



Google Drive: Share Files Online | ADTA\_5770 - Google Drive | QASearch\_adta5770\_ALL\_GOO | QASearch\_adta5770\_ALL\_PHASES | QASearch\_adta5770\_coding\_lectures.ipynb

File Edit View Insert Runtime Tools Help

Commands + Code + Text

```
ask("what is annotator diversity?")
> Entering new RetrievalQA chain...
> Entering new StuffDocumentsChain chain...
> Entering new LLMChain chain...
Prompt after formatting:
SYSTEM: You are an intelligent assistant helping the users with their questions on research papers.
Question: What is annotator diversity?
Strictly use ONLY the following pieces of context to answer the question at the end. Think step-by-step and then answer.
Do not try to make up an answer:
- If the answer to the question cannot be determined from the context alone, say "I cannot determine the answer to that."
- If the context is empty, just say "I do not know the answer to that."
=====
According to the survey results, a significant majority of respondents (75%) considered diversity to be a somewhat to extremely influential factor (>=3 on the Likert scale) in the 4.2 Approaches to Annotator Diversity Below, we capture the varied perspectives to annotator diversity among the participants in our study, from some who considered diversity as irrelevant to others who saw it as a key factor in ensuring high-quality data. Annotator Diversity in Data Practices The desire to control ambiguity and complexity in data annotations extended to addressing annotator subjectivities, which were regularly framed as a form of 'bias' manifested in di
```

The screenshot shows a Jupyter Notebook interface within Google Colab. The notebook title is "QASearch\_adta5770\_coding\_lectures.ipynb". The code cell contains the following Python code:

```
filters = {
    "namespace": "document_name",
    "allow_list": ["detecting-news-headline-hallucinations-with-explanations.pdf"]
}
ask("what are video localized narratives?", filters=filters)
```

Below the code, the output shows the execution steps:

```
> Entering new RetrievalQA chain...
> Entering new StuffDocumentsChain chain...
> Entering new LLMChain chain...
Prompt after formatting:
SYSTEM: You are an intelligent assistant helping the users with their questions on research papers.
Question: What are video localized narratives?
Strictly Use ONLY the following pieces of context to answer the question at the end. Think step-by-step and then answer.
```

Do not try to make up an answer:
- If the answer to the question cannot be determined from the context alone, say "I cannot determine the answer to that."
- If the context is empty, just say "I do not know the answer to that."

=====

Google Drive: Share Files Online | ADTA\_5770 - Google Drive | QASearch\_adta5770\_ALL\_GOO | QASearch\_adta5770\_ALL\_PHASES | QASearch\_adta5770\_coding\_lectures.ipynb

colab.research.google.com/drive/1yIrKCUUeq2aJXZMsdyVX935l6wQ8E18I#scrollTo=U1Ycy478VSgA

PRO File Edit View Insert Runtime Tools Help

Commands + Code + Text

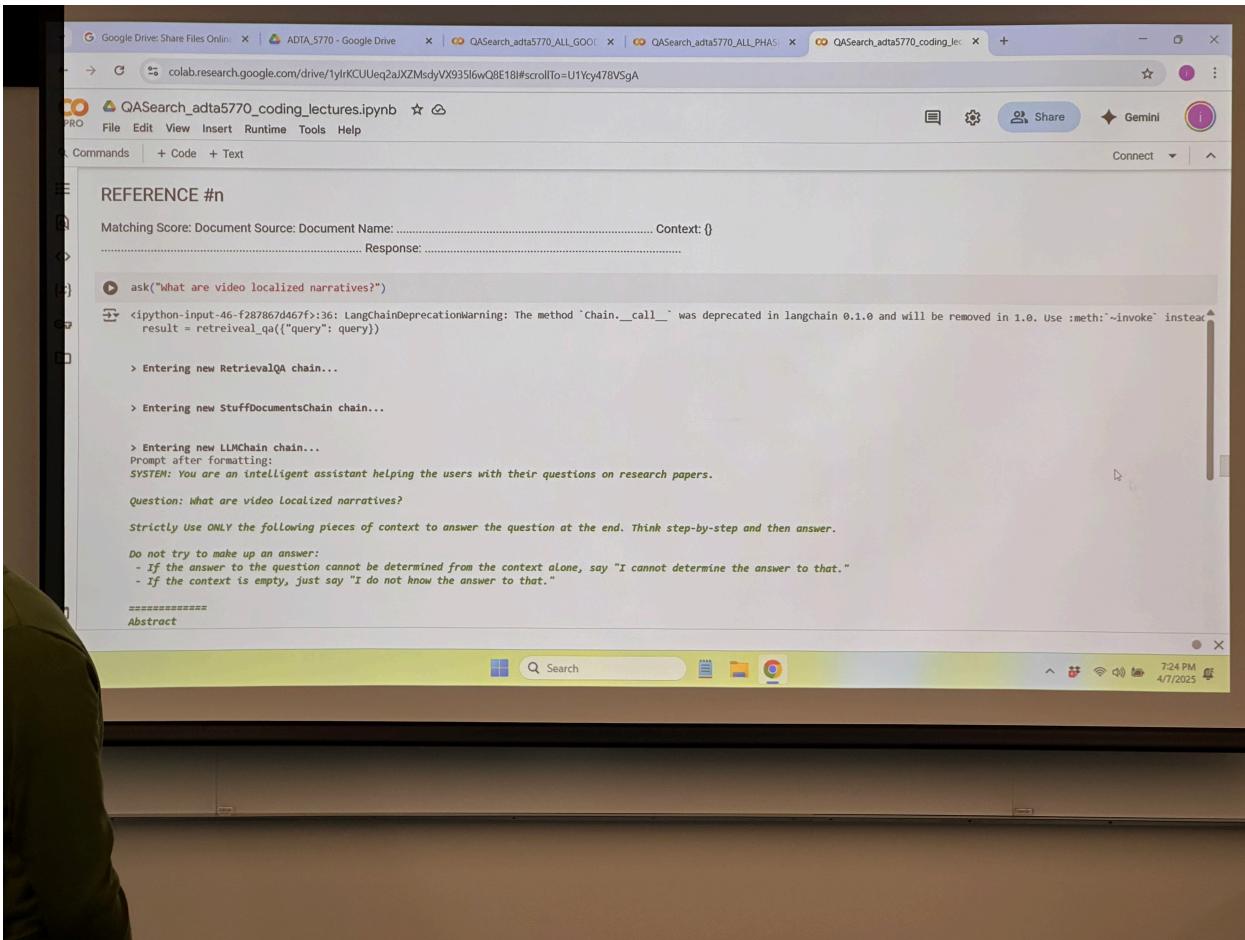
Enable VERBOSE

Enable verbose logging for debugging and troubleshooting the chains which includes the complete prompt to the LLM

```
[ ] # Enable for troubleshooting
retrieval_qa.combine_documents_chain.verbose = True
retrieval_qa.combine_documents_chain.llm_chain.verbose = True
retrieval_qa.combine_documents_chain.llm_chain.llm.verbose = True
```

Utility function used to format the response

```
def formatter(result):
    print(f"Query: {result['query'])")
    print("." * 80)
    if "source_documents" in result.keys():
        for idx, ref in enumerate(result["source_documents"]):
            print("-" * 80)
            print(f"REFERENCE #{idx}")
            print("." * 80)
            if "score" in ref.metadata:
                print(f"Matching Score: {ref.metadata['score'])")
            if "source" in ref.metadata:
                print(f"Document Source: {ref.metadata['source'])")
            if "document_name" in ref.metadata:
                print(f"Document Name: {ref.metadata['document_name'])")
            print("." * 80)
            print(f"Content: \n{ref.page_content}")
    print("." * 80)
```



The screenshot shows a Jupyter Notebook interface in Google Colab. The title bar indicates the notebook is titled "QASearch\_adta5770\_coding\_lectures.ipynb". The code cell contains the following Python script:

```
print("-" * 80)
if "score" in ref.metadata:
    print(f"Matching Score: {ref.metadata['score']}")
if "source" in ref.metadata:
    print(f"Document Source: {ref.metadata['source']}")
if "document_name" in ref.metadata:
    print(f"Document Name: {ref.metadata['document_name']}")
print("." * 80)
print(f"Content: \n{wrap(ref.page_content)}")
print("." * 80)
print(f"Response: {wrap(result['result'])}")
print("." * 80)

def wrap(s):
    return "\n".join(textwrap.wrap(s, width=120, break_long_words=False))

def ask(
    query,
    qa=retrieveal_qa,
    k=NUMBER_OF_RESULTS,
    search_distance=SEARCH_DISTANCE_THRESHOLD,
    filters={}
):
    retrieveal_qa.retriever.search_kwargs["search_distance"] = search_distance
    retrieveal_qa.retriever.search_kwargs["k"] = k
    retrieveal_qa.retriever.search_kwargs["filters"] = filters
    result = retrieveal_qa({"query": query})
    return formatter(result)
```

The screenshot shows a Jupyter Notebook interface with the following content:

```
File Edit View Insert Runtime Tools Help
Commands + Code + Text
Run QA chain on sample questions
Following are sample questions: To run a query
--> RetrievalQA chain takes the user question.
--> Next, RetrievalQA chain calls the retriever to fetch top k semantically similar texts from the Matching Engine Index (vector store). --> Then
RetrievalQA chain passes to the LLM as part of the prompt.
--> The final prompt sent to the LLM looks like this format:
SYSTEM: {system}
=====
{context}
Question: {question} Helpful Answer: where:
system: Instructions for LLM on how to respond to the question based on the context context: Semantically similar text (a.k.a snippets)
retrieved from the vector store question: question posed by the user The response returned from the LLM includes both the response and
references that lead to the response. This way the response from LLM is always grounded to the sources. Here we have formatted the
response as:
Question: {question}
REFERENCE #n
Matching Score: Document Source: Document Name: ..... Context: {}
```

The screenshot shows a Jupyter Notebook interface with the following details:

- Title Bar:** Google Drive: Share Files Online | ADTA\_5770 - Google Drive | QAsearch\_adta5770\_ALL\_GOO | QAsearch\_adta5770\_ALL\_PHASE | QAsearch\_adta5770\_coding\_lectures.ipynb
- File Menu:** File Edit View Insert Runtime Tools Help
- Toolbar:** Commands + Code + Text
- Section:** Configure a RetrievalQA chain
- VIP NOTES:**
  - > A Vertex AI vector search index, i.e., a vector database, set up as a Vertex AI vector store can be configured as a RetrievalQA.
  - > RetrievalQA is a Vertex AI vector search tool specifically used for Q&A Search
- Code Block:**

```
# Uses LLM to synthesize results from the search index.  
# Use Vertex AI PaLM Text API for LLM  
retrieval_qa = RetrievalQA.from_chain_type(  
    llm=llm,  
    chain_type="stuff",  
    retriever=retriever,  
    return_source_documents=True,  
    verbose=True,  
    chain_type_kwargs={  
        "prompt": PromptTemplate(  
            template=template,  
            input_variables=["context", "question"],  
        ),  
    },  
)
```
- Section:** Enable VERBOSE
- Description:** Enable verbose logging for debugging and troubleshooting the chains which includes the complete prompt to the LLM

The screenshot shows a Jupyter Notebook interface with the following details:

- Title Bar:** Google Drive: Share Files Online | ADTA\_5770 - Google Drive | QASearch\_adta5770\_ALL\_GOOD | QASearch\_adta5770\_ALL\_PHASES | QASearch\_adta5770\_coding\_lectures
- Header:** colab.research.google.com/drive/1yIrKCUUeq2aJXZMsdyVX935l6wQ8E18I#scrollTo=U1Ycy478VSgA
- Toolbar:** PRO File Edit View Insert Runtime Tools Help Commands + Code + Text
- Section:** Configure a RetrievalQA chain
- VIP NOTES:**
  - ) A Vertex AI vector search index, i.e., a vector database, set up as a Vertex AI vector store can be configured as a RetrievalQA.
  - ) RetrievalQA is a Vertex AI vector search tool specifically used for Q&A Search
- Code Block:**

```
# Uses LLM to synthesize results from the search index.  
# Use Vertex AI PaLM Text API for LLM  
retrieveal_qa = RetrievalQA.from_chain_type(  
    llm=llm,  
    chain_type="stuff",  
    retriever=retriever,  
    return_source_documents=True,  
    verbose=True,  
    chain_type_kwargs={  
        "prompt": PromptTemplate(  
            template=template,  
            input_variables=["context", "question"],  
        ),  
    },  
)
```
- Section:** Enable VERBOSE
- Description:** Enable verbose logging for debugging and troubleshooting the chains which includes the complete prompt to the LLM

The screenshot shows a Google Colab notebook interface. The title bar indicates the file is named 'QASearch\_adta5770\_coding\_lectures.ipynb'. The main content area contains two sections:

- Customize a prompt template**:

```
[ ] template = """SYSTEM: You are an intelligent assistant helping the users with their questions on research papers.  
Question: {question}  
===== {context} =====  
Question: {question}  
Helpful Answer:"""
```
- Configure a RetrievalQA chain**:

**VIP NOTES:**

  - ) A Vertex AI vector search index, i.e., a vector database, set up as a Vertex AI vector store can be configured as a RetrievalQA.
  - ) RetrievalQA is a Vertex AI vector search tool specifically used for Q&A Search

```
[ ] # Uses LLM to synthesize results from the search index.
```

The screenshot shows a Google Colab notebook interface. The title bar indicates the file is 'QASearch\_adta5770\_coding\_lectures.ipynb'. The menu bar includes File, Edit, View, Insert, Runtime, Tools, Help, Commands, + Code, and + Text. The main content area displays the following code:

```
# Create chain to answer questions
# Number of responses that the system tries to search for
NUMBER_OF_RESULTS = 10

# The distance (similarity search) used as a threshold to search for responses
SEARCH_DISTANCE_THRESHOLD = 0.6

# Expose index to the retriever
retriever = vsvectordb.as_retriever(
    search_type="similarity",
    search_kwargs={
        "k": NUMBER_OF_RESULTS,
        "search_distance": SEARCH_DISTANCE_THRESHOLD,
    },
    filters=None,
)

[ ] template = """SYSTEM: You are an intelligent assistant helping the users with their questions on research papers.
```

-) A Vertex AI vector search index, i.e., a vector database, set up as a Vertex AI vector store can be configured as a RetrievalQA.  
-) RetrievalQA is a Vertex AI vector search tool specifically used for Q&A Search

First, need to configure a retriever

-) Retriever is used to retrieve response from the vector search index (or vector store)

```
# Create chain to answer questions
# Number of responses that the system tries to search for
NUMBER_OF_RESULTS = 10

# The distance (similarity search) used as a threshold to search for responses
SEARCH_DISTANCE_THRESHOLD = 0.6

# Expose index to the retriever
retriever = vsvectordb.as_retriever(
    search_type="similarity",
    search_kwargs={
        "K": NUMBER_OF_RESULTS,
        "search_distance": SEARCH_DISTANCE_THRESHOLD,
    },
    filters=None,
)
```

Customize a prompt template

```
[ ] template = """SYSTEM: You are an intelligent assistant helping the users with their questions on research papers.
```

Google Drive: Share Files Online | ADTA\_5770 - Google Drive | QASearch\_adta5770\_ALL\_GOO! | QASearch\_adta5770\_ALL\_PHASES | QASearch\_adta5770\_coding\_lectures

colab.research.google.com/drive/1yIrKCUUeq2ajXZMsdyVX935l6wQ8E18I#scrollTo=U1Ycy478VSgA

QASearch\_adta5770\_coding\_lectures.ipynb

File Edit View Insert Runtime Tools Help

Commands + Code + Text

-) A Vertex AI vector search index, i.e., a vector database, set up as a Vertex AI vector store can be configured as a RetrievalQA.  
-) RetrievalQA is a Vertex AI vector search tool specifically used for Q&A Search

First, need to configure a retriever

-) Retriever is used to retrieve response from the vector search index (or vector store)

```
# Create chain to answer questions

# Number of responses that the system tries to search for
NUMBER_OF_RESULTS = 10

# The distance (similarity search) used as a threshold to search for responses
SEARCH_DISTANCE_THRESHOLD = 0.6

# Expose index to the retriever
retriever = vsvectordb.as_retriever(
    search_type="similarity",
    search_kwargs={
        "k": NUMBER_OF_RESULTS,
        "search_distance": SEARCH_DISTANCE_THRESHOLD,
    },
    filters=None,
)
```

Customize a prompt template

```
[ ] template = """SYSTEM: You are an intelligent assistant helping the users with their questions on research papers.
```

The screenshot shows a Jupyter Notebook interface with the following details:

- Title Bar:** Google Drive: Share Files Online, ADTA\_5770 - Google Drive, QASearch\_adta5770\_ALL\_GOO!, QASearch\_adta5770\_ALL\_PHASE!, QASearch\_adta5770\_coding\_lectures.ipynb
- File Menu:** File, Edit, View, Insert, Runtime, Tools, Help
- Toolbar:** Commands, + Code, + Text, Share
- Section Header:** Q&A SEARCH
- Section:** Create an LLM based on an existing LLM
- Code:**

```
[ ] # Create a model to perform Q&A Search
# This model is based on one of pre-trained LLM

from langchain_google_vertexai import VertexAI

llm = VertexAI(
    model_name ="gemini-2.5-pro-exp-03-25",
    max_output_tokens=8192,
    temperature=0.2,
    top_p=0.8,
    top_k=40,
    verbose=True,
)
```
- Section:** Configure an instance of RetrievalQA
- VIP NOTES:**
  - ) A Vertex AI vector search index, i.e., a vector database, set up as a Vertex AI vector store can be configured as a RetrievalQA.
  - ) RetrievalQA is a Vertex AI vector search tool specifically used for Q&A Search
- Text:** First, need to configure a retriever

The screenshot shows a Google Colab notebook interface. The title bar indicates the file is 'QASearch\_adta5770\_coding\_lectures.ipynb'. The menu bar includes File, Edit, View, Insert, Runtime, Tools, Help, Commands, + Code, and + Text. The main content area has a section titled 'VIP NOTES:' which contains the following text:

ACCESS THE EXISTING VECTOR SEARCH INDEX TO PERFORM Q&A SEARCH

Validate semantic search with Matching Engine is working

```
# Test whether search from vector store is working
vsvectordb.similarity_search("what are video localized narratives?", k=2)
```

[Document(metadata={'source': 'gs://adta5770-bucket-1/docs/pdfs', 'document\_name': 'connecting-vision-and-language-with-video-localized-narratives.pdf', 'page\_content': 'Abstract\nWe propose Video Localized Narratives, a new form of multimodal video annotations connecting vision and language. In the original [40], annotators speak and move their mouse simultaneously on an image, thus grounding each word with a mouse trace segment. However, this is challenging protocol empowers annotators to tell the story of a video with localized Narratives, capturing even complex events involving multiple actors interacting with several passive objects. We annotated 20k videos of the OVIS, UVO, and OOPS datasets, totalling 1.7M words. Based on this data, we also construct new benchmark narrative grounding and video question answering tasks, and provide reference results from strong baseline models. Our annotations are available at https://.../localized-narratives/.'}), Document(metadata={'source': 'gs://adta5770-bucket-1/docs/pdfs', 'document\_name': 'connecting-vision-and-language-with-video-localized-narratives.pdf', 'page\_content': 'Vision and language is a very active research area, which experienced much exciting progress recently [1,32,42,43, 52]. At the heart of many connecting still images to captions, while grounding some of the words in the caption to regions in the image, made with a variety of protocols [6,21,24,26, recent Localized Narratives (IMLN) [40] offer a particularly attractive solution: the annotators describe an image with their voice while simultaneously regions they are describing. Speaking is natural and efficient, resulting in long captions that describe the whole scene. Moreover, the synchronization between pointer yields dense visual grounding for every word. Yet, still images only show one instant in time. Annotating videos would be even more interesting, as featuring a flow of events involving multiple actors and objects interacting with each other.'})]

Q&A SEARCH

**VIP NOTES:**

ACCESS THE EXISTING VECTOR SEARCH INDEX TO PERFORM Q&A SEARCH

Validate semantic search with Matching Engine is working

```
# Test whether search from vector store is working
svvectordb.similarity_search("What are video localized narratives?", k=2)
```

```
[Document(metadata={'source': 'gs://adta5770-bucket-1/docs/pdfs', 'document_name': 'connecting-vision-and-language-with-video-localized-narratives.pdf', 'chunk': 1145}, page_content='Abstract\n\nWe propose Video Localized Narratives, a new form of multimodal video annotations connecting vision and language. In the original Localized [40], annotators speak and move their mouse simultaneously on an image, thus grounding each word with a mouse trace segment. However, this is challenging on a video protocol empowers annotators to tell the story of a video with Localized Narratives, capturing even complex events involving multiple actors interacting with each other and passive objects. We annotated 20k videos of the OTIS, IVO, and Dops datasets, totalling 1.7M words. Based on this data, we also construct new benchmarks for narrative grounding and video question answering tasks, and provide reference results from strong baseline models. Our annotations are available at https://google.github.io/localized-narratives/.'), Document(metadata={'source': 'gs://adta5770-bucket-1/docs/pdfs', 'document_name': 'connecting-vision-and-language-with-video-localized-narratives.pdf', 'chunk': 1146}, page_content='Vision and language is a very active research area, which experienced much exciting progress recently [1,32,42,43, 52]. At the heart of many development recent Localized Narratives (LmNs) [40] offer a particularly attractive solution: the annotators describe an image with their voice while simultaneously moving the regions they are describing. Speaking is natural and efficient, resulting in long captions that describe the whole scene. Moreover, the synchronization between the voice pointer yields dense visual grounding for every word. Yet, still images only show one instant in time. Annotating videos would be even more interesting, as they should feature a flow of events involving multiple actors and objects interacting with each other.'))]
```

**Q&A SEARCH**

Google Drive: Share Files Online | ADTA\_5770 - Google Drive | QASearch\_adta5770\_ALL\_GOO! | QASearch\_adta5770\_ALL\_PHASES | QASearch\_adta5770\_coding\_lectures.ipynb | +

colab.research.google.com/drive/1ylrKCUUeq2aXZMsdyVX935l6wQ8E18l#scrollTo=U1Ycy478VSgA

QASearch\_adta5770\_coding\_lectures.ipynb ☆ ↗

File Edit View Insert Runtime Tools Help

ands | + Code + Text

Add documents as embeddings in Matching Engine as index

The document chunks are transformed as embeddings (vectors) using Vertex AI Embeddings API and added to the index with [streaming index update](#). With Streaming Updates, you can update and query your index within a few seconds.

The original document text is stored on Cloud Storage bucket had referenced by id.

```
[ ] # Store docs as embeddings in Matching Engine/Vector Search index
# It may take a while since API is rate limited
texts = [doc.page_content for doc in doc_splits]
metadata = [doc.metadata for doc in doc_splits]
```

Add embeddings to the vector store

OTE:  
Depending on the volume and size of documents, this step may take time.

```
[ ] for i in range(0, len(texts), 1000):
    vsvectordb.add_texts(texts[i : i + 1000], metadata[i : i + 1000])
```

**VIP NOTES:**

ACCESS THE EXISTING VECTOR SEARCH INDEX TO PERFORM Q&A SEARCH

The screenshot shows a Google Colab notebook interface. The title bar says 'QAsearch\_adta5770\_coding\_lectures.ipynb'. The menu bar includes File, Edit, View, Insert, Runtime, Tools, Help, Commands, + Code, and + Text. The main content area has a section titled 'Add documents as embeddings in Matching Engine as index'. It contains a note about using Vertex AI Embeddings API and streaming index update. Below is a code snippet:

```
[ ] # Store docs as embeddings in Matching Engine/Vector Search index
# It may take a while since API is rate limited
texts = [doc.page_content for doc in doc_splits]
metadata = [doc.metadata for doc in doc_splits]
```

There is also a section 'Add embeddings to the vector store' with a note about volume and size. A code snippet for adding texts to a database is shown:

```
for i in range(0, len(texts), 1000):
    vsvectordb.add_texts(texts[i : i + 1000], metadata[i : i + 1000])
```

A section titled 'VIP NOTES:' is present with the heading 'ACCESS THE EXISTING VECTOR SEARCH INDEX TO PERFORM Q&A SEARCH'.

colab.research.google.com/drive/1ylrKCUUeq2aJXZMsdyVX935l6wQ8E18I#scrollTo=U1Ycy478VSgA

QASearch\_adta5770\_coding\_lectures.ipynb

File Edit View Insert Runtime Tools Help

Commands + Code + Text

**VIP NOTES:**

-) if an existing bucket has been created before for previous code runs, MUST delete it first

```
[ ] # Create a GCS bucket to store the Matching Engine/Vector Search index
! set -x && gsutil mb -p $PROJECT_ID -l us-central1 gs://$VSVDB_EMBEDDING_DIR
→ + gsutil mb -p ai-dl-vertexai-examples -l us-central1 gs://adta5770-vsvdb-bucket
Creating gs://adta5770-vsvdb-bucket/...

[ ] # Specify the embedding model used for embedding
embeddings = VertexAIEmbeddings(model_name="text-embedding-005")

▶ # Embeddings API integrated with LangChain
# Create a Vector Search vector database, vsvectordb
# It is actually a GCP Vertex AI vector store
# initialize the vector store

vsvectordb = VectorSearchVectorStore.from_components(
    project_id=PROJECT_ID,
    region=VSVDB_REGION,
    gcs_bucket_name=f"gs://{VSVDB_EMBEDDING_DIR}".split("/")[2],
    embedding=embeddings,
    index_id=index.name,
    endpoint_id=index_endpoint.name,
    stream_update=True,
)
```

Search

Google Drive: Share Files Online x | ADTA\_5770 - Google Drive x | QASearch\_adta5770\_ALL\_GOOD x | QASearch\_adta5770\_ALL\_PHASES x | QASearch\_adta5770\_coding\_Lec x +

colab.research.google.com/drive/1ylrKCUUeq2aJZMsdyVX935l6wQ8E18I#scrollTo=U1Ycy478VSgA

QASearch\_adta5770\_coding\_lectures.ipynb ☆ ↗

File Edit View Insert Runtime Tools Help

hands | + Code + Text

# of documents = 2080

▼ PHASE 4 START

Configure Matching Engine/Vector Search as Vector Store

-) Initialize Matching Engine/Vector Search vector store with text embeddings mode

▼ SET UP to access existing Matching Engine/Vector Search index (GCP: Vector Store)

```
[ ] # Define credentials
# VSDB: Vector Search Vector Database (Vector Database: GCP Vector Store)

VSDB_REGION = "us-central1"
VSDB_INDEX_NAME = "adta5770-vsdb-index" # @param {type:"string"}
VSDB_EMBEDDING_DIR = "adta5770-vsdb-bucket" # @param {type:"string"}
VSDB_DIMENSIONS = 768 # when using Vertex AI text-embedding-002
```

VSDB\_INDEX\_NAME: "adta5770-vsdb-index"

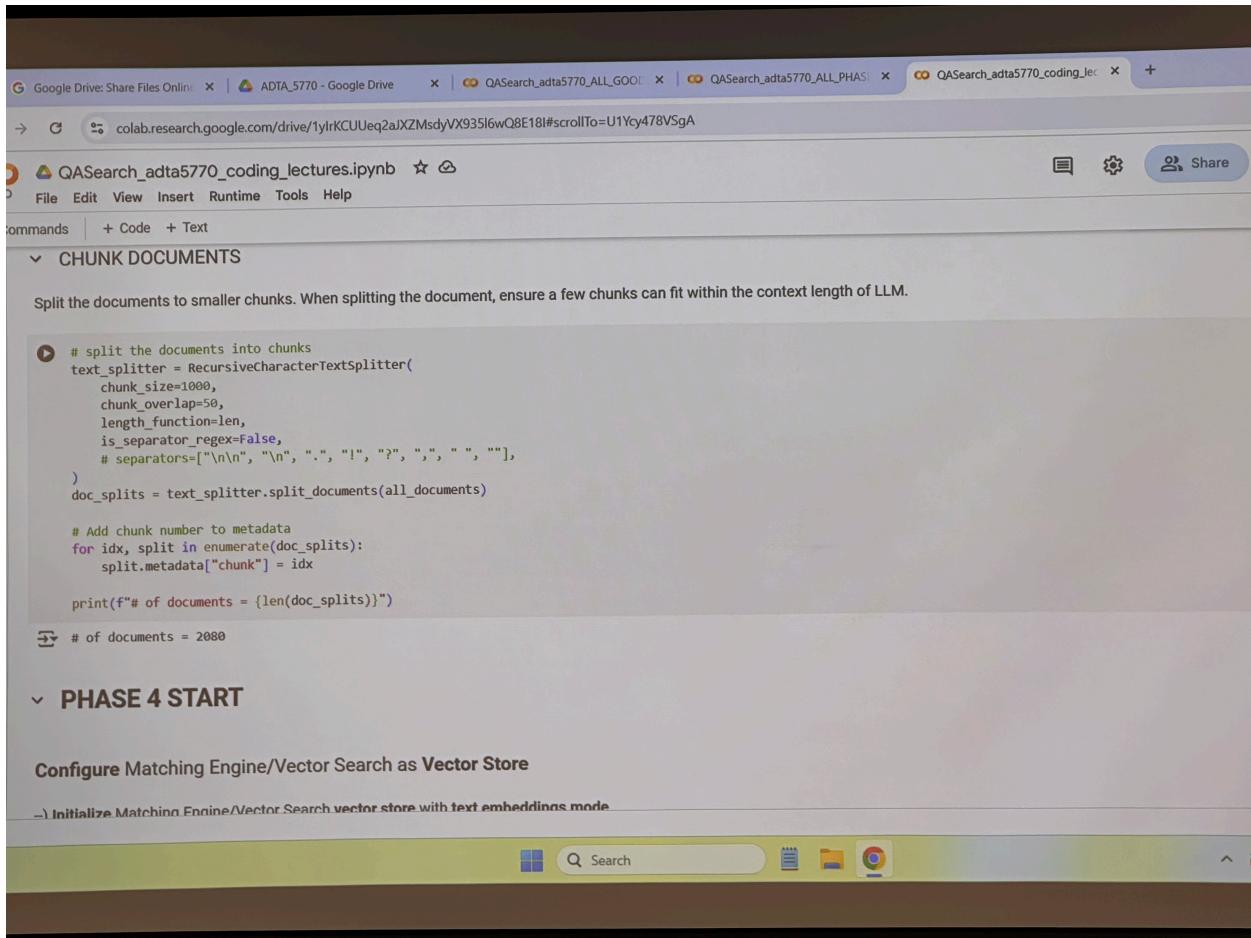
VSDB\_EMBEDDING\_DIR: "adta5770-vsdb-bucket"

▼ Create a GCS bucket to store vector search index

VIP NOTES:

-) if an existing bucket has been created before for previous code runs, MUST delete it first

```
[# Create a GCS bucket to store the Matching Engine/Vector Search index
! set -x && gsutil mb -p $PROJECT_ID -l us-central1 gs://$VSVD8_EMBEDDING_DIR
+ gsutil mb -p ai-dl-vertexai-examples -l us-central1 gs://adta5770-vsvdb-bucket
Creating gs://adta5770-vsvdb-bucket/...
] # Specify the embedding model used for embedding
embeddings = VertexAIEmbeddings(model_name="text-embedding-005")
# Embeddings API integrated with LangChain
# Create a Vector Search vector database, vsvectordb
# It is actually a GCP Vertex AI vector store
# initialize the vector store
vsvectordb = VectorSearchVectorStore.from_components(
    project_id=PROJECT_ID,
    region=VSVD8_REGION,
    gcs_bucket_name=f"gs://{(VSVD8_EMBEDDING_DIR).split('/')[2]}",
    embedding=embeddings,
    index_id=index.name,
    endpoint_id=index_endpoint.name,
    stream_update=True,
)
```



```
# split the documents into chunks
text_splitter = RecursiveCharacterTextSplitter(
    chunk_size=1000,
    chunk_overlap=50,
    length_function=len,
    is_separator_regex=False,
    # separators=[r'\n\n', r'\n', '.', '!', '?', ',', '\"', '\"', '\"'],
)
doc_splits = text_splitter.split_documents(all_documents)

# Add chunk number to metadata
for idx, split in enumerate(doc_splits):
    split.metadata["chunk"] = idx

print(f"# of documents = {len(doc_splits)}")

# of documents = 2080
```

Split the documents to smaller chunks. When splitting the document, ensure a few chunks can fit within the context length of LLM.

## CHUNK DOCUMENTS

Configure Matching Engine/Vector Search as Vector Store

→ Initialize Matching Engine/Vector Search vector store with text embeddings mode

The screenshot shows a Google Colab notebook interface. The title bar reads "QASearch\_adta5770\_coding\_lectures.ipynb". The code cell contains the following Python script:

```
from langchain_google_community import GCSDirectoryLoader
from langchain_google_community import GCSFileLoader

# Make sure these variables are set in your environment or defined here
PROJECT_ID = "ai-dl-vertexai-examples" # Replace with your Project ID if not using env var
GCS_BUCKET_NAME = "adta5770-bucket-1" # Replace with your bucket name (e.g., the value of GCS_BUCKET_DOCS)
FOLDER_PREFIX = "docs/pdfs/" # Replace with the folder prefix within the bucket (e.g., the value of folder_prefix)
# Ensure FOLDER_PREFIX ends with a '/' if it's targeting a directory structure

# --- Initialization ---
print(f"Processing documents from gs://{GCS_BUCKET_NAME}/{FOLDER_PREFIX}")
storage_client = storage.Client(project=PROJECT_ID)
bucket = storage_client.bucket(GCS_BUCKET_NAME)

# --- Load Documents ---
all_documents = []
blobs = bucket.list_blobs(prefix=FOLDER_PREFIX) # List blobs matching the prefix

for blob in blobs:
    print(str(blob))

    # Skip directories/folders if represented as blobs, and ensure it's a PDF
    if blob.name.endswith("/") or not blob.name.lower().endswith(".pdf"):
        continue

    print(f"Loading document: {blob.name}")
```

The screenshot shows a Google Colab notebook interface. The title bar reads "QASearch\_adta5770\_coding\_lectures.ipynb". The code cell contains the following Python script:

```
# Skip directories/folders if represented as blobs, and ensure it's a PDF
if blob.name.endswith("/") or not blob.name.lower().endswith(".pdf"):
    continue

print(f" Loading document: {blob.name}")

# Use GCSFileLoader for each PDF blob
loader = GCSFileLoader(
    project_name=PROJECT_ID, bucket=GCS_BUCKET_NAME, blob=blob.name
)

# Load documents (GCSFileLoader often returns one Document per page)
documents_from_blob = loader.load()

===== NEW NEW NEW

# --- Metadata Enhancement (similar to original logic) ---
# Derive document name from the blob name
document_name = blob.name.split("/")[-1]

# Derive doc source prefix and suffix to match original logic
# Note: The original code's 'source' derivation seemed to point to the
#       "folder" rather than the file. This replicates that.
#       A more common approach might be to use the blob's full gs:// path
#       as the source.
doc_source_prefix = f"gs://{GCS_BUCKET_NAME}"
# Get the directory path containing the blob
doc_source_suffix = "/".join(blob.name.split("/")[:-1])
source = f"{doc_source_prefix}/{doc_source_suffix}" # Folder path as source
```

The screenshot shows a Jupyter Notebook interface in Google Colab. The notebook title is "QASearch\_adta5770\_coding\_lectures.ipynb". The code cell contains Python code for processing PDFs, specifically for extracting documents from a Google Cloud Storage bucket and adding metadata. The output cell displays the logs of the PDFMiner library, which are numerous "WARNING" messages indicating missing CropBox information for each page of the loaded PDFs. The browser tabs at the top show other open notebooks and Google Drive.

```
doc_source_prefix = f"gs://{GCS_BUCKET_NAME}"
# Get the directory path containing the blob
doc_source_suffix = "/".join(blob.name.split("/")[0:-1])
source = f"{doc_source_prefix}/{doc_source_suffix}" # Folder path as source

for document in documents_from_blob:
    # Add derived metadata to each document (page)
    document.metadata["source"] = source
    document.metadata["document_name"] = document_name

# GCSFileLoader might add other useful metadata like 'page' automatically
all_documents.extend(documents_from_blob)

# The 'all_documents' list now contains LangChain Document objects,
# likely one per page from all loaded PDFs.
print(f"# of document pages loaded (pre-chunking) = {len(all_documents)}")
```

```
<Blob: adta5770-bucket-1, docs/pdfs/differentially-private-continual-releases-of-streaming-freq-moment-estimations.pdf, 1743568051399383>
Loading document: docs/pdfs/development-of-a-ml-model-for-sonographic-assessment-of-gestational-age.pdf
<Blob: adta5770-bucket-1, docs/pdfs/differentially-private-continual-releases-of-streaming-freq-moment-estimations.pdf, 1743568051399383>
Loading document: docs/pdfs/differentially-private-continual-releases-of-streaming-freq-moment-estimations.pdf
WARNING:pdfminer.pdfpage:CropBox missing from /Page, defaulting to MediaBox
```

```
resource name: projects/446492160130/locations/us-central1/indexEndpoints/2235188392012808192
INDEX ENDPOINT ID or name: 2235188392012808192

▼ PHASE 3 START

} ▼ Add Document Embeddings to Matching Engine - Vector Store

□ ▼ Ingest PDF files

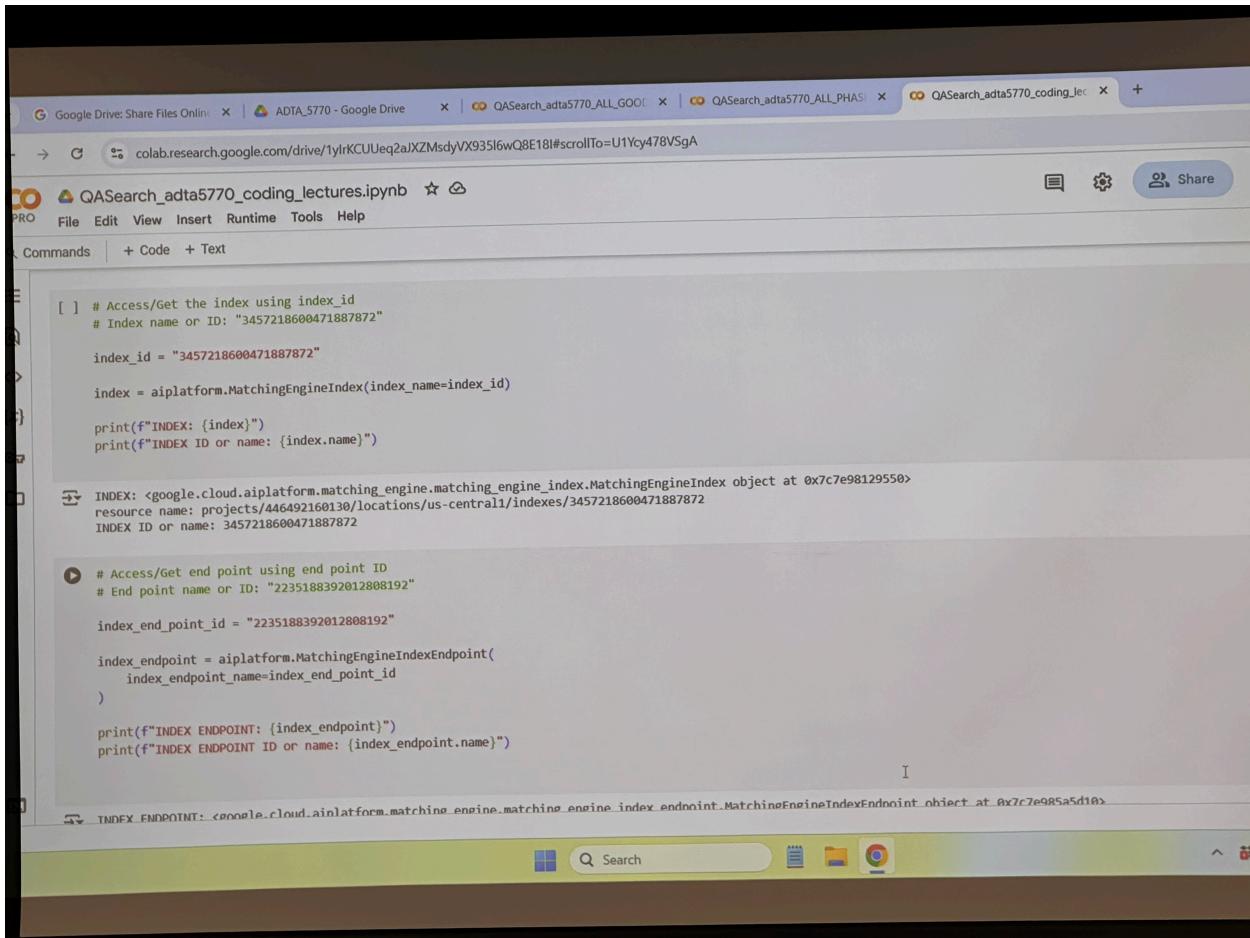
VIP NOTES:
→) A bucket has been created in GCS
→) PDFs have been uploaded into the bucket

▶ # For GCS: Google Cloud Storage access and loading
# MUST import again here!
from google.cloud import storage

client = storage.Client()

for blob in client.list_blobs('adta5770-bucket-1', prefix='docs/pdfs'):
    print(str(blob))

<Blob: adta5770-bucket-1, docs/pdfs/, 1743567844880661>
<Blob: adta5770-bucket-1, docs/pdfs/a-human-centered-approach-to-developer-productivity.pdf, 1743568042270002>
<Blob: adta5770-bucket-1, docs/pdfs/a-mixed-methods-approach-to-understanding-user-trust-after-voice-assistant-failures.pdf, 1743568043092578>
```



```
[ ] # Access/Get the index using index_id
# Index name or ID: "3457218600471887872"

index_id = "3457218600471887872"
index = aiplatformMatchingEngineIndex(index_name=index_id)

print(f"INDEX: {index}")
print(f"INDEX ID or name: {index.name}")

INDEX: <google.cloud.aiplatform.matching_engine.matching_engine_index.MatchingEngineIndex object at 0x7c7e98129550>
resource name: projects/446492160130/locations/us-central1/indexes/3457218600471887872
INDEX ID or name: 3457218600471887872

[ ] # Access/Get end point using end point ID
# End point name or ID: "2235188392012808192"

index_end_point_id = "2235188392012808192"
index_endpoint = aiplatformMatchingEngineEndpoint(
    index_endpoint_name=index_end_point_id
)

print(f"INDEX ENDPOINT: {index_endpoint}")
print(f"INDEX ENDPOINT ID or name: {index_endpoint.name}")

INDEX ENDPOINT: <google.cloud.aiplatform.matching_engine.matching_engine_index.endpoint.MatchingEngineIndexEndpoint object at 0x7c7e985a5d10>
```

The screenshot shows a Jupyter Notebook interface in Google Colab. The notebook title is `QASearch_adta5770_coding_lectures.ipynb`. The code cell contains the following Python code:

```
[ ] # First, get the list of indexes that has been created for the vector search system
list_indexes = aiplatformMatchingEngineIndex.list()
print(f"List of indexes: {list_indexes}")

# Access/Get the index using index_id
# Index name or ID: "3457218600471887872"
index_id = "3457218600471887872"
```

The output of the first cell is:

```
>List of indexes: [

The notebook also includes a section titled "GET & SET index_id and end_point_id" with notes and a code example.


```

```
from google.cloud import storage
# IMPORT FROM LANGCHAIN
import langchain
from langchain.chains import RetrievalQA
# from langchain.document_loaders import GCSDirectoryLoader
# from langchain.document_loaders import GCSFileLoader
from langchain.prompts import PromptTemplate
from langchain.text_splitter import RecursiveCharacterTextSplitter
from langchain_community.document_loaders import PyPDFLoader
from langchain_text_splitters import RecursiveCharacterTextSplitter
# IMPORT FROM GCP: VERTEX AI & LANGCHAIN API
from langchain_google_vertexai import VertexAI
from langchain_google_vertexai import VertexAIEMBEDDINGS
from langchain_google_vertexai import (
    VectorSearchVectorStore,
    VectorSearchVectorStoreDatastore,
)
# IMPORT OTHERS
import textwrap
```

The screenshot shows a Jupyter Notebook interface with the following code content:

```
[ ] from google.cloud import storage  
[ ] # IMPORT FROM LANGCHAIN  
import langchain  
from langchain.chains import RetrievalQA  
# from langchain.document_loaders import GCSDirectoryLoader  
# from langchain.document_loaders import GCSFileLoader  
from langchain.prompts import PromptTemplate  
from langchain.text_splitter import RecursiveCharacterTextSplitter  
from langchain_community.document_loaders import PyPDFLoader  
from langchain_text_splitters import RecursiveCharacterTextSplitter  
[ ] from langchain_google_community import GCSDirectoryLoader  
from langchain_google_community import GCSFileLoader  
[ ] # IMPORT FROM GCP: VERTEX AI & LANGCHAIN API  
from langchain_google_vertexai import VertexAI  
from langchain_google_vertexai import VertexAIEmbeddings  
from langchain_google_vertexai import (  
    VectorSearchVectorStore,  
    VectorSearchVectorStoreDatastore,  
)  
[ ] # IMPORT OTHERS  
import textwrap
```

```
import IPython
app = IPython.Application.instance()
app.kernel.do_shutdown(True)

{'status': 'ok', 'restart': True}

import sys
if "google.colab" in sys.modules:
    from google.colab import auth
    auth.authenticate_user()

# IMPORT FROM GCP: VERTEX AI
from google.cloud import aiplatform
import vertexai
from google.cloud.aiplatform.matching_engine.matching_engine_index_endpoint import (
    Namespace,
    NumericNamespace,
)
# For GCS: Google Cloud Storage access and loading
from google.cloud import storage

# IMPORT FROM LANGCHAIN
import langchain
```

The screenshot shows a Jupyter Notebook interface in Google Colab. The code cell contains the following command:

```
!pip uninstall torch torchvision torchaudio
```

With the note: "# If using CUDA, install the correct PyTorch package.  
# Replace 'cu118' with your CUDA version from below:  
!pip install torch torchvision torchaudio --index-url https://download.pytorch.org/whl/cu118"

The terminal output shows the process of uninstalling these packages. It starts by listing the packages to be removed, then asks for confirmation ("Proceed (Y/n)? Y"). It then lists the files being uninstalled, followed by another confirmation ("Proceed (Y/n)? Y"). Finally, it lists the files being removed again.

```
Proceed (Y/n)? Y
Successfully uninstalled torch-2.6.0+cu124
Found existing installation: torchvision 0.21.0+cu124
Uninstalling torchvision-0.21.0+cu124:
Would remove:
/usr/local/lib/python3.11/dist-packages/torchvision/*
/usr/local/lib/python3.11/dist-packages/torchvision.libs/libcudart.41118559.so.12
/usr/local/lib/python3.11/dist-packages/torchvision.libs/libjpeg.1c1cb009.so.8
/usr/local/lib/python3.11/dist-packages/torchvision.libs/libnvjpeg.02bd700.so.12
/usr/local/lib/python3.11/dist-packages/torchvision.libs/libpng16.0364a1db.so.16
/usr/local/lib/python3.11/dist-packages/torchvision.libs/libsharpyuv.5c41a003.so.0
/usr/local/lib/python3.11/dist-packages/torchvision.libs/libwebrtc.54aa0d2a.so.7
/usr/local/lib/python3.11/dist-packages/torchvision.libs/libz.d13a2644.so.1
/usr/local/lib/python3.11/dist-packages/torchvision/*
Proceed (Y/n)? Y
Successfully uninstalled torchvision-0.21.0+cu124
Found existing installation: torchaudio 2.6.0+cu124
Uninstalling torchaudio-2.6.0+cu124:
Would remove:
/usr/local/lib/python3.11/dist-packages/torchaudio-2.6.0+cu124.dist-info/*
/usr/local/lib/python3.11/dist-packages/torchaudio/*
/usr/local/lib/python3.11/dist-packages/torio/*
Proceed (Y/n)? Y
```

The screenshot shows a Google Colab notebook interface. The title bar indicates the file is 'QASearch\_adta5770\_coding\_lectures.ipynb'. The notebook content includes:

- A code cell with the following Python code:

```
[ ] # Automatically restart kernel after installs so that your environment can access the new packages
import IPython

app = IPython.Application.instance()
app.kernel.do_shutdown(True)

{'status': 'ok', 'restart': True}
```
- A section titled "UNINSTALL UNCOMPATIBLE PYTORCH-RELATED MODULES --> CANNOT LOADING PDF FILES INTO LOADER".
- A section titled "REINSTALL COMPATIBLE PYTORCH-RELATED MODULES".
- A section titled "VIP NOTES" with the following bullet points:
  - ) MUST MANUALLY ENTER 'Y' TO ANSWER EACH PROMPT
    - ) First, to manually enter 'Y', use mouse to click next to the right of the prompt "Proceed (Y/n)?"
    - ) A text field will show up
    - ) Type into the text field 'Y'
    - ) Press ENTER

```
[ ] # Automatically restart kernel after installs so that your environment can access the new packages
import IPython

app = IPython.Application.instance()
app.kernel.do_shutdown(True)

{'status': 'ok', 'restart': True}

# Install: Utilities libraries needed for Q&A Search

# Dependencies required by Unstructured PDF loader
! sudo apt -y -qq install tesseract-ocr libtesseract-dev
! sudo apt-get -y -qq install poppler-utils

#pip install --user --upgrade unstructured pdf2image pytesseract pdfminer.six
#pip install --user --upgrade pillow-heif opencv-python unstructured-inference pikepdf pypdf

# For Matching Engine (NEW NAME: Vector Search) integration dependencies (default embeddings)
#pip install --user --upgrade tensorflow_hub tensorflow_text

# For loading PDF documents in Google Cloud Storage into LangChain Loader
#pip install --user --upgrade pi-heif

The following additional packages will be installed:
libarchive-dev liblibleptonica-dev tesseract-ocr-eng tesseract-ocr-osd
The following NEW packages will be installed:
libarchive-dev liblibleptonica-dev libtesseract-dev tesseract-ocr-
tesseract-ocr-eng tesseract-ocr-osd
0 upgraded, 6 newly installed, 0 to remove and 30 not upgraded.
Need to get 8,560 kB of archives.
After this operation, 31.6 MB of additional disk space will be used.
```