi



- You no longer need to use the terminal to download a model. Simply navigate to the “Settings --> Models” section in the WebUI, and select the specific model you want to download.

The image consists of two side-by-side screenshots. The left screenshot shows the Ollama WebUI homepage with a sidebar containing a 'New Chat' button, a workspace section, and a search bar. A model named 'llama3:latest (4.3GB)' is listed. The right screenshot shows the 'Settings' page with a sidebar containing 'General', 'Advanced', 'Models' (which is highlighted), 'External', 'Interface', 'Audio', 'Chats', 'Account', and 'About'. The 'Models' section contains fields for 'Enter model tag (e.g. mistral:7b)', 'Delete a model', 'Upload a GGUF model (Experimental)', and 'Click here to select'.

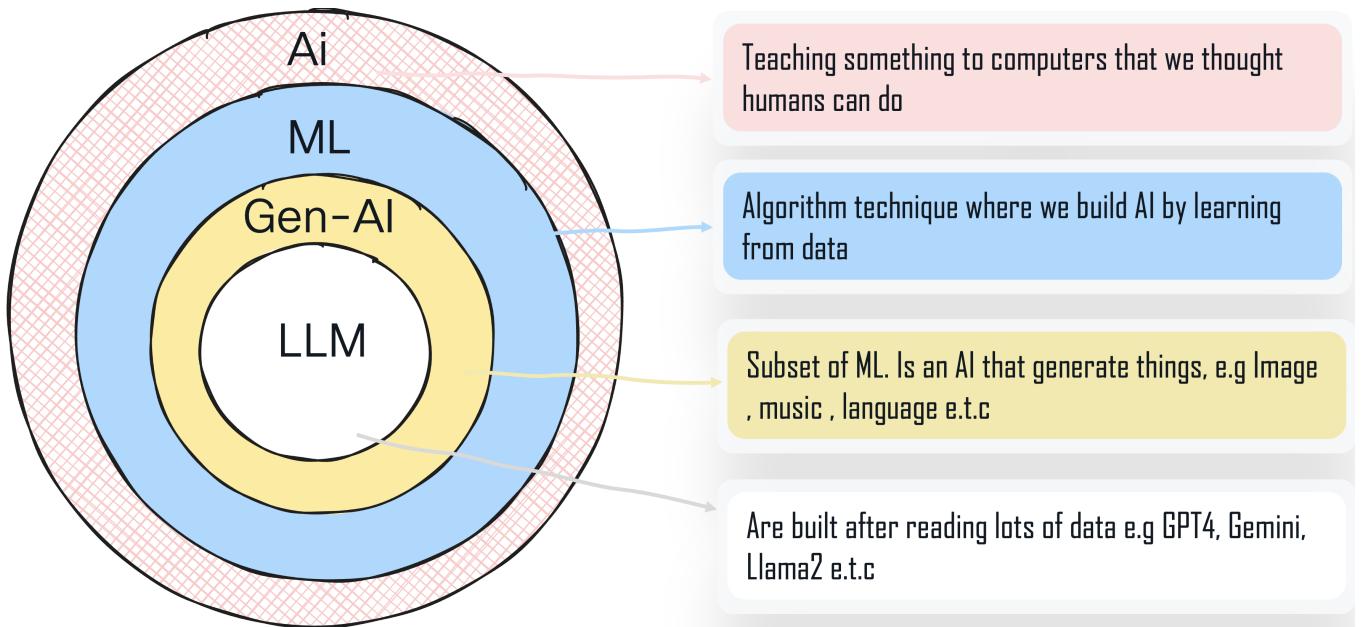
- After installing, all the installed models will be displayed in the UI.



Note: Using LangChain with Ollama allows you to leverage Ollama's capabilities within the LangChain ecosystem. You can find more information [here](#)

1.3 AI/ML Revolution Unveiled

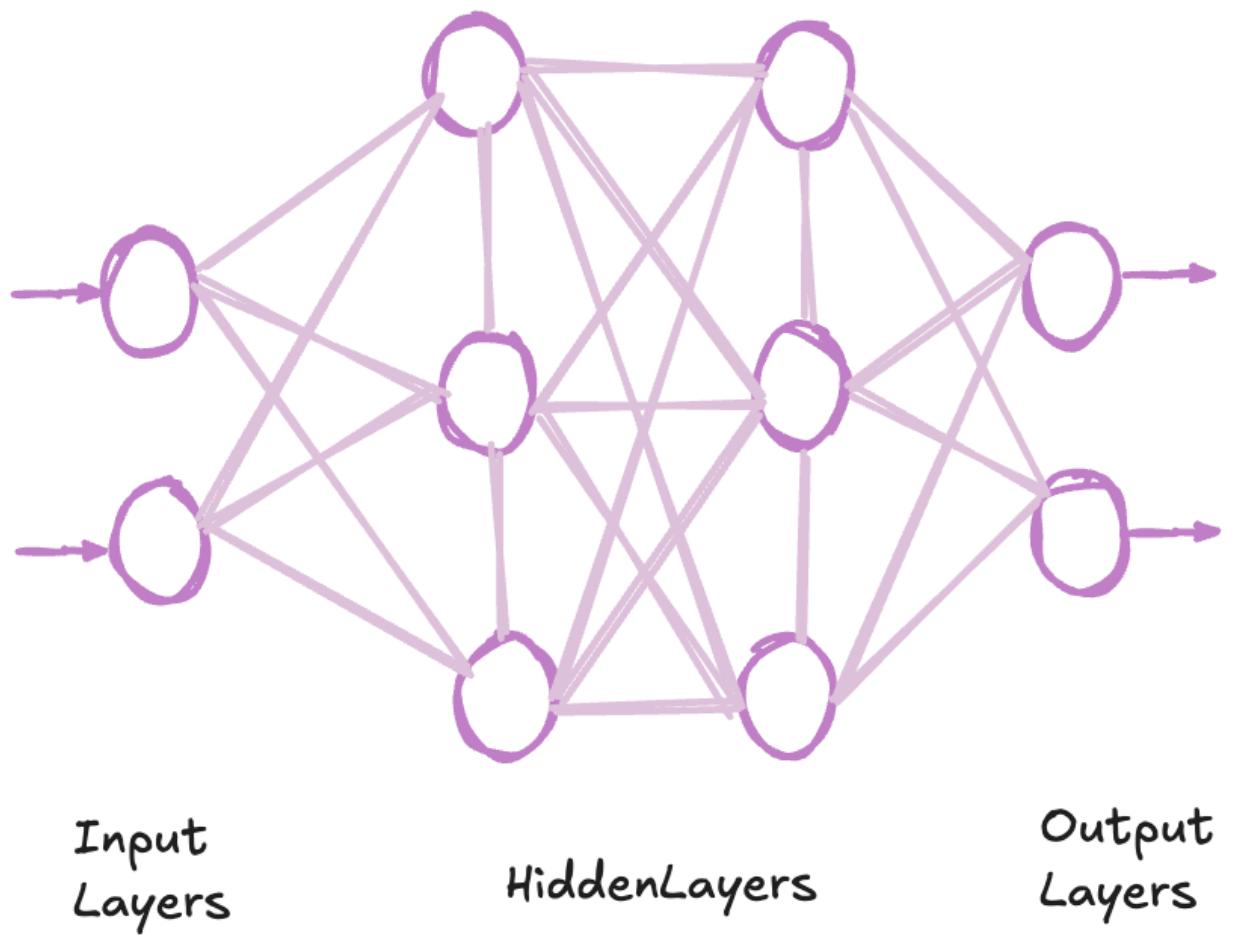
1.3.1 Different Types of AI



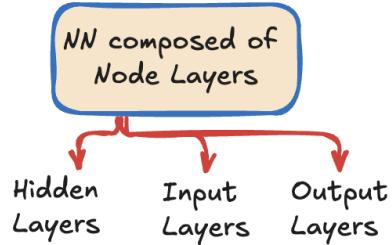
Remember: Subset of AI is ML and kind of ML is Gen-AI . LLM e.g GPT are subset of Gen-AI

1.3.2 Neural Network

Artificial Neural Network



Note: The way nodes are connected called weights



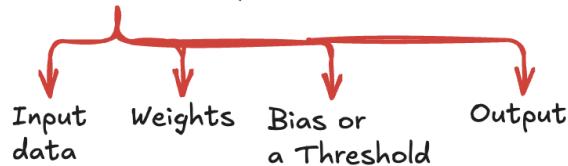
These neural network reflects the behaviour of human brain allowing computers to recognize patterns and solve common problems

Each neuron has its own Linear Regression Model

Linear Regression is the mathematical model that is used to predict future events

The weights of the connections between the nodes determine how much influence each input has on output

Each **Node** is composed of

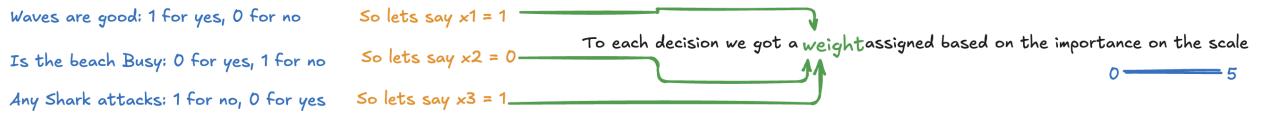


Data Passed between Nodes called FFN

The data is passed in the neural network from one layer to another which is known as feed forward propagation

Example: Lets see how a single node in NN looks like when deciding if you wanna go surfing. The decision to go or not will be our predictive outcome or lets call it our $y\hat{}$

Assume there are 3 factors influencing our decision:



Lets assign weights to our decisions:

So Waves are good we will assign a weight $w_1 = 5$

Beach not busy not of great importance, $w_2 = 2$

No Shark attacks, lets give it a weight $w_3 = 4$

Overall Result:

$$\begin{array}{lll} x_1 = 1 & x_2 = 0 & x_3 = 1 \\ w_1 = 5 & w_2 = 2 & w_3 = 4 \end{array}$$

We can plug these values into our formulae to get the desired output:

$$y\hat{=} (1*5) + (0*2) + (1*4)$$

$$y\hat{=} 9 - 3 = 6$$

|
Threshold or Bias

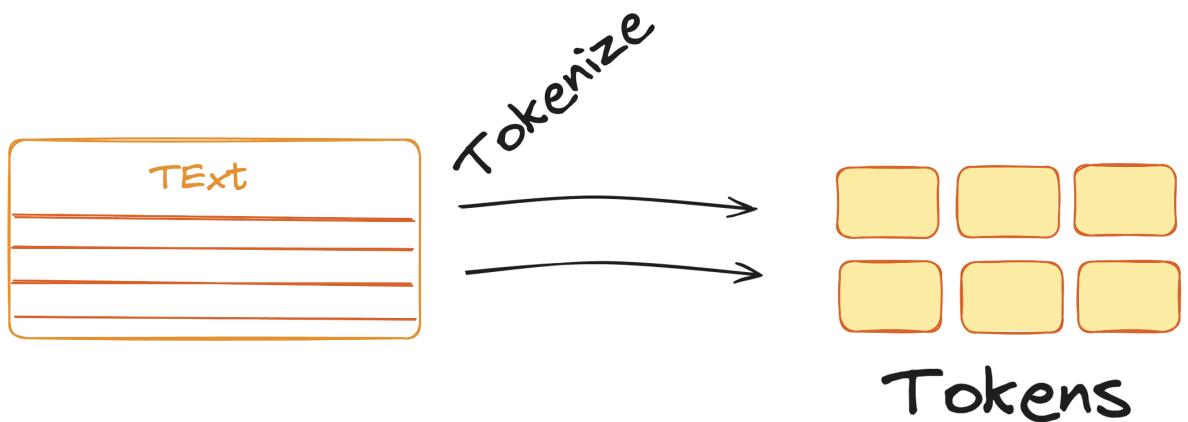
As 6 is greater than 0 so the output of this node is 1. We are going surfing

We can always adjust weights to get different outcomes

Note: Neural network rely on training data to learn and improve their accuracy over time. We used supervised learning to train the algorithm. There are multiple types of Neural network other than the Feed forward network that we have defined here. Example: CNN: Convolutional Neural Network - Unique architecture for identifying patterns like image recognition, or RNN: Recurrent neural network, that uses time series data to make prediction about future events like sales forecasting.

1.4 Tokenization

What is Tokenization and why its needed?



- Language models have a limit on how much text they can handle at once, known as their context window. While these limits are expanding, research shows that LLM's often perform better when provided with less, but more relevant, information. However, selecting the most relevant information is straightforward for humans but challenging for computers.

A common approach to manage large amounts of data is to break it down into smaller, more manageable parts a process often referred to as Tokenization or Chunking. Tokenization is a key step in this process, where raw text is divided into smaller units, called tokens, which can then be processed by a neural network.

Tokenization

Different ways to Tokenize

Word-Based

Welcome

to

Cisco

Character

W e l c o m e

t o c i s c o

Sentence

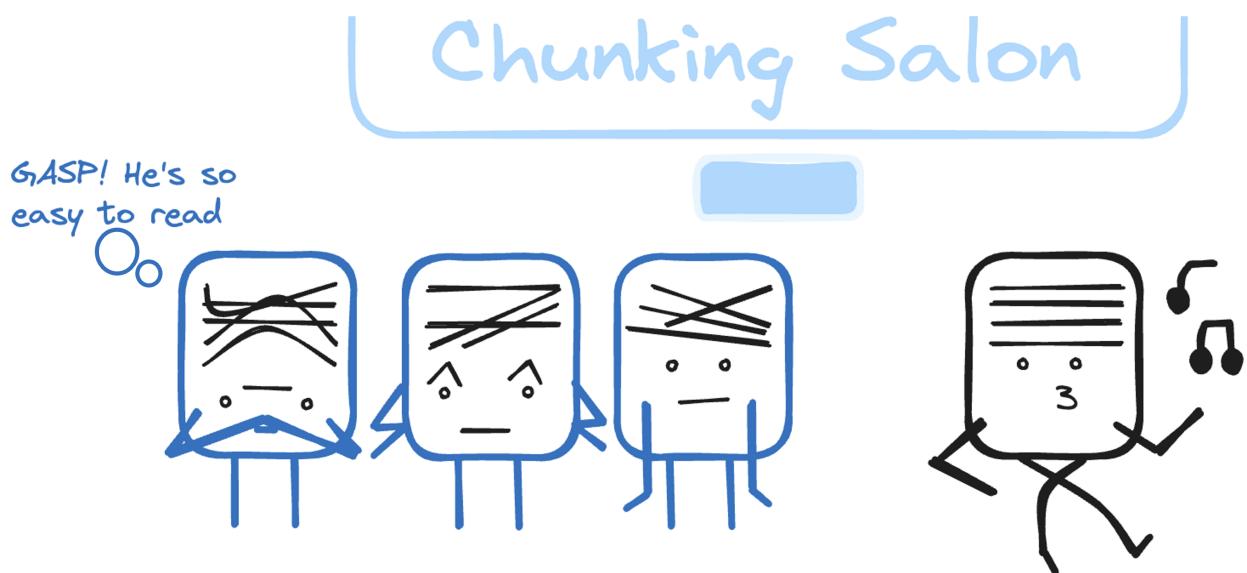
Welcome to Cisco.

In order to do this you need to pick a chunk strategy. Just to name a few:

- Word-Based: A simple and straightforward method that most of us would propose is to use word-based tokens, splitting the text by spaces.
- Character based tokenization: Individual words are considered as tokens . Lot of computing resources needed as now e.g for a 3 word sentence where you might need 3 tokens now 15 – 20 tokens needed
- Sentence Based: We need a .(fullstop) for it to work

Note: We've all heard of GPT and OpenAI. They utilize a tokenization method called Byte Pair Encoding (BPE), which is a middle ground between word-based and character-based tokenization. In BPE, words are broken down into smaller character sequences that the model encountered during training, allowing it to make informed predictions.

1.4.1 Why we need Tokenization or Text Splitting?



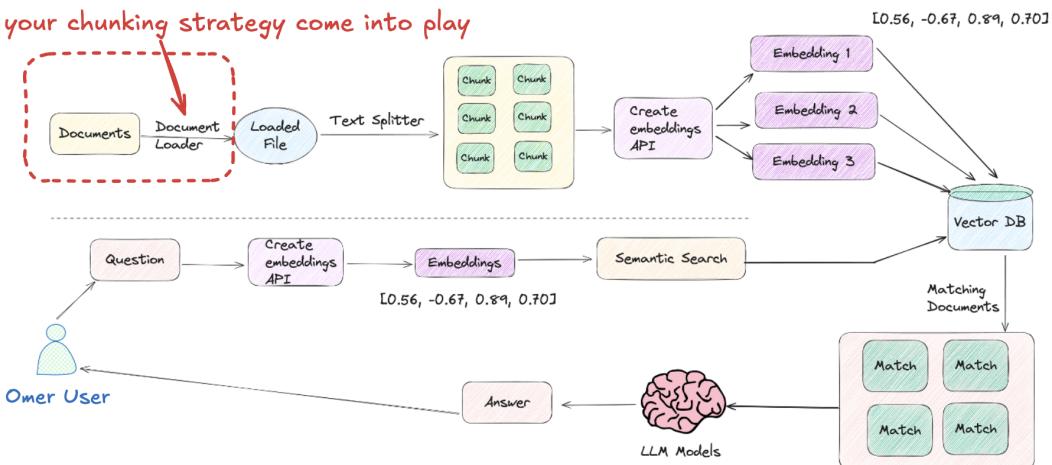
Historically, applications perform better when they are provided with your own data. However, you can't feed unlimited data to your LLMs due to two key limitations:

- Context window limit: LLMs have a finite context window for processing data.
- Signal-to-noise ratio (SNR): LLMs perform better when the SNR is high, meaning the information provided is useful, relevant, and clear. Clear, unambiguous instructions help the model deliver more accurate and detailed results, while ambiguous or complex input can lead to less accurate or incomplete outputs.

As noted, chunking or splitting refers to breaking your data into smaller, manageable pieces.

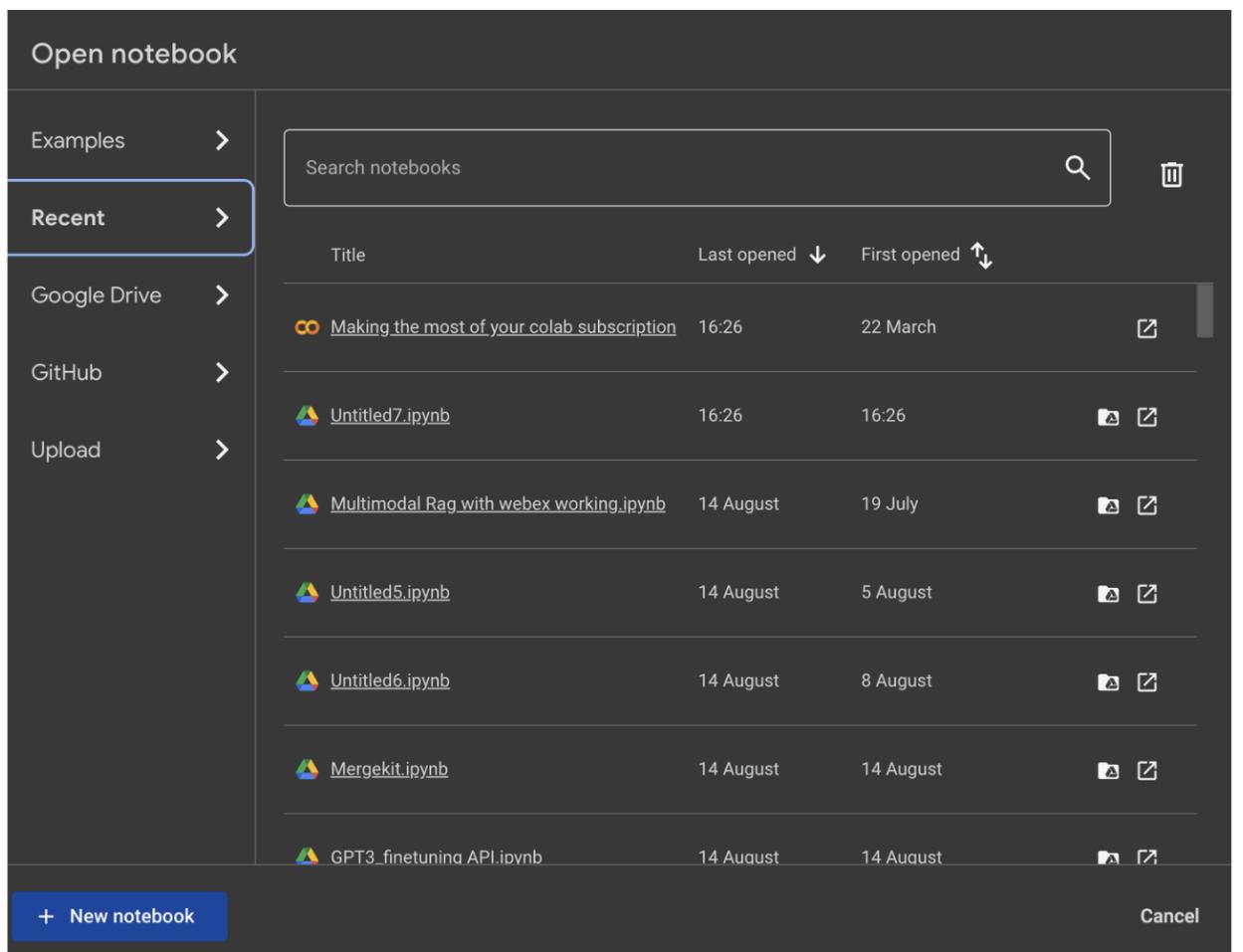
AI Full Stack Frame-work

That's where your chunking strategy come into play

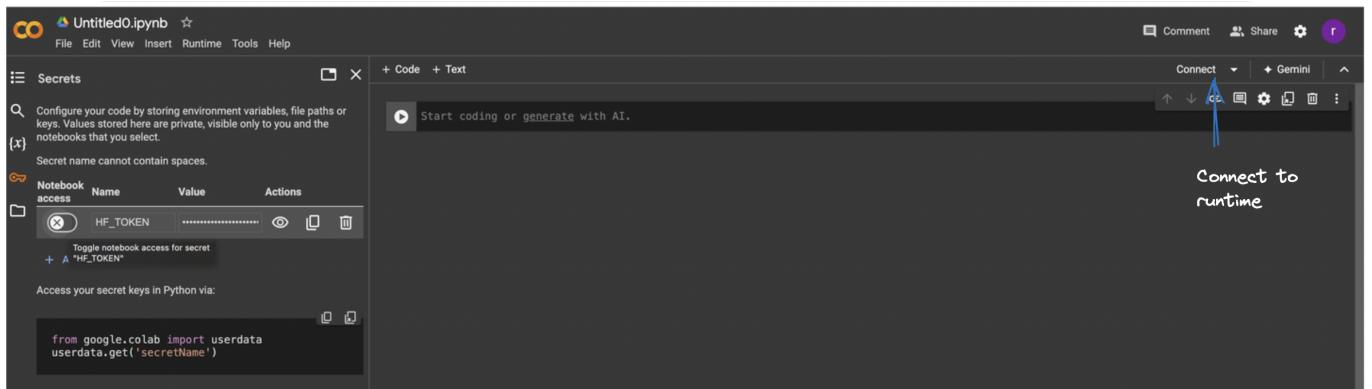


1.4.2 Task 1: Log into the Lab Environment

- Open Google Colab and create a new notebook. Click on "File" > "New notebook". Please refer to the following section to create Google Colab account.



- Make sure you are connected to a runtime. For this task, you can use the CPU as the runtime environment.



Reminder: Whenever you want to copy the code in Google Colab and run it, be sure to click on + Code to add a new code cell.



Reminder: Click the play button to the left of the code, or use the keyboard shortcut "Command/Ctrl+Enter" while the cell is selected.



Manual testing for Chunking

We will create Chunks of 35 characters, so first 35 characters as chunk 1 , next 35 characters chunk2 and so on

```

1  # Manual Splitting
2  text = """WebexOne is an annual in-person and virtual event that takes place over four days. It's an event focused on AI collaboration and customer
3  experience.
4  It features a range of activities such as insightful breakout sessions, technical training courses, hands-on labs, inspiring keynotes, epic
5  entertainment, a solutions showcase and expo, customer awards, meet the experts, 1:1 executive meetings, a partner program, and more!"""
6
7  # Create a list that will hold your chunks
8  chunks = []
9
10 chunk_size = 35
11
12 # Run through the a range with the length of your text and iterate every chunk_size you want
13 for i in range(0, len(text), chunk_size):
14     chunk = text[i:i + chunk_size]
15     chunks.append(chunk)

print(chunks)

```

WebexOne is an annual in-person and virtual event that takes place over four days. It's an event focused on AI collaboration and customer experience. It features a range of activities such as insightful breakout sessions, technical training courses, hands-on labs, inspiring keynotes, epic entertainment, a solutions showcase and expo, customer awards, meet the experts, 1:1 executive meetings, a partner program, and more!

The screenshot shows a 'Character Splitter' interface. At the top, there is a dropdown menu labeled 'Splitter: Character Splitter' with a gear icon. Below it are two sliders: 'Chunk Size' set to 35 and 'Chunk Overlap' set to 0. To the right of these sliders is a button labeled 'Upload .txt'. Below the sliders, the text 'Total Characters: 422', 'Number of chunks: 13', and 'Average chunk size: 32.5' is displayed.

WebexOne is an annual in-person and virtual event that takes place over four days. It's an event focused on AI collaboration and customer experience. It features a range of activities such as insightful breakout sessions, technical training courses, hands-on labs, inspiring keynotes, epic entertainment, a solutions showcase and expo, customer awards, meet the experts, 1:1 executive meetings, a partner program, and more!

Note: The text contained a total of 422 characters, divided into 13 chunks. Problem with the above chunking technique is that it got split at 'r'. How we know when to chunk . Before we look into that. Lets look into Langchain Splitter

LangChain Text Splitter

- Let's retrieve the langchain library from the Python Package. In the example below, we'll configure the chunk_overlap parameter, which ensures that our chunks are blended together—meaning the end of chunk 1 will overlap with the beginning of chunk 2.

```
1 !pip install langchain
```

```
Collecting langchain
  Downloading langchain-0.2.14-py3-none-any.whl.metadata (7.1 kB)
Requirement already satisfied: PyYAML>=5.3 in /usr/local/lib/python3.10/dist-packages (from langchain) (6.0.2)
Requirement already satisfied: SQLAlchemy<3,>=1.4 in /usr/local/lib/python3.10/dist-packages (from langchain) (2.0.32)
Requirement already satisfied: aiohttp<4.0.0,>=3.8.3 in /usr/local/lib/python3.10/dist-packages (from langchain) (3.10.2)
Requirement already satisfied: async-timeout<5.0.0,>=4.0.0 in /usr/local/lib/python3.10/dist-packages (from langchain) (4.0.3)
Collecting langchain-core<0.3.0,>=0.2.32 (from langchain)
  Downloading langchain_core-0.2.33-py3-none-any.whl.metadata (6.2 kB)
Collecting langchain-text-splitters<0.3.0,>=0.2.0 (from langchain)
  Downloading langchain_text_splitters-0.2.2-py3-none-any.whl.metadata (2.1 kB)
Collecting langsmith<0.2.0,>=0.1.17 (from langchain)
  Downloading langsmith-0.1.99-py3-none-any.whl.metadata (13 kB)
Requirement already satisfied: numpy<2,>=1 in /usr/local/lib/python3.10/dist-packages (from langchain) (1.26.4)
Requirement already satisfied: pydantic<3,>=1 in /usr/local/lib/python3.10/dist-packages (from langchain) (2.8.2)
Requirement already satisfied: requests<,>=2 in /usr/local/lib/python3.10/dist-packages (from langchain) (2.32.3)
Collecting tenacity!=8.4.0,<9.0.0,>=8.1.0 (from langchain)
  Downloading tenacity-8.5.0-py3-none-any.whl.metadata (1.2 kB)
Requirement already satisfied: aiohappyeyeballs>=2.3.0 in /usr/local/lib/python3.10/dist-packages (from aiohttp<4.0.0,>=3.8.3->langchain) (2.3.5)
Requirement already satisfied: aiosignal<1.2.1 in /usr/local/lib/python3.10/dist-packages (from aiohttp<4.0.0,>=3.8.3->langchain) (1.3.1)
Requirement already satisfied: attrs>=17.3.0 in /usr/local/lib/python3.10/dist-packages (from aiohttp<4.0.0,>=3.8.3->langchain) (24.2.0)
Requirement already satisfied: frozenlist>=1.1.1 in /usr/local/lib/python3.10/dist-packages (from aiohttp<4.0.0,>=3.8.3->langchain) (1.4.1)
Requirement already satisfied: multidict<7.0,>=4.5 in /usr/local/lib/python3.10/dist-packages (from aiohttp<4.0.0,>=3.8.3->langchain) (6.0.5)
Requirement already satisfied: yarl<2.0,>=1.0 in /usr/local/lib/python3.10/dist-packages (from aiohttp<4.0.0,>=3.8.3->langchain) (1.9.4)
Collecting jsonpatch<2.0,>=1.33 (from langchain-core<0.3.0,>=0.2.32->langchain)
  Downloading jsonpatch-1.33-py2.py3-none-any.whl.metadata (3.0 kB)
Requirement already satisfied: packaging<25,>=23.2 in /usr/local/lib/python3.10/dist-packages (from langchain-core<0.3.0,>=0.2.32->langchain) (24.1)
Requirement already satisfied: typing_extensions>=4.7 in /usr/local/lib/python3.10/dist-packages (from langchain-core<0.3.0,>=0.2.32->langchain) (4.12.2)
```

```
1 text = """WebexOne is an annual in-person and virtual event that takes place over four days. It's an event focused on AI collaboration and customer
2 experience.
It features a range of activities such as insightful breakout sessions, technical training courses, hands-on labs, inspiring keynotes, epic
entertainment, a solutions showcase and expo, customer awards, meet the experts, 1:1 executive meetings, a partner program, and more!"""

```

```
1 from langchain.text_splitter import CharacterTextSplitter
2
3 text_splitter = CharacterTextSplitter(chunk_size=35, chunk_overlap=3, separator='', strip_whitespace=False) # separate means you split by character
4
5 a = text_splitter.create_documents([text])
6
7 print(a)
```

Note: Separators are characters sequence you wanna split on.

WebexOne is an annual in-person and virtual event that takes place over four days. It's an event focused on AI collaboration and customer experience. It features a range of activities such as insightful breakout sessions, technical training courses, hands-on labs, inspiring keynotes, epic entertainment, a solutions showcase and expo, customer awards, meet the experts, 1:1 executive meetings, a partner program, and more!

The screenshot shows a text splitting interface with the following settings:

- Splitter: Character Splitter
- Chunk Size: 35
- Chunk Overlap: 3

Statistics displayed:

- Total Characters: 461
- Number of chunks: 14
- Average chunk size: 32.9

The text "WebexOne is an annual in-person and virtual event that takes place over four days. It's an event focused on AI collaboration and customer experience. It features a range of activities such as insightful breakout sessions, technical training courses, hands-on labs, inspiring keynotes, epic entertainment, a solutions showcase and expo, customer awards, meet the experts, 1:1 executive meetings, a partner program, and more!" is shown, divided into two chunks labeled "Chunk 1" and "Chunk 2". A red arrow labeled "overlap" points to the boundary between the two chunks, indicating the overlap specified in the settings.

Recursive Charector Text Splitting

In the previous example, when we used Character Splitting, we split the text based on a fixed number of characters. Specifically, we divided the text into chunks of 35 characters each.

However, with Recursive Text Splitting (RTS), the process is more dynamic and considers the physical structure of the text to determine the appropriate chunk size.

Here's how RTS works:

- Instead of relying on a static number of characters, RTS examines the structure of the document.
- It begins by identifying the largest logical divisions, such as paragraphs, and splits the text at each double newline (indicating paragraph breaks).
- If any of these chunks are still too large, RTS will then move to the next level of separation, which is single newlines (often indicating sentences or list items).
- If necessary, it will continue to break down the text further, using spaces and eventually individual characters as separators.

With RTS, you don't need to manually specify the chunk size by character count. You simply pass your text to the RTS process, and it will intelligently determine how to split the text based on its structure, resulting in chunks that are logically organized and more contextually meaningful.

- "\n\n" - Double new line, or most commonly paragraph breaks
- "\n" - New lines
- " " - Spaces
- "" - Characters

```

1 text = """WebexOne is an annual in-person and virtual event that takes place over four days. It's an event focused on AI collaboration and customer
2 experience.
It features a range of activities such as insightful breakout sessions, technical training courses, hands-on labs, inspiring keynotes, epic
entertainment, a solutions showcase and expo, customer awards, meet the experts, 1:1 executive meetings, a partner program, and more!"""

```

```

1 from langchain.text_splitter import RecursiveCharacterTextSplitter
2 text_splitter = RecursiveCharacterTextSplitter(chunk_size = 35, chunk_overlap = 0)
3 # how many chunks we have
4 print(len(text_splitter.create_documents([text])))
5
6 text_splitter.create_documents([text])

```

Note: We avoid splitting in the middle of words by using spaces as one of the separators, ensuring that each chunk ends on a complete word. This is crucial because RCS (Recursive Character Splitting) helps maintain the context within sentences by keeping related words together.

WebexOne is an annual in-person and virtual event that takes place over four days. It's an event focused on AI collaboration and customer experience. It features a range of activities such as insightful breakout sessions, technical training courses, hands-on labs, inspiring keynotes, epic entertainment, a solutions showcase and expo, customer awards, meet the experts, 1:1 executive meetings, a partner program, and more!

Upload .txt

Splitter: Recursive Character Text Splitter

Chunk Size: 35

Chunk Overlap: 0

Total Number of Chunks: 14

Total Characters: 422

Average chunk size: 30.1

RTS maintains words together

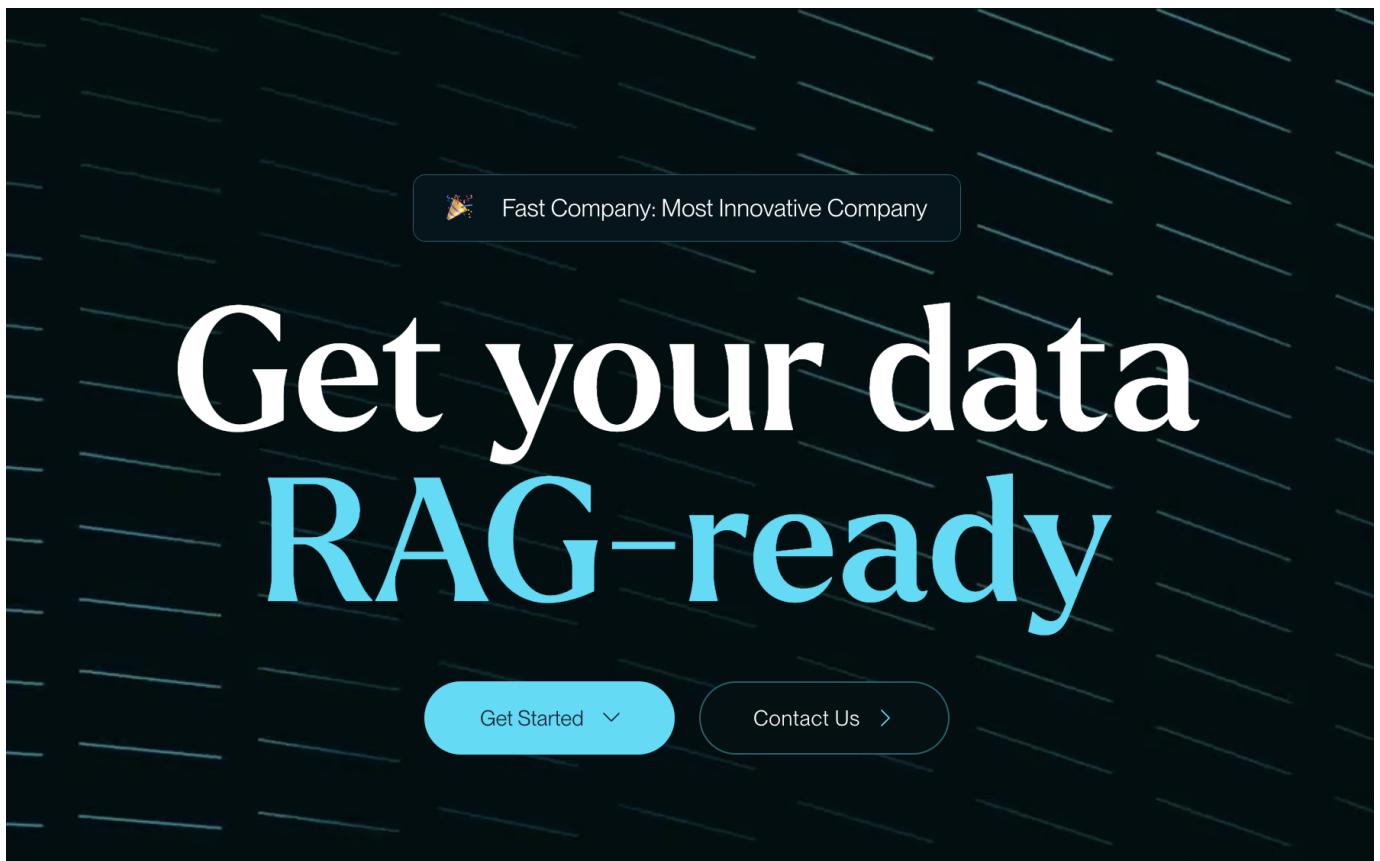
Note: If you're new to AI, I would personally recommend starting with Recursive Text Splitting (RTS).

Document Level Splitting

So far, we've been working with splitting regular documents. But what if we have markdown files or PDF or Python documentation? There's a better way to handle those cases, and that's where specialized document splitting comes into play.

PDF WITH TABLE

PDFs often contain tables and other structured data that can be challenging to split accurately using character-based methods. For PDFs, it's important to extract and chunk all elements, including tables, effectively. We'll accomplish this using the Unstructured library, which is specifically designed for handling such tasks. If you have a large collection of PDFs, Unstructured is an excellent tool to manage them efficiently.



- Install relevant libraries including Unstructured

```
1 !pip install scikit-learn langchain_community unstructured[all-docs] unstructured pdfminer pdfminer.six pdf2image pillow_heif opencv-python
unstructured_inference pytesseract unstructured_pytesseract python-dotenv openai
```

```
Requirement already satisfied: scikit-learn in /usr/local/lib/python3.10/dist-packages (1.3.2)
Requirement already satisfied: langchain_community in /usr/local/lib/python3.10/dist-packages (0.2.12)
Requirement already satisfied: unstructured in /usr/local/lib/python3.10/dist-packages (0.15.5)
Requirement already satisfied: pdfminer in /usr/local/lib/python3.10/dist-packages (20191125)
Requirement already satisfied: pdfminer.six in /usr/local/lib/python3.10/dist-packages (20231228)
Requirement already satisfied: pdf2image in /usr/local/lib/python3.10/dist-packages (1.17.0)
Requirement already satisfied: pillow_heif in /usr/local/lib/python3.10/dist-packages (0.18.0)
Requirement already satisfied: opencv-python in /usr/local/lib/python3.10/dist-packages (4.10.0.84)
Requirement already satisfied: unstructured_inference in /usr/local/lib/python3.10/dist-packages (0.7.36)
Requirement already satisfied: pytesseract in /usr/local/lib/python3.10/dist-packages (0.3.13)
Requirement already satisfied: unstructured_pytesseract in /usr/local/lib/python3.10/dist-packages (0.3.13)
Requirement already satisfied: numpy<2.0,>=1.17.3 in /usr/local/lib/python3.10/dist-packages (from scikit-learn) (1.26.4)
Requirement already satisfied: scipy<1.5.0 in /usr/local/lib/python3.10/dist-packages (from scikit-learn) (1.13.1)
Requirement already satisfied: joblib<1.1.1 in /usr/local/lib/python3.10/dist-packages (from scikit-learn) (1.4.2)
Requirement already satisfied: threadpoolctl<2.0.0 in /usr/local/lib/python3.10/dist-packages (from scikit-learn) (3.5.0)
Requirement already satisfied: PyYAML<5.3 in /usr/local/lib/python3.10/dist-packages (from langchain_community) (6.0.2)
Requirement already satisfied: SQLAlchemy<3,>=1.4 in /usr/local/lib/python3.10/dist-packages (from langchain_community) (2.0.32)
Requirement already satisfied: aiohttp<4.0.0,>=3.8.3 in /usr/local/lib/python3.10/dist-packages (from langchain_community) (3.10.2)
Requirement already satisfied: dataclasses-json<0.7,>=0.5.7 in /usr/local/lib/python3.10/dist-packages (from langchain_community) (0.6.7)
Requirement already satisfied: langchain<0.3.0,>=0.2.13 in /usr/local/lib/python3.10/dist-packages (from langchain_community) (0.2.14)
Requirement already satisfied: langchain-core<0.3.0,>=0.2.30 in /usr/local/lib/python3.10/dist-packages (from langchain_community) (0.2.33)
Requirement already satisfied: langsmith<0.2.0,>=0.1.0 in /usr/local/lib/python3.10/dist-packages (from langchain_community) (0.1.99)
```

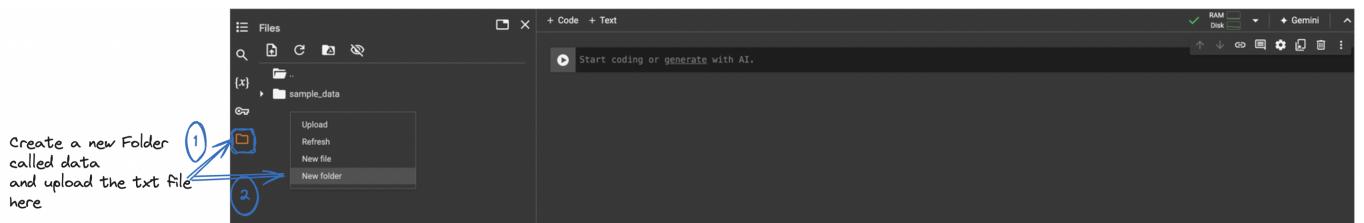
```
1 !apt-get install -y poppler-utils && apt-get install -y tesseract-ocr
```

```

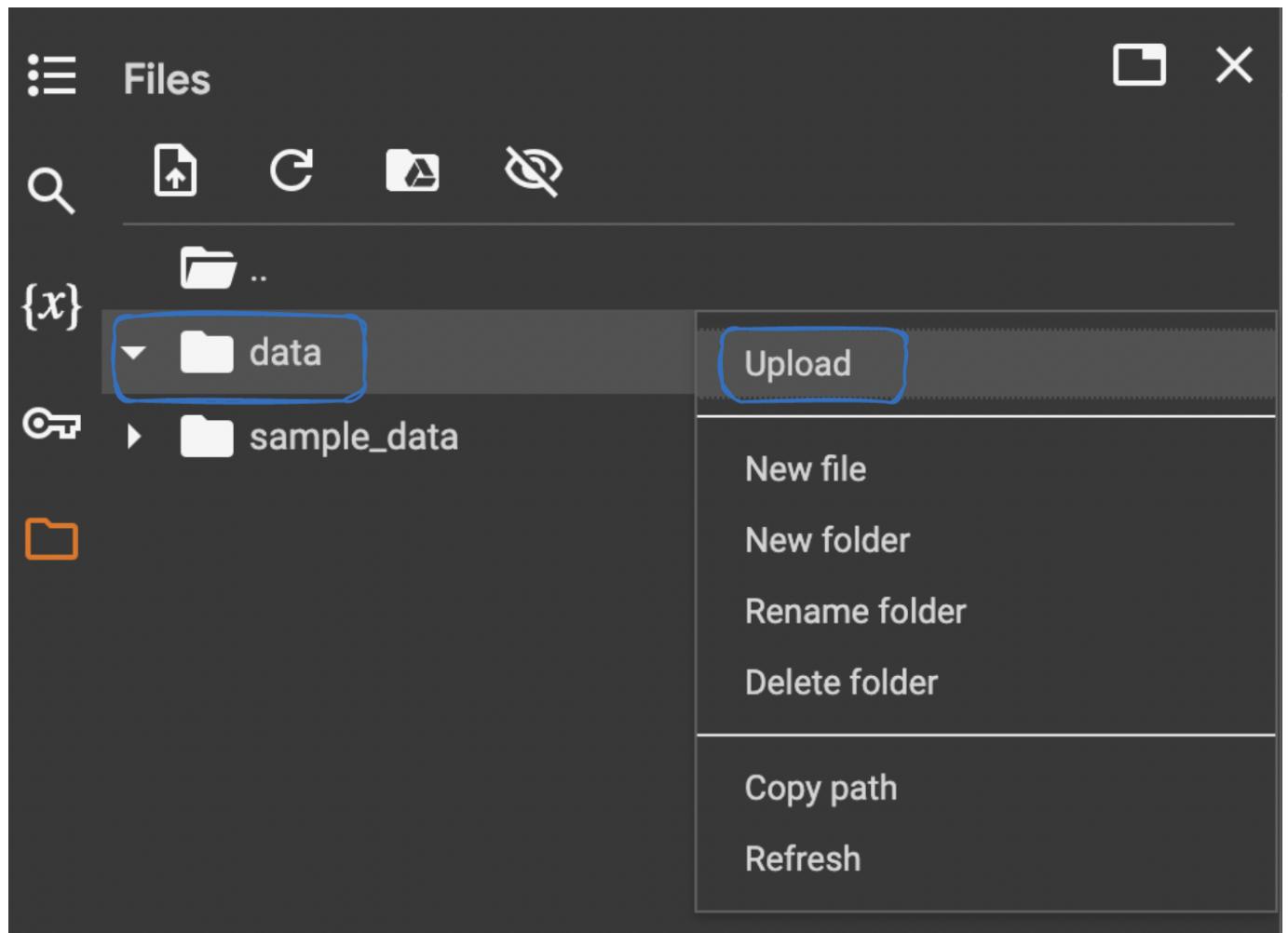
Reading package lists... Done
Building dependency tree... Done
Reading state information... Done
The following NEW packages will be installed:
  poppler-utils
0 upgraded, 1 newly installed, 0 to remove and 45 not upgraded.
Need to get 186 kB of archives.
After this operation, 696 kB of additional disk space will be used.
Get:1 http://archive.ubuntu.com/ubuntu jammy-updates/main amd64 poppler-utils amd64 22.02.0-2ubuntu0.5 [186 kB]
Fetched 186 kB in 1s (228 kB/s)
Selecting previously unselected package poppler-utils.
(Reading database ... 123594 files and directories currently installed.)
Preparing to unpack .../poppler-utils_22.02.0-2ubuntu0.5_amd64.deb ...
Unpacking poppler-utils (22.02.0-2ubuntu0.5) ...
Setting up poppler-utils (22.02.0-2ubuntu0.5) ...
Processing triggers for man-db (2.10.2-1) ...

```

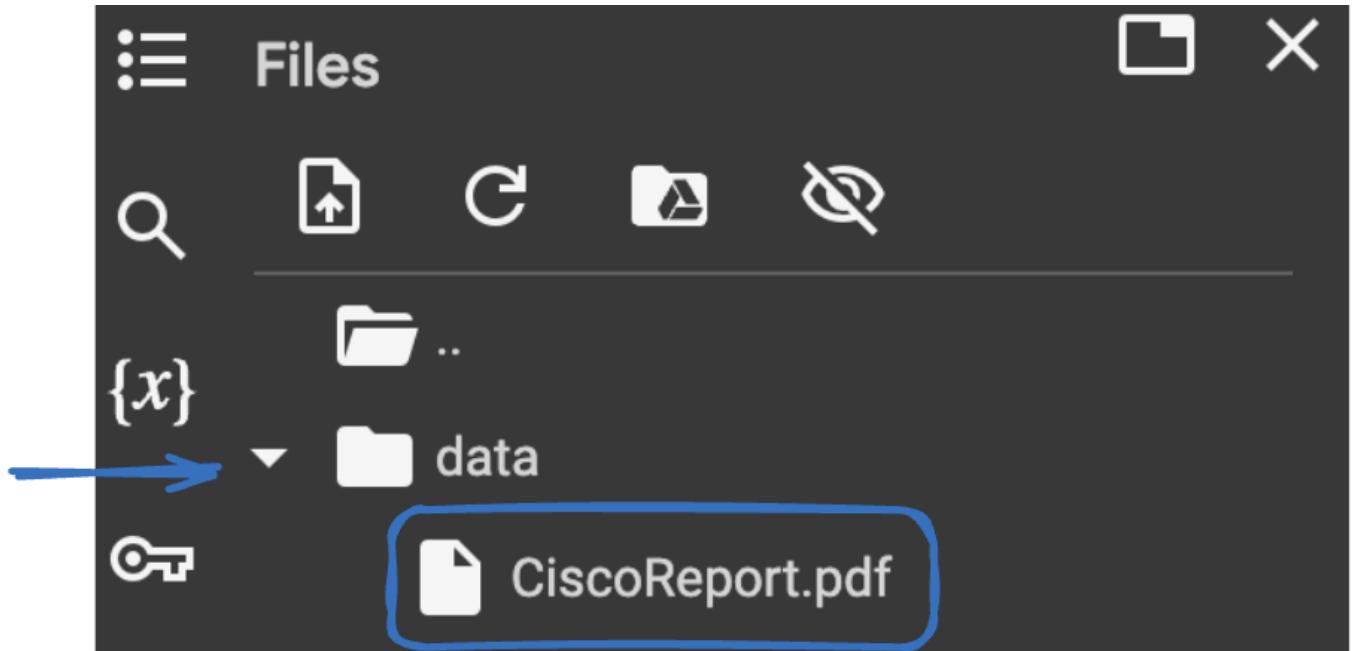
- Let's load our PDF files into Google Colab. For this example, we can use the Cisco Financial Results. Please Download the file here as we will be using in the next step.
- Within Google Colab, Click on Folder and create a new folder called "data"



- Click on [...], select Upload



- Choose your CiscoReport.pdf file and click Open



- We'll use the following code from the Unstructured library to demonstrate how tables can be extracted

```
1 import os
2 from unstructured.partition.pdf import partition_pdf
3 from unstructured.staging.base import elements_to_json
4
5 # Let's load up our PDF and then partition it.
6 filename = "/content/data/CiscoReport.pdf"# Use relative path since the file is in the same directory
7
8 # Extracts the elements from the PDF
9 elements = partition_pdf(
10     filename=filename,
11     extract_images_in_pdf=True,
12     strategy="hi_res",
13     infer_table_structure=True,
14     hi_res_model_name="yolox"
15 )
16
17 # Let's look at our elements
18 print(elements)
```

```
config.json: 100% [1.47k/1.47k [00:00<00:00, 76.5kB/s]
model.safetensors: 100% [115M/115M [00:01<00:00, 125MB/s]
model.safetensors: 100% [46.8M/46.8M [00:00<00:00, 117MB/s]

<unstructured.documents.elements.Title at 0x7ba52cc9c5e0>,
<unstructured.documents.elements.ListItem at 0x7ba52cc9d7b0>,
<unstructured.documents.elements.ListItem at 0x7ba52cc9d3f0>,
<unstructured.documents.elements.Title at 0x7ba52cc9caf0>,
<unstructured.documents.elements.ListItem at 0x7ba4c0997ca0>,
<unstructured.documents.elements.Text at 0x7ba4c09943a0>,
<unstructured.documents.elements.NarrativeText at 0x7ba52cc9ebc0>, ← Regular Text
<unstructured.documents.elements.NarrativeText at 0x7ba52cc9f370>,
<unstructured.documents.elements.NarrativeText at 0x7ba52cc9d720>,
<unstructured.documents.elements.Title at 0x7ba52cc9de10>,
<unstructured.documents.elements.Text at 0x7ba4c0997f10>,
<unstructured.documents.elements.Text at 0x7ba4c0994460>,
<unstructured.documents.elements.Text at 0x7ba4c0994520>,
<unstructured.documents.elements.Table at 0x7ba4c0994160>,
<unstructured.documents.elements.NarrativeText at 0x7ba4c09944f0>,
<unstructured.documents.elements.Title at 0x7ba4c09946a0>,
<unstructured.documents.elements.Text at 0x7ba4c09947f0>,
<unstructured.documents.elements.Text at 0x7ba4c0994730>,
<unstructured.documents.elements.Text at 0x7ba4c0994430>,
<unstructured.documents.elements.Table at 0x7ba4c0995150>, ← Table
<unstructured.documents.elements.Title at 0x7ba4c0994790>,
```

- Lets grab element 37 as its the one where our table is

```
1     elements[37].metadata.text_as_html
```

OUTPUT

```
'<table><thead><tr><th>Revenue</th><th>$8.51 - $8.53 Billion</th><th>$34.5 - $34.7 Billion</th><th>Y/Y Growth</td><td>-10%</td><td>-10%</td></tr></thead><tbody><tr><td>FX Impact</td><td>no impact</td><td>no impact</td><td>GAAP Operating Margin</td><td>-11.4%</td></tr><tr><td>Non-GAAP Operating Margin'</td><td></td><td>28.0%</td><td>GAAP Earnings per Share?</td><td>$0.79 - $0.80</td><td>$2.67 - $2.69</td></tr><tr><td>Non-GAAP Earnings per Share()</td><td>$1.89 - $1.90</td><td>$7.41 - $7.43</td><td>Operating Cash Flow Growth (Y/Y)</td><td>-17%</td></tr><tr><td>Current Remaining Performance Obligation Growth (Y/Y)</td><td>-10%</td><td></td><td></td><td>16%
```

- Tables are straightforward for humans to read, but they aren't as easy for language models to interpret. Language models are typically trained on HTML tables, so when you pass HTML-formatted tables to an LLM, it will better understand the structure and content. You can paste the HTML into an HTML viewer to see how it looks.

HTML Viewer

```
<table><tbody><tr><td>Revenue</td><td>$ 13.6 billion</td><td>Net Income</td><td>$ 0.97</td></tr></tbody></table>
```

HTML Code

Output

	Revenue	Net Income	Diluted Earnings per Share (EPS)
	\$ 13.6 billion	\$ 2.2 billion	\$ 0.97

Our First table in the pdf

Note: That's how you handle tables within a PDF

Multi-Modal (Text + images)

What if there are images within a PDF or other documents? How can you extract them? We'll use the Unstructured library again to handle this.

EXTRACT IMAGES WITHIN PDF

- Install relevant libraries

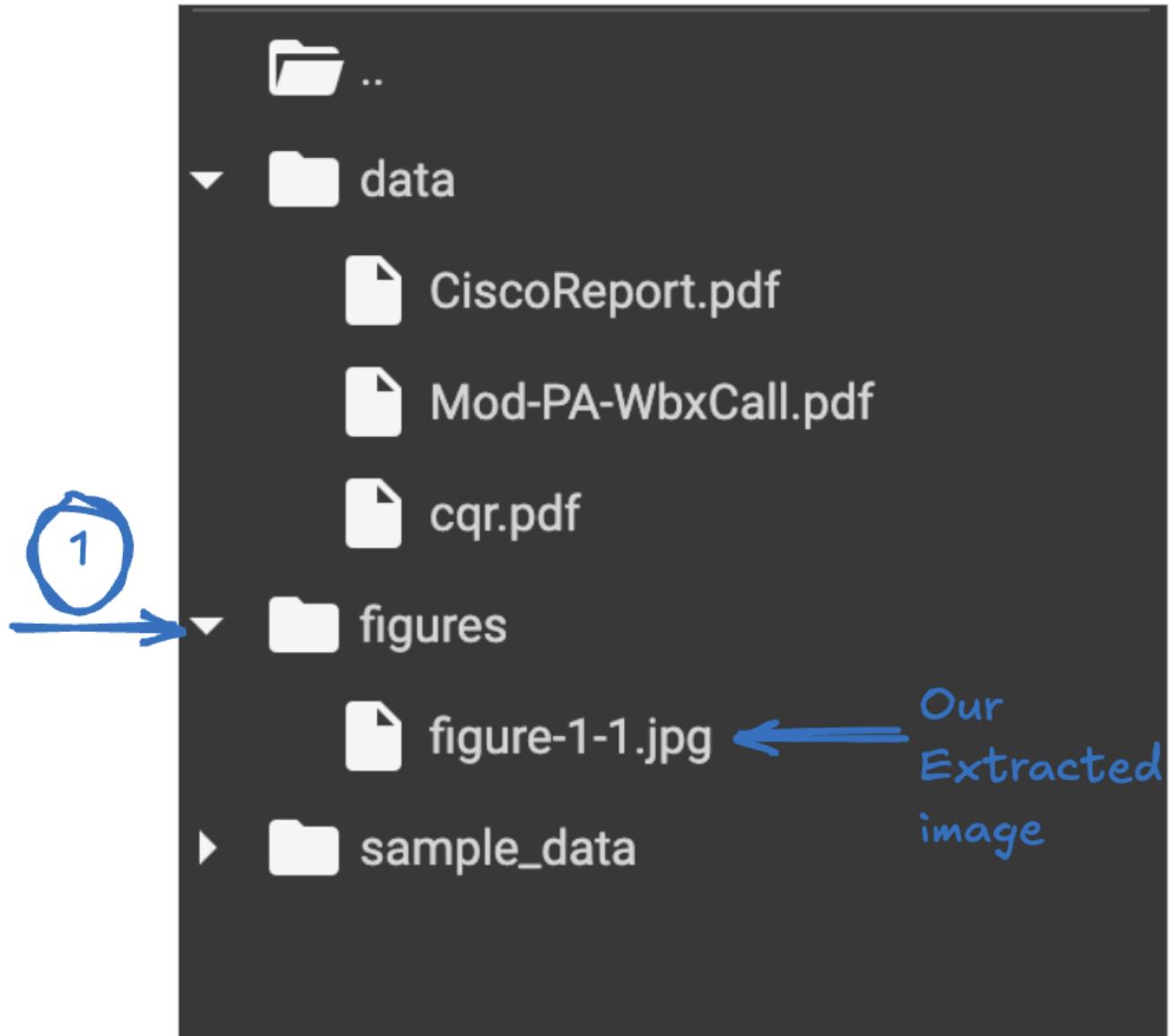
```

1  from typing import Any
2  from pydantic import BaseModel
3  from unstructured.partition.pdf import partition_pdf

1  filepath = "/content/data/CiscoReport.pdf"
2  # Get elements
3  raw_pdf_elements = partition_pdf(
4      filename=filepath,
5
6      # Using pdf format to find embedded image blocks
7      extract_images_in_pdf=True,
8
9      # Use layout model (YOLOX) to get bounding boxes (for tables) and find titles
10     infer_table_structure=True,
11
12     # Specifies the strategy to be used for chunking the text. In our case, it will chunk the text based on titles found in the document.
13     chunking_strategy="by_title",
14
15     # Sets a hard limit on the number of characters allowed in each chunk.
16     max_characters=4000,
17     # This function will attempt to create a new chunk after every 3,800 characters, allowing some flexibility while chunking.
18     new_after_n_chars=3800,
19     # If a chunk has fewer than 2,000 characters, the function will
20     # attempt to combine it with neighboring text blocks to create a larger, more meaningful chunk.
21     combine_text_under_n_chars=2000
22 )

```

Note: You'll notice that a folder called 'figures' is created, where all the extracted images are stored. To make these images more useful, we can generate embeddings for them, which can later be used for semantic search. Typically, embedding models are specialized—they either handle text or images, but not both. This means the vector lengths won't align, and using different models for text and image embeddings can complicate similarity searches. However, the CLIP model can generate embeddings for both images and text, making it a versatile option for this purpose. To simplify things, we can first create a text summary for each image, and then generate embeddings based on those summaries. This approach allows us to leverage text-based semantic search while still incorporating the visual information from the images.



Note: After extracting images from a PDF or any other source, you can use the below method to send them to a multimodal model like GPT-4o to understand and generate insights about the images. This approach allows you to combine textual and visual data for a richer understanding of the content.

SET OPENAI TOKEN - MULTIMODAL

Note: First, create an account from the [OpenAI official website](#).

- Create a new project API key by browsing to API Keys web page. Select Create new secret key. The API key is automatically generated. Save the APi Key as we will be using it in the later steps .

Save your key

Please save this secret key somewhere safe and accessible. For security reasons, **you won't be able to view it again** through your OpenAI account. If you lose this secret key, you'll need to generate a new one.

JWZs13myG1i0gbLGkU9XHe6wEAhEsRwm8gxjX6eYwf8A

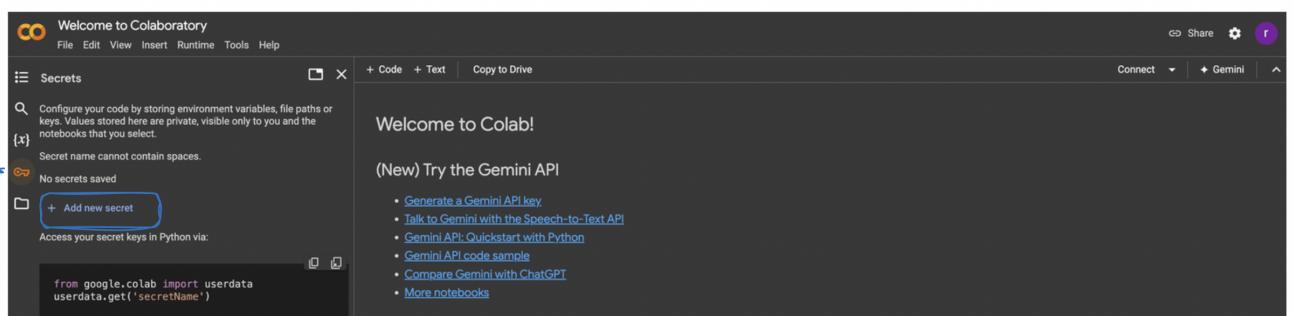
Copy

Permissions

Read and write API resources

Done

- Within your existing Google Colab notebook navigate to the new “Secrets” section in the sidebar.



- Click on “Add a new secret.” Enter the name example: OPENAI_API_KEY and value of the secret(the API key created above). Note: The name is permanent once set.
- The list of secrets is global across all your notebooks.
- Use the “Notebook access” toggle to grant or revoke access to a secret for each notebook.

Secrets

Configure your code by storing environment variables, file paths or keys. Values stored here are private, visible only to you and the notebooks that you select.

Secret name cannot contain spaces.

Notebook access	Name	Value	Actions		
<input type="button" value="X"/>	HF_TOKEN	<input type="button" value=""/>	<input type="button" value=""/>	<input type="button" value=""/>
<input checked="" type="button" value=""/>	OPENAI_API_KE	<input type="button" value=""/>	<input type="button" value=""/>	<input type="button" value=""/>

+ Add new secret

Create Gemini API key

Access your secret keys in Python via:

```
from google.colab import userdata
userdata.get('secretName')
```

Lets import and load the environment variables

```
1 from langchain.chat_models import ChatOpenAI
2 from langchain.schema.messages import HumanMessage
3 import os
4 from dotenv import load_dotenv
5 from PIL import Image
6 import base64
7 import io
8 load_dotenv()
```

Incorporating Secrets into Your Code - We will use it later in our lab

```
1 import os
2 from google.colab import userdata
3 os.environ['OPENAI_API_KEY'] = userdata.get('OPENAI_API_KEY')
4 llm = ChatOpenAI(model="gpt-4o")
```

Function to Convert Image to Base64:

```
1 # Function to convert image to base64
2 def image_to_base64(image_path):
3     with Image.open(image_path) as image:
4         buffered = io.BytesIO()
5         image.save(buffered, format=image.format)
6         img_str = base64.b64encode(buffered.getvalue())
7         return img_str.decode('utf-8')
8
9 image_str = image_to_base64("/content/figures/figure-1-1.jpg")
```

Note: This is the image that was extracted from our data and saved in the figures folder.

```
1     print(image_str)
```

Output: Our base64-encoded string

Initializing the Multimodal and sending the Image to GPT-4o for Analysis:

```

1  chat = ChatOpenAI(model="gpt-4o", max_tokens=1024)
2
3  msg = chat.invoke(
4      [
5          HumanMessage(
6              content=[
7                  {"type": "text", "text": "Please give a summary of the image provided. Be descriptive"},
8                  {"type": "image_url", "image_url": {
9                      "url": f"data:image/jpeg;base64,{image_str}"
10                 }
11             },
12         ]
13     ]
14 )
15 )

```

- Retrieving the Response:

```
1  msg.content
```

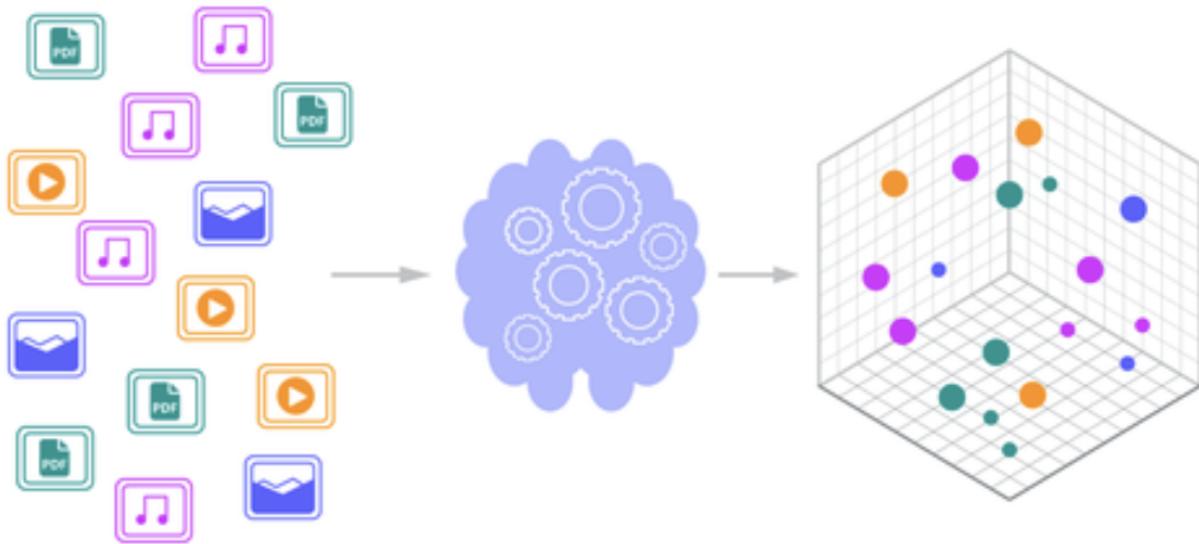
Output: The response from the model is stored in `msg.content`, which will contain the descriptive summary of the image.

The image is the logo of Cisco Systems, Inc., commonly known as Cisco. The logo features the company name "cisco" written in lowercase letters in a distinct, modern sans-serif font. Above the text, there is a series of vertical bars of varying heights, which is stylized to resemble the Golden Gate Bridge. The logo is rendered in a light blue color. The design is clean, professional, and easily recognizable, reflecting Cisco's identity as a major tech company specializing in networking hardware, telecommunications equipment, and other high-technology services and products.

Note: The description of your image may vary depending on the image retrieved.

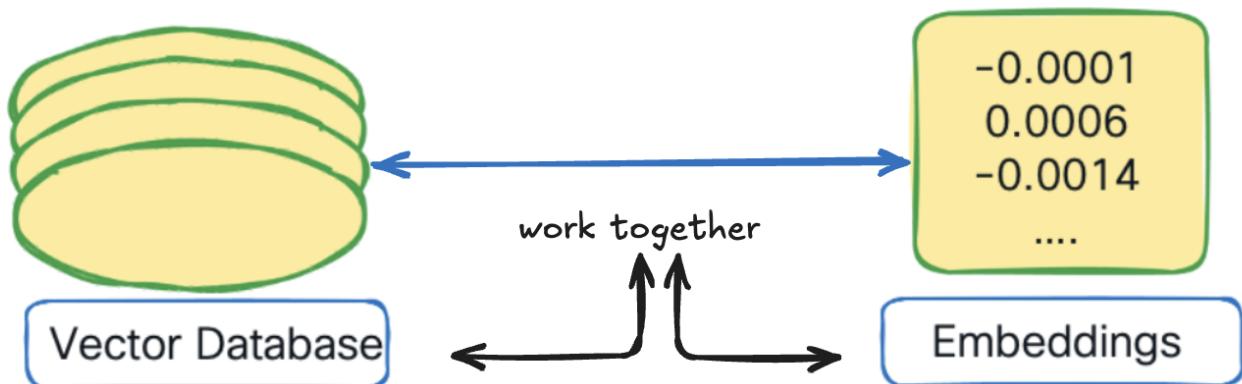
1.5 Task 4: Embeddings and Vector Database

- Vector databases are special databases meant to store content such as PDFs, Blogs, Word Documents, Images, Audio Files, and even Videos as embedding vectors, allowing for semantic based retrieval.



- In the process of a large language model, embedding is generally the first step to convert discrete tokens (like words or subwords) into dense vector representations, which will allow the rest of the network to do the math necessary to predict the next word. Embeddings and vector databases are essential if you are building any kind of AI product.

Embeddings and Vector Database



1.5.1 What are embeddings?

Embeddings are basically data.

Embeddings and Vector databases

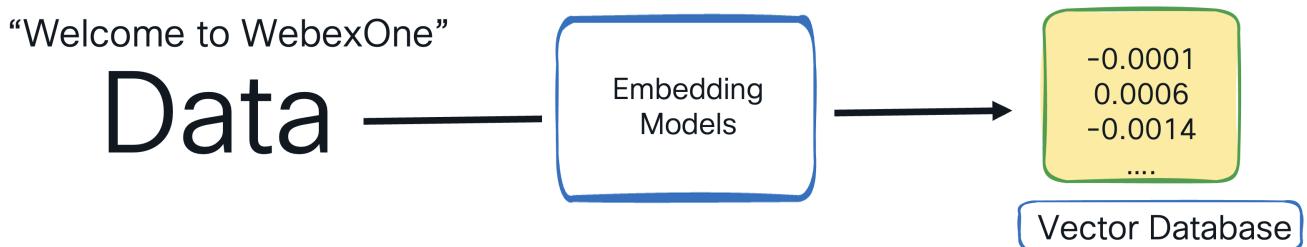
What are Embeddings?



That are converted into array of numbers called vectors that contain pattern of relationship, and can be used for similarity search.

Embeddings and Vector databases

What are Embeddings?



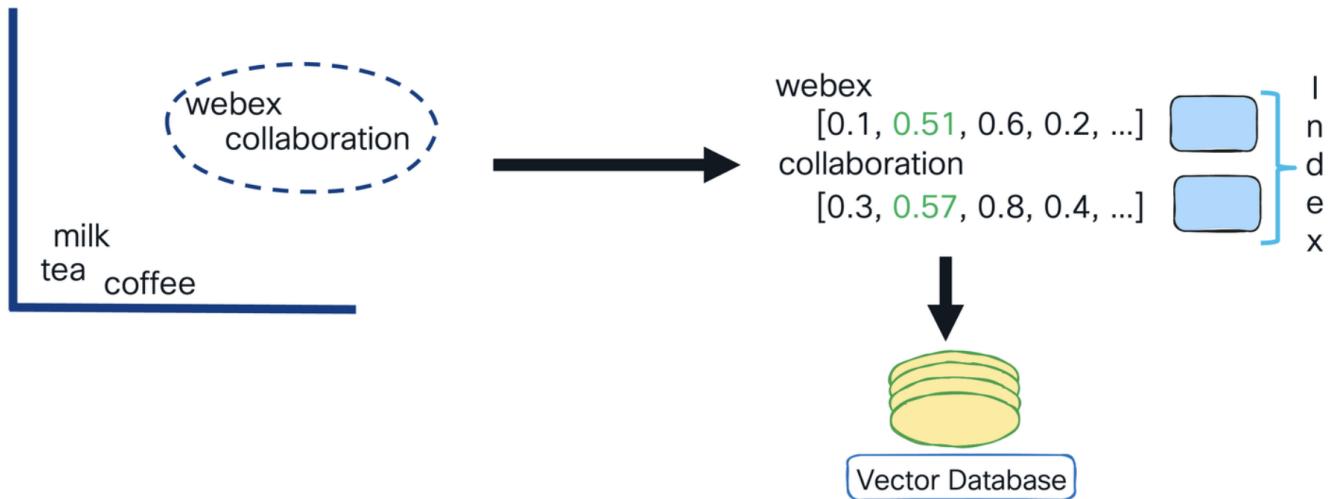
Let me explain this concept using a 2D graph. In this graph, words like "Webex" and "collaboration" are often used in similar context, so their embeddings essentially vector representations are positioned close to each other. These vectors are numerical arrays that computers can easily interpret.

The advantage of using vectors is that you can find similar items by calculating the distances between their vectors, a process known as nearest neighbor search. However, simply storing these embeddings isn't enough. Performing queries across thousands of vectors can be incredibly slow. To overcome this, these vectors also need to be indexed in the vector database.

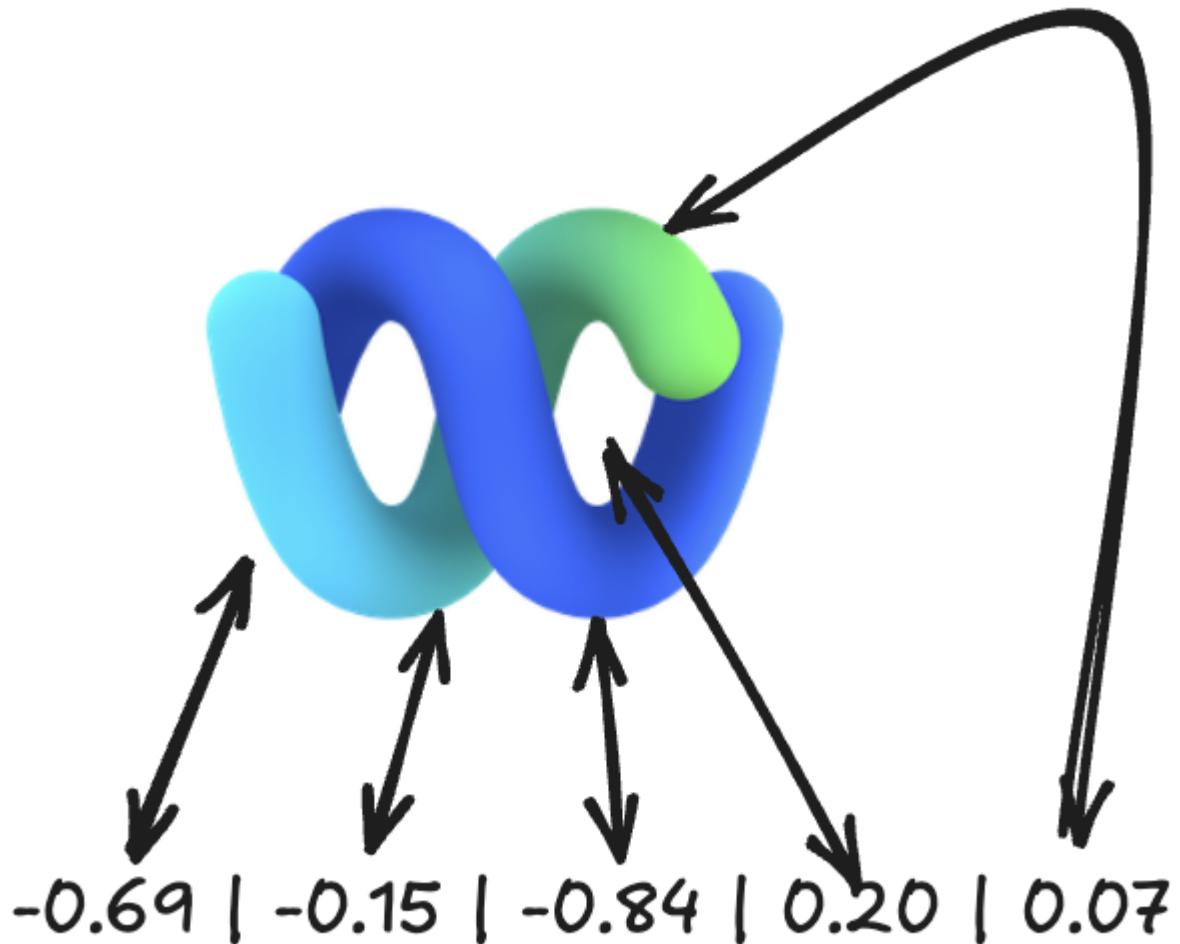
Note: In reality vectors can have hundred of dimensions

Embeddings and Vector databases

What are Embeddings?



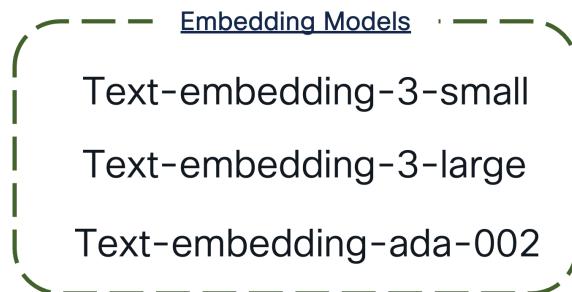
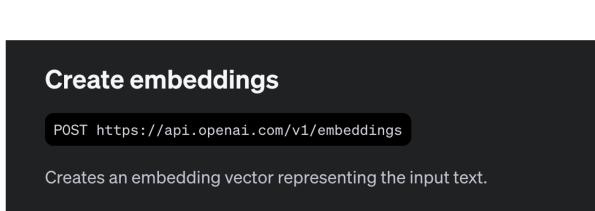
Similarly, images are also broken into vectors, which are arrays of numbers that machines can process. Once these embeddings (both for words and images) are generated, they are stored in a specialized database known as a vector database.



There are numerous embedding models available, such as Google's Word2Vec, CLIP (Contrastive Language-Image Pretraining), and even those provided by OpenAI, which offer excellent capabilities for generating embeddings. However, the challenge is that these models don't include tools for storing and managing those embeddings. That's where vector databases become essential.

Embeddings and Vector Database

OpenAI - Embeddings



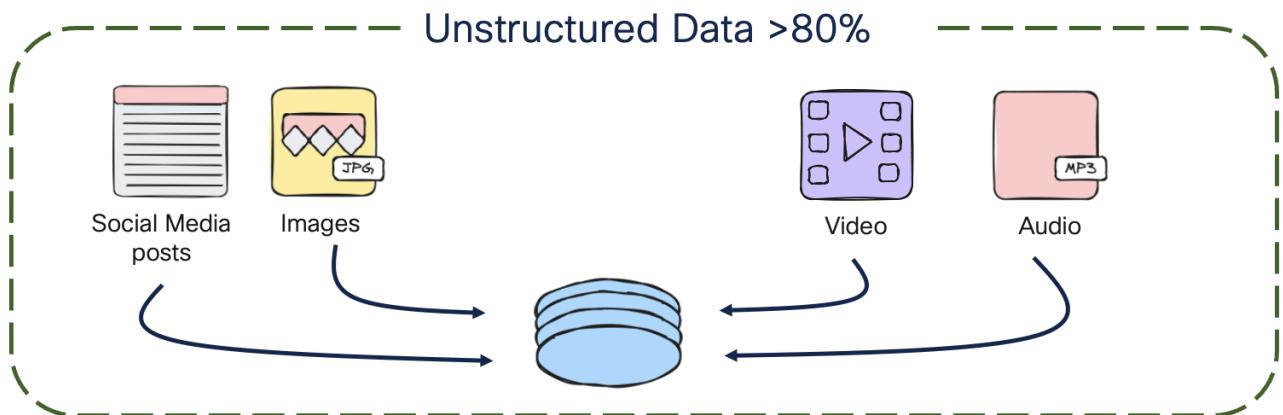
Note: In this lab session, we will be using OpenAI embeddings. More info can be found [here](#)

1.5.2 Why Vector Databases when we have relational DB's?

Around 80% of the data we encounter is unstructured, including social media posts, images, and videos. This type of data doesn't easily fit into traditional relational databases, which is where vector databases come into play.

Embeddings and Vector databases

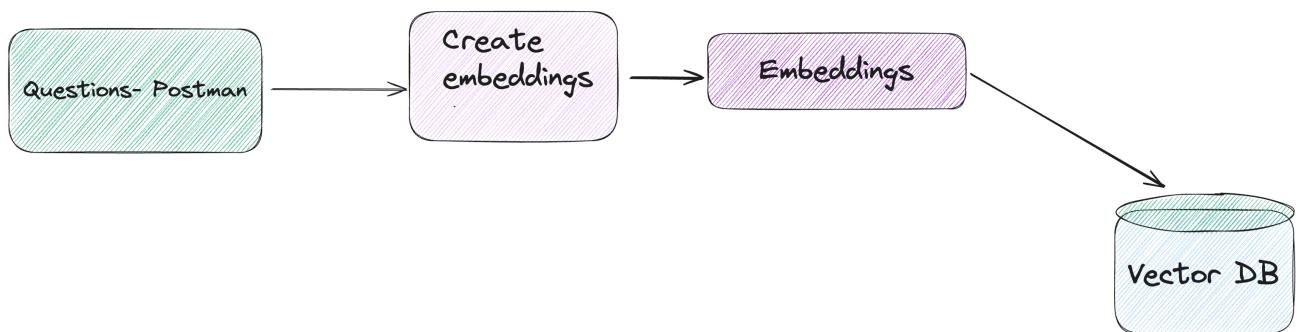
Why Vector Database?



1.5.3 Understanding working of Embeddings and Vector Databases

Practical Example of Embedding Techniques Using Postman

Embeddings and Vector databases High Level Overview



SET OPENAI TOKEN

Note: If you've already obtained the token in previous steps, you can skip this section

Note: First, create an account from the [OpenAI official website](#).

- Create a new project API key by browsing to API Keys web page. Select Create new secret key. The API key is automatically generated. Save the APi Key as we will be using it in the later steps .

Save your key

Please save this secret key somewhere safe and accessible. For security reasons, **you won't be able to view it again** through your OpenAI account. If you lose this secret key, you'll need to generate a new one.

JWZs13myG1i0gbLGkU9XHe6wEAhEsRwm8gxjX6eYwf8A

Copy

Permissions

Read and write API resources

Done

Note: We will use the OpenAI key in Postman to generate embeddings.

LOGIN TO POSTMAN

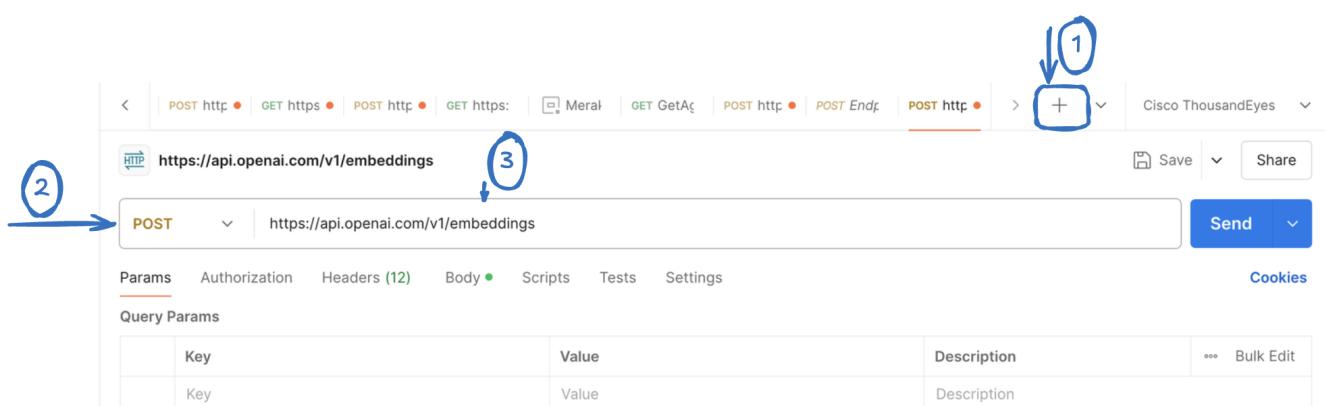
If you haven't installed Postman on your machine yet, you can download it from the following link

As we will be using embedding APIs to create an embedding vector that represents our input text, more information can be found at the following link

The POST request we will be using:

```
https://api.openai.com/v1/embeddings
```

* Open Postman, create a new request by pressing +. Select POST as your request and enter the request url



- Click on the Headers tab and enter the Following creds

```
Authorization: Bearer <replace with your openAI API key>
Content-Type : application/json
```

HTTP <https://api.openai.com/v1/embeddings>

POST https://api.openai.com/v1/embeddings

Params Authorization Headers (12) Body Scripts Tests Settings Cookies

Headers 10 hidden

Key	Value	Description	Bulk Edit	Presets
<input checked="" type="checkbox"/> Authorization	Bearer sk-Xm2PdTmYBhBCVf0dvpOhT3BlkFJ1ht...			
<input checked="" type="checkbox"/> Content-Type	application/json			
Key	Value	Description		

Enter your API key and Content-Type

- In this lab, we will be using the text-embedding-ada-002 model, but feel free to use any other embedding model of your choice.

Embedding models

OpenAI offers two powerful third-generation embedding model (denoted by -3 in the model ID). You can read the embedding v3 [announcement blog post](#) for more details.

Usage is priced per input token, below is an example of pricing pages of text per US dollar (assuming ~800 tokens per page):

MODEL	~ PAGES PER DOLLAR	PERFORMANCE ON MTEB EVAL	MAX INPUT
text-embedding-3-small	62,500	62.3%	8191
text-embedding-3-large	9,615	64.6%	8191
text-embedding-ada-002	12,500	61.0%	8191

- In Postman, click on the "Body" tab, enter the following information, and then press "Send."

```
{
  "input": "What is WebexOne 2024? WebexOne is an annual in-person and virtual event that takes place over four days. It's an event focused on AI collaboration and customer experience, and it features a range of activities such as insightful breakout sessions, technical training courses, hands-on labs, inspiring keynotes, epic entertainment, a solutions showcase and expo, customer awards, meet the experts, 1:1 executive meetings, a partner program, and more!",
  "model": "text-embedding-ada-002"
}
```

HTTP <https://api.openai.com/v1/embeddings>

POST https://api.openai.com/v1/embeddings

Params Authorization Headers (12) **Body** Scripts Tests Settings Cookies

Body none form-data x-www-form-urlencoded raw binary GraphQL JSON

Beautify

Input

Model

Send

- You will receive a 200 OK message, confirming success, and you'll notice that our text has been converted into embeddings (vectors). In the upcoming steps, I will show you on how to manually save this information in a vector database.

```

1 {
2   "object": "list",
3   "data": [
4     {
5       "object": "embedding",
6       "index": 0,
7       "embedding": [
8         -0.0042485874,
9         -0.02296605,
10        0.011213377,
11        0.01114761,
12        -0.027280405,
13        0.0023528363,
14        -0.004350527,
15        0.0043768343,
16        -0.00068562734,
17        -0.011660597,
18        0.010397859,
19        0.018730614,
20        -0.0158763,
21        -0.013469206,
22        0.0026997605,
      ]
    }
  ]
}

```

Embeddings

Note: Text embedding models, convert text into numerical data(embeddings) that represent text meaning.

INSERTING VALUES IN VECTOR DATABASE

We have a variety of databases available. In the upcoming step, I will demonstrate how to use SingleStore as a vector database and save our embeddings there.

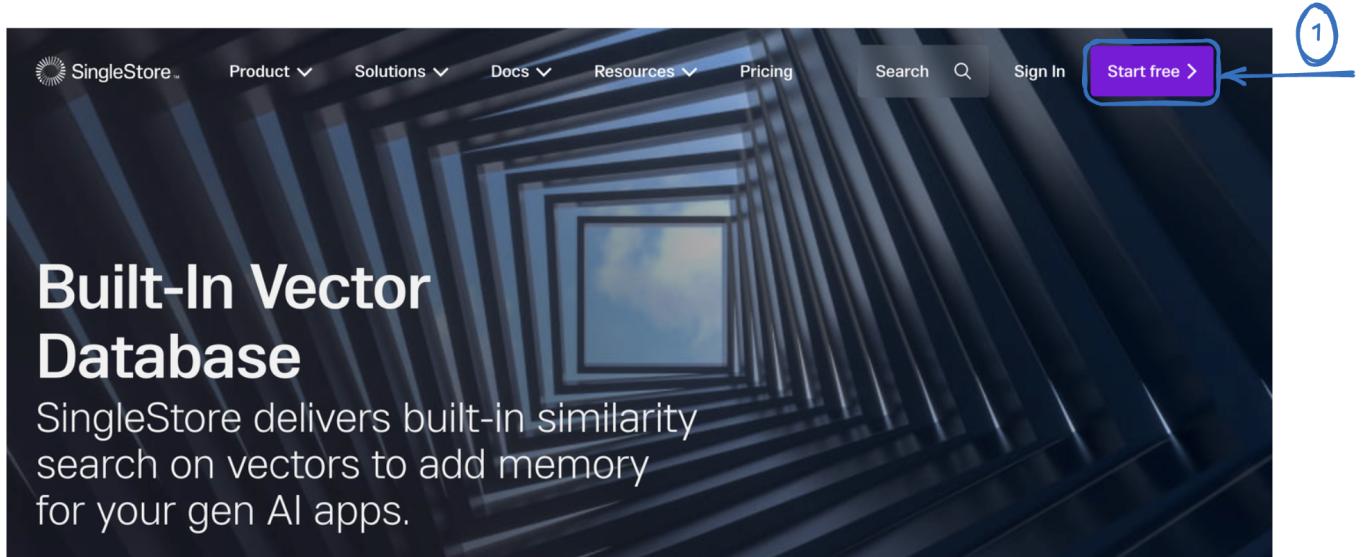
Note: This manual step is simply to illustrate how embeddings are stored in vector databases.

Embeddings and Vector Database

A few examples of Vector Database?



Note: We will create a Free account on SingleStore since it provides some credits to set up a trial account.



- I'm choosing to create an account using Google, but feel free to use any other method that works for you.



Try it free



Sign up with Google

Sign up with Microsoft

OR

First Name

Last Name

Enter your first name

Company Email Address

This email address will be used as your username

Continue →

Already have an account? [Sign In](#)

- Select Continue

 Sign in with Google



Sign in to SingleStore

By continuing, Google will share your name, email address, language preference, and profile picture with SingleStore. See SingleStore's [Privacy Policy](#) and [Terms of Service](#).

You can manage Sign in with Google in your [Google Account](#).

[Cancel](#) [Continue](#)

1



Add your details

Please provide your details to receive a personalised experience and make the most of your trial.

JOB TITLE *

Developer / Software En... ▾

COUNTRY *

United Kingdom of Grea... ▾

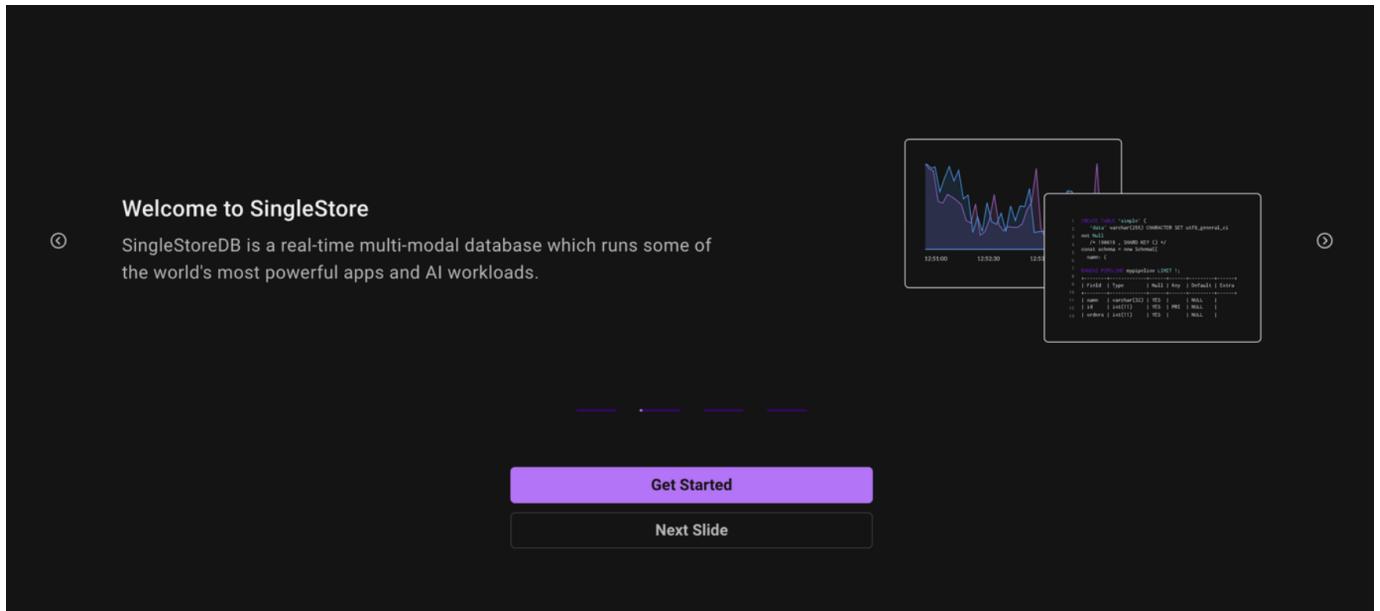
Company Name *

Omer

Company Email Address

omer.ilias@boldbetz.com

Continue >



- Click on New Deployment

1

- Give your workspace a name, keep all other settings at their default, and click "Next."

The screenshot shows the configuration page for creating a new workspace. A blue circle labeled '1' points to the 'Workspace Group Name' field, which contains 'OpenAI VectorDB'. A blue circle labeled '2' points to the 'Next' button at the bottom right.

Standard Workspace

Starter PREVIEW FREE
Ideal for starting projects and prototyping with our Shared Edition

Standard
Get full functionality and performance of SingleStore deployed on dedicated compute

1 Workspace Group Name
OpenAI VectorDB
A Workspace Group enables you to share data and access settings between multiple databases and workspaces. It is good practice to label a group as either: dev, prod or staging.

2 Cloud Provider
 AWS
 GCP
 Azure

3 Region
US East 1 (N. Virginia)
Want to try SingleStore in another region? Contact Us.

Advanced Settings Optional
Deploy In Two Availability Zones
Enable deployment across two Availability Zones, incurring a 1.3x increase in credit cost.

Summary
Workspace Group: OpenAI VectorDB
Provider and Region: AWS US East 1 (N. Virginia)
Workspace Name: —

Compute
Size: —
Credits per Hour: —
Credits Available: 60 CR
Free Trial Credits Applied: 60 Credits (US\$216). More will be added upon completion of the onboarding checklist.

Next

- Then, leave everything as is and click on "Create Workspace."

Create Workspace

Workspaces provide compute resources to run operations on a database. After you set up your first workspace you can create databases and attach them to this workspace to perform different types of operations.

1 Workspace Details

Workspace Name
workspace-1

Choose a name that will help identify what type of workloads will be connected to this workspace

2 Size

S-00			Edit
vCPUs	RAM	Credits per hour	
2	16 GB	0.25 CR	

Compute consumes a fixed number of credits per hour, and can be scaled up or down to increase or decrease performance. [Learn about credits](#) or [contact us to prepay for a discount](#).

Storage is separate from Compute
Storage is billed separately from compute charges and are based on the average total number of GB used per month.

Settings

SingleStore Kai Enable the MongoDB API and endpoint for your workspace.

Deployment Type
Updates roll out to non-production deployments prior to production.
Production

Auto Suspend

Summary

Workspace Group
OpenAI VectorDB

Provider and Region
AWS US East 1 (N. Virginia)

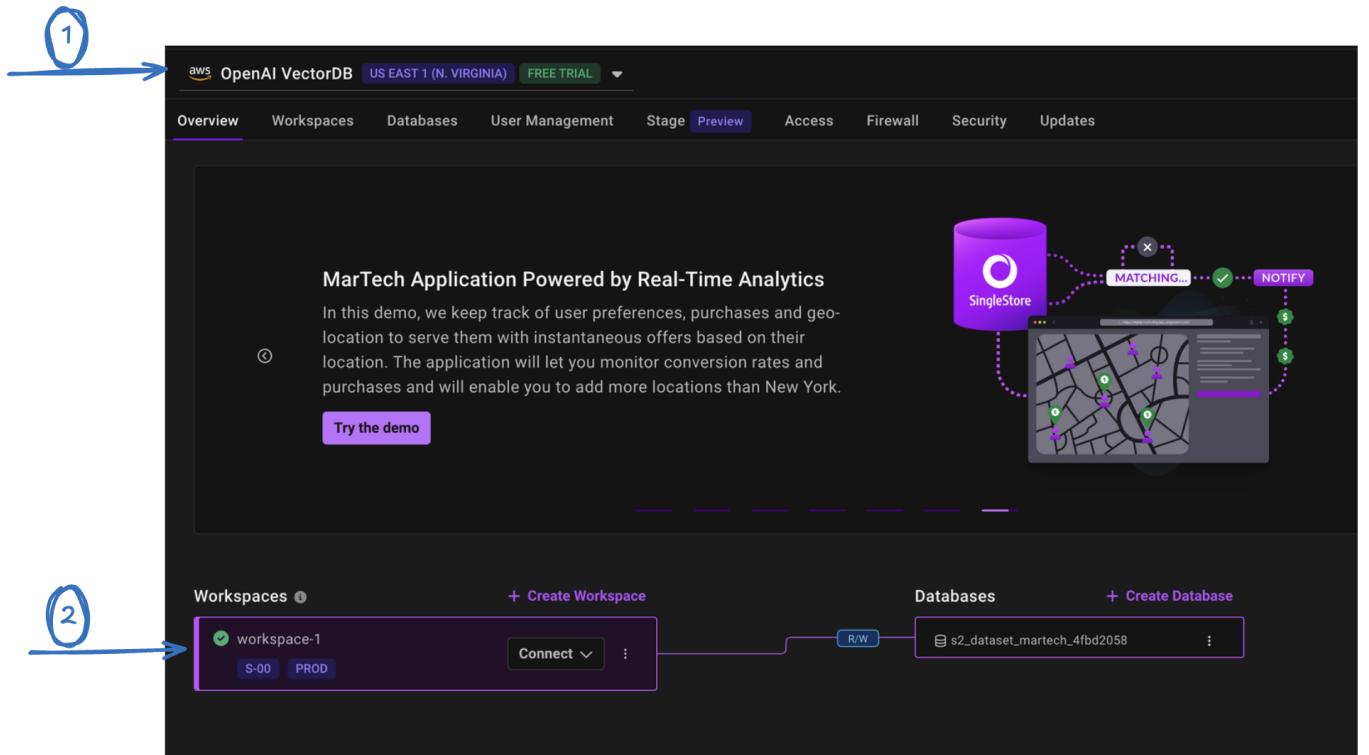
Workspace Name
workspace-1

Compute

Size	S-00
Credits per Hour	0.25 CR / hour
Credits Available	60 CR

Free Trial Credits Applied: 60 Credits (US\$216).
More will be added upon completion of the onboarding checklist.

Create Workspace



- Now we can create a database

aws OpenAI VectorDB US EAST 1 (N. VIRGINIA) FREE TRIAL

Overview Workspaces Databases User Management Stage Preview Access Firewall Security Updates

MarTech Application Powered by Real-Time Analytics
In this demo, we keep track of user preferences, purchases and geo-location to serve them with instantaneous offers based on their location. Ingest live and historical data from different sources into your workspace. We support AWS S3, GCS, Azure Blob and Kafka.

Ingest your data through pipelines
Ingest live and historical data from different sources into your workspace. We support AWS S3, GCS, Azure Blob and Kafka.

Try the demo

Workspaces (1) **+ Create Workspace**

- workspace-1** S-00 PROD
- Connect** ⋮
- R/W**

Databases (1) **+ Create Database**

- s2_dataset_martech_4fb02058** ⋮

Create Database

Choose a new database name
WebexOneDB

Attach To
workspace-1

Attachment Type
Read & Write **Submit**

Pricing (Average GB per Month)
\$0.023

Cancel **Create Database**

Welcome to Workspaces

Workspaces are dedicated compute and memory resources. They are separate from databases giving flexibility on how you connect storage and compute resources.

Workspaces

- workspace-1 (S-00 PROD)

Databases

- s2_dataset_martech_4fbcd2058
- WebexOneDB

- Click on Database tab

Databases

MarTech Application Powered by Real-Time Analytics

In this demo, we keep track of user preferences, purchases and geo-location to serve them with instantaneous offers based on their location. The application will let you monitor conversion rates and purchases and will enable you to add more locations than New York.

Try the demo

Workspaces

- workspace-1 (S-00 PROD)

Databases

- s2_dataset_martech_4fbcd2058
- WebexOneDB

No tables or data yet

Name	Status	Table Count	Memory Usage	Index Disk Usage
WebexOneDB	ATTACHED	0	0 B	0 B
s2_dataset_martech_4fbdb2058	ATTACHED	12	91 MB	20 MB
information_schema	ATTACHED	198	-	-

- Return to the home screen and navigate to Develop > Data Studio. Open the SQL Editor.

1

2

- Be sure to select your workspace and database.

Your Workspace Your Database

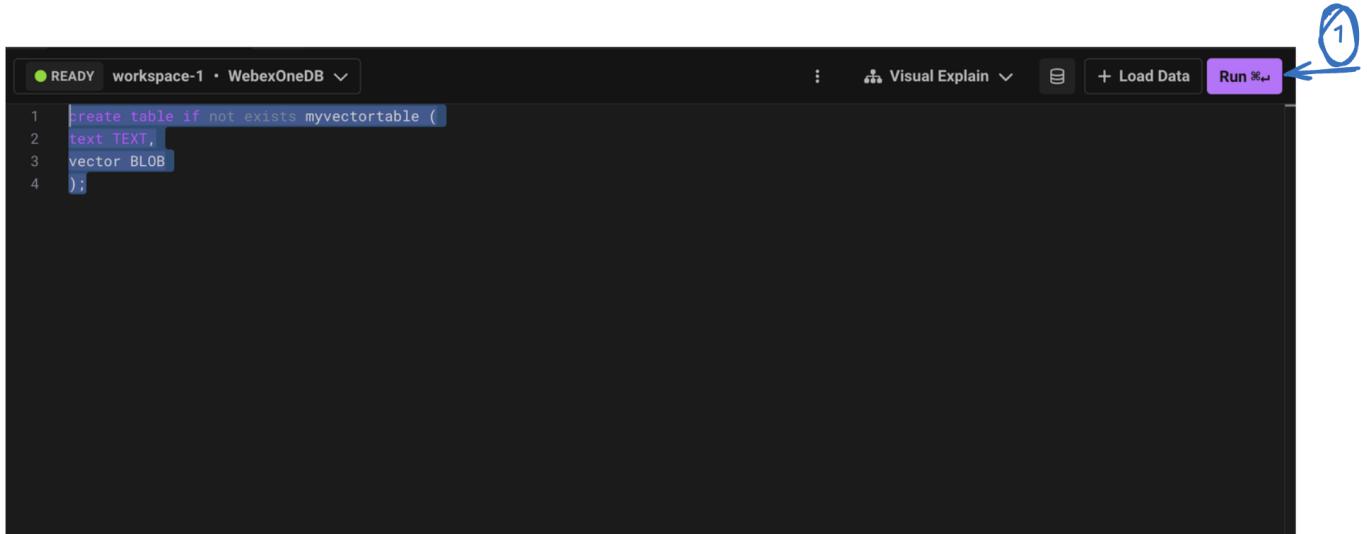
READY workspace-1 • WebexOneDB

Connection Database

workspace-1 WebexOneDB

- Run a SQL command to create a table - Press Run

```
create table if not exists myvectortable (
text TEXT,
vector BLOB
);
```



```

1 create table if not exists myvectortable (
2     text TEXT,
3     vector BLOB
4 );

```

- You can navigate to the Databases tab and verify that the table has been created.



Name	Status	Table Count	Memory Usage	Index Disk Usage	Disk Usage	Partition Count
WebexOneDB	ATTACHED	1	0 B	0 B	0 B	8

- Let's copy the embeddings we generated earlier using Postman so that we can insert them into our database.
- To proceed, copy the input text and all the embedding values, including the square brackets [], as shown in the image. These will be used for insertion into our database.

Copy Input

```

1 {
2   "input": "What is WebexOne 2024? WebexOne is an annual in-person and virtual event that takes place over four days. It's an event focused on AI collaboration and customer experience, and it features a range of activities such as insightful breakout sessions, technical training courses, hands-on labs, inspiring keynotes, epic entertainment, a solutions showcase and expo, customer awards, meet the experts, 1:1 executive meetings, a partner program, and more!",
3   "model": "text-embedding-ada-002"

```

Copy the vectors

```

1 {
2   "object": "list",
3   "data": [
4     {
5       "object": "embedding",
6       "index": 0,
7       "embedding": [
8         -0.0042485874,
9         -0.02296605,
10        0.011213377,
11        0.01114761,
12        -0.027280405,
13        0.0023528363,
14        -0.004350527,
15        0.0043768343,
16        -0.00068562734,
17        -0.011660597,
18        0.010397859,
19        0.018730614,
20        -0.0158763,
21        -0.013469206,
22        0.0026997605,

```

- Let's head back to our SQL editor and insert the values into our database. Once done press "Run"

```
insert into myvectortable (text ,vector) values ("your_input_text", JSON_ARRAY_PACK("your_embeddings"))
```

In this query:

```
* Replace "your_input_text" with the input values from Postman.
* Replace your_embeddings with the embeddings you copied earlier, including the square brackets [].

insert into myvectortable (text ,vector) values ("What is WebexOne 2024? WebexOne is an annual in-person and virtual event that takes place over four days. It's an event focused on AI collaboration and customer experience, and it features a range of activities such as insightful breakout sessions, technical training courses, hands-on labs, inspiring keynotes, epic entertainment, a solutions showcase and expo, customer awards, meet the experts, 1:1 executive meetings, a partner program, and more!", JSON_ARRAY_PACK([
-0.0042485874,
-0.02296605,
0.011213377,
0.01114761,
-0.027280405,
0.0023528363,
-0.004350527,
0.0043768343,
-0.00068562734,
-0.011660597,
0.010397859,
0.018730614,
-0.0158763,
-0.013469206,
0.0026997605,
-0.00003768792,
0.014850326,
-0.03183152,
0.0026175508,
-0.0459321,
-0.0055803815,
-0.008247258,
-0.0117313211,
0.0006169824,
-0.017112732,
-0.002885554,
-0.004113764,
-0.01348236,
```

[")
] "))

- After running the above command, you will see that our table now contains both the input text and the corresponding embeddings.

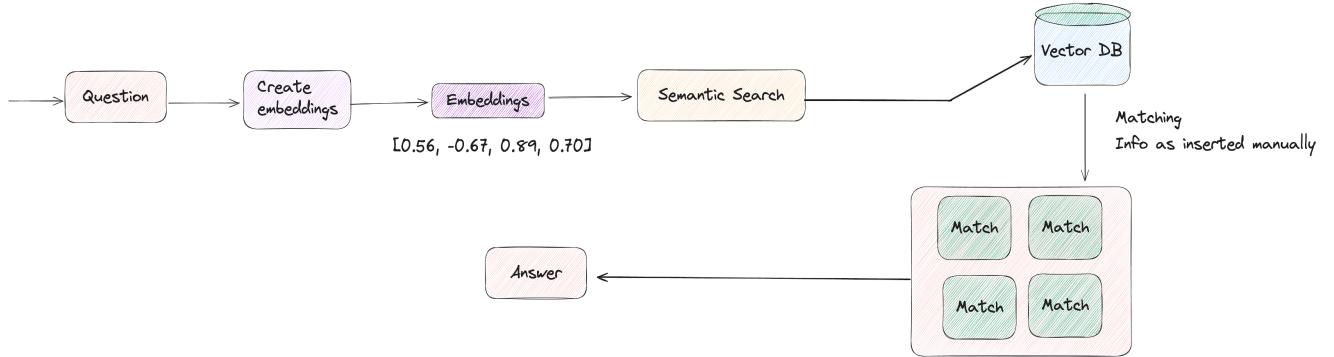
OpenAI VectorDB > WebexOneDB > myvectortable

Columns 2 Indexes 0 Sample Data 1 SQL

text	vector
What is WebexOne 2024? WebexOne is an annual in-person and virtual event that takes place over four days. It's an ...	[REDACTED]

RETRIEVING VALUES FROM THE VECTOR DATABASE

High Level Overview



- Searching the vector database is quite simple. Example: we want to find anything related to WebexOne. To do this, we'll create embeddings for our query, "Is WebexOne an annual event?" and then search the vector database to find matches against the existing embeddings.
- Let's open Postman and create an embedding for our question. Be sure to click on "Send."

The screenshot shows a Postman interface with a blue arrow pointing from step 1 to the request body and another arrow pointing from step 2 to the first embedding in the response.

```

1 {
  ...
  "input": "Is WebexOne an annual event?",
  ...
  "model": "text-embedding-ada-002"
}

```

```

1 {
  "object": "list",
  "data": [
    {
      "object": "embedding",
      "index": 0,
      "embedding": [
        -0.0015118581,
        -0.021267666,
        0.023597047,
        -0.0072638104,
        -0.02286653,
        0.005909599,
        -0.026119394,
        -0.010268575,
        0.0052169873,
        -0.009551842,
        0.024065679,
        0.020178782,
        0.015464887,
        0.016815653,
        -0.016333235,
        -0.0006641838,
        0.018069934,
        -0.010564916,
        0.007911626,
        -0.026367495,
        -0.0019417254,
        -0.005547787,
        -0.012604848,
        0.0003715036,
        -0.013962504,
        0.011729606,
        -0.00078737224,
        -0.011881223,
        0.007815143,
        0.009992908,
        ...
      ]
    }
  ]
}

```

- Let's head over to the SQL editor and run a query for our search.

```

select text,dot_product(vector,JSON_ARRAY_PACK("[your_embeddings]")) as score
from myvectortable
order by score desc
limit 5;

```

* In this query: * Replace your_embeddings with the embeddings you copied earlier, including the square brackets [] for the question.

```

select text,dot_product(vector,JSON_ARRAY_PACK([
  -0.0015118581,
  -0.021267666,
  0.023597047,
  -0.0072638104,
  -0.02286653,
  0.005909599,
  -0.026119394,
  -0.010268575,
  0.0052169873,
  -0.009551842,
  0.024065679,
  0.020178782,
  0.015464887,
  0.016815653,
  -0.016333235,
  -0.0006641838,
  0.018069934,
  -0.010564916,
  0.007911626,
  -0.026367495,
  -0.0019417254,
  -0.005547787,
  -0.012604848,
  0.0003715036,
  -0.013962504,
  0.011729606,
  -0.00078737224,
  -0.011881223,
  0.007815143,
  0.009992908,
  ...
])

```

```

0.020316616,
-0.0032580325,
-0.0066401153,
0.0097310245,
-0.0028652078,
-0.008145943,
-0.0074843434,
.....
]) as score
from myvectortable
order by score desc
limit 5;

```

The screenshot shows a SQL Editor interface with a dark theme. At the top, there's a toolbar with icons for file, edit, and database management. The main area is a code editor containing a SQL query. The query is a SELECT statement using the dot_product function to calculate similarity between a query vector and database vectors, ordered by score and limited to 5 rows. Below the code editor is a results table with two columns: 'text' and 'score'. The first row of results is highlighted. At the bottom of the results table, there's a message log and several export options: 'Save as CSV', 'Save to Stage', and 'Copy to Clipboard'. The status bar at the bottom indicates the query took 411 ms and returned 1 row.

text	score
What is WebexOne 2024? WebexOne is an annual in-person and virtual event that takes place over four days. It's an event focused on AI collaboration and customer ex...	0.8964636921882629

Note: You'll be able to see the success of the vector database—higher scores indicate better matches for the answer.

1.5.4 Conclusion

In summary, vector databases allow LLMs to have long-term memory. In this section, we explored how to use embeddings and vector databases by generating embeddings with the OpenAI API through Postman. After setting up a free SingleStore account, embeddings were stored in a vector database. The process included creating an embedding for a query, searching the database, and confirming that it effectively retrieved relevant information. This demonstrated how vector databases can efficiently manage and query embeddings, making it easier to find relevant information based on the data stored.

1.6 Task 5: Context Windows and Retrieval Augmented Generation (RAG)

What is Context window

- The context window is the limit on how much text the model can keep in mind and process at once. It includes your messages, the model's replies, and any other text in the conversation. If the conversation gets too long and goes beyond this limit, the model might start to forget what was said earlier.
- Example
- For GPT-3.5, the context window limit is approximately 8,000 tokens (one token typically represents about 4 characters on average.). If the input text along with previous interactions (if any) exceeds the limit of LLM, the model may not be able to see or remember parts of the text that fall outside this window. This can result in errors or the inability to refer back to earlier parts of the conversation or document.

Context Window

Model Context Window = 8K

The screenshot shows the ChatGPT 3.5 interface. At the top, there is a blue header bar with the text "ChatGPT 3.5". Below this, the main content area displays a large list of academic papers, likely from a search or retrieval system. The list includes references such as [25] Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words with subword units. arXiv preprint arXiv:1508.07909, 2016. and [29] Ilya Sutskever, Oriol Vinyals, and Quoc VV Le. Sequence to sequence learning with neural networks. In Advances in Neural Information Processing Systems, pages 3104–3112, 2014. The list continues with many other references. At the bottom of the list, there is a small number "11" followed by a blue circular icon with a dot. Below the list, there is a red error message box containing the text: "The message you submitted was too long, please shorten it and submit something shorter." Below the error message, it says "There was an error generating a response". At the very bottom of the interface, there is a green button labeled "Regenerate".

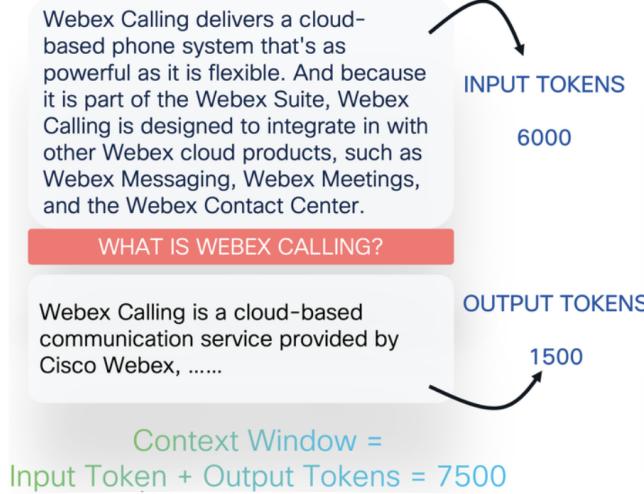
Outside of context window →

Note: A token can be as small as a piece of a word or as large as a word itself, depending on the language and complexity of the text

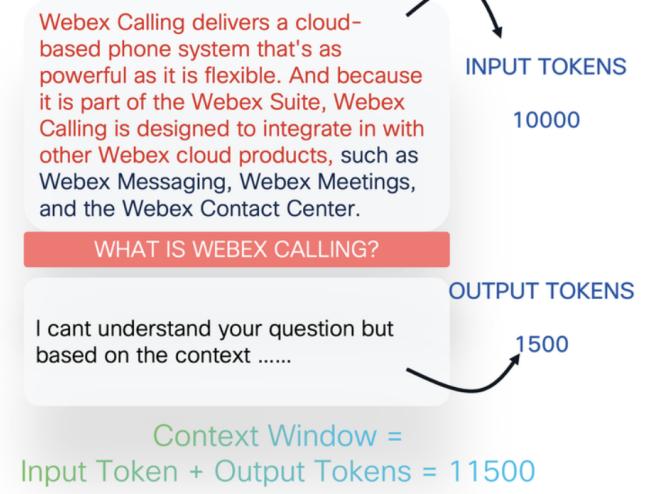
- Consider a model with a context window limit of 8,000 tokens. If we input 6,000 tokens and the model generates a response of 1,500 tokens, the total token count is 7,500. This is within the model's context window, so all information is processed correctly. Conversely, if the input increases to 10,000 tokens and the response remains at 1,500 tokens, the total becomes 11,500 tokens. This exceeds the model's 8,000-token limit. Consequently, the model might lose context or be unable to respond properly as the excess information falls outside its context window. This demonstrates the importance of managing input size to stay within the model's processing capabilities.

Context Window

Model Context Window = 8K



Model Context Window = 8K



1.6.1 How Context window work

- In the below example , multiple potential predicted words exist, but the model chooses the words based on the surrounding context and their probability. However, if the context window size increases and the word "Webex" moves out of the window, the model may lose critical context. As a result, the next predicted word might be incorrect, as the model can no longer reference "Webex" to inform its predictions accurately. This highlights the importance of keeping key information within the context window to maintain the accuracy of the model's predictions.
- Additionally, the temperature setting in the model can influence the type of content it generates. Lower temperatures are useful for summarizing or generating more deterministic and concise outputs, while higher temperatures encourage creativity, making the model more suitable for writing stories or poems. Adjusting the temperature allows you to control the balance between predictability and creativity in the model's output.

1	2	3	4	5	6	7	8	Predicted word
					Webex	Calling	is	cloud
				Webex	Calling	is	cloud	based
			Webex	Calling	is	cloud	based	telephony
		Webex	Calling	is	cloud	based	telephony	solution
	Webex	Calling	is	cloud	based	telephony	solution	that
Webex	Calling	is	cloud	based	telephony	solution	that	offers
Webex	Calling	is	cloud	based	telephony	solution	that	offers
Webex	Calling	is	cloud	based	telephony	solution	that	enterprise

How can we address the challenges related to the context window?

- That's where RAG comes into play

1.6.2 Retrieval Augmented Generation (RAG)

Lets try to understand what RAGs are all about before we can look how they can solve context window issues

RAG = Retrieval Augmented Generation

RAG = Retrieval Augmented Generation

Lets focus on the **Generation** part. Generation basically means responding to user query. For example, you might ask the model whether DX or Navigators are compatible with supporting ThousandEyes.

Retrieval Augmented Generation

Generation



Response to user
Query also known as
Prompt



Can have some
undesirable behavior

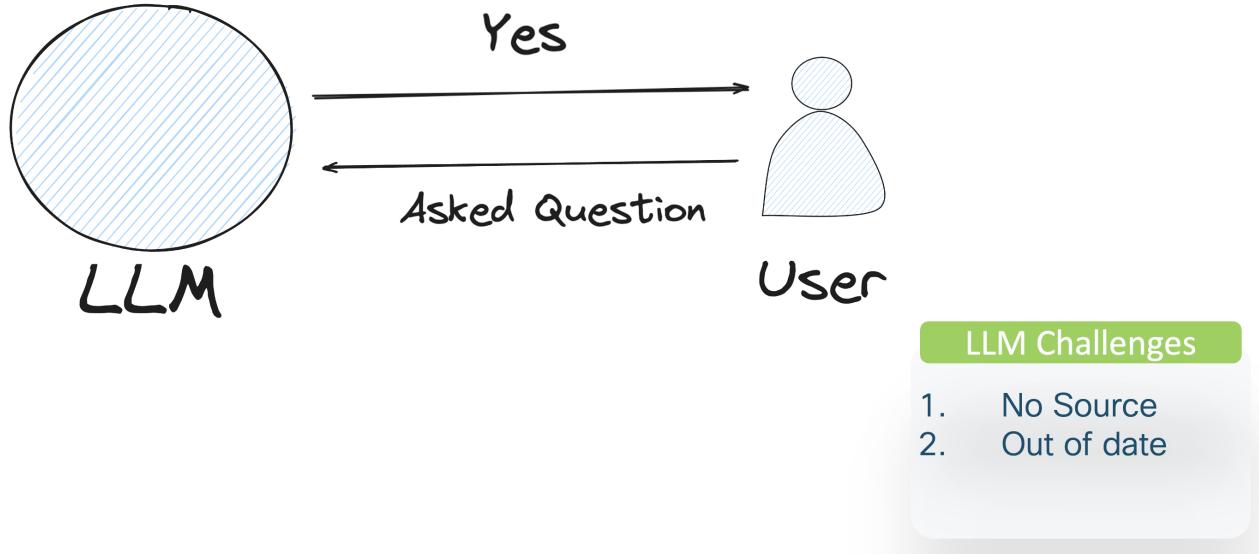


MTR?

ThousandEyes?

Source

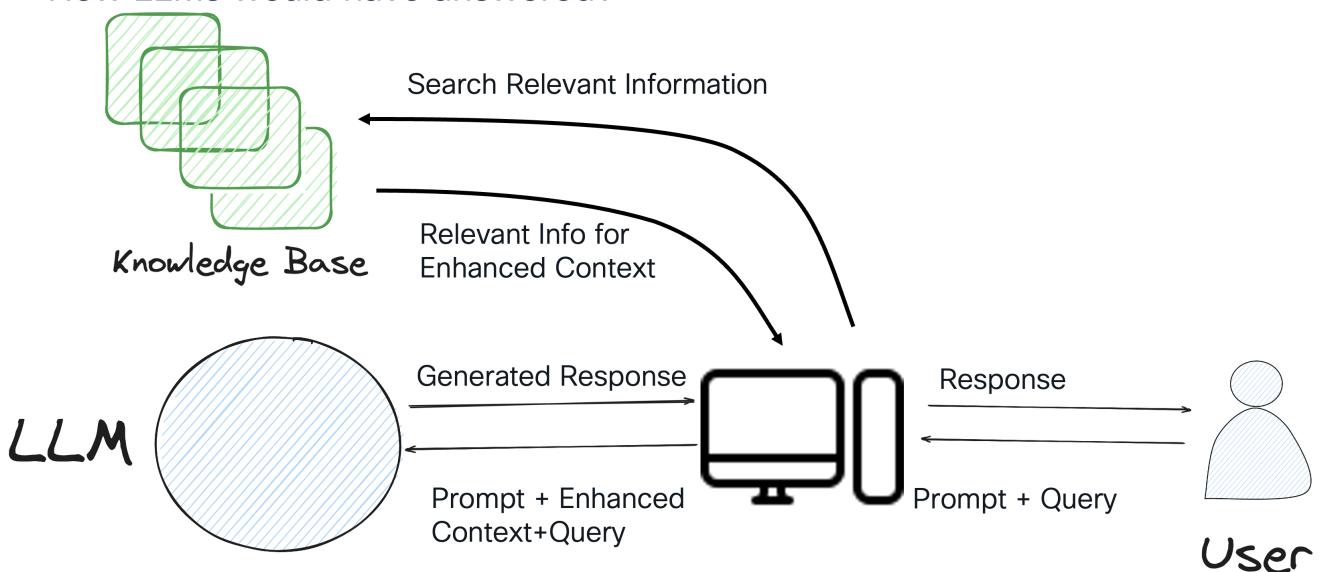
That's an excellent question. LLMs are trained on vast amounts of historical data, so they may recognize that DX and Navigators are well-regarded Cisco products. As a result, the LLM might confidently respond that you can install ThousandEyes on those devices. However, the issue is that the response may lack a verifiable source to back up that claim or the information can be out-of-date.



What if, instead of just asking the question, I had first consulted a reputable source and only then informed the user that ThousandEyes can't be installed on those devices? This is where Retrieval-Augmented Generation comes into play. By integrating a knowledge base—such as PDFs, CSVs, images, and other resources—into the model, we ensure that any time a question is asked, the response is based on the most current and accurate information available.

Retrieval Augmented Generation

How LLMs would have answered?



1.6.3 How can RAG address the limitations we've seen with the context window?

To address the context window limitation, we can use an embedding model to convert relevant text into embeddings and store them in a vector database. When a user submits a query, we transform it into an embedding and compare it to the stored embeddings to find the closest match.

Once we identify the match, we retrieve that specific chunk of text to respond to the query. This approach effectively overcomes the context window limitation by ensuring that our responses are accurate and contextually relevant.

RAG

Webex Calling delivers a cloud-based phone system that's as powerful as it is flexible. And because it is part of the Webex Suite, Webex Calling is designed to integrate in with other Webex cloud products, such as Webex Messaging, Webex Meetings, and the Webex Contact Center.

vector embeddings

```
[ -0.011466463, -0.02045446, -0.012214276, -0.005793769, -0.00562284, 0.012370961, -0.0033829627, 0.017833555, 0.0056691333, -0.033388063, 0.027633464, 0.011195826, 0.018887615, -0.010939433, -0.015568751, 0.009906739]
```



Compare embedded query to embedded text and find the closest match
Return closest chunk of text into prompt

User Query to embeddings

What is Webex Calling?

[-0.011466463, -0.02045446, -0.012214276, -0.005793769]



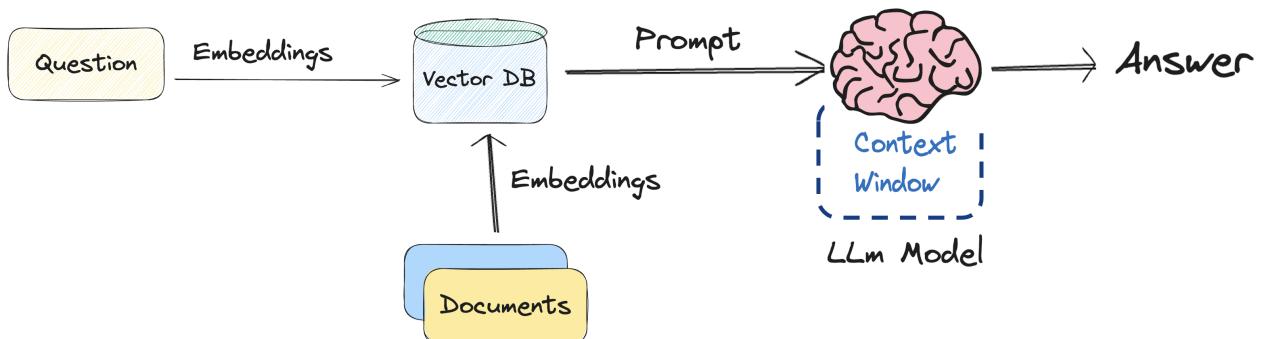
Webex Calling service provided by Cisco Webex,

What is Webex Calling?

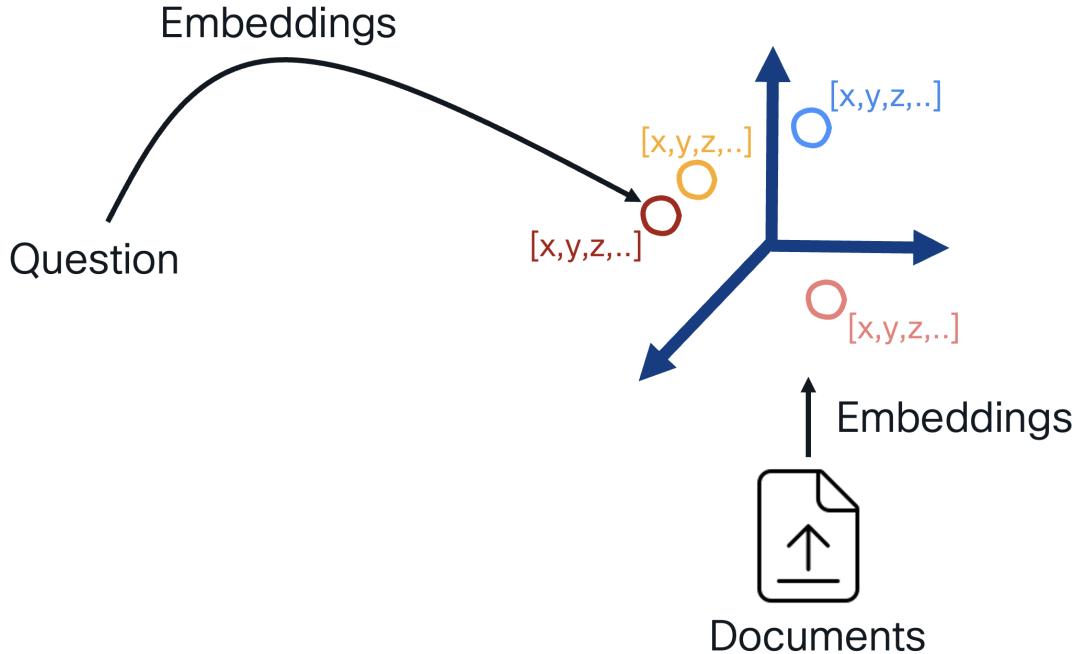
Webex Calling is a cloud-based communication service provided by Cisco Webex,

Answer the user query with fetched content

Now you can see how RAG can address some of the challenges with LLMs. Instead of re-training or fine-tuning the model with new information, we can augment our data sources to retrieve the most up-to-date information. Additionally, LLMs are now instructed to prioritize primary source data before generating responses, making them less likely to hallucinate and encouraging more accurate and reliable outputs.



1.6.4 What's Happening Under the Hood



We normally start by taking multiple documents and converting them into embeddings. Imagine these embeddings as points in a 3D space, where each document is projected based on its semantic meaning or content. Documents located near each other in this space contain similar semantic information, which is crucial when we search for or retrieve information.

Now, when we receive a query, we also convert it into an embedding (numerical value). We then perform a similarity search in this 3D space, looking for the documents that are closest to our query's embedding. These nearby documents are likely to contain the most relevant information.

Once we've identified these matching documents, we retrieve the relevant sections (or "splits") and feed all this information into the LLM's context window, allowing it to generate a well-informed response.

Large context window size. Do we still need RAG?

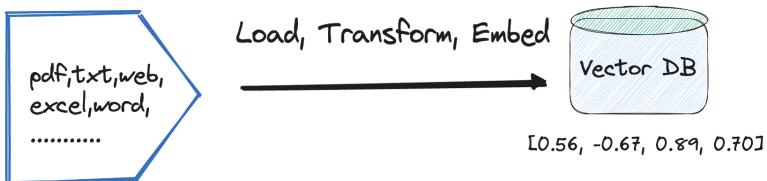
We've seen models like Gemini Pro with a context window size of millions, allowing them to retain more information. So, why use RAG? The answer is that RAG remains highly beneficial for sourcing real-time data, fact-checking, and accessing external knowledge that isn't contained within the context window.

1.6.5 Lets build a quick RAG application

While LLMs (Large Language Models) like GPT-4, Llama, Gemini possess impressive general knowledge, they don't inherently have access to your data. To connect LLMs with your own data sources, we can use frameworks like [LangChain](#). These frameworks enable us to leverage our own data by breaking down documents into smaller chunks and storing them as embeddings in vector databases.

This approach allows us to build applications that can effectively use language models. Example: When a user asks a question, we convert it into embeddings, perform a semantic search to find the most relevant answers, and then send both the question and the retrieved information to the LLM to generate a response.

In this section, we'll briefly introduce LangChain, but in the upcoming chapters, we'll dive deeper to gain a more comprehensive understanding of how it works.

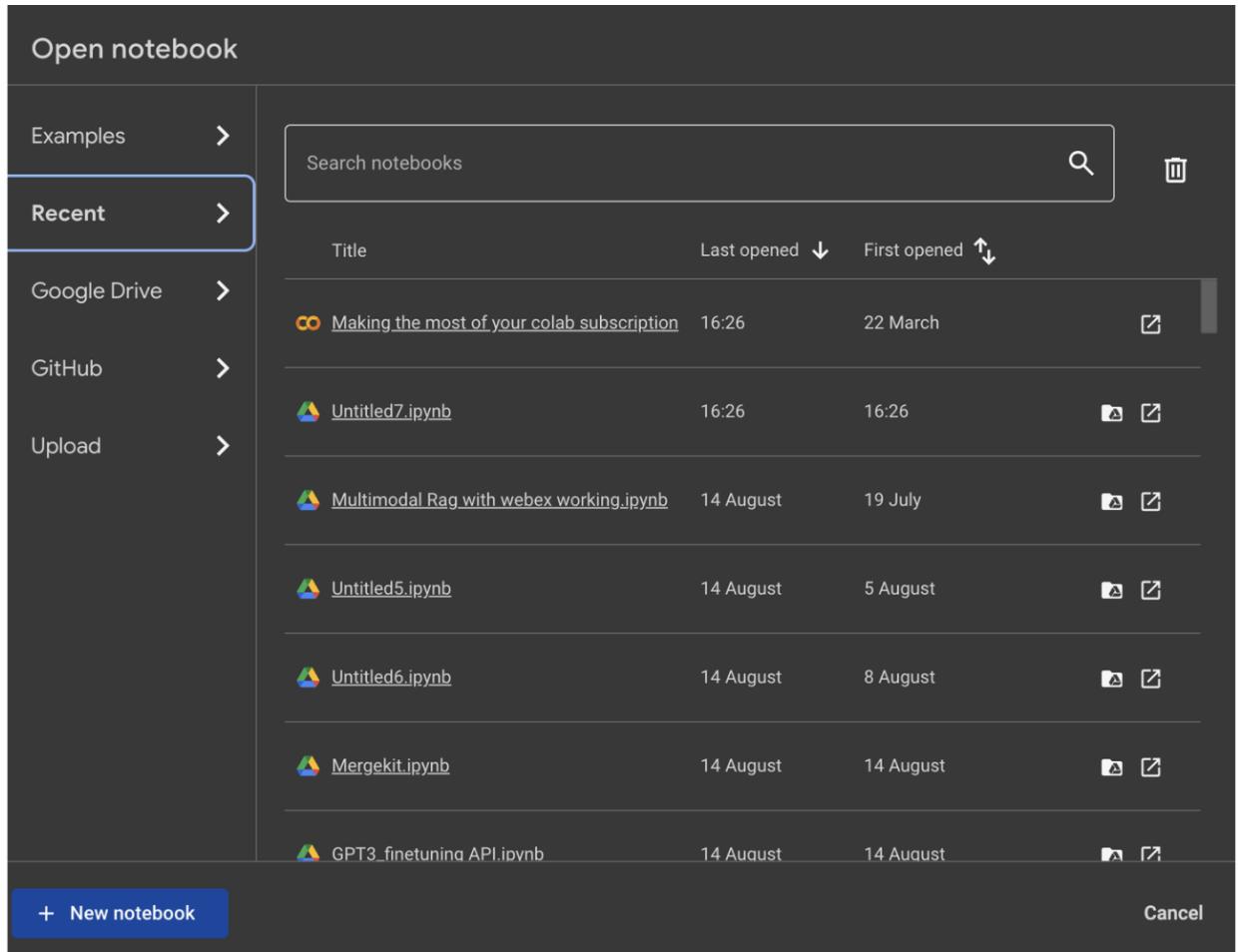


- Step 1: Load data source also called Data ingestion
- Step 2: Transform, where we break data into small chunks
- Step 3: Convert Chunks into vectors also called Embeddings
- Step 4: Save in Vector Database

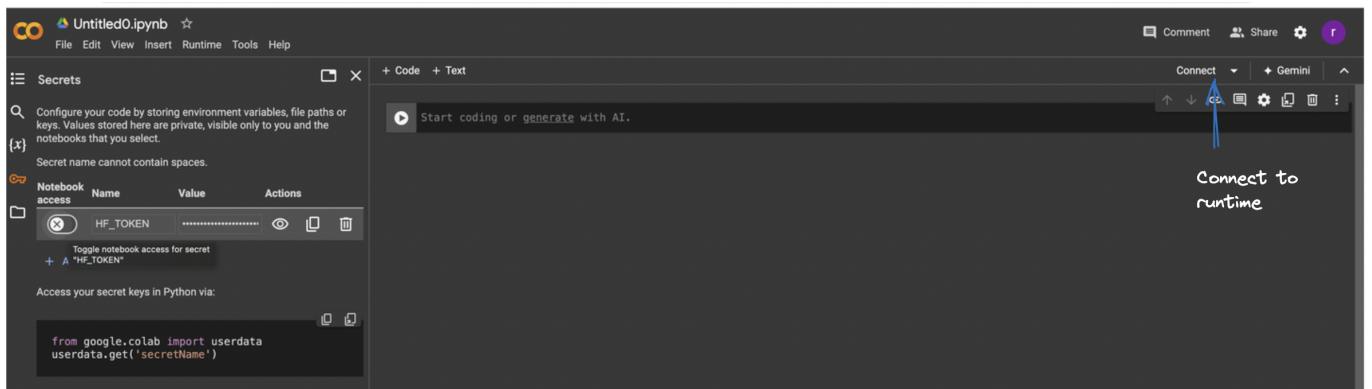
} **Entire RAG pipeline**

Note: In this section, we will be using the Cisco DECT 6800 Deployment Guide. For the purposes of this lab, I've modified the original PDF by removing a few pages and saving it as a new file. The PDF we are using in this lab can be downloaded from [here](#)

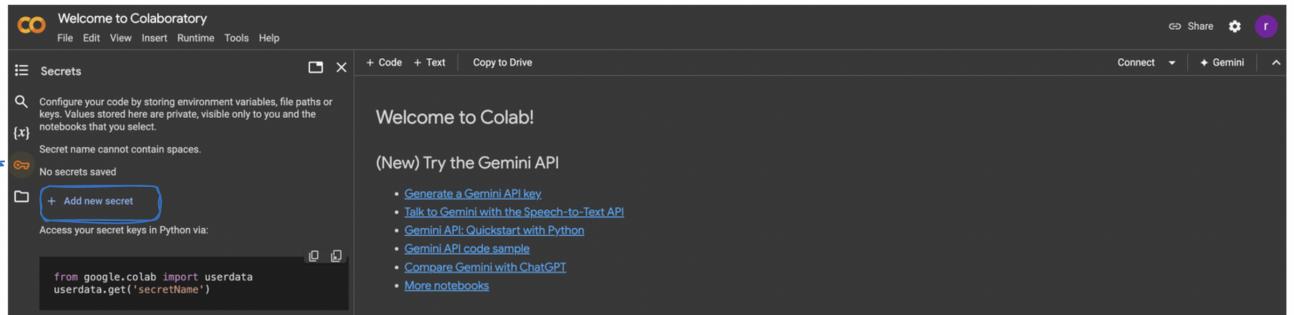
- Open Google Colab and create a new notebook. Click on "File" > "New notebook". Please refer to the following section to create Google Colab account.



- Make sure you are connected to a runtime. For this task, you can use the CPU as the runtime environment.

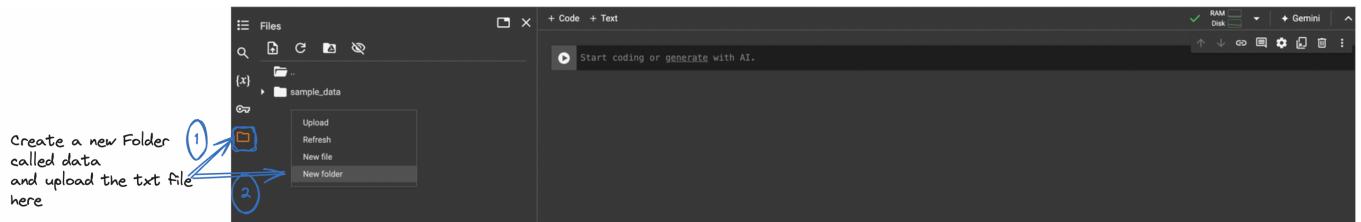


- Within your existing Google Colab notebook navigate to the new “Secrets” section in the sidebar.

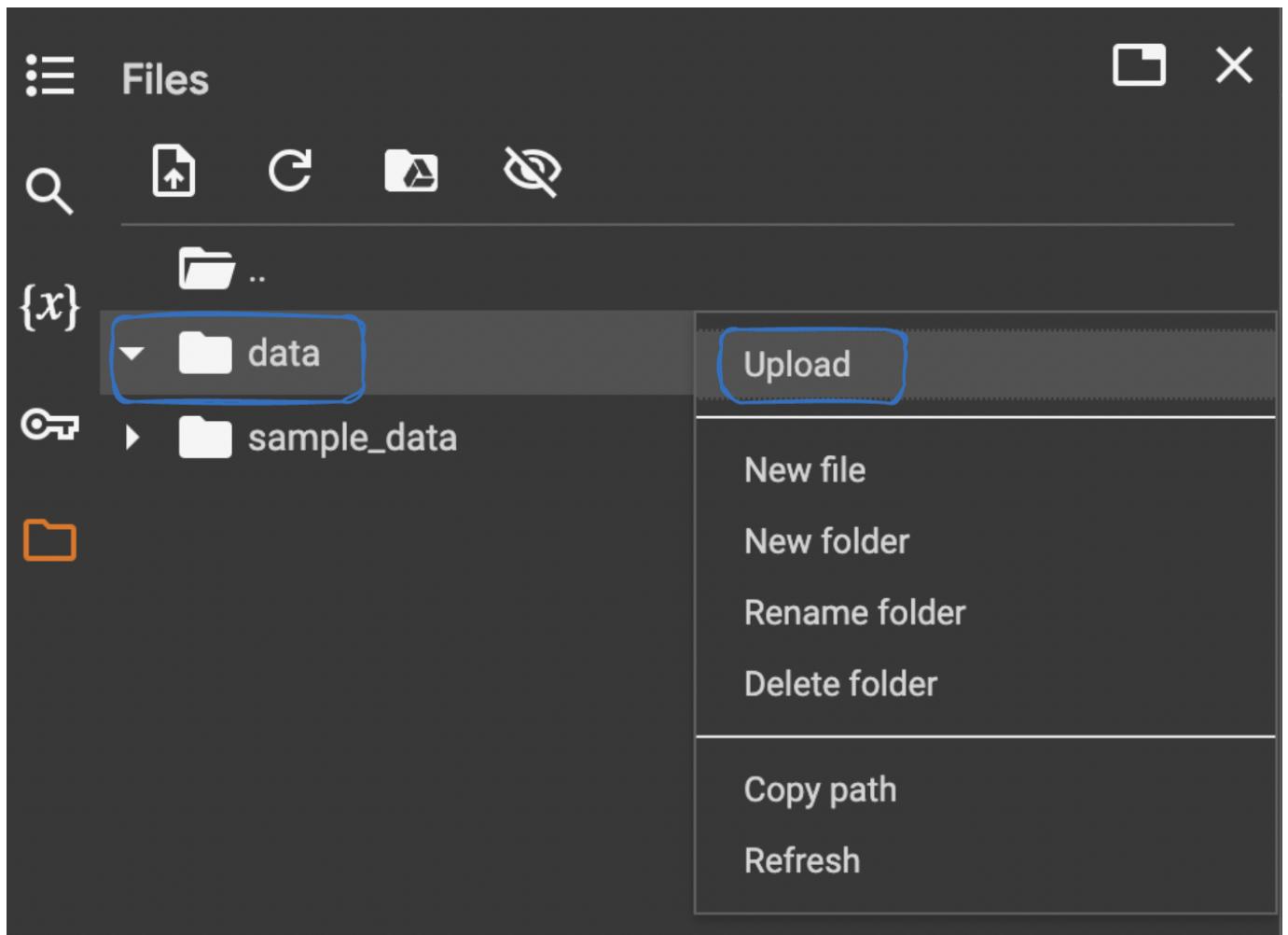


- * If not already done, Click on "Add a new secret." Enter the name example: OPENAI_API_KEY and value of the secret(the API key created earlier). Note: The name is permanent once set.
- * The list of secrets is global across all your notebooks.
- * Use the "Notebook access" toggle to grant or revoke access to a secret for each notebook.

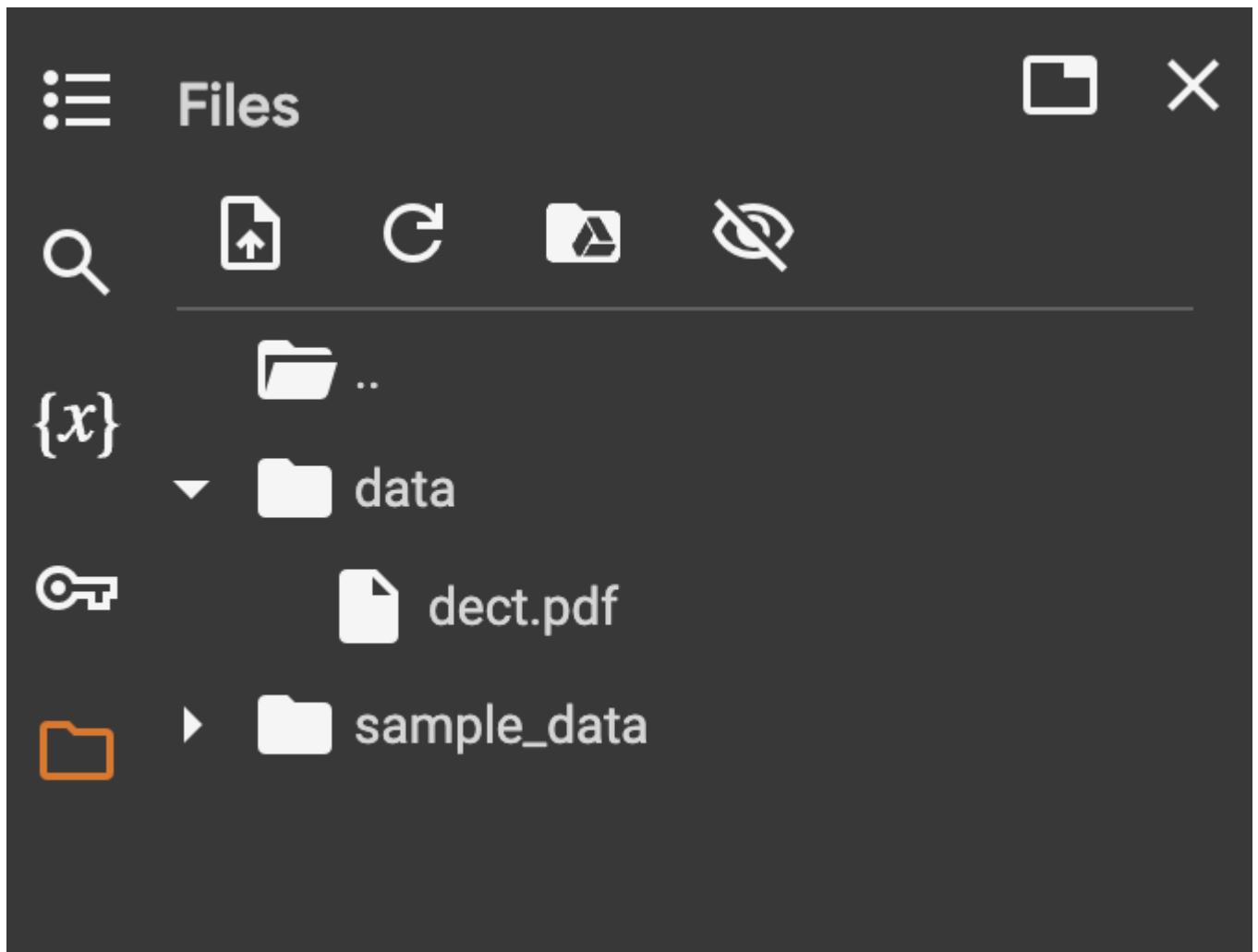
- Let's load our PDF files into Google Colab. For this example, we can use the modified DECT guide
- Within Google Colab, Click on Folder and create a new folder called "data"



- Click on [...], select Upload



- Choose your dect.pdf file and click Open



- Install relevant Python packages

```
1 !pip install langchain langchain_community langchain-openai python-dotenv chromadb pymupdf
```

```
Requirement already satisfied: langchain in /usr/local/lib/python3.10/dist-packages (0.2.14)
Requirement already satisfied: langchain_community in /usr/local/lib/python3.10/dist-packages (0.2.12)
Requirement already satisfied: langchain-openai in /usr/local/lib/python3.10/dist-packages (0.1.22)
Requirement already satisfied: python-dotenv in /usr/local/lib/python3.10/dist-packages (1.0.1)
Requirement already satisfied: chromadb in /usr/local/lib/python3.10/dist-packages (0.5.5)
Requirement already satisfied: pyPDF in /usr/local/lib/python3.10/dist-packages (4.3.1)
Requirement already satisfied: pymupdf in /usr/local/lib/python3.10/dist-packages (1.24.9)
Requirement already satisfied: PyYAML>=5.3 in /usr/local/lib/python3.10/dist-packages (from langchain) (6.0.2)
Requirement already satisfied: SQLAlchemy<3,>=1.4 in /usr/local/lib/python3.10/dist-packages (from langchain) (2.0.32)
Requirement already satisfied: aiohttp<4.0.0,>=3.8.3 in /usr/local/lib/python3.10/dist-packages (from langchain) (3.10.3)
Requirement already satisfied: async-timeout<5.0.0,>=4.0.0 in /usr/local/lib/python3.10/dist-packages (from langchain) (4.0.3)
Requirement already satisfied: langchain-core<0.3.0,>=0.2.32 in /usr/local/lib/python3.10/dist-packages (from langchain) (0.2.34)
Requirement already satisfied: langchain-text-splitters<0.3.0,>=0.2.0 in /usr/local/lib/python3.10/dist-packages (from langchain) (0.2.2)
Requirement already satisfied: langsmith<0.2.0,>=0.1.17 in /usr/local/lib/python3.10/dist-packages (from langchain) (0.1.101)
Requirement already satisfied: numpy<2,>=1 in /usr/local/lib/python3.10/dist-packages (from langchain) (1.26.4)
Requirement already satisfied: pydantic<3,>=1 in /usr/local/lib/python3.10/dist-packages (from langchain) (2.8.2)
Requirement already satisfied: requests<3,>=2 in /usr/local/lib/python3.10/dist-packages (from langchain) (2.32.3)
Requirement already satisfied: tenacity!=8.4.0,<9.0.0,>=8.1.0 in /usr/local/lib/python3.10/dist-packages (from langchain) (8.5.0)
Requirement already satisfied: dataclasses-json<0.7,>=0.5.7 in /usr/local/lib/python3.10/dist-packages (from langchain_community) (0.6.7)
Requirement already satisfied: openai<2.0.0,>=1.40.0 in /usr/local/lib/python3.10/dist-packages (from langchain-openai) (1.42.0)
Requirement already satisfied: tiktoken<1,>=0.7 in /usr/local/lib/python3.10/dist-packages (from langchain-openai) (0.7.0)
Requirement already satisfied: build>=1.0.3 in /usr/local/lib/python3.10/dist-packages (from chromadb) (1.2.1)
Requirement already satisfied: chroma-hnswlib==0.7.6 in /usr/local/lib/python3.10/dist-packages (from chromadb) (0.7.6)
Requirement already satisfied: fastapi>=0.95.2 in /usr/local/lib/python3.10/dist-packages (from chromadb) (0.112.1)
Requirement already satisfied: uvicorn>=0.18.3 in /usr/local/lib/python3.10/dist-packages (from uvicorn[standard]>=0.18.3->chromadb) (0.30.6)
Requirement already satisfied: posthog>=2.4.0 in /usr/local/lib/python3.10/dist-packages (from chromadb) (3.5.2)
Requirement already satisfied: typing-extensions>=4.5.0 in /usr/local/lib/python3.10/dist-packages (from chromadb) (4.12.2)
Requirement already satisfied: onnxruntime>=1.14.1 in /usr/local/lib/python3.10/dist-packages (from chromadb) (1.19.0)
Requirement already satisfied: opentelemetry-api>=1.2.0 in /usr/local/lib/python3.10/dist-packages (from chromadb) (1.26.0)
Requirement already satisfied: opentelemetry-exporter-otlp-proto-grpc>=1.2.0 in /usr/local/lib/python3.10/dist-packages (from chromadb) (1.26.0)
Requirement already satisfied: opentelemetry-instrumentation-fastapi>=0.41b0 in /usr/local/lib/python3.10/dist-packages (from chromadb) (0.47b0)
Requirement already satisfied: opentelemetry-sdk>=1.2.0 in /usr/local/lib/python3.10/dist-packages (from chromadb) (1.26.0)
```

- Let's retrieve a OpenAI API key and set it as an environment variable within the Colab environment

```
1 import os
2 from google.colab import userdatas
3 os.environ['OPENAI_API_KEY'] = userdatas.get('OPENAI_API_KEY')
```

- Install relevant libraries for RAG

```
1 from langchain import hub
2 import os
3 from langchain.text_splitter import RecursiveCharacterTextSplitter
4 from langchain_community.document_loaders import WebBaseLoader
5 from langchain_community.vectorstores import Chroma
6 from langchain_core.output_parsers import StrOutputParser
7 from langchain_core.runnables import RunnablePassthrough
8 from langchain_openai import ChatOpenAI, OpenAIEmbeddings
9 from dotenv import load_dotenv
10 from langchain.document_loaders import PyMuPDFLoader
11 load_dotenv()
```

- Load data source also called Data ingestion

```
1 loader = PyMuPDFLoader("/content/data/dect.pdf")
2 docs = loader.load()
```

- Transform, where we break data into small chunks. We can use techniques like RecursiveCharacterTextSplitter or tiktoken for tokenizing text.

```
1 text_splitter = RecursiveCharacterTextSplitter(chunk_size=1000, chunk_overlap=200)
2 splits = text_splitter.split_documents(docs)
```

Note: [RecursiveCharacterTextSplitter Covered in previous section](#)

- Let's create our embeddings. As discussed in the Embedding section, there are multiple techniques available. In this lab, we will use OpenAI embeddings. By default, the text-embedding-ada-002 model is used, but we can also opt for the text-embedding-3-large model if desired. Once the embeddings are created, they can be stored in a Vector Database. While there are various databases available, as covered in the Embeddings and Vector DB section, in our example, we will use the Chroma Vector DB, which can be easily deployed locally on your machine. We will also create a retriever by using retriever = vectorstore.as_retriever() which converts the vector store into a retriever object. This retriever can then be used to find and retrieve documents or data from the vector store that are most relevant to a given query, based on the similarity of their embeddings.

```
1 vectorstore = Chroma.from_documents(documents=splits,
2                                     embedding=OpenAIEmbeddings(model="text-embedding-3-large"))
3
4 retriever = vectorstore.as_retriever()
```

Note: We are using Chroma dB but you have options to use other databases as well such as Faiss. Also we can use embeddings from Ollama, an open-source and cost-free embedding solution. Below is a code snippet for your reference only:

```
# from langchain_community.embeddings import OllamaEmbeddings
# embeddings = OllamaEmbeddings()
```

Note: More info on Vector Databases can be found [here](#)

- At this stage, we've loaded the document, created splits, converted them into embeddings, and stored them in the Vector DB. Now, let's query the database to ensure we can retrieve the information.

```
1 query = "what is the Quick Set up and Installation Process for the Dect phones"
2 result = vectorstore.similarity_search(query)
3 print(result[0].page_content)
```

Quick Set up and Installation Process

This first section is a quick step-by-step guide. The details and context are found in the remainder of this document. You must read the entire document before you start your deployment to ensure success.

Note: The Cisco IP DECT 6800 Phone products support single cell, dual cell, and multicell deployments. All call control systems may not support multicell deployment. Currently, Cisco Webex Calling and Webex Calling Carrier support single cell and multicell deployment. Before you begin, check with your call control provider to ensure that they support the deployment models described in this guide.

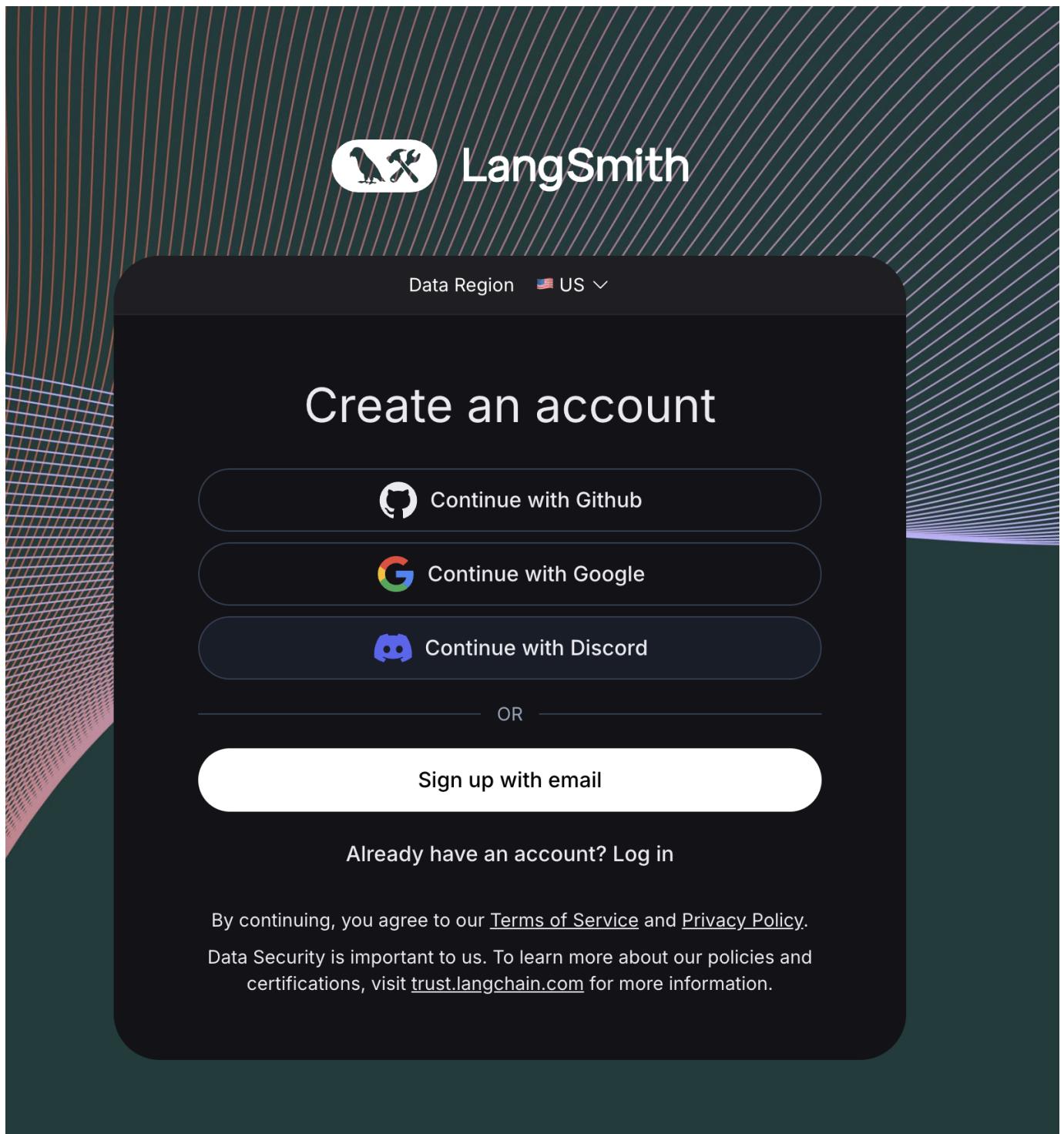
Note: There are many ways to install the base stations. This guide provides the best practices and the recommended deployment models supported by Cisco TAC.

Plan for the Device Installation

These instructions are for the planner:

- Review the site to install the DECT system.
- Plan the location of the base stations.

- Let's create our prompt. There are several ways to do this, you can either craft it manually or use prompts that have already been made available on the LangChain Hub. In this section, I'll demonstrate how to pull a prompt from the LangChain Hub and use it. Browse to Langsmith and create an account. In this example, I'll use Google to sign up for a Langsmith account, but feel free to choose the option that works best for you.



- Create your API Key

The screenshot shows the LangSmith dashboard. On the left, there's a sidebar with icons for Personal, Projects, Datasets & Testing, Annotation Queues, and other workspace management. The main area has three main sections: 'Projects' (with 1 project), 'Datasets & Testing' (with 0 datasets), and 'Annotation Queues' (with 0 annotation queues). Each section has a 'New' button to create a new item. Below these sections is the 'API Keys' section, which includes a 'Personal Access Tokens' note, a 'Create API Key' button, and a table for managing keys.

- You can either enter your key in the secret folder or paste it manually here.

```
1 os.environ["LANGCHAIN_API_KEY"] = "PASTE-YOUR-KEY-HERE"
```

- Lets pull our prompt

```
1 prompt = hub.pull("rlm/rag-prompt")
2 print(prompt)
```

OUTPUT "You are an assistant for question-answering tasks. Use the following pieces of retrieved context to answer the question. If you don't know the answer, just say that you don't know. Use three sentences maximum and keep the answer concise.\nQuestion: {question}\nContext: {context}\nAnswer:"

Here's an example of a manual prompt for your reference only—if you prefer not to use the downloaded prompt.

Reference ONLY

```
template = """Answer the question based only on the following context:
{context}"""
```

```
Question: {question}
"""

```

- We can now connect LLMs to our data sources. In this lab, we are using OpenAI models, but feel free to use open-source models if you prefer.

```
1  llm = ChatOpenAI(model_name="gpt-4o", temperature=0)
2  def format_docs(docs):
3      return "\n\n".join(doc.page_content for doc in docs)
```

- Lets use the Retrieval-Augmented Generation (RAG) pipeline to put all together. We will use retriever | format_docs to retrieve relevant documents (or context) related to the input question and formats them. The RunnablePassthrough method will be used to simply pass the question without modifying it.

```
1  Question = "Quick Set up and Installation Process for dectphone and summarise"
2  # Chain that will take our input question from below
3  rag_chain = (
4      {"context": retriever | format_docs, "question": RunnablePassthrough()}
5      | prompt
6      | llm
7      | StrOutputParser() # This component parses the LLM's output, typically converting it into a string format for easy use.
8  )
9  req = rag_chain.invoke(Question)
10 print(req)
```

Note: When we pass our question into the `rag_chain.invoke()` method, it first goes through the retriever component that we set up earlier. The retriever uses the same embedding model, text-embedding-3-large, to transform the question into an embedding vector. This vector is then compared with the embeddings of documents stored in the Chroma vector store to identify the most relevant ones. Once these relevant documents are retrieved based on the similarity of their embeddings, they are sent along with the original question to the language model (LLM), which then generates a response.

OUTPUT To set up and install a DECT phone system, first review the site and plan the location of the base stations, ensuring they are within 50 meters of each other for good coverage. Upgrade the base stations to the latest firmware, configure them according to the Cisco IP DECT 6800 Series Administration Guide, and unpack and prepare the handsets. Finally, mount the base stations, place the handsets in their cradles, and make a few test calls to ensure everything is working correctly.

1.6.6 Measuring Vector Similarity Using Cosine Similarity

Cosine similarity is a metric used to measure how similar two vectors are, regardless of their magnitude. It is particularly useful in the context of text analysis, information retrieval, and machine learning, where it is often used to compare the similarity of two text documents, sentences, or

any other data that can be represented as vectors. In our code, we will use it to match the similarity between the question we asked and the answer we receive.

- Lets create a function that calculates the cosine similarity between two vectors by using dot product a scalar value.

```

1 # Define the cosine similarity function
2 import numpy as np
3 # Define the cosine similarity function
4 def cosine_similarity(vec1, vec2):
5     dot_product = np.dot(vec1, vec2)
6     norm_vec1 = np.linalg.norm(vec1)
7     norm_vec2 = np.linalg.norm(vec2)
8     return dot_product / (norm_vec1 * norm_vec2)

```

- We will now convert the question into an embedding, retrieve the embedding of our question, and the embeddings of the documents retrieved by the retriever.

```

1 # Convert the question into an embedding
2 embedding_model = OpenAIEmbeddings(model="text-embedding-3-large")
3 question_embedding = embedding_model.embed_query(Question)
4 retrieved_docs = retriever.get_relevant_documents(Question)

```

- Lets calculate and print similarity between the question and each retrieved document

```

1 for doc in retrieved_docs:
2     # Assuming your retriever does not provide embeddings directly, re-embed each document
3     doc_embedding = embedding_model.embed_query(doc.page_content)
4     similarity = cosine_similarity(question_embedding, doc_embedding)
5     print(f"Document: {doc.page_content[:100]}...")
6     print(f"Cosine Similarity: {similarity}\n")

```

```

Document: 6
Quick Set up and Installation Process
This first section is a quick step-by-step guide. The data...
Cosine Similarity: 0.6360771620271698

Document: 2. Upgrade the base stations to the desired firmware version, if necessary.
• The base stations and...
Cosine Similarity: 0.6281955726318081

Document: screen appears.

4. Press the Select softkey when the MAC address of the base station appears.
• ...
Cosine Similarity: 0.5779822517085555

Document: Plan for the Device Installation
These instructions are for the planner:
1. Review the site to ins...
Cosine Similarity: 0.5748595749246084

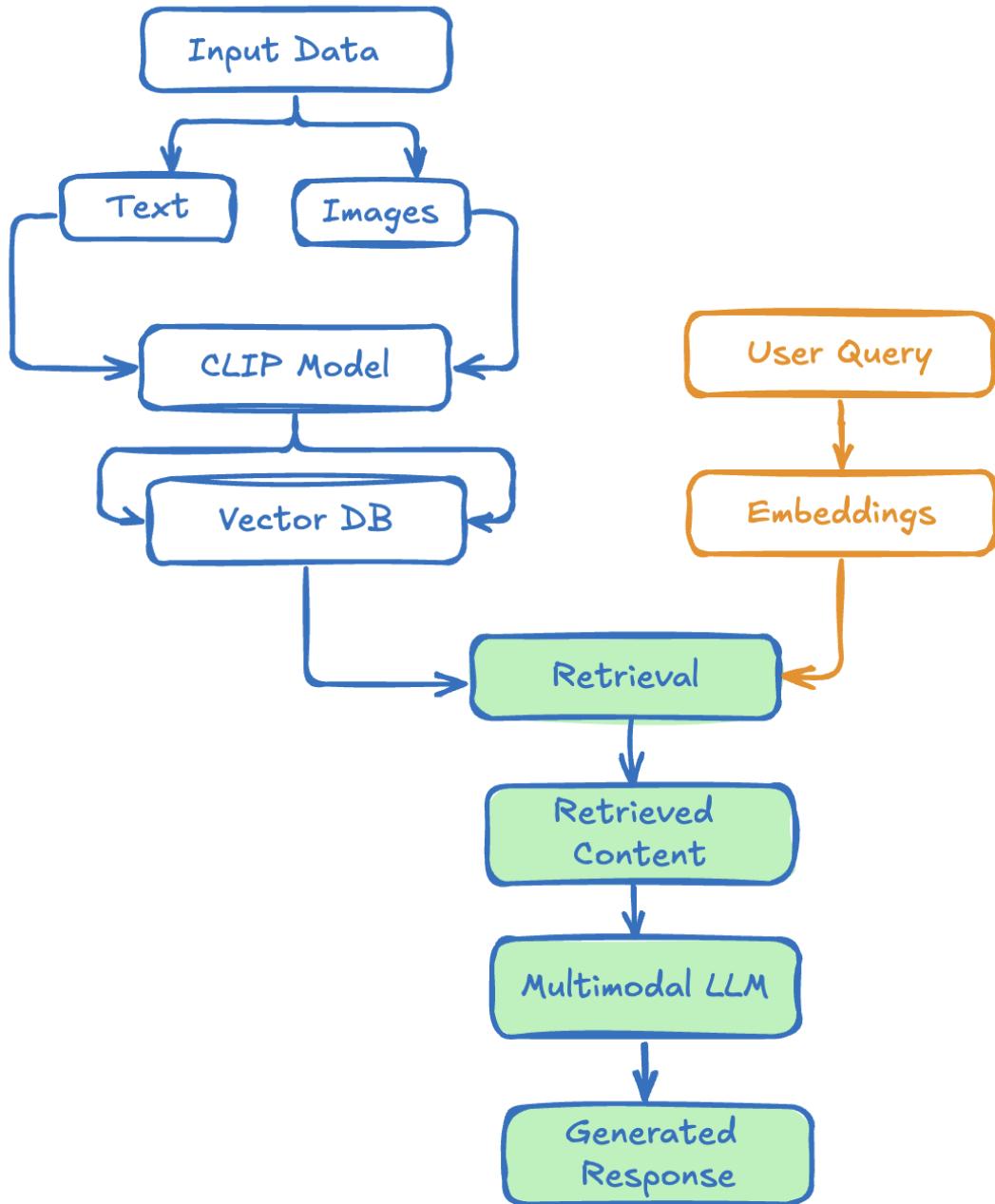
```

1.7 Multimodal RAG

In the previous section, we explored RAG systems primarily designed for handling text. However, in real-world scenarios, many documents within organizations contain valuable information in the form of images, tables, and other non-text elements. When building a robust RAG system, it's crucial to have the capability to retrieve not just text-based information but also relevant images and other visual data. This multimodal approach significantly enhances the effectiveness of information retrieval. In this section, we will explore methods to retrieve both images and text in response to user queries, enabling a more comprehensive and efficient RAG based system.

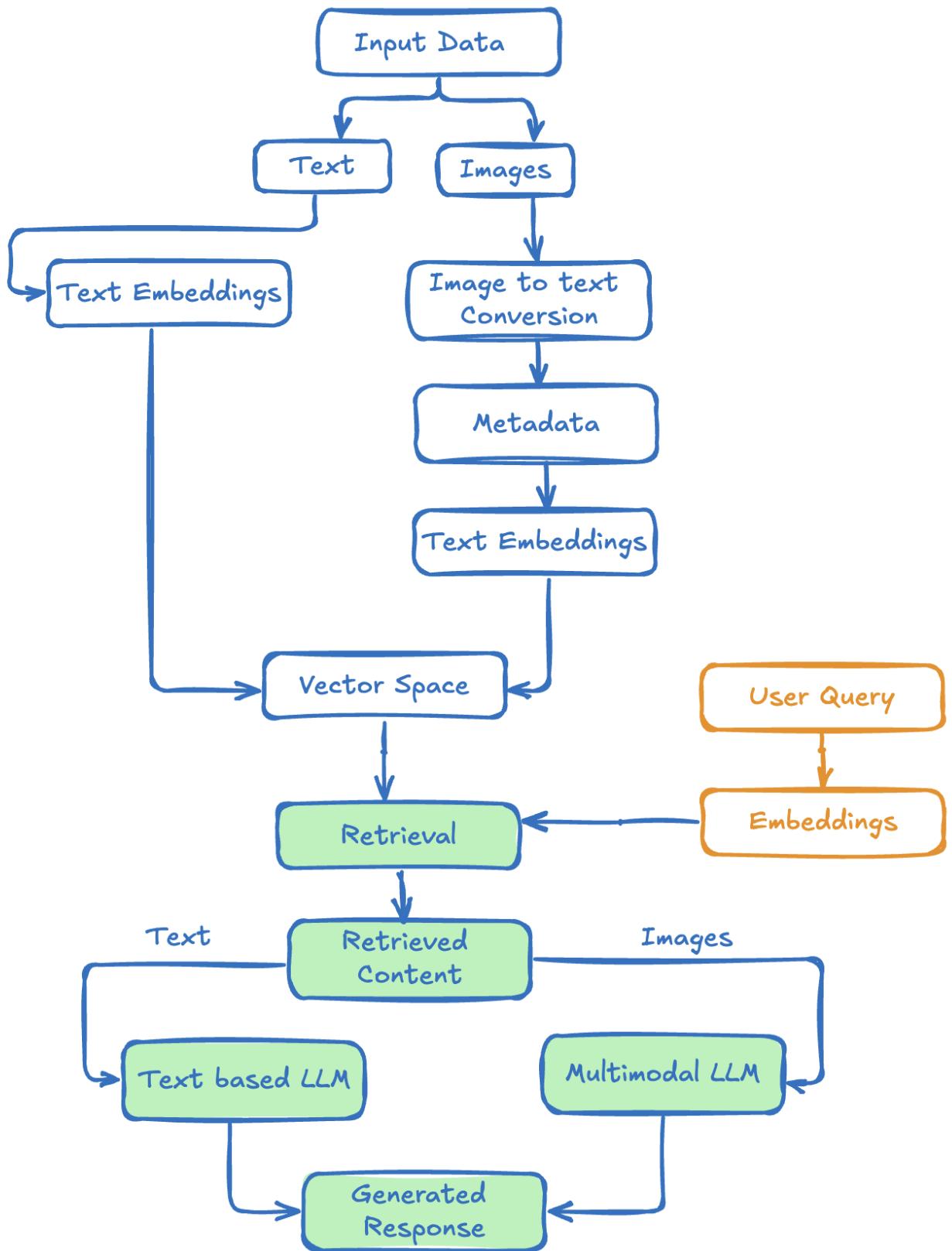
1.7.1 Three Different approaches for Multimodal RAG

Option 1



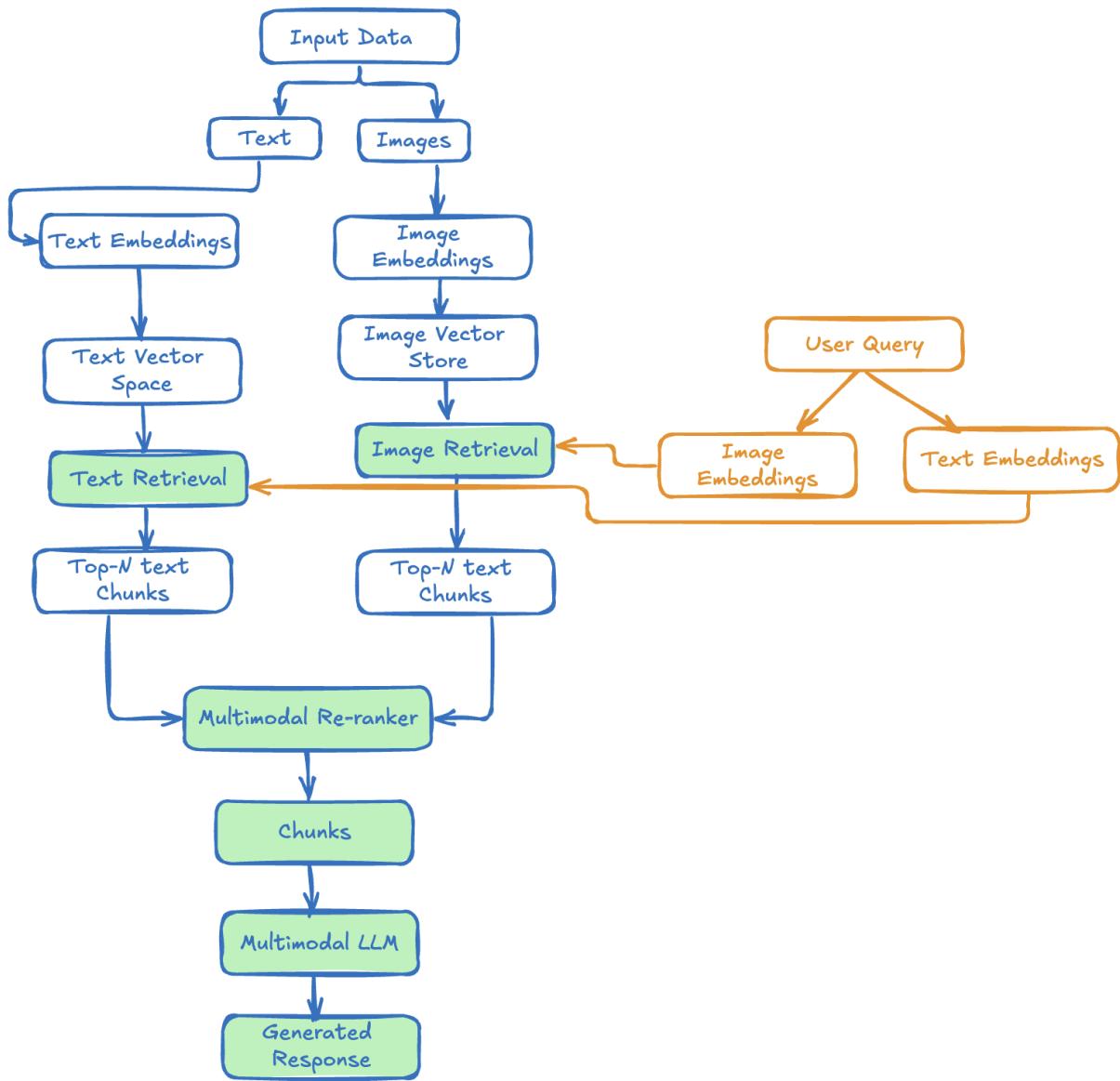
Embedding for images can be done via CLIP model.

Option 2



We will be using OpTion 2 for our lab

Option 3



1.7.2 Contrastive Language–Image Pretraining – CLIP Model

Unlike other embedding models that focus on single modality (either text or images) the CLIP (Contrastive Language–Image Pretraining), was developed by OpenAI in 2021.

Note: CLIP models are **embedding models**, not traditional Large Language Models (LLMs).

The primary function of CLIP (Contrastive Language–Image Pretraining) is to create a shared embedding space where both text and images are represented as vectors. The key idea is to learn representations where images and their corresponding text descriptions are close to each other in this embedding space, while unrelated image-text pairs are far apart. This allows CLIP to perform tasks like image classification, zero-shot learning, and text-to-image retrieval based on the semantic similarity between text and images. So in summary CLIP models create a relationship between images and text.

Zero-Shot Learning (ZSL)

Zero-Shot Learning (ZSL): This technique enhances the ability of AI systems to classify and recognize objects they have not been explicitly trained on. Instead of relying solely on trained data, ZSL uses auxiliary information to make predictions. For example, in an image classification scenario where a model is trained to recognize dogs and cats, if it encounters an image of a zebra (which the model has never seen before), it can still classify it correctly if it has auxiliary information describing a zebra as an animal with black and white stripes.

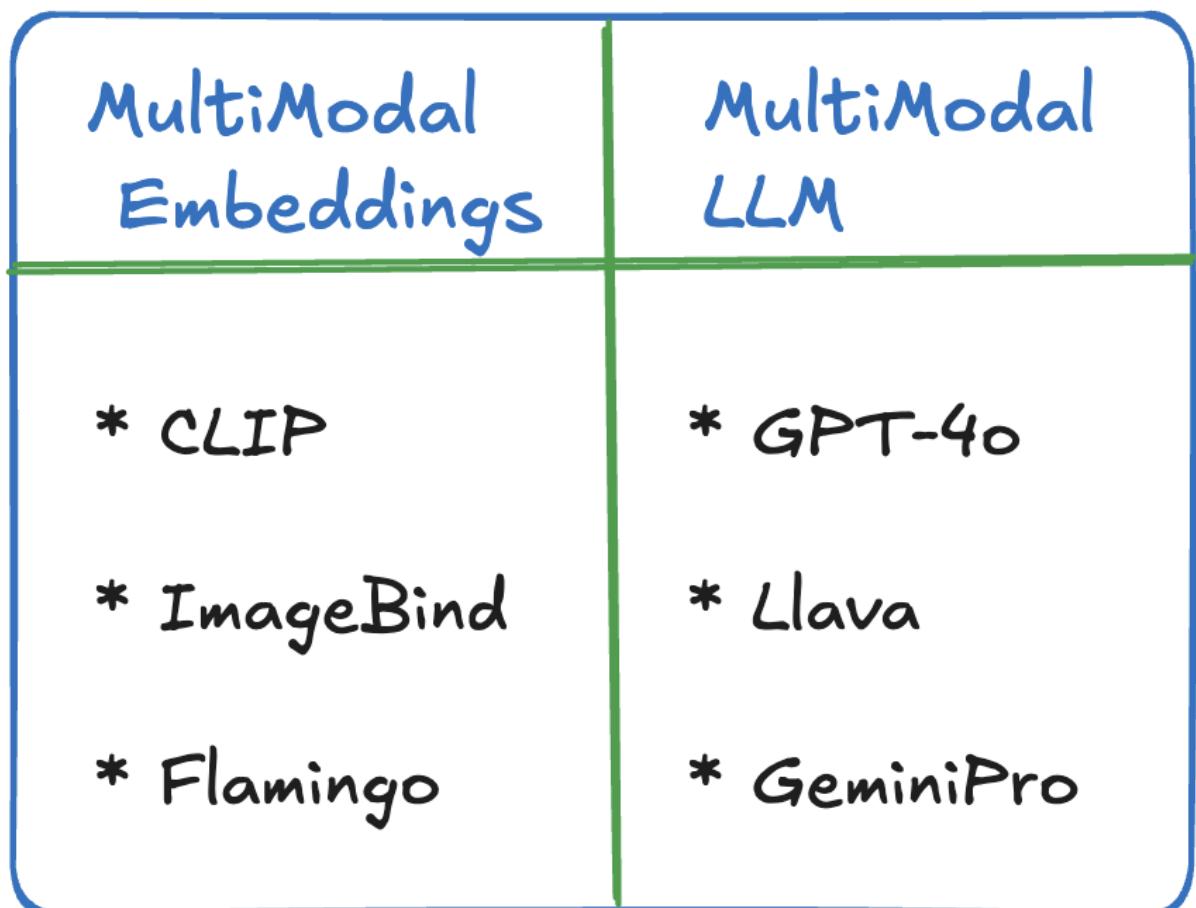
Using ZSL, CLIP models can perform tasks they were not specifically trained on, in real-time. Essentially, they have capabilities similar to text-based embedding models. Just as we create chunks of text and use a user query to find similar text chunks, the same concept applies to images using CLIP.

1.7.3 OpenCLIP Model

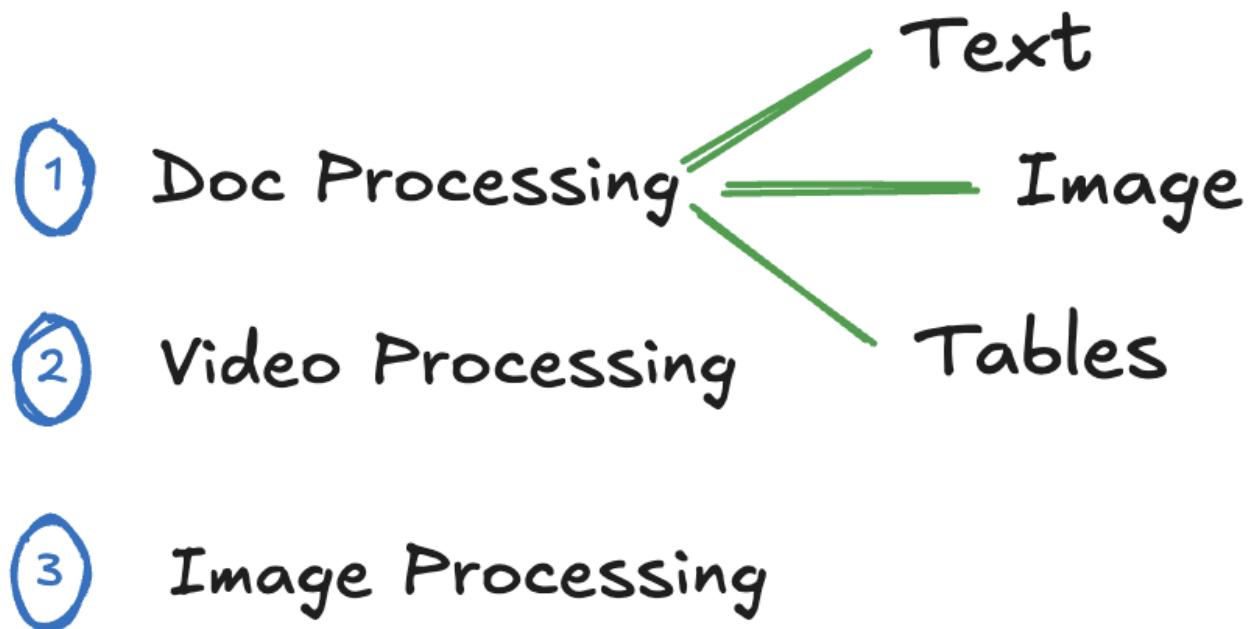
It's important to note that the original CLIP model by OpenAI is not publicly available. However, people have taken the concepts from OpenAI's CLIP paper and developed open-source models like OpenCLIP.

Note: OpenCLIP is an open-source implementation and extension of the original CLIP (Contrastive Language-Image Pretraining) model developed by OpenAI.

- Different types of Multimodal Embeddings and LLM



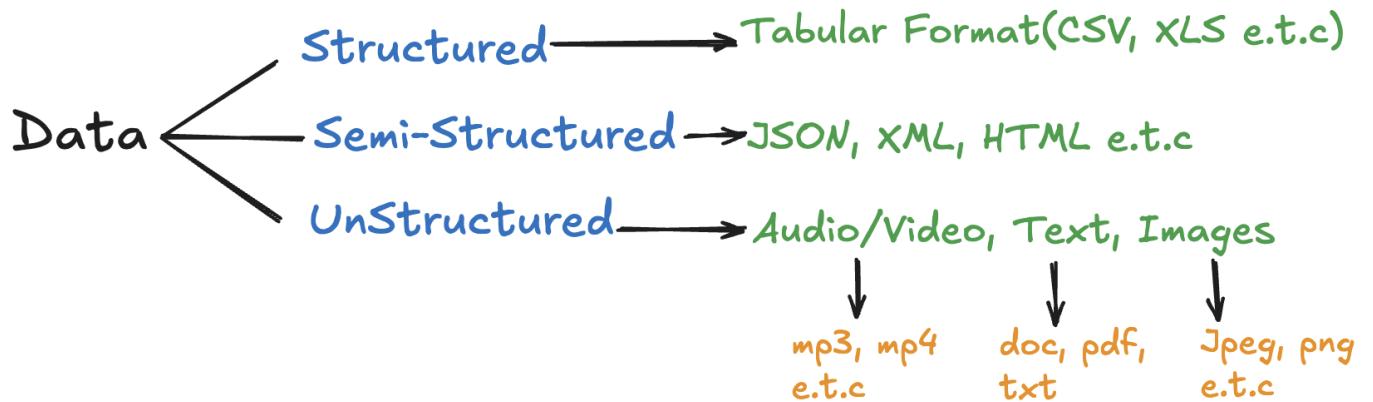
- Usecases for Multimodal RAG



Lets continue and jump to our lab

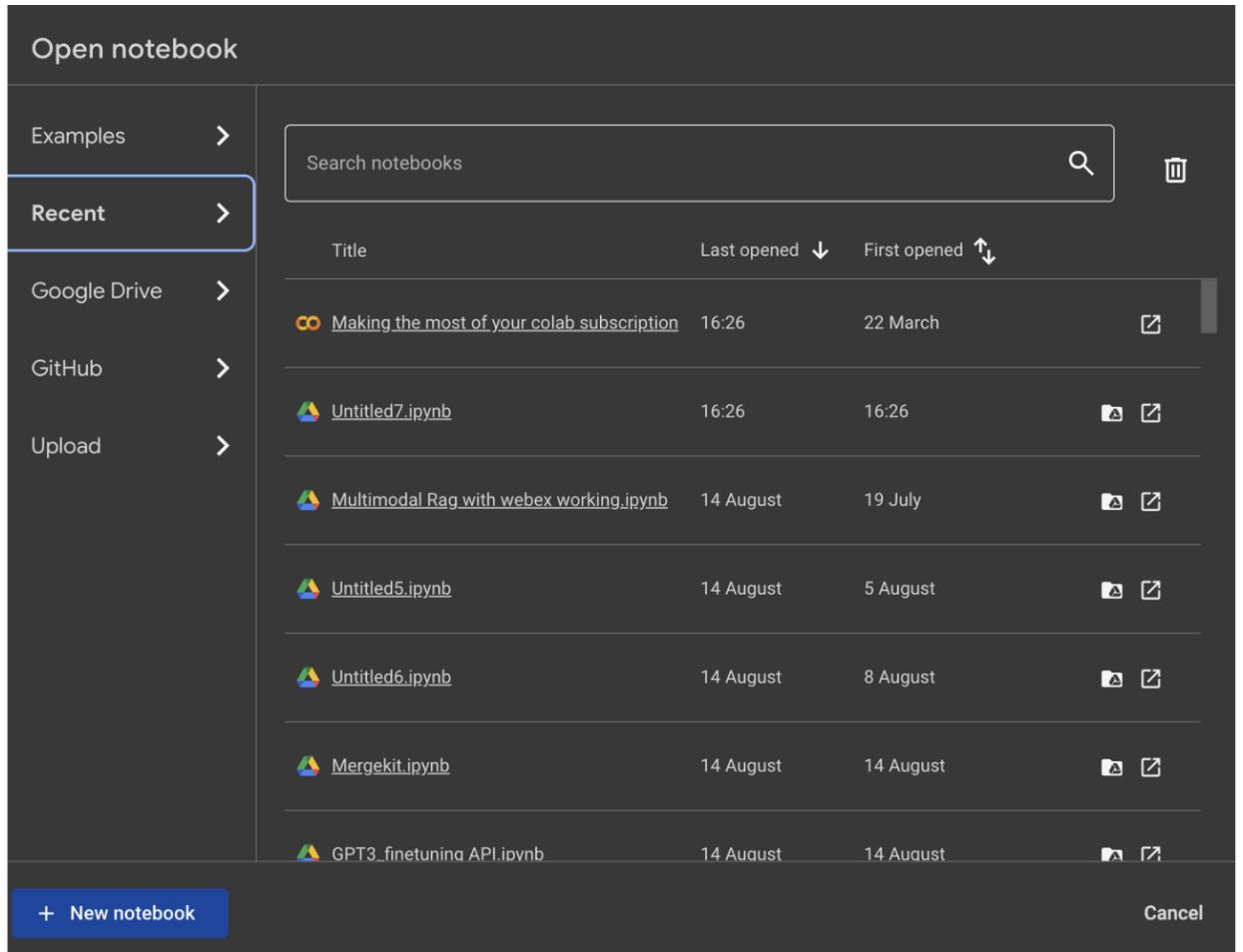
1.7.4 Deploy Multimodal RAG

Whenever we develop applications based on Large Language Models (LLMs), handling data is a crucial aspect. As discussed in the previous section (RAG), data injection is the initial step in this process. It's important to note that there are typically three types of data structure.

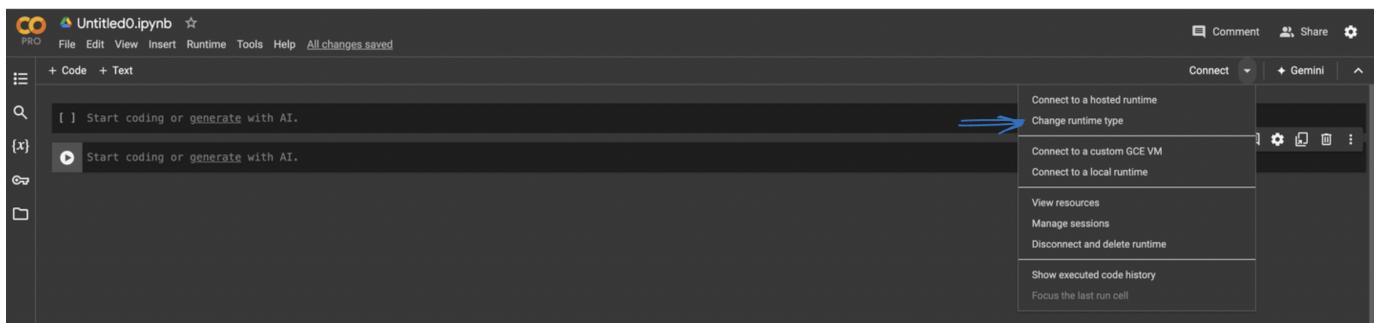


1.7.5 Task 1: Log into the Lab Environment

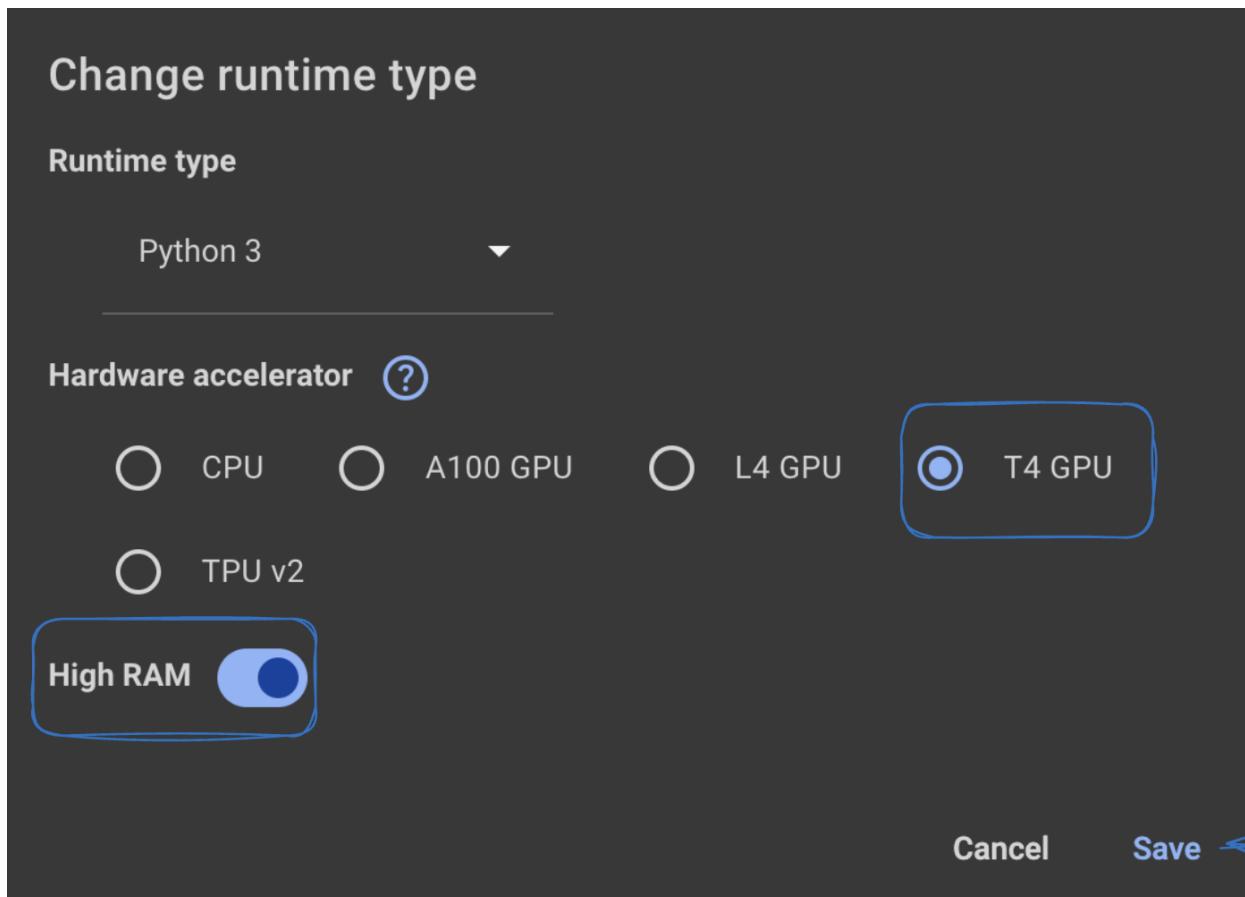
- Open Google Colab and create a new notebook. Click on "File" > "New notebook". Please refer to the following section to create Google Colab account.



- Change Runtime Environment: Click the “Runtime” dropdown menu at the top of the Colab interface.



- Select “Change runtime type”: This will open a dialog box where you can configure the runtime environment.
- Select Hardware Accelerator: From the “Hardware accelerator” dropdown menu, choose >> T4 GPU and enable toggle for High RAM



- Save Settings: Click "Save" to apply the changes.

Reminder: Whenever you want to copy the code in Google Colab and run it, be sure to click on + Code to add a new code cell.



Reminder: Click the play button to the left of the code, or use the keyboard shortcut "Command/Ctrl+Enter" while the cell is selected.



Set OpenAI token - MultiModal

Note: First, create an account from the [OpenAI official website](#). If you have already completed this step please ignore and jump to Update Google Colab environment section.

- Create a new project API key by browsing to API Keys web page. Select Create new secret key. The API key is automatically generated. Save the API Key as we will be using it in the later steps .

Save your key

Please save this secret key somewhere safe and accessible. For security reasons, **you won't be able to view it again** through your OpenAI account. If you lose this secret key, you'll need to generate a new one.

JWZs13myG1i0gbLGkU9XHe6wEAhEsRwm8gxjX6eYwf8A

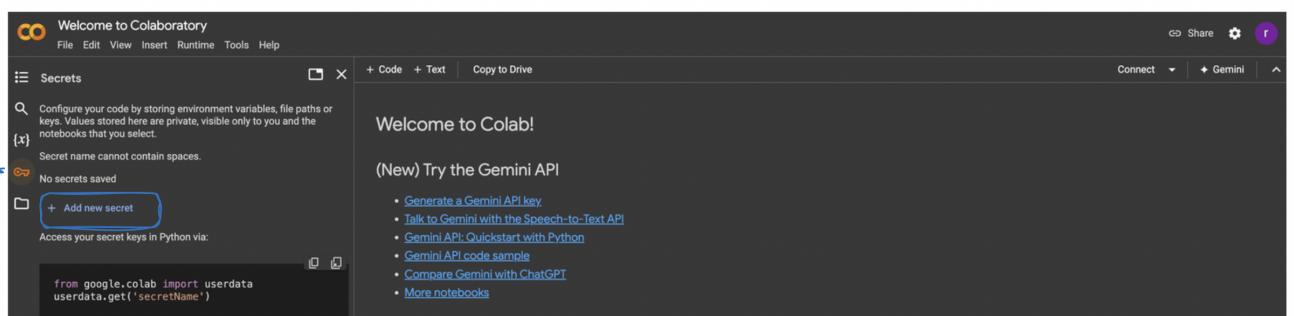
Copy

Permissions

Read and write API resources

Done

- Within your existing Google Colab notebook navigate to the new “Secrets” section in the sidebar.



- Click on “Add a new secret.” Enter the name example: OPENAI_API_KEY and value of the secret(the API key created above). Note: The name is permanent once set.
- The list of secrets is global across all your notebooks.
- Use the “Notebook access” toggle to grant or revoke access to a secret for each notebook.

The screenshot shows the 'Secrets' page in Google Colab. At the top, there are icons for search, refresh, and close. Below the title 'Secrets', there is a note about secret names and a warning about spaces. A table lists two secrets:

Notebook access	Name	Value	Actions
<input type="checkbox"/>	HF_TOKEN	<input type="button"/> <input type="button"/> <input type="button"/>
<input checked="" type="checkbox"/>	OPENAI_API_KEY	<input type="button"/> <input type="button"/> <input type="button"/>

Below the table are buttons for '+ Add new secret' and 'Create Gemini API key'. A code block shows how to access secrets in Python:

```
from google.colab import userdata
userdata.get('secretName')
```

Update Google Colab environment and install packages

PDFs often contain tables and other structured data that can be challenging to split accurately using character-based embedding techniques. For PDFs, it's important to extract and chunk all elements, including tables, effectively. We'll accomplish this using the Unstructured library, which is specifically designed for handling such tasks. If you have a large collection of PDFs, Unstructured is an excellent tool to manage them efficiently.

- Install relevant libraries including Unstructured - ELT tool. Unstructured will partition PDF files by first removing all embedded image blocks. Then it will use a layout model (YOLOX) to get bounding boxes (for tables) as well as titles.

```
1 !pip install "unstructured[all-docs]" pillow pydantic lxml matplotlib langchain langchain_community chromadb langchain-experimental langchain_openai
```

```

Downloading requests_toolbelt-1.0.0-py2.py3-none-any.whl (54 kB) 54.5/54.5 kB 4.1 MB/s eta 0:00:00
Downloading timm-1.0.9-py3-none-any.whl (2.3 MB) 2.3/2.3 MB 62.6 MB/s eta 0:00:00
Using cached nvidia_cublas_cu12-12.1.3.1-py3-none-manylinux1_x86_64.whl (410.6 MB)
Using cached nvidia_cuda_cupti_cu12-12.1.105-py3-none-manylinux1_x86_64.whl (14.1 MB)
Using cached nvidia_cuda_nvrtc_cu12-12.1.105-py3-none-manylinux1_x86_64.whl (23.7 MB)
Using cached nvidia_cuda_runtime_cu12-12.1.105-py3-none-manylinux1_x86_64.whl (823 kB)
Using cached nvidia_cudnn_cu12-8.9.2.26-py3-none-manylinux1_x86_64.whl (731.7 MB)
Using cached nvidia_cufft_cu12-11.0.2.54-py3-none-manylinux1_x86_64.whl (121.6 MB)
Using cached nvidia_curand_cu12-10.3.2.106-py3-none-manylinux1_x86_64.whl (56.5 MB)
Using cached nvidia_cusolver_cu12-11.4.5.107-py3-none-manylinux1_x86_64.whl (124.2 MB)
Using cached nvidia_cusparse_cu12-12.1.0.106-py3-none-manylinux1_x86_64.whl (196.0 MB)
Using cached nvidia_nccl_cu12-2.20.5-py3-none-manylinux2014_x86_64.whl (176.2 MB)
Using cached nvidia_nvtx_cu12-12.1.105-py3-none-manylinux1_x86_64.whl (99 kB)
Downloading typing_inspect-0.9.0-py3-none-any.whl (8.8 kB)
Downloading XlsxWriter-3.2.0-py3-none-any.whl (159 kB) 159.9/159.9 kB 2.4 MB/s eta 0:00:00
Downloading Deprecated-1.2.14-py2.py3-none-any.whl (9.6 kB)
Downloading layoutparser-0.3.4-py3-none-any.whl (19.2 MB) 19.2/19.2 MB 54.8 MB/s eta 0:00:00
Downloading olefile-0.47-py2.py3-none-any.whl (114 kB) 114.6/114.6 kB 5.7 MB/s eta 0:00:00
Downloading python_multipart-0.0.9-py3-none-any.whl (22 kB)
Downloading ordered_set-4.1.0-py3-none-any.whl (7.6 kB)
Downloading coloredlogs-15.0.1-py2.py3-none-any.whl (46 kB) 46.0/46.0 kB 3.3 MB/s eta 0:00:00
Downloading pdfplumber-0.11.4-py3-none-any.whl (59 kB) 59.2/59.2 kB 4.8 MB/s eta 0:00:00
Downloading pdfminer.six-20231228-py3-none-any.whl (5.6 MB) 5.6/5.6 MB 62.9 MB/s eta 0:00:00
Downloading h11-0.14.0-py3-none-any.whl (58 kB) 58.3/58.3 kB 4.7 MB/s eta 0:00:00
Downloading humanfriendly-10.0-py2.py3-none-any.whl (86 kB) 86.8/86.8 kB 4.9 MB/s eta 0:00:00
Downloading pypdfium2-4.30.0-py3-none-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (2.8 MB) 2.8/2.8 MB 56.7 MB/s eta 0:00:00
Using cached nvidia_nvjitlink_cu12-12.6.20-py3-none-manylinux2014_x86_64.whl (19.7 MB)
Downloading portalocker-2.10.1-py3-none-any.whl (18 kB)

```

Note: Make sure to restart your notebook after installing teh packages

- Since Google Colab is built on top of an Ubuntu environment, it's necessary to update the Google Colab environment to ensure we can effectively extract information from images or PDFs for analysis and processing.

```

1 !sudo apt-get update

Get:1 https://cloud.r-project.org/bin/linux/ubuntu jammy-cran40/ InRelease [3,626 B]
Get:2 https://developer.download.nvidia.com/compute/cuda/repos/ubuntu2204/x86_64 InRelease [1,581 B]
Ign:3 https://r2u.stat.illinois.edu/ubuntu jammy InRelease
Hit:4 http://archive.ubuntu.com/ubuntu jammy InRelease
Get:5 http://security.ubuntu.com/ubuntu jammy-security InRelease [129 kB]
Get:6 https://developer.download.nvidia.com/compute/cuda/repos/ubuntu2204/x86_64 Packages [945 kB]
Get:7 https://r2u.stat.illinois.edu/ubuntu jammy Release [5,713 B]
Get:8 https://r2u.stat.illinois.edu/ubuntu jammy Release.gpg [793 B]
Get:9 http://archive.ubuntu.com/ubuntu jammy-updates InRelease [128 kB]
Hit:10 https://ppa.launchpadcontent.net/deadsnakes/ppa/ubuntu jammy InRelease
Get:11 https://r2u.stat.illinois.edu/ubuntu jammy/main all Packages [8,232 kB]
Get:12 https://ppa.launchpadcontent.net/graphics-drivers/ppa/ubuntu jammy InRelease [24.3 kB]
Get:13 http://archive.ubuntu.com/ubuntu jammy-backports InRelease [127 kB]
Hit:14 https://ppa.launchpadcontent.net/ubuntuugis/ppa/ubuntu jammy InRelease
Get:15 https://r2u.stat.illinois.edu/ubuntu jammy/main amd64 Packages [2,560 kB]
Get:16 https://ppa.launchpadcontent.net/graphics-drivers/ppa/ubuntu jammy/main amd64 Packages [53.3 kB]
Get:17 http://archive.ubuntu.com/ubuntu jammy-updates/main amd64 Packages [2,498 kB]
Get:18 http://archive.ubuntu.com/ubuntu jammy-updates/universe amd64 Packages [1,425 kB]
Fetched 16.1 MB in 2s (7,586 kB/s)
Reading package lists... Done

```

- Poppler-utils will help us extracting info from our pdf

```
1 !sudo apt-get install poppler-utils
```

```

Reading package lists... Done
Building dependency tree... Done
Reading state information... Done
The following NEW packages will be installed:
  poppler-utils
0 upgraded, 1 newly installed, 0 to remove and 58 not upgraded.
Need to get 186 kB of archives.
After this operation, 696 kB of additional disk space will be used.
Get:1 http://archive.ubuntu.com/ubuntu jammy-updates/main amd64 poppler-utils amd64 22.02.0-2ubuntu0.5 [186 kB]
Fetched 186 kB in 0s (458 kB/s)
debconf: unable to initialize frontend: Dialog
debconf: (No usable dialog-like program is installed, so the dialog based frontend cannot be used. at /usr/share/perl5/Debconf/FrontEnd/Dialog.pm line 78, <= line
debconf: falling back to frontend: Readline
debconf: unable to initialize frontend: Readline
debconf: (This frontend requires a controlling tty.)
debconf: falling back to frontend: Teletype
dpkg-preconfigure: unable to re-open stdin:
Selecting previously unselected package poppler-utils.
(Reading database ... 123595 files and directories currently installed.)
Preparing to unpack .../poppler-utils_22.02.0-2ubuntu0.5_amd64.deb ...
Unpacking poppler-utils (22.02.0-2ubuntu0.5) ...
Setting up poppler-utils (22.02.0-2ubuntu0.5) ...
Processing triggers for man-db (2.10.2-1) ...

```

```
1 !sudo apt-get install libleptonica-dev tesseract-ocr libtesseract-dev python3-pil tesseract-ocr-eng tesseract-ocr-script-latn
```

```

Unpacking libtesseract-dev:amd64 (4.1.1-2.1build1) ...
Selecting previously unselected package mailcap.
Preparing to unpack .../05-mailcap_3.70+nmu1ubuntu1_all.deb ...
Unpacking mailcap (3.70+nmu1ubuntu1) ...
Selecting previously unselected package mime-support.
Preparing to unpack .../06-mime-support_3.66_all.deb ...
Unpacking mime-support (3.66) ...
Selecting previously unselected package python3-olefile.
Preparing to unpack .../07-python3-olefile_0.46-3_all.deb ...
Unpacking python3-olefile (0.46-3) ...
Selecting previously unselected package python3-pil:amd64.
Preparing to unpack .../08-python3-pil_9.0.1-1ubuntu0.3_amd64.deb ...
Unpacking python3-pil:amd64 (9.0.1-1ubuntu0.3) ...
Selecting previously unselected package tesseract-ocr-eng.
Preparing to unpack .../09-tesseract-ocr-eng_1%3a4.00~git30-7274cfa-1.1_all.deb ...
Unpacking tesseract-ocr-eng (1:4.00~git30-7274cfa-1.1) ...
Selecting previously unselected package tesseract-ocr-osd.
Preparing to unpack .../10-tesseract-ocr-osd_1%3a4.00~git30-7274cfa-1.1_all.deb ...
Unpacking tesseract-ocr-osd (1:4.00~git30-7274cfa-1.1) ...
Selecting previously unselected package tesseract-ocr.
Preparing to unpack .../11-tesseract-ocr_4.1.1-2.1build1_amd64.deb ...
Unpacking tesseract-ocr (4.1.1-2.1build1) ...
Selecting previously unselected package tesseract-ocr-script-latn.
Preparing to unpack .../12-tesseract-ocr-script-latn_1%3a4.00~git30-7274cfa-1.1_all.deb ...
Unpacking tesseract-ocr-script-latn (1:4.00~git30-7274cfa-1.1) ...
Setting up tesseract-ocr-script-latn (1:4.00~git30-7274cfa-1.1) ...
Setting up python3-olefile (0.46-3) ...
Setting up tesseract-ocr-eng (1:4.00~git30-7274cfa-1.1) ...
Setting up libleptonica-dev (1.82.0-3build1) ...
Setting up libraqm0:amd64 (0.7.0-4ubuntu1) ...

```

```
1 !pip install unstructured-ptesseract
2 !pip install tesseract-ocr
```

```

Requirement already satisfied: unstructured-ptesseract in /usr/local/lib/python3.10/dist-packages (0.3.13)
Requirement already satisfied: packaging>=21.3 in /usr/local/lib/python3.10/dist-packages (from unstructured-ptesseract) (24.1)
Requirement already satisfied: Pillow>=8.0.0 in /usr/local/lib/python3.10/dist-packages (from unstructured-ptesseract) (10.4.0)
Collecting tesseract-ocr
  Downloading tesseract-ocr-0.0.1.tar.gz (33 kB)
  Preparing metadata (setup.py) ... done
Requirement already satisfied: cython in /usr/local/lib/python3.10/dist-packages (from tesseract-ocr) (3.0.11)
Building wheels for collected packages: tesseract-ocr
  Building wheel for tesseract-ocr (setup.py) ... done
    Created wheel for tesseract-ocr: filename=tesseract_ocr-0.0.1-cp310-cp310-linux_x86_64.whl size=169769 sha256=ce48d78777bf67501e8e60975bf1e6bfb0b3d729b7d1d2b8bb1
    Stored in directory: /root/.cache/pip/wheels/bb/f3/5c231ecbb80a1fe33204ff3021d99b54ef6daf6f8099311b8
Successfully built tesseract-ocr
Installing collected packages: tesseract-ocr
Successfully installed tesseract-ocr-0.0.1

```

- In Google Colab, I noticed some challenges with the NLTK library, so let's set the path correctly to resolve these issues.

```
1 import os
2 # Set the NLTK_DATA environment variable to the correct path
3 os.environ['NLTK_DATA'] = '/root/nltk_data'
```

- Lets make sure we have nltk version 3.9.1 installed

```
1 pip install --upgrade nltk
```

```

Requirement already satisfied: nltk in /usr/local/lib/python3.10/dist-packages (3.8.1)
Collecting nltk
  Downloading nltk-3.9.1-py3-none-any.whl.metadata (2.9 kB)
Requirement already satisfied: click in /usr/local/lib/python3.10/dist-packages (from nltk) (8.1.7)
Requirement already satisfied: joblib in /usr/local/lib/python3.10/dist-packages (from nltk) (1.4.2)
Requirement already satisfied: regex>=2021.8.3 in /usr/local/lib/python3.10/dist-packages (from nltk) (2024.5.15)
Requirement already satisfied: tqdm in /usr/local/lib/python3.10/dist-packages (from nltk) (4.66.5)
Downloading nltk-3.9.1-py3-none-any.whl (1.5 MB)
   1.5/1.5 MB 24.8 MB/s eta 0:00:00
Installing collected packages: nltk
  Attempting uninstall: nltk
    Found existing installation: nltk 3.8.1
    Uninstalling nltk-3.8.1:
      Successfully uninstalled nltk-3.8.1
Successfully installed nltk-3.9.1

```

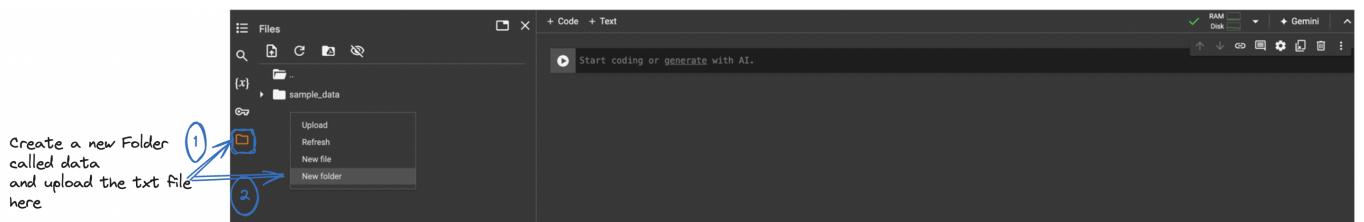


PDF with table

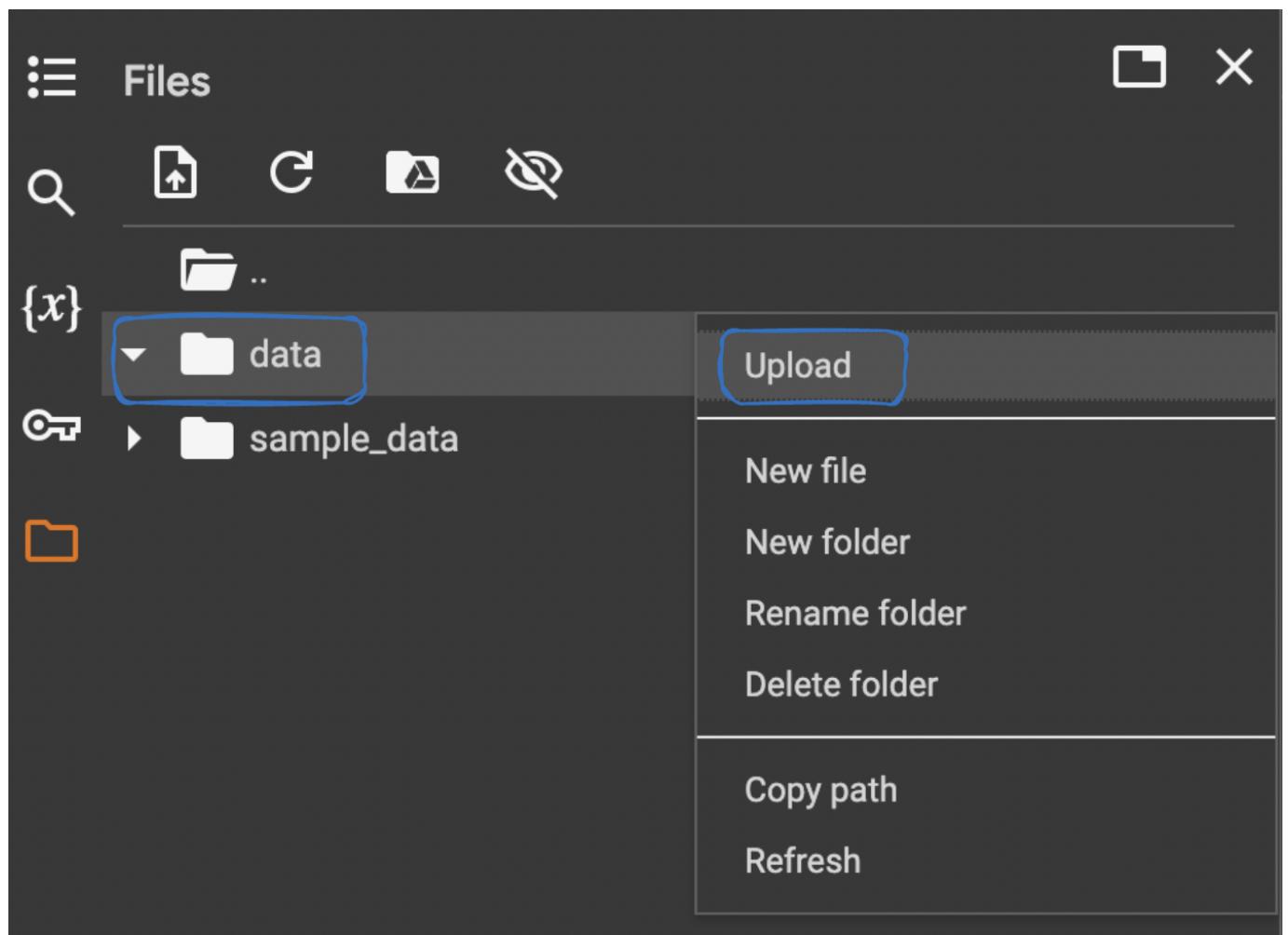
- Let's load our PDF files into Google Colab. For this example, we can use the article titled "Webex Customer Experience Essentials" from the Webex Help Center. You can also download the modified article here as we will be using in the next step.

Note: To save time and speed up processing in this lab, I have modified the CX-Essentials.pdf file by reducing it to only a few pages. You can download this version from the modified link provided above.

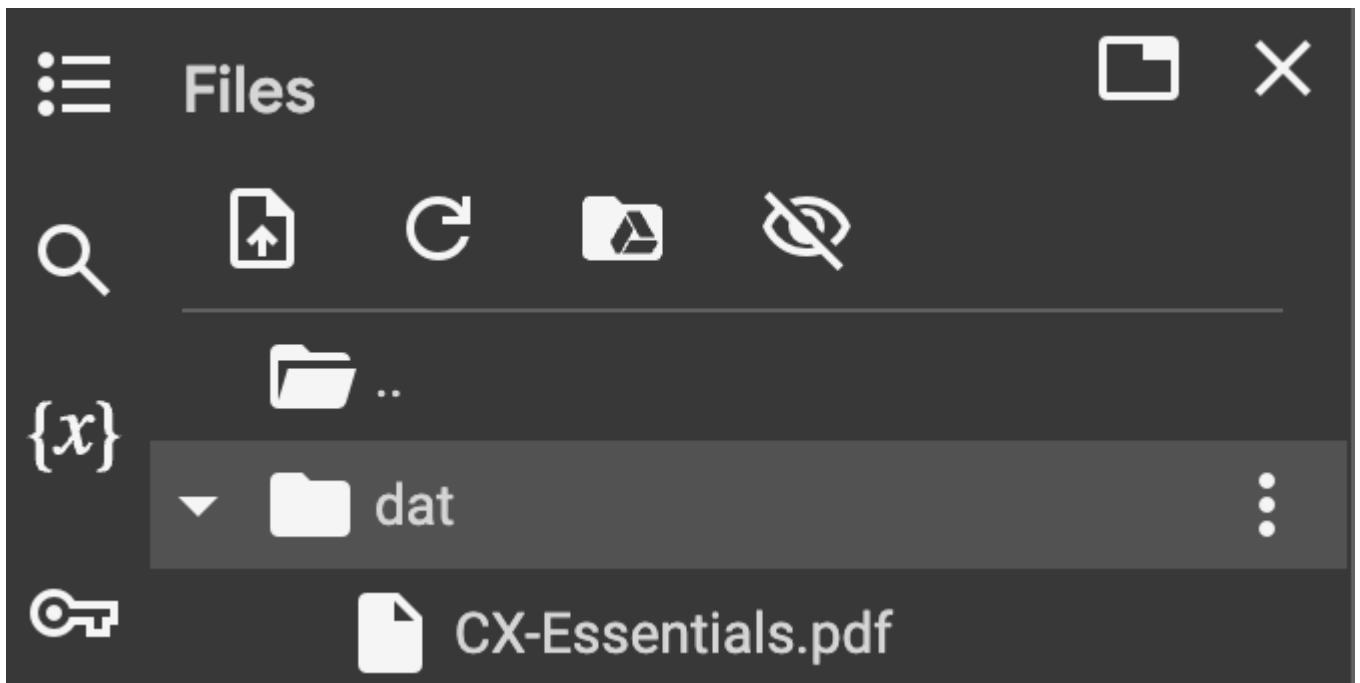
- Within Google Colab, Click on Folder and create a new folder called "dat"



- Click on [...], select Upload



- Choose your CX-Essentials.pdf file and click Open



- Let's process our PDF document by splitting it into smaller, manageable chunks based on titles, extracting images, and handling text in a way that ensures the chunks are neither too large nor too small. We will store the extracted elements (text, images, etc.) in the pdf_elements variable, and images are saved in the specified path(image_path). To understand more about unstructured partition_pdf click [here](#)

Unstructured will partition PDF files by first extracting embedded images if specified. It then processes the document's layout, dividing the content into structured elements such as titles, tables, and paragraphs based on the layout and text structure. For more detailed layout analysis, it may use advanced models like YOLOX to identify and extract bounding boxes for elements like tables and titles. This allows for the precise extraction and organization of text and images, making the content more suitable for further analysis and processing.

```

1  from unstructured.partition.pdf import partition_pdf
2  image_path = "./content/images"
3  pdf_elements = partition_pdf(
4      filename="/content/dat/CX-Essentials.pdf",
5      chunking_strategy="by_title",
6      extract_images_in_pdf=True,
7      # extract_image_block_types=["Image", "Table"],
8      max_characters=3000, # Sets the maximum number of characters
9      new_after_n_chars=2800, # Character threshold after which a new chunk will start. Ensures chunks are created before the max_characters limit is
10     reached
11     combine_text_under_n_chars=2000, #If chunk of text is smaller than 2000 characters, it should be combined with the following chunk to create a more
12     substantial piece
13     image_output_dir_path=image_path
14 )

```



Note: This process may take a few minutes to complete as we are extracting information from the PDF.

- You'll also notice that a folder has been created containing all the extracted images.

