# *Project Progress Report: Information Retrieval using Large Language Models through Question Answering for Machine Learning and Artificial Intelligence Research Papers*

*( CS 410:  Text Information Systems , Fall 2023)*

*(Team: Lexical Wizard)*

| Team : Lexical Wizards | | | |
|---|---|---|---|
| # | Name | ID | Email |
| 1 | Kumar, Amrit | amritk2 | amritk2@illinois.edu |
| 2. | Gattu, Sudha Mrudula | sudhamg2 | sudhamg2@illinois.edu |
| 3. | Madhavan, Siddharth | sm120 | sm120@illinois.edu |
| 4. | Awachat, Ishika | awachat2 | awachat2@illinois.edu |

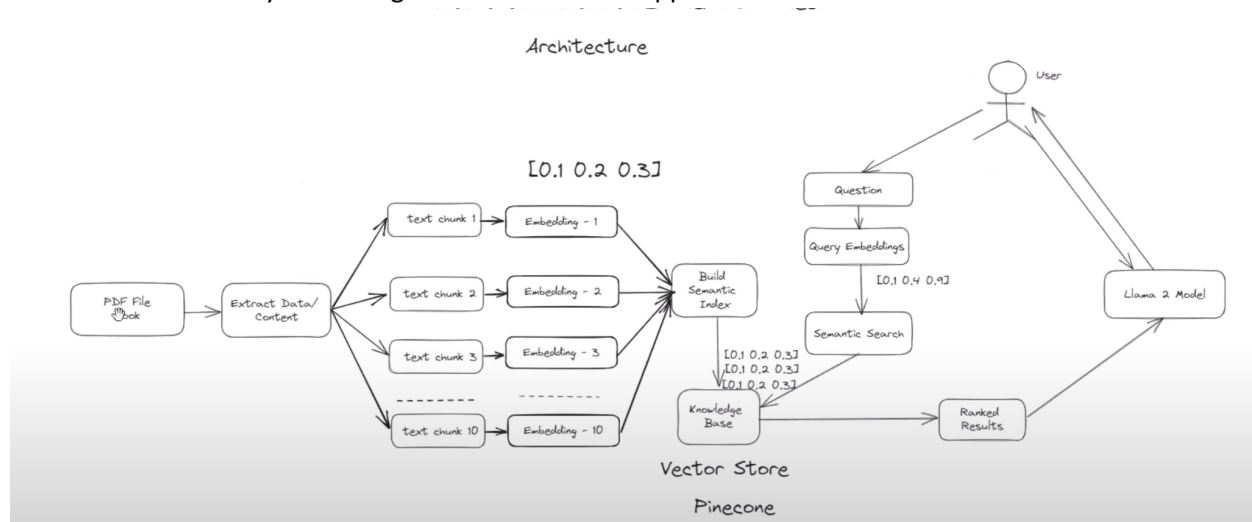| Update History | | | | |
|---|---|---|---|---|
| # | Version | Date | Author | Update Details |
| 1 | draft | 2023-11-18 | Amrit Kumar | Draft Document |
| 2 | draft | 2023-11-18 | Siddharth Madhavan | Text embedding section |
| 3 | draft | 2023-11-19 | Amrit Kumar | System Architecture and Virtual Store Setup |
| 4 | draft | 2023-11-19 | Sudha Mrudula | Web Scraping Section |
| 5 | draft | 2023-11-19 | Amrit Kumar, Ishika Awachat | Interactive Chat window section |
| 6 | draft | 2023-11-19 | Amrit Kumar | Summary section and Overall progress |
| 7 | v1 | 2023-11-19 | Siddharth Madhavan | Edits and formatting |

## Table of Contents

## Introduction:

The rapid proliferation of research papers in the fields of Machine Learning and Artificial Intelligence has led to a substantial increase in the volume of textual information. This surge poses a significant challenge for knowledge seekers who wish to access specific information within this vast corpus. To address this issue, we propose a comprehensive project titled "Information Retrieval using Large Language Models through Question Answering." This project aims to leverage the power of modern NLP models to extract, index, and retrieve information from research papers in the domain of Machine Learning and Artificial Intelligence.

The aim of this document is to share the progress on this project implementation. This project has following main modules,

- System Design/ Architecture
- Web scraping:
- Text Chunk generation/ embedding
- Storing embedded data in Index Vector store.
- Information Retrieval and chat conversation using Large Language Model
- Interactive Chat interface

## System Design/Architecture

- This Large Language Model powered application is intended to perform information retrieval from the loaded document/web and provide the personalized answer of user questions through a chat window .
- We are intending to use LLAMA2 (large language model) for the implementation .
- Below is the system design for **LexicalWizard** application



Architecture

- **Flow**:
  - **Preprocessing**
    - The data (PDF/ Web Data) for a specific domain will be loaded
    - Data from the source will be extracted and splitted.
    - Splitted data (chunks) will be converted into **embeddings** .

- Embeddings will be saved into a Indexed Vector Database i.e. Pinecone
- Indexed Vector Store will be acting as a "Knowledge Base" to answer user queries.
  - **Answering User queries**:
    - When user posts his/her question, embeddings will be created for the provided query
    - Similarity search will be performed in the Pinecone vector store using the query embeddings.
    - Similarity Search will return **TOP RANKED** chunks from the index vector store.
    - The r**anked results** will be sent to the **LLAMA2** model along with the user/system prompts to answer the asked queries
    - Answers from the LLAMA2 model will be displayed in the chat window as an answer.
    - For further questions, previous chat conversations will be passed as context to continue the conversation.

## Web Scraping:

1. Which tasks have been completed?
   a. Created a baseline template in which the program accepts a user query and conducts searches on the DuckDuckGo search engine to identify potential and relevant websites that address the query.
   b. Integrated the DDGS (DuckDuckGoSearch) api to retrieve search results with a given query
   c. After retrieving the relevant website for the given query, the program saves the search results into a text file
2. Which tasks are pending?
   a. Rather than fetching websites and urls, implementing the feature to where the program also retrieves PDF documents
   b. Integrating the url lib package to read and parse the urls
   c. Upon parsing the URLs, the program retrieves the most optimal answer for the given query.
3. Any challenges?
   a. Identifying the most suitable website/PDF document from the search results

## Text Chunk generation/ embedding

1. Which tasks have been completed?
   a. We have integrated hugging face api to get text embeddings
   b. We have successfully split a pdf into separate lines and have used chuck_overlap to retain context of that paragraph
   c. We are able to get some questions answered based on a sample hugging face Question and Answer module (proof of concept)
2. Which tasks are pending?
   a. Incorporating the logic to take on multiple pdfs
3. Any challenges?
   a. Improving text embeddings to capture the context for multiple pdfs
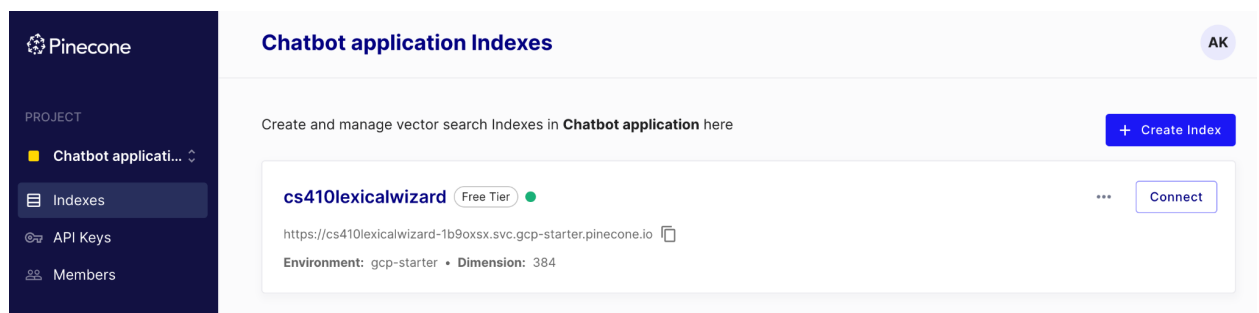
## Storing the text embedding in index database

We are using Pinecone Vector Database to store the vector embeddings.
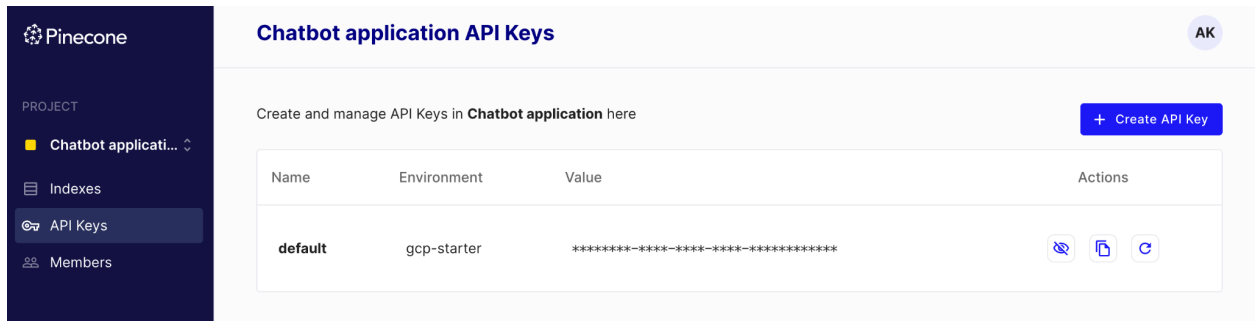
### Why Vector Database

- Applications that involve large language models, generative AI, and semantic search rely on vector embeddings, a type of data that represents semantic information.
- This information allows AI applications to gain understanding and maintain a long-term memory that they can draw upon when executing complex tasks.
- Vector databases like Pinecone offer optimized storage and querying capabilities for embeddings.
- Reference: https://docs.pinecone.io/docs/overview

### Configuring Pinecone Database

- Pinecone vector database will store embeddings for knowledge source chunk embeddings and the embedding of the user questions.
- Setup a index database in Pinecone io: https://app.pinecone.io/



- Generate/Copy the `PINECONE_API_KEY` , `PINECONE_API_ENV` values, this will be needed for connecting to the pinecone vector database programmatically.

## Using Pinecone Vector database programmatically

- Package to be installed : "pinecone-client"

- Load Data

```
64  # #**Step 3: Load the Data**
65  loader = PyPDFLoader("content/The-Field-Guide-to-Data-Science.pdf")
66
67  data = loader.load()
```

- Split the text into chunks

```
76  # #**Step 4: Split the Text into Chunks**
77  text_splitter=RecursiveCharacterTextSplitter(chunk_size=500, chunk_overlap=0)
78
79  docs=text_splitter.split_documents(data)
```

- Setup the PineCone environment

```
92  # #**Step 5: Setup the Environment**
93  os.environ["HUGGINGFACEHUB_API_TOKEN"] = "********************"
94  PINECONE_API_KEY = os.environ.get('PINECONE_API_KEY', '********************')
95  PINECONE_API_ENV = os.environ.get('PINECONE_API_ENV', 'gcp-starter')
```

- Download the embeddings

```
99   # #**Step 6: Downlaod the Embeddings**
100
101
102  embeddings=HuggingFaceEmbeddings(model_name='sentence-transformers/all-MiniLM-L6-v2')
103
```

- Initialize the Pinecone Vector Database

```
105     # #**Step 7: Initializing the Pinecone**
106
107     pinecone.init(
108         api_key=PINECONE_API_KEY,  # find at app.pinecone.io
109         environment=PINECONE_API_ENV  # next to api key in console
110     )
111     index_name = "cs410lexicalwizard" # put in the name of your pinecone index here
112
```

- Create embeddings for the chunks created above

```
114     # #**Step 8: Create Embeddings for Each of the Text Chunk**
115     docsearch=Pinecone.from_texts([t.page_content for t in docs], embeddings, index_name=index_name)
116
117     # # If pinecone index are already present, it can be resused as below
118
119     #docsearch = Pinecone.from_existing_index(index_name, embeddings)
120
```

- Perform Similarity Search

```
122     # #**Step 9: Similarity Search**
123     query="YOLOv7 outperforms which models"  # Sample query to be searched in the loaded doc/PDF
124
125     docs=docsearch.similarity_search(query)
126     print(f"docs :{docs}")
127
```

**Console Output:**

```
LLAMA2_CS410_Project
  docs length: 383
  docs 0: page_content='THE FIELD  GUIDE   \n        to   DATA  SCIENCE\n© COPYRIGHT 2013 BOOZ ALLEN HAMILTON INC. ALL RIGHTS RESERVED.' metadata={'source': 'content/
  docs retrieved after Similarity Search:
  [Document(page_content='Online models have both advantages and disadvantages. /T_hey \ndynamically evolve over time, meaning they only require a single \ndeployment in
```

- Verify the index database if embeddings are loaded or not in Pinecone portal

## indexes from the loaded data



## Usage Metrics

**Chatbot application Indexes**                                                    AK

| METRIC | DIMENSIONS | POD TYPE | HOST |
|--------|-----------|----------|------|
| cosine | 384 | starter | https://cs410lexicalwizard-1b9oxsx.svc.gcp-starter.pinecone.io |

| PROVIDER | REGION | ENVIRONMENT | VECTOR COUNT |
|----------|--------|-------------|--------------|
| ☁ GCP | 🇺🇸 Iowa (us-central1) | gcp-starter | 1149 |

BROWSER    METRICS

Past  5 min  15 min  **4 hours**  12 hours  1 day  2 days  1 week

REQUESTS ⓘ                                          Operation   query  upsert  update  delete  fetch

```
0.05/sec
0.04/sec
0.03/sec                                                            11/19/23
0.02/sec                                                            10:56:15 AM
0.01/sec                                                            | UPSERT    0.043 /SEC
0/sec
     9:27 AM      9:44 AM      9:58 AM     10:13 AM     10:27 AM     10:41 AM     10:56 AM
```

● QUERY   ● UPSERT   ● UPDATE   ● DELETE   ● FETCH

## Information Retrieval and chat conversation using Large Language Model

This module's development is in progress.

```python
143    from langchain.llms import LlamaCpp
144    from langchain.callbacks.manager import CallbackManager
145    from langchain.callbacks.streaming_stdout import StreamingStdOutCallbackHandler
146    from huggingface_hub import hf_hub_download
147    from langchain.chains.question_answering import load_qa_chain
```

```python
177    model_name_or_path = "TheBloke/Llama-2-13B-chat-GGML"
178    model_basename = "llama-2-13b-chat.ggmlv3.q5_1.bin"  # the model is in bin format
179
180    model_path = hf_hub_download(repo_id=model_name_or_path, filename=model_basename)
181
```

```
190     n_gpu_layers = 40
191     n_batch = 256
192     # Loading model,
193     llm = LlamaCpp(
194         model_path=model_path,
195         max_tokens=256,
196         n_gpu_layers=n_gpu_layers,
197         n_batch=n_batch,
198         callback_manager=callback_manager,
199         n_ctx=1024,
200         verbose=False,
201     )
202
203     chain=load_qa_chain(llm, chain_type="stuff")
204
205     query="YOLOv7 outperforms which models"
206     docs=docsearch.similarity_search(query)
207     chain.run(input_documents=docs, question=query)
```

1. Which tasks have been completed?:
   a. Using LLAMA2 implement the Information Retrieval module.
   b. LLAMA2 is one of the most popular open source large language model implementation provided by **META**
   c. Installed LLAMA2 as per the instructions https://ai.meta.com/llama/
   d. Sample module has been implemented , code testing is in progress.
2. Which tasks are pending?
   a. As per the current progress, getting exceptions while loading the LlamaCpp model in local, which we are working on.
   b. We are working on resolving the issue.

```
docs retrieved after Similarity Search:
 [Document(page_content='Online models have both advantages and disadvantages. /T_hey \ndynamically evolve over time, meaning they only require a single \ndeployment
Traceback (most recent call last):
  File "/Users/amritpandey/Desktop/Study/workspace/MSDS/CS410_Text_Processing_System_Fall_2023/Project/CS410TeamProject/lexical_wizard_cs410_project/LLAMA2_CS410_Pr
    llm = LlamaCpp(
  File "/Users/amritpandey/Desktop/Study/workspace/MSDS/CS410_Text_Processing_System_Fall_2023/Project/CS410TeamProject/lexical_wizard_cs410_project/venv/lib/python
    super().__init__(**kwargs)
  File "/Users/amritpandey/Desktop/Study/workspace/MSDS/CS410_Text_Processing_System_Fall_2023/Project/CS410TeamProject/lexical_wizard_cs410_project/venv/lib/python
    raise validation_error
pydantic.v1.error_wrappers.ValidationError: 1 validation error for LlamaCpp
__root__
  Could not load Llama model from path: /Users/amritpandey/.cache/huggingface/hub/models--TheBloke--Llama-2-13B-chat-GGML/snapshots/3140827b4dfcb6b562cd87ee3d7f0710
Exception ignored in: <function _LlamaModel.__del__ at 0x1860c1a60>
Traceback (most recent call last):
  File "/Users/amritpandey/Desktop/Study/workspace/MSDS/CS410_Text_Processing_System_Fall_2023/Project/CS410TeamProject/lexical_wizard_cs410_project/venv/lib/python
  File "/Users/amritpandey/Desktop/Study/workspace/MSDS/CS410_Text_Processing_System_Fall_2023/Project/CS410TeamProject/lexical_wizard_cs410_project/venv/lib/python
RuntimeError: could not find io module state (interpreter shutdown?)
Exception ignored in: <function _LlamaContext.__del__ at 0x1860d39d0>
```

3. Any challenges?
   a. No blocker as of now.



```
190     n_gpu_layers = 40
191     n_batch = 256
192     # Loading model,
193     llm = LlamaCpp(
194         model_path=model_path,
195         max_tokens=256,
196         n_gpu_layers=n_gpu_layers,
197         n_batch=n_batch,
198         callback_manager=callback_manager,
199         n_ctx=1024,
200         verbose=False,
201     )
202
203     chain=load_qa_chain(llm, chain_type="stuff")
204
205     query="YOLOv7 outperforms which models"
206     docs=docsearch.similarity_search(query)
207     chain.run(input_documents=docs, question=query)
```

1. Which tasks have been completed?:
   a. Using LLAMA2 implement the Information Retrieval module.
   b. LLAMA2 is one of the most popular open source large language model implementation provided by **META**
   c. Installed LLAMA2 as per the instructions https://ai.meta.com/llama/
   d. Sample module has been implemented , code testing is in progress.
2. Which tasks are pending?
   a. As per the current progress, getting exceptions while loading the LlamaCpp model in local, which we are working on.
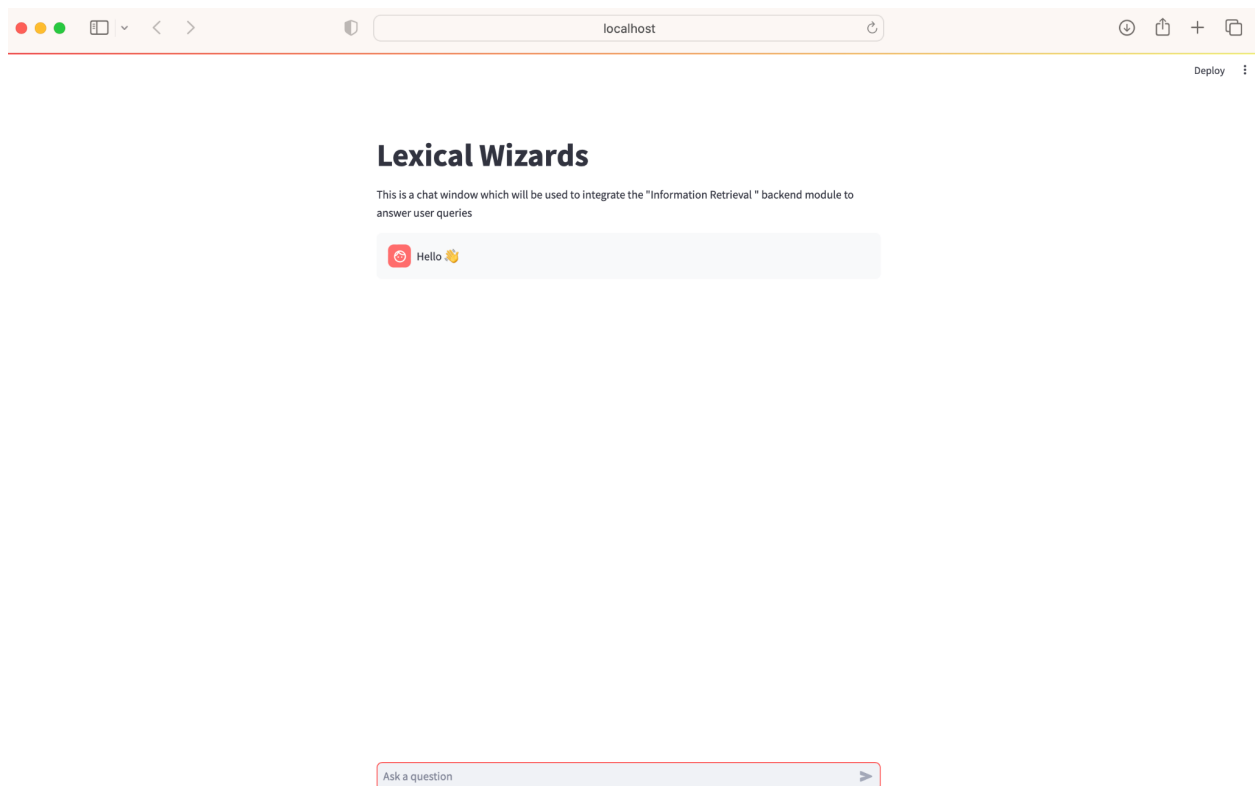   b. We are working on resolving the issue.

```
docs retrieved after Similarity Search:
 [Document(page_content='Online models have both advantages and disadvantages. /T_hey \ndynamically evolve over time, meaning they only require a single \ndeploymen
Traceback (most recent call last):
  File "/Users/amritpandey/Desktop/Study/workspace/MSDS/CS410_Text_Processing_System_Fall_2023/Project/CS410TeamProject/lexical_wizard_cs410_project/LLAMA2_CS410_Pr
    llm = LlamaCpp(
  File "/Users/amritpandey/Desktop/Study/workspace/MSDS/CS410_Text_Processing_System_Fall_2023/Project/CS410TeamProject/lexical_wizard_cs410_project/venv/lib/python
    super().__init__(**kwargs)
  File "/Users/amritpandey/Desktop/Study/workspace/MSDS/CS410_Text_Processing_System_Fall_2023/Project/CS410TeamProject/lexical_wizard_cs410_project/venv/lib/python
    raise validation_error
pydantic.v1.error_wrappers.ValidationError: 1 validation error for LlamaCpp
__root__
  Could not load Llama model from path: /Users/amritpandey/.cache/huggingface/hub/models--TheBloke--Llama-2-13B-chat-GGML/snapshots/3140827b4dfcb6b562cd87ee3d7f0710
Exception ignored in: <function _LlamaModel.__del__ at 0x1860c1a60>
Traceback (most recent call last):
  File "/Users/amritpandey/Desktop/Study/workspace/MSDS/CS410_Text_Processing_System_Fall_2023/Project/CS410TeamProject/lexical_wizard_cs410_project/venv/lib/python
  File "/Users/amritpandey/Desktop/Study/workspace/MSDS/CS410_Text_Processing_System_Fall_2023/Project/CS410TeamProject/lexical_wizard_cs410_project/venv/lib/python
RuntimeError: could not find io module state (interpreter shutdown?)
Exception ignored in: <function _LlamaContext.__del__ at 0x1860d39d0>
```

3. Any challenges?
   a. No blocker as of now.

## Interactive Chat interface

1. Which tasks have been completed?
    a. We have been able to implement a simple chat window using Streamlit in a python script.
    b. This window will be act as an user interface where an user can post domain specific questions and the application will post relevant answer based on Information retrieval using LLAMA2 (large language model)
    c. Tested with posting simple texts as shown in the screenshot below, the window works as expected.
2. Which tasks are pending?
    a. Integration with the backend Information retrieval system (LLAMA2 powered) API
        i. Specifically sending user input to retrieval system and then displaying the fetched information for user
    b. Conversation context: to remember previous conversations.
3. Any challenges?
    a. No challenges so far, besides potential challenges involving integration.

## Summary:

Overall , our progress is good so far. We have been able to finish following items:
- System Design
- Tech stack: Here are few important packages, libraries
  - Large Language Model: LLAMA2
  - DuckDuckGo search: for web scraping
  - Vector Store Database: Pinecone
  - UI Interface: Streamlit
  - Other important libraries/APIs: Hugging Face embeddings, langchain, pyloader etc
- **Proof of Concepts**: We have been able to implement POC to ensure different modules are implementable and narrow down the tech stack system design. Here are the proof of concepts that we have worked on:
  - **Text chunk generation/ embeddings**:
    - Able to implement a template code for loading the PDF , converting into embeddings and implementing a QA module, which answers user queries based on the user query.
    - This was a very simple implementation to prove the viability of the project goal.

  - **Web Scraping**:
    - Able to programmatically scrape information from the web pages using DuckDuckGo search
    - The scrapped information will be embedded in the vector store to enhance the knowledge base and provide rich context for the system to provide answers to user queries in addition to the context created through the PDF load.
  - **Vector Store Database Using Pinecone**:
    - Successfully able to set up the Pinecone vector database to store vector embeddings as indexes for quick retrieval and persistent memory.
    - Successfully implemented the code to perform:
      - saving logic for embedded chunks
      - Ranking logic: Similarity Search
  - **Chat window**:
    - Chatbot user interface for user conversation with the system.
- Pending Items:
  - Proof of concept:
    - Chat conversation LLAMA 2 implementation (IN PROGRESS currently)
  - Integrating different modules together
  - Code cleanup and reviews