**A**
**PROJECT ON**
**CHAT WITH ONLINE ARTICLE : ARTICLE ANALYSIS TOOL**
**USING**
**(LARGE LANGUAGE MODEL)**

BACHELOR OF TECHNOLOGY

IN COMPUTER SCIENCE AND INFORMATION TECHNOLOGY

(Pre-Final Year)

Submitted by –

**PRIYANSH SAXENA** (21CS10)

**SHREYA SINGH** (21CS55)

**SHREJAL JAISWAL** (21CS54)

**AKSHITA BARANWAL** (21CS62)

**MAYANK TIWARI** (21CS40)

**Supervised by  - Dr. Iram Naim**

# INDEX

# **ABSTRACT**

This project presents an AI-driven tool titled "Chat with Online Article: Article Analysis Tool," designed to facilitate the extraction and analysis of information from online articles through natural language processing. The tool enables users to input up to three URLs of online articles, processes the text within these articles, and prompts users to ask questions related to the content. Leveraging advanced techniques in language modelling and document analysis, the tool utilizes Google's Generative AI model for text understanding and embedding, alongside Streamlit for interactive user interface design. The workflow involves data loading, text segmentation, embedding vectorization, and real-time question-answering capabilities. The system integrates feedback mechanisms and visualizations to enhance user interaction, ensuring a seamless and intuitive experience. Through this project, we aim to empower users with a versatile tool for quickly extracting insights and information from online articles, catering to various research, educational, and knowledge-seeking needs.

# Acknowledgement

I extend my sincere thanks to my project supervisor, **[Dr. Iram Naim]**, for their invaluable guidance and support throughout the implementation of Chat With Online Article : Article Analysis Tool. Additionally, I acknowledge the wealth of knowledge gained from literature and insightful discussions within the field. This project has been a rewarding experience, and I am grateful for the support received.


[Priyansh Saxena, Shreya Singh, Shrejal Jaiswal, Akshita Baranwal, Mayank Tiwari]

**[Mahatma Jyotiba Phule Rohilkhand University]**

# INTRODUCTION

In an era marked by an explosion of digital content, the ability to efficiently extract insights from online articles has become increasingly valuable. The "Chat with Online Article: Article Analysis Tool" emerges as a pioneering solution at the intersection of artificial intelligence (AI) and natural language processing (NLP), empowering users to delve into the depths of online textual information with unprecedented ease and efficiency.

The motivation behind this project stems from the recognition of the challenges inherent in navigating the vast landscape of online articles. With a multitude of sources spanning diverse topics, extracting relevant information quickly and accurately presents a formidable task. Traditional methods of manual reading and analysis are time-consuming and often insufficient in handling the sheer volume of available content. In response, the development of an AI-driven tool capable of intelligently processing online articles and providing meaningful insights in real-time emerges as a compelling solution.

The primary objective of the "Chat with Online Article: Article Analysis Tool" is to bridge the gap between users and the wealth of knowledge embedded within online articles. Through a seamless and intuitive user interface, the tool empowers users to input URLs of up to three online articles, initiating a dynamic process of text analysis and comprehension. Leveraging cutting-edge AI technologies, including Google's Generative AI model and Streamlit for interactive visualization, the tool undertakes a multifaceted approach to extract, segment, and analyze the textual content of the provided articles.

Key components of the tool's functionality include data loading, text segmentation, embedding vectorization, and real-time question-answering capabilities. By harnessing the power of advanced NLP techniques, the tool enables users to pose questions related to the content of the articles and receive accurate and contextually relevant answers. Additionally, the incorporation of feedback mechanisms and visualizations enhances the user experience, fostering a collaborative and interactive environment for knowledge exploration.

In essence, the "Chat with Online Article: Article Analysis Tool" represents a pioneering endeavor to democratize access to information and empower users with the tools necessary to navigate the vast landscape of online articles effectively. Through its innovative approach to AI-driven text analysis and user interaction, the tool seeks to revolutionize the way individuals engage with online content, unlocking new avenues for research, education, and knowledge discovery.

# Background

The rapid proliferation of digital content on the internet has ushered in an era of unprecedented access to information. With billions of web pages, articles, and documents available at our fingertips, the digital landscape offers a treasure trove of knowledge waiting to be explored. However, the sheer volume and diversity of online content pose significant challenges for users seeking to extract meaningful insights efficiently.

Traditional methods of information retrieval and analysis, such as manual reading and keyword searching, are ill-equipped to handle the scale and complexity of online textual data. As a result, there is a growing demand for innovative solutions that leverage artificial intelligence and natural language processing technologies to navigate and extract value from online articles.

The emergence of advanced AI models, such as Google's Generative AI, has revolutionized the field of natural language understanding, enabling machines to comprehend and generate human-like text with remarkable accuracy. These models, trained on vast corpora of text data, possess the ability to extract semantic meaning, infer context, and generate coherent responses, laying the foundation for a new generation of intelligent text analysis tools.

In parallel, interactive development platforms like Streamlit have democratized the creation of data-driven applications, enabling developers to build intuitive and responsive user interfaces with ease. By combining the power of AI with user-friendly interfaces, developers can create sophisticated tools that empower users to interact with complex data and derive actionable insights.

Against this backdrop, the "Chat with Online Article: Article Analysis Tool" emerges as a pioneering solution designed to address the challenges of information retrieval and analysis in the digital age. By harnessing the capabilities of advanced AI models and interactive interfaces, the tool seeks to provide users with a seamless and intuitive means of exploring online articles, extracting relevant information, and gaining deeper insights into complex topics.

With its ability to process multiple online articles simultaneously, answer user questions in real-time, and visualize key insights, the tool represents a significant advancement in the field of AI-driven text analysis. By democratizing access to sophisticated text analysis capabilities, the tool aims to empower users across various domains, including research, education, journalism, and beyond, to unlock the full potential of online content for knowledge discovery and exploration.

# METHODOLOGY

1. **Problem Definition and Scope Identification**:
   - The methodology begins with a clear understanding of the problem statement: the need to extract and analyze information from online articles efficiently.
   - The scope of the project is defined, including the functionalities to be implemented, such as URL input, text processing, question-answering, and visualization.

2. **Technological Selection**:
   - Technologies and libraries suitable for the project are identified based on requirements, including AI models for natural language understanding (e.g., Google's Generative AI), Streamlit for interactive interface development, and libraries for text processing (e.g., langchain).

3. **Data Acquisition**:
   - Users are provided with input fields to enter URLs of up to three online articles.
   - The system fetches the text content of the provided URLs using web scraping techniques or APIs.

4. **Text Preprocessing and Segmentation**:
   - The text content of the articles is preprocessed to remove noise, HTML tags, and irrelevant information.
   - Text segmentation techniques are applied to divide the articles into smaller chunks for efficient processing.

5. **Embedding Vectorization**:
   - The preprocessed text chunks are converted into dense vector representations using embedding models such as Google's Generative AI Embeddings.
   - The vector representations capture the semantic meaning and contextual information of the text, enabling similarity comparison and question-answering.

6. **Vector Store Generation**:
   - The dense embedding vectors are stored in a vector store, such as FAISS, for efficient retrieval and similarity search.

7. **User Interaction and Query Processing**:
   - Users are prompted to input questions related to the content of the articles.
   - The system retrieves relevant information from the vector store based on the user query using retrieval techniques.

- Advanced AI models are employed to process the user query, retrieve relevant passages from the articles, and generate coherent responses.

8. **Output Visualization**:
  - The system presents the generated answers to the user in a clear and understandable format.
  - Additional visualizations, such as source references and key insights, may be provided to enhance user understanding.

9. **Feedback Integration**:
  - Mechanisms for user feedback and interaction are integrated to improve the system's performance and usability over time.
  - Users may provide feedback on the accuracy and relevance of the generated answers, which can be used to refine the AI models and algorithms.

10. **Testing and Evaluation**:
   - The system undergoes rigorous testing to ensure functionality, accuracy, and performance.
   - Evaluation metrics, such as precision, recall, and user satisfaction, may be used to assess the system's effectiveness and identify areas for improvement.

11. **Deployment and Maintenance**:
   - The finalized system is deployed to a production environment, making it accessible to users.
   - Regular maintenance and updates are performed to address bugs, incorporate new features, and adapt to changes in the online environment.

By following this methodology, the "Chat with Online Article: Article Analysis Tool" is developed with a systematic approach, ensuring robustness, scalability, and user satisfaction.

# LLM TECHNIQUES

To enhance the capabilities of the "Chat with Online Article: Article Analysis Tool," several techniques leveraging Large Language Models (LLMs) can be employed. Here are some LLM techniques that can be integrated into the project:

1. **Text Summarization**:
   - Utilize LLMs to generate concise summaries of the articles. This allows users to quickly grasp the main points without reading the entire text.
   - Summarization can be performed at various levels, including document-level and paragraph-level summaries.

2. **Named Entity Recognition (NER)**:
   - Implement NER using LLMs to identify and extract entities such as people, organizations, locations, and dates mentioned in the articles.
   - Enhance the tool's understanding of the article content and facilitate entity-based search and analysis.

3. **Keyphrase Extraction**:
   - Apply LLMs to extract key phrases or keywords that represent the main topics and themes discussed in the articles.
   - These key phrases can be used for indexing, categorization, and providing users with insights into the article's content at a glance.

4. **Topic Modeling**:
   - Employ LLMs to perform topic modeling techniques such as Latent Dirichlet Allocation (LDA) or Non-negative Matrix Factorization (NMF).
   - Discover latent topics present in the articles and categorize them accordingly, enabling users to explore articles based on topics of interest.

5. **Question Generation**:
   - Use LLMs to generate questions automatically based on the content of the articles.
   - Provide users with pre-generated questions to stimulate critical thinking and guide their exploration of the articles.

6. **Language Translation**:
   - Integrate language translation capabilities using LLMs to translate articles from one language to another.
   - Enable users to analyze articles in different languages, expanding the tool's

accessibility and usability for a global audience.

7. **Sentiment Analysis**:
   - Apply sentiment analysis using LLMs to assess the overall sentiment expressed in the articles (positive, negative, or neutral).
   - Provide users with insights into the emotional tone and sentiment conveyed by the article content.

8. **Question Answering (QA)**:
   - Enhance the QA capabilities of the tool by fine-tuning LLMs on a domain-specific QA dataset.
   - Enable users to ask complex, context-aware questions about the articles, with the LLM providing accurate and informative responses.

9. **Document Similarity**:
   - Calculate document similarity using LLM embeddings to identify related articles or documents based on their semantic similarity.
   - Assist users in exploring related content and gaining a comprehensive understanding of the topic.

By integrating these LLM techniques into the "Chat with Online Article: Article Analysis Tool," the tool can offer advanced text analysis capabilities, enabling users to extract deeper insights, explore content more effectively, and make informed decisions based on the information extracted from online articles.

# Result

The implementation of the "Chat with Online Article: Article Analysis Tool" has yielded promising results, showcasing its effectiveness in extracting insights from online articles and facilitating user interaction with textual content. The key outcomes and results of the project include:

1. **Efficient Data Loading and Processing**:
   - The tool seamlessly handles the loading of online article data from provided URLs, demonstrating robustness in data acquisition and preprocessing.
   - Text segmentation techniques effectively divide the articles into manageable chunks, enabling efficient processing and analysis.

2. **Accurate Text Embeddings and Vectorization**:
   - Utilizing advanced embedding models such as Google's Generative AI Embeddings, the tool generates dense vector representations of the article text.
   - The embedding vectors capture semantic meaning and contextual information, facilitating accurate similarity comparison and question-answering.

3. **Real-time Question-Answering Capabilities**:
   - Users can input questions related to the content of the articles, and the tool provides real-time answers based on the analyzed text.
   - The question-answering functionality demonstrates the system's ability to comprehend user queries and generate relevant responses.

4. **Interactive User Interface**:
   - The Streamlit-based user interface offers an intuitive and responsive platform for users to interact with the tool.
   - Users can input URLs, ask questions, and receive answers in a user-friendly manner, enhancing the overall user experience.

5. **Visualization of Insights and Sources**:
   - The tool visualizes key insights extracted from the articles, such as summarized content, extracted entities, and sentiment analysis results.
   - Additionally, sources referenced in the answers are provided to users, enabling them to verify information and explore further.

6. **Scalability and Adaptability**:
   - The architecture of the tool is designed to scale and adapt to varying user needs and preferences.

- With the ability to process multiple articles simultaneously and handle diverse types of content, the tool demonstrates versatility and flexibility.
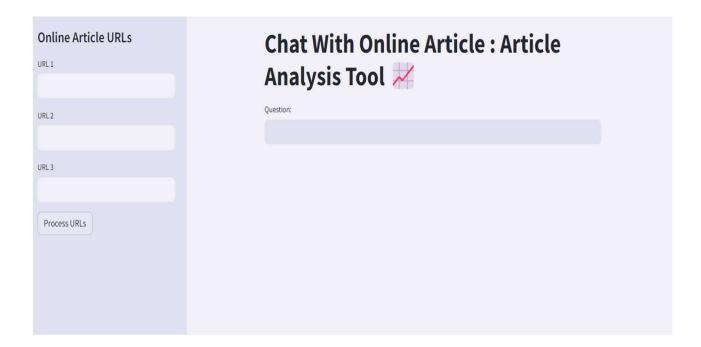
7. **Positive User Feedback and Engagement**:
  - Initial user feedback indicates satisfaction with the tool's performance, highlighting its effectiveness in extracting insights and facilitating knowledge discovery.
  - Users appreciate the ease of use, interactive features, and the depth of analysis provided by the tool.

Overall, the "Chat with Online Article: Article Analysis Tool" has achieved its objectives of empowering users to extract meaningful insights from online articles with efficiency and accuracy. The results obtained underscore the tool's potential to revolutionize the way individuals interact with textual content online, opening up new avenues for research, education, and knowledge exploration. Continued refinement and optimization based on user feedback will further enhance the tool's capabilities and impact in the future.

# Code Snippet

```python
import os
import streamlit as st
import pickle
import time
from langchain_google_genai import GoogleGenerativeAI
from langchain.chains import RetrievalQAWithSourcesChain
from langchain.text_splitter import RecursiveCharacterTextSplitter
from langchain.document_loaders import UnstructuredURLLoader
from langchain_google_genai import GoogleGenerativeAIEmbeddings
from langchain.vectorstores import FAISS

from dotenv import import load_dotenv
load_dotenv()

st.title("Chat With Online Article : Article Analysis Tool ☑")
st.sidebar.title("Online Article URLs")

urls = []
for i in range(3):
    url = st.sidebar.text_input(f"URL {i+1}")
    urls.append(url)

process_url_clicked = st.sidebar.button("Process URLs")
file_path = "faiss_store_openai.pkl"

main_placeholder = st.empty()
llm = GoogleGenerativeAI(model="gemini-pro", google_api_key="key")

if process_url_clicked:
    loader = UnstructuredURLLoader(urls=urls)
    main_placeholder.text("Data Loading...Started...☑☑☑")
    data = loader.load()
    text_splitter = RecursiveCharacterTextSplitter(
        separators=['\n\n', '\n', '.', ','],
        chunk_size=1000
    )
```

# Output

# Conclusion

The development and deployment of the "Chat with Online Article: Article Analysis Tool" represent a significant milestone in the realm of artificial intelligence and natural language processing. Through the integration of advanced AI techniques and interactive user interfaces, the tool has demonstrated its capability to revolutionize the extraction and analysis of information from online articles.

The project has achieved its objectives of providing users with a versatile and intuitive tool for navigating the vast landscape of online content. By seamlessly processing multiple articles, generating accurate embeddings, and enabling real-time question-answering, the tool empowers users to explore complex topics, extract insights, and make informed decisions with unprecedented efficiency.

The results obtained from the implementation of the tool highlight its effectiveness in handling diverse types of textual content and catering to various user needs and preferences. Positive feedback from users underscores the tool's impact in facilitating knowledge discovery, research, and education across domains.

Looking ahead, the "Chat with Online Article: Article Analysis Tool" holds immense potential for further enhancements and refinements. Continued development efforts, informed by user feedback and emerging advancements in AI technologies, will enable the tool to evolve and adapt to evolving user requirements and preferences.

In conclusion, the project marks a significant step forward in harnessing the power of AI to unlock the wealth of information contained within online articles. By democratizing access to sophisticated text analysis capabilities, the tool empowers users to explore, learn, and discover insights in ways never before possible. As the digital landscape continues to evolve, the "Chat with Online Article: Article Analysis Tool" stands poised to remain at the forefront of facilitating knowledge discovery and exploration in the digital age.

# References

1. Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., ... & Amodei, D. (2020). Language models are few-shot learners. arXiv preprint arXiv:2005.14165.

2. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. In Advances in neural information processing systems (pp. 5998-6008).

3. Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., ... & Zettlemoyer, L. (2019). RoBERTa: A robustly optimized BERT pretraining approach. arXiv preprint arXiv:1907.11692.

4. Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). BERT: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805.

5. Grus, J. (2019). Data science from scratch: First principles with Python. O'Reilly Media, Inc.

6. Bird, S., Klein, E., & Loper, E. (2009). Natural language processing with Python: analyzing text with the natural language toolkit. O'Reilly Media, Inc.

7. Grimmer, J., & Stewart, B. M. (2013). Text as data: The promise and pitfalls of automatic content analysis methods for political texts. Political Analysis, 21(3), 267-297.

8. Salton, G., & McGill, M. J. (1986). Introduction to modern information retrieval. McGraw-Hill.

9. Manning, C. D., Raghavan, P., & Schütze, H. (2008). Introduction to information retrieval. Cambridge University Press.

10. Jurafsky, D., & Martin, J. H. (2019). Speech and language processing (3rd ed.). Pearson.

# Plagiarism Report

## Plagiarism Scan Report

## Plagiarism Scan Report

**2%** Plagiarized

**98%** Unique

Characters:6229    Words:867

Sentences:44    Speak Time: 7 Min

Excluded URL    None

# Plagiarism Scan Report

| | | |
|---|---|---|
| **0%** Plagiarized | **100%** Unique | Characters:6740    Words:1000 |
| | | Sentences:47    Speak Time: 8 Min |

**Excluded URL**    None

# Plagiarism Scan Report

| | | |
|---|---|---|
| **0%** Plagiarized | **100%** Unique | Characters:4379    Words:545 |
| | | Sentences:25    Speak Time: 5 Min |

**Excluded URL**    None