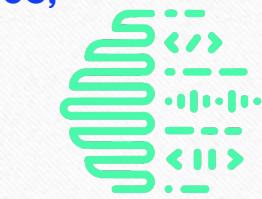




TEAM DETAILS AND PROBLEM STATEMENT



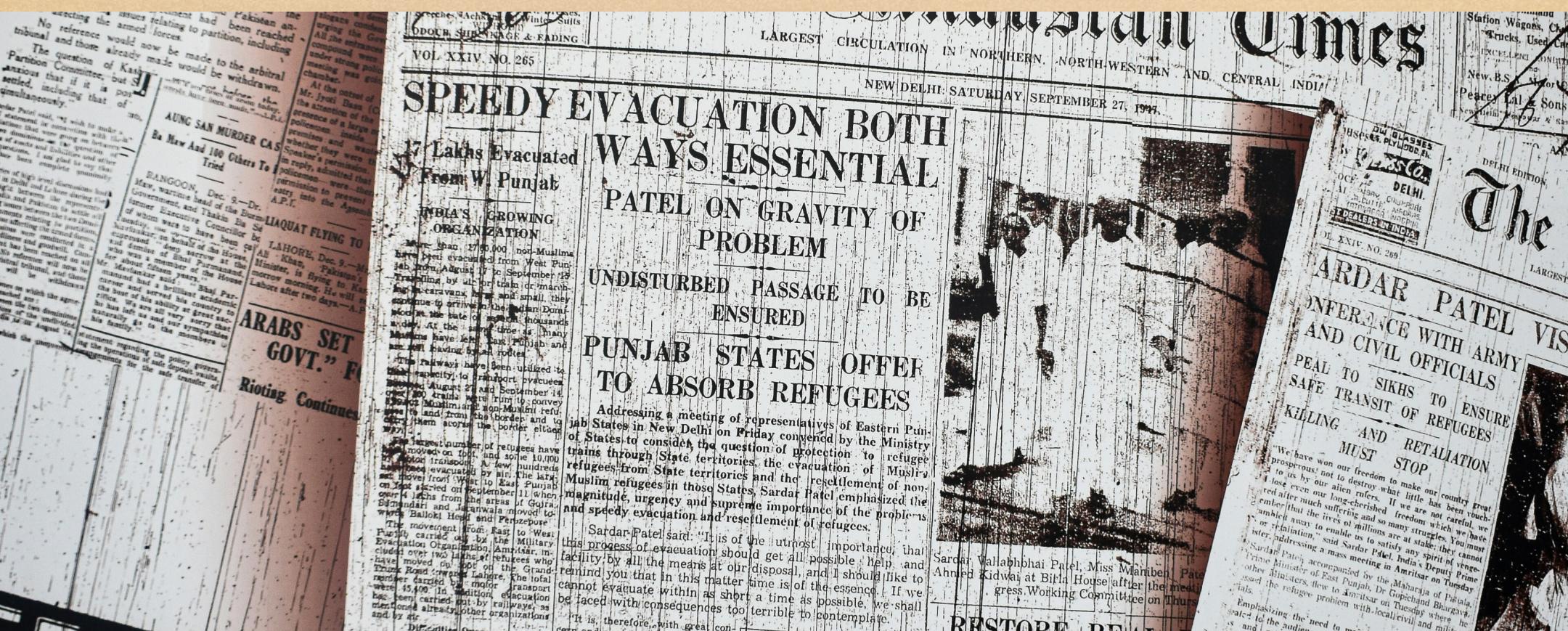
- Problem Statement : Global Article Summarization and Link Analysis Platform
- Team Name : Resc.ai
- Team Leader Details (Name, Phone Number, Email) : Devesh mishra,+916306962393, mishradevesh5518@gmail.com
- Institution Name : Bennett university
- Course Enrolled : Btech cse



OPEN HACK
SHAPE TOMORROW

Problem statement

- Develop an interactive platform using advanced language processing to filter and summarize articles based on user inputs.
 - Users struggle to extract insights from the vast amount of available information, necessitating concise, relevant summaries.
 - The platform should be customizable and analyze relationships between articles to enhance understanding.



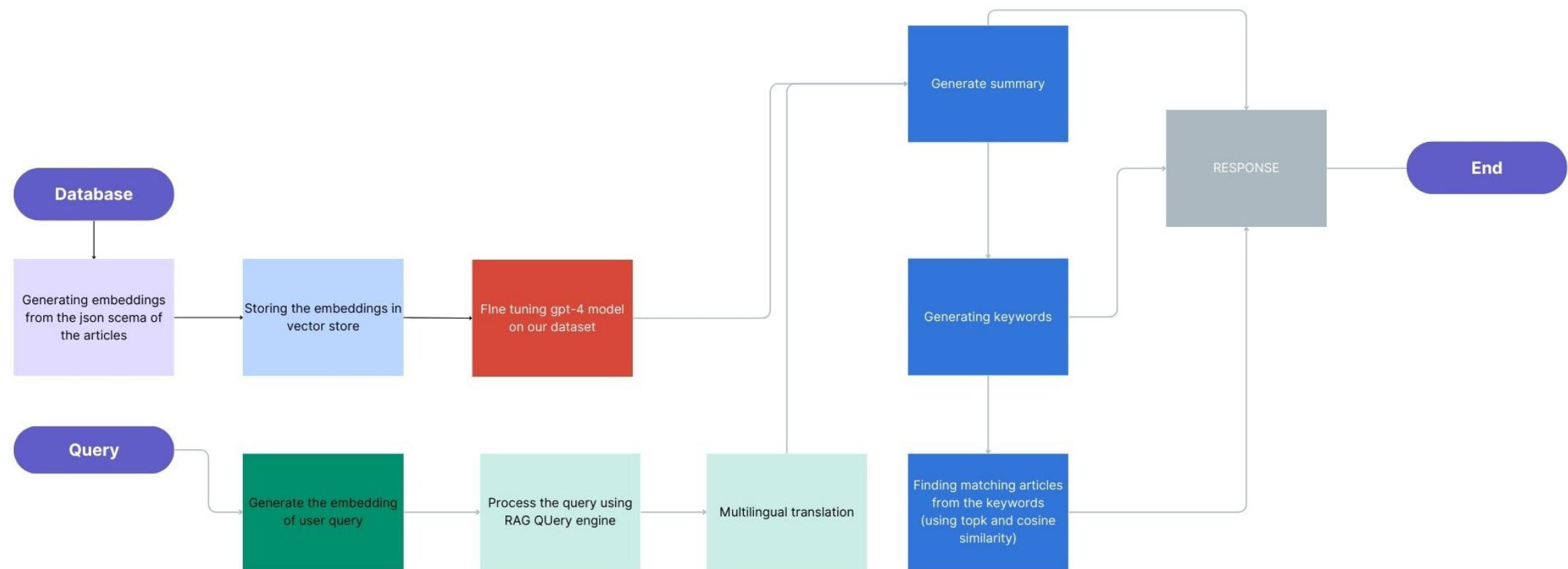
Outline the goals and objective of your approach and what you are aiming to achieve.

- 1. Develop a newspaper article summarizer using a custom dataset and fine-tuned models.**
- 2. Integrate LlamaIndex for efficient article indexing and retrieval.**
- 3. Utilize GPT-4 embeddings to enhance the summarization process.**
- 4. Ensure the summarizer produces concise and relevant summaries across diverse topics.**
- 5. Evaluate summarizer performance based on accuracy, coherence, and relevance metrics.**
- 6. Implement a user-friendly interface for inputting queries and accessing article summaries.**
- 7. Enable the summarizer to suggest related articles based on summary content, promoting user engagement.**

Describe your Methodology

1. **Data Collection and Preprocessing:** Gather diverse newspaper articles and preprocess them to remove noise and standardize formatting.
2. **Model Selection:** Choose GPT-4 as the language model for fine-tuning due to its advanced language processing capabilities.
3. **Fine-tuning:** Adapt the selected language model to the task of article summarization taking reference from news-article-dataset.json then creating our own dataset of 10k+ news articles for better description
4. **LlamaIndex Integration:** Integrate LlamaIndex for efficient indexing and recursive retrieval with response synthesizer using topk similarity for better finetune result.
5. **Embedding Generation:** Generate embeddings for articles using GPT-4-0613 7B model, enhancing summarization quality and semantic understanding.
6. **Evaluation and User Interface:** Assess summarizer performance using relevant metrics and develop a user-friendly interface for easy access to multilingual article summaries and related suggestions.

Relevant image related to your idea.



Mention your technology stack

1. **LLAMA Index**
 - **VectorStoreIndex**
 - **CustomQueryEngine**
2. **OpenAI GPT-4**
3. **Langchain**
4. **Google Translator API**
5. **deep_translator**
6. **rouge_score**
 - **NLTK (Natural Language Toolkit)**

```

main.py | 43     record = json.loads(line)
        | 44
        | 45     # Extract link and headline from the record
        | 46     links.append(record['link'])
        | 47     headlines.append(record['headline'])
        |
        | 48     # Print the extracted lists
        | 49     print("Links:", links[0:50])
        | 50     print("Headlines:", headlines[0:50])
        |
        | 51 def create_index():
        | 52
        | 53     llm = OpenAI(temperature=0.1, model="gpt-4")
        | 54     service_context = ServiceContext.from_defaults(llm=llm)
        |
        | 55     documents = SimpleDirectoryReader("./data").load_data()
        | 56     index = VectorStoreIndex.from_documents(documents, service_context=service_context, show_progress=True)
        |
        | 57     index.storage_context.persist(persist_dir="./content/Emed2")
        |
        | 58 def processing_without_Custom_QE():
        | 59
        | 60     storage_context = StorageContext.from_defaults(persist_dir="./content/drive/MyDrive/Emed2")
        | 61
        | 62     index = load_index_from_storage(storage_context)
        |
        | 63     query_engine = index.as_query_engine(LLM=ChatOpenAI())
        |
        | 64     #TESTCASE 1-
        | 65     response = query_engine.query("give me article about Over 4 Million Americans Roll Up Sleeves For Omicron-Targeted COVID Boosters in 500 words")
        | 66     # response = llm.complete(prompt)
        | 67     # print(response.text)
        | 68     print(response)
        |
        | 69 PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS
        |
        | D:\Resc\client>git remote add origin https://github.com/devesh-mishra13/Resc_client.git
        | error: remote origin already exists.
        |
        | D:\Resc\client>push -u origin main
        | Enumerating objects: 5, done.
        | Counting objects: 100% (5/5), done.
        | Delta compression using up to 16 threads
        | Compressing objects: 100% (3/3), done.
        | Writing objects: 100% (3/3), 307 bytes | 307.00 KiB/s, done.
        | Total 3 (delta 2), reused 0 (delta 0), pack-reused 0
        | remote: Resolving deltas: 100% (2/2), completed with 2 local objects.
        | To https://github.com/devesh-mishra13/prohun.git
        | d50ae5c...50ae5c main > main
        | branch 'main' set up to track 'origin/main'.
        |
        | D:\Resc\client>

```



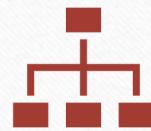
```

main.py | 133     def extract_news_desc():
        | 134         record = json.loads(line)
        | 135
        | 136         # Extract link and headline from the record
        | 137         links.append(record['link'])
        | 138         headlines.append(record['headline'])
        |
        | 139         # Print the extracted lists
        | 140         print("Links:", links[0:50])
        | 141         print("Headlines:", headlines[0:50])
        |
        | 142     class RAGStringQueryEngine(CustomQueryEngine):
        | 143         """RAG String Query Engine."""
        |
        | 144         retriever: BaseRetriever
        | 145         response_synthesizer: BaseSynthesizer
        | 146         llm: OpenAI
        | 147         q_prompt: PromptTemplate
        |
        | 148         def custom_query(self, query_str: str):
        | 149             nodes = self.retriever.retrieve(query_str)
        |
        | 150             context_str = "\n\n".join([n.node.get_content() for n in nodes])
        | 151             response = self.llm.complete(
        | 152                 qa_prompt.format(context_str=context_str, query_str=query_str, n_words=n_words)
        | 153             )
        |
        | 154             return str(response)
        |
        | 155         retriever = VectorIndexRetriever(
        | 156             index=index,
        | 157             similarity_top_k=1,
        | 158         )
        |
        | 159         # configure response synthesizer
        | 160         response_synthesizer = get_response_synthesizer(
        | 161             response_mode="tree_summarize",
        | 162         )
        |
        | 163         llm = OpenAI(model="gpt-3.5-turbo")
        |
        | 164         query_engine = RAGStringQueryEngine(
        | 165             retriever=retriever,
        | 166         )
        |
        | 167 PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS
        |
        | D:\Resc\client>git remote add origin https://github.com/devesh-mishra13/Resc_client.git
        | error: remote origin already exists.
        |
        | D:\Resc\client>push -u origin main
        | Enumerating objects: 5, done.
        | Counting objects: 100% (5/5), done.
        | Delta compression using up to 16 threads
        | Compressing objects: 100% (3/3), done.
        | Writing objects: 100% (3/3), 307 bytes | 307.00 KiB/s, done.
        | Total 3 (delta 2), reused 0 (delta 0), pack-reused 0
        | remote: Resolving deltas: 100% (2/2), completed with 2 local objects.
        | To https://github.com/devesh-mishra13/prohun.git
        | d50ae5c...50ae5c main > main
        | branch 'main' set up to track 'origin/main'.
        |
        | D:\Resc\client>

```

GitHub backend link
with summaries
GitHub frontend link

- <https://github.com/devesh-mishra13/Resc>
https://github.com/devesh-mishra13/Resc_client



Describe your use cases.

1. Article Summarization
2. Keyword Extraction
3. Related Article Suggestions
4. Customized Summarization Length
5. Language Support
6. Interactive User Interface
7. Educational and Research Purposes



Future Work.

1. Incorporating User Feedback Mechanisms for Improved Summaries
2. Integrating Semantic Understanding for Better Keyword Extraction
3. Enhancing Link Analysis for More Relevant Article Suggestions
4. Expanding Language Support to Include More Languages
5. Developing Mobile Applications for Accessibility
6. Exploring Integration with Voice Assistants for Hands-free Interaction