# CMPE255 Team Project: Research Paper Summarization

Team Name: Bay Area Rockers
Team Members: Shawn Chumbar, Dhruval Shah, and Sajal Agarwal
Project Title: Research Paper Summarization

The project plan is transcribed below as markdown for ease of reading.

## Project Description

Given the vast number of academic papers published, we aim to provide concise summaries to help researchers quickly grasp the essence of research papers. We will summarize abstracts from a collection of papers on a specific topic to provide an overview of the entire field or detect emerging research areas.

## Background

In recent years, the academic world has witnessed an unprecedented surge in the number of published research papers.
Data from the International Association of Scientific, Technical, and Medical Publishers indicates that over 2.5 million new research articles are published annually, a number that continues to grow. This proliferation of academic content presents a significant challenge for researchers and scholars in keeping abreast of developments within their respective fields.

## Objectives

- **Develop an Advanced Summarization System**: Create a sophisticated system capable of summarizing academic paper abstracts, utilizing the latest advancements in natural language processing (NLP) and machine learning.

- **Provide Concise and Accurate Summaries**: Generate summaries that are both brief and precise, enabling researchers to quickly understand the core content of research papers.

- **Offer an Overview of Field Developments**: Aggregate and synthesize information from multiple papers to present an overarching view of current trends and developments in specific academic fields.

- **Identify Emerging Research Areas**: Use summarization techniques to detect and highlight new and evolving areas of research within the academic literature.

- **Reduce Time and Effort in Literature Reviews**: Aim to significantly lower the time and effort required for conducting literature reviews, thereby increasing efficiency in academic research.

## Methodology

1. Data Collection and Preprocessing
   a. **Corpus Gathering**: Collect a large dataset of academic papers, particularly focusing on abstracts.

b. **Text Preprocessing**: Clean and preprocess the text data, which includes tokenization, removing stop words, stemming or lemmatization, and handling special characters or equations found in academic texts.

2. Natural Language Processing Techniques
   a. **Text Representation**: Convert text data into a format understandable by machine learning models, using techniques like TF-IDF, word embeddings (Word2Vec, GloVe), or more advanced embeddings from models like GPT.

   b. **Named Entity Recognition (NER)**: Identify and categorize key terms and entities in the text for identifying key concepts, authors, or research terms.

3. Summarization Approaches
   a. **Extractive Summarization**: Develop models to identify and extract key sentences or phrases directly from the text. This includes techniques like sentence ranking based on relevance, clustering, or graph-based models.

   b. **Abstractive Summarization**: Implement models that generate new text that summarizes the original content. Use advanced language models like ChatGPT-4 to understand and generate human-like text.

4. Machine Learning and Deep Learning Models
   a. **Supervised Learning Models**: Use labeled datasets (if available) to train models to generate summaries.

   b. **Fine-tuning Pre-trained Models**: Utilize pre-trained language models and fine-tune them for better performance in summarization tasks.

5. Evaluation Metrics
   a. **ROUGE Scores**: Use Recall-Oriented Understudy for Gisting Evaluation (ROUGE) metrics to evaluate the quality of the summaries against reference summaries.

   b. **BLEU Score**: use BLEU scores to evaluate the quality of generated summaries.

6. Visualization and Analysis
   a. **Trend Analysis**: Implement methods to analyze and visualize key trends and patterns in the summarized content, such as topic modeling or sentiment analysis.

   b. **Interactive Dashboards**: Create interactive tools or dashboards for users to explore and visualize the summaries and underlying data, enhancing the usability of the summaries.

## Technologies to be Used

- **Machine Learning/Artificial Intelligence**: For summarization of papers. This may include Natural LanguageProcessing.

- **Cloud Computing**: For scalable data storage, real-time analytics, and collaborative planning.

- **Advanced Analytics**: For data visualization, trend analysis, and insight generation regarding research papers.

## Expected Outcomes

- **Development of an Effective Summarization Tool**: A robust system capable of generating accurate and concise summaries of academic paper abstracts, enhancing understanding and accessibility of complex research content.

- **Significant Time Savings in Research**: A notable reduction in the time and effort required for literature reviews, allowing researchers and students to stay updated with less effort.

- **Insights into Academic Trends and Themes**: Ability to identify and analyze prevailing trends, emerging research areas, and key themes in various academic fields.

- **Advancements in NLP and ML Techniques**: Contributions to the field of natural language processing and machine learning, especially in the context of processing and summarizing specialized academic text.

- **Scalability and Adaptability of the Model**: Demonstrated effectiveness of the model across different academic disciplines and its potential for future expansion and customization.

- **Potential for Academic Publications and Open-Source Contributions**: Opportunities for publishing findings in academic journals and contributing to the open-source community.

- **Real-World Application and Commercial Potential**: Demonstrating practical utility in academic settings and exploring possibilities for commercialization as a research tool.

## Project Deliverables

Please see below for a list of deliverables that we plan to complete for this project.

1. **Project Plan:** Document detailing the project's scope, objectives, methodology, and expected outcomes.

2. **Project Report:** Document detailing what we learned from performing this project. This document will also include the results of our analysis and any conclusions we draw from them, as well as any related charts and plots.

3. **README.md:** File which contains the project plan and details about deliverables. This file will also document how to use the Google Colab File, and details about the data used in the project.

4. **Google Colab Project File:** A Google Colab Notebook File containing the code for the project.

5. **Dataset:** Dataset containing the data that was used for this project. This dataset will be uploaded to the GitHub repository.

## Conclusion

This project aims to transform academic research review by developing a system for efficiently summarizing academic paper abstracts using advanced NLP and machine learning technologies. Our goal is to provide a tool that not only condenses content but also enhances understanding and accessibility of complex research for a wide audience. We envision this tool as a catalyst for research efficiency, fostering deeper insight and cross-disciplinary collaboration. Ultimately, this initiative represents a step towards democratizing scientific knowledge and advancing global research and education.