

Research Paper Summarization and Retrieval Assistant

- Suriya Narayanan Rajavel
- Swetha Gendlur Nagarajan

Agenda

- Project Motivation & Overview
- System Architecture
- Datasets & Preprocessing
- Model Architecture
- Results & Performance
- Future Work & Applications

Underlying Thief?

In today's rapidly evolving academic landscape, researchers face an overwhelming challenge: ****information overload****. With thousands of research papers published daily across numerous disciplines, staying current and discovering relevant work has become increasingly difficult.

Introducing Our Solution: The Research Paper Summarization and Retrieval Assistant



Project Motivation:

Time-Intensive Research



Researchers invest countless hours reading and summarizing academic papers, diverting time from critical analysis and innovation.

Difficulty in Discovering Relevant Work



With the exponential growth of research publications, identifying relevant papers across vast repositories has become increasingly challenging.

Gap in Research Discovery Tools



Existing tools often lack the ability to provide concise summaries or pinpoint the most relevant studies, slowing down the research process.

Opportunity with Advanced LLMs:



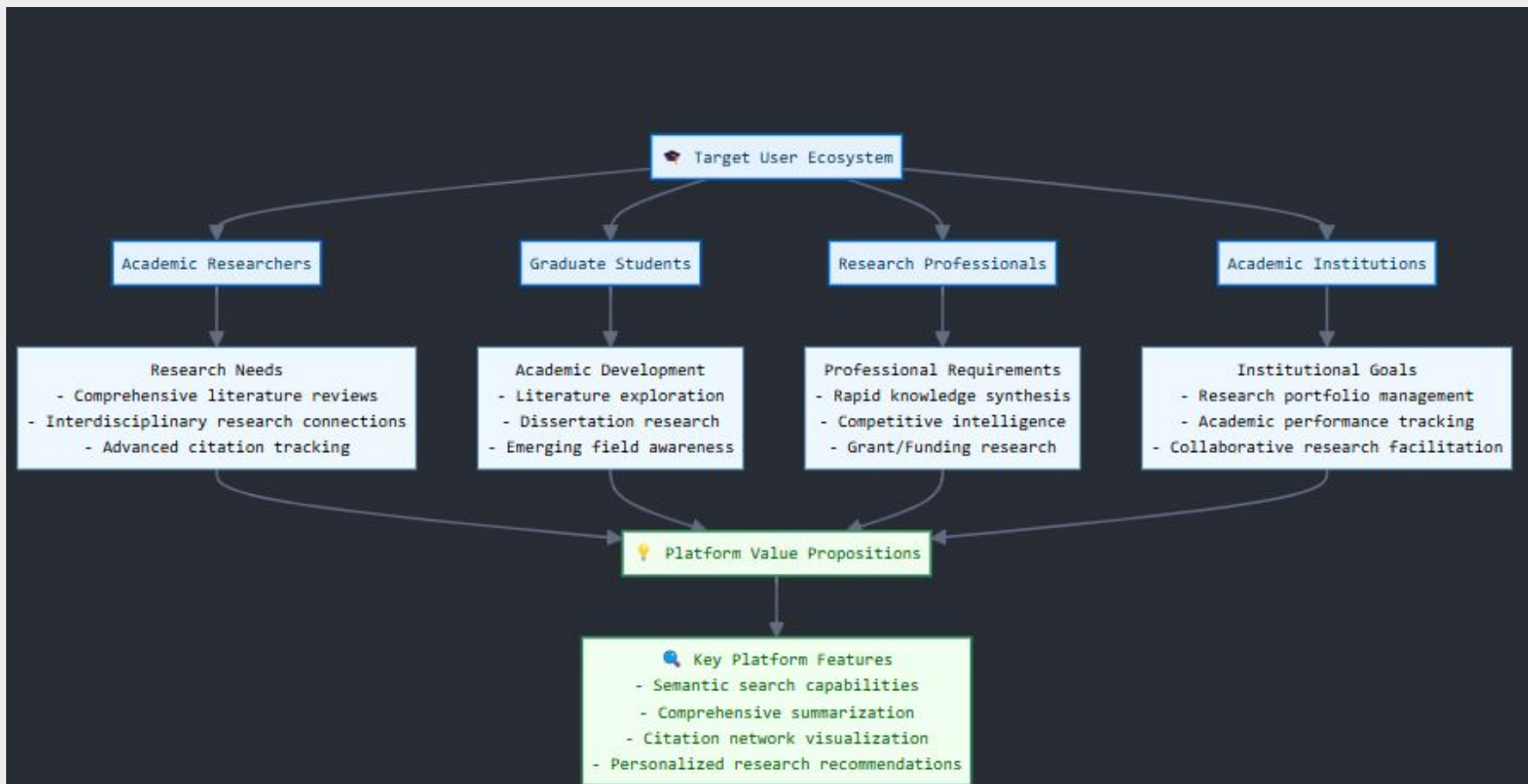
Large language models (LLMs) offer a transformative solution by automating paper summarization and enabling efficient retrieval, enhancing research productivity.

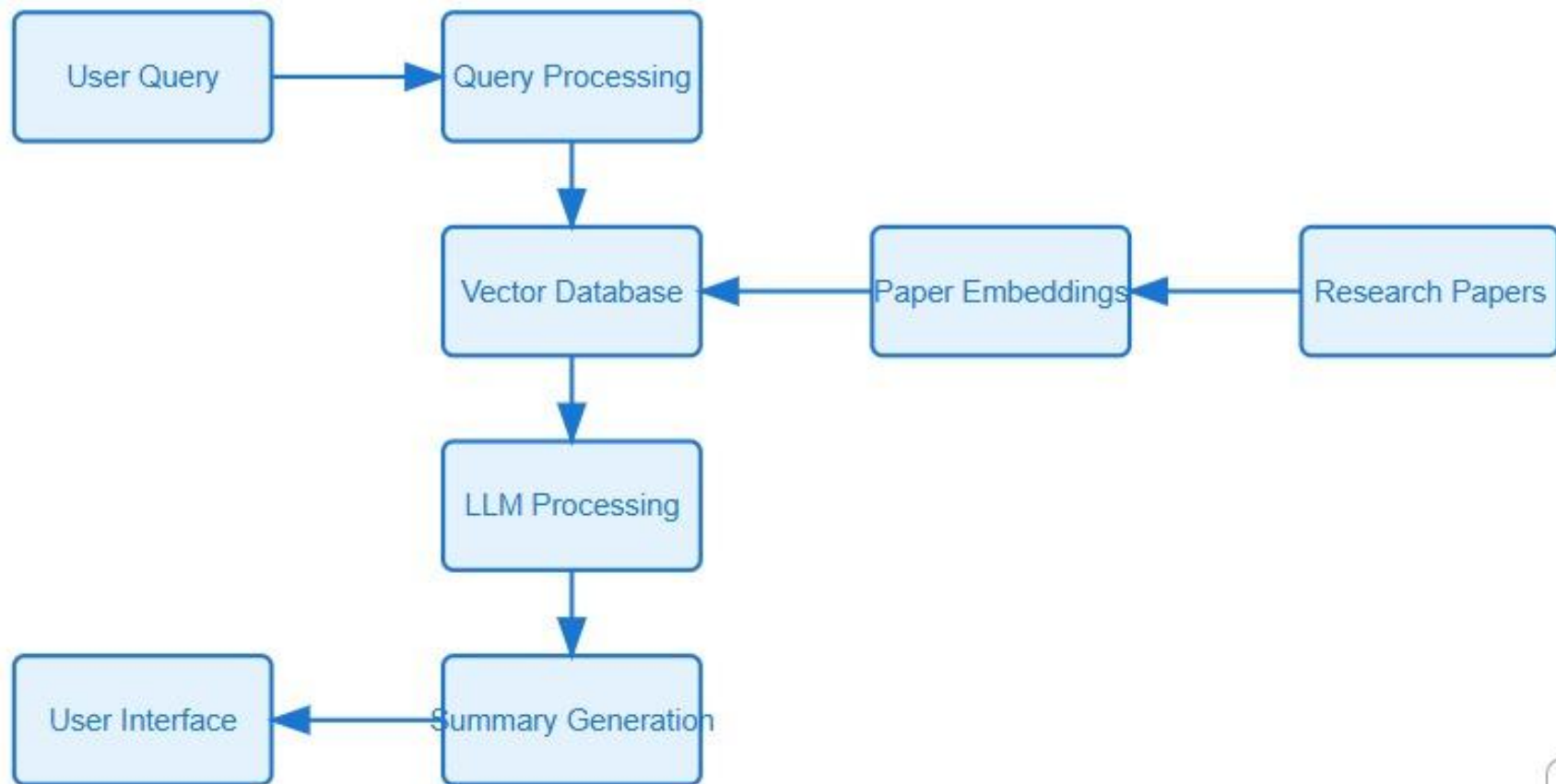
Project Overview:

Objective

Build an **AI-powered assistant** that:

- Provides high-quality summaries of research papers
- Helps discover related academic work
- Enhances research efficiency through intelligent retrieval





Source of the Papers

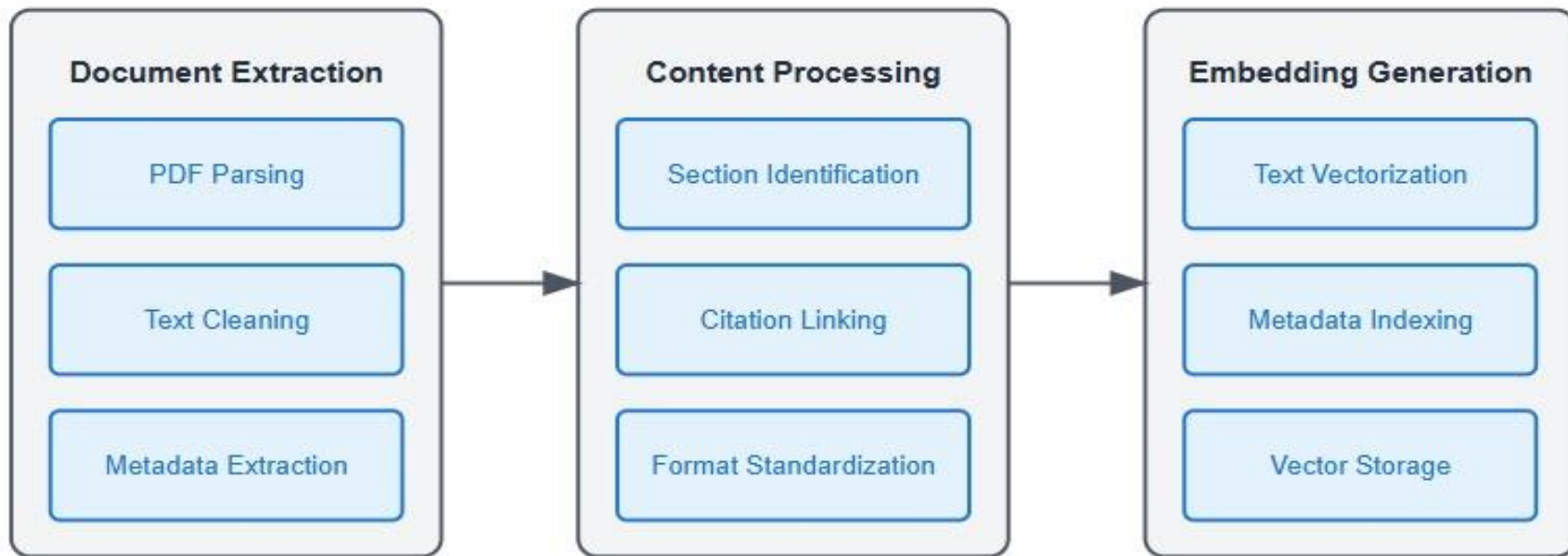
Primary Sources:

- arXiv Research Papers (Computer Science)
- Semantic Scholar Open Research Corpus

Characteristics

- Rich academic content across disciplines
- Structured metadata
- Open access availability
- Comprehensive coverage

Data Preprocessing Pipeline



Data Preprocessing Pipeline

- **Document Extraction**

PDF parsing techniques to extract raw text and relevant metadata from research papers.

Ensure the text is cleaned by removing non-informative elements

- **Content Processing**

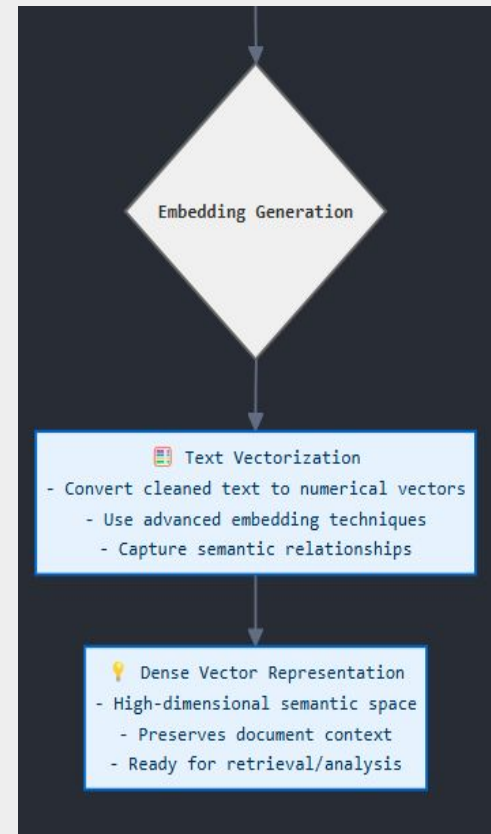
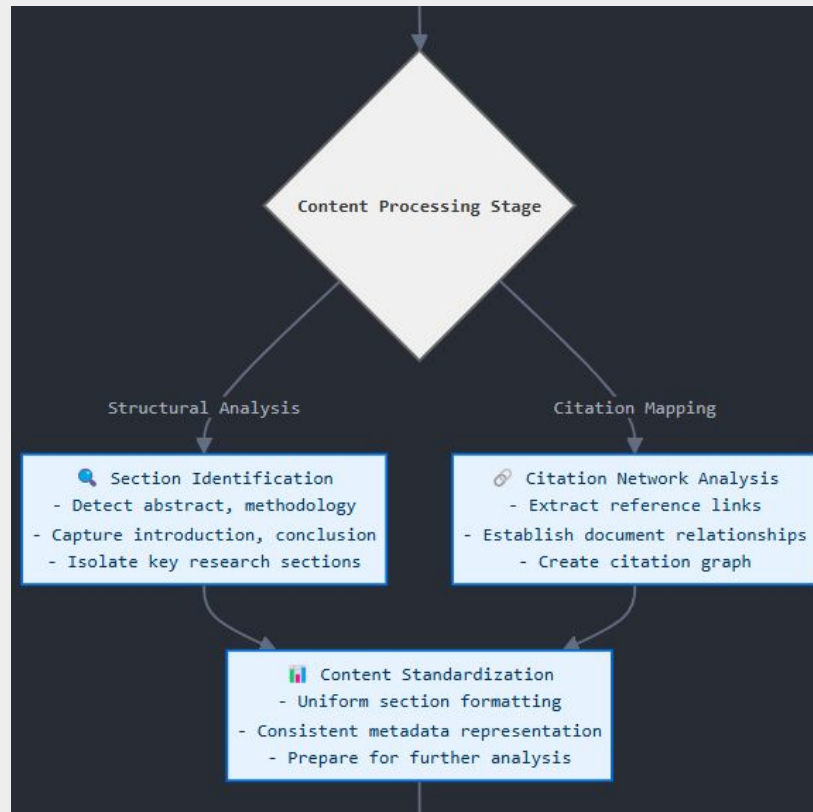
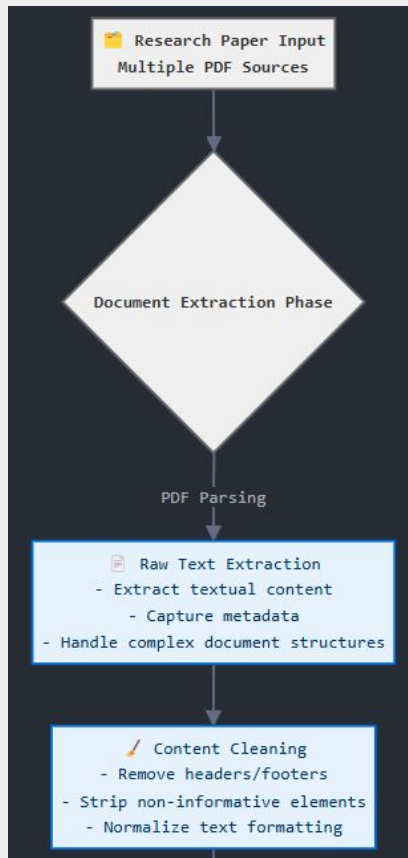
Identify key sections of papers for targeted analysis.

Link citations to establish relationships between documents and standardize formatting for consistency.

- **Embedding Generation**

Transform cleaned text into dense numerical vectors using advanced text vectorization techniques.

Data Preprocessing Pipeline: Detailed Overview



Technical Implementation: Core Components

- **Large Language Models**

Utilize GPT-4 for generating accurate and concise summaries of research papers. The model is fine-tuned specifically for academic content to improve relevance and quality.

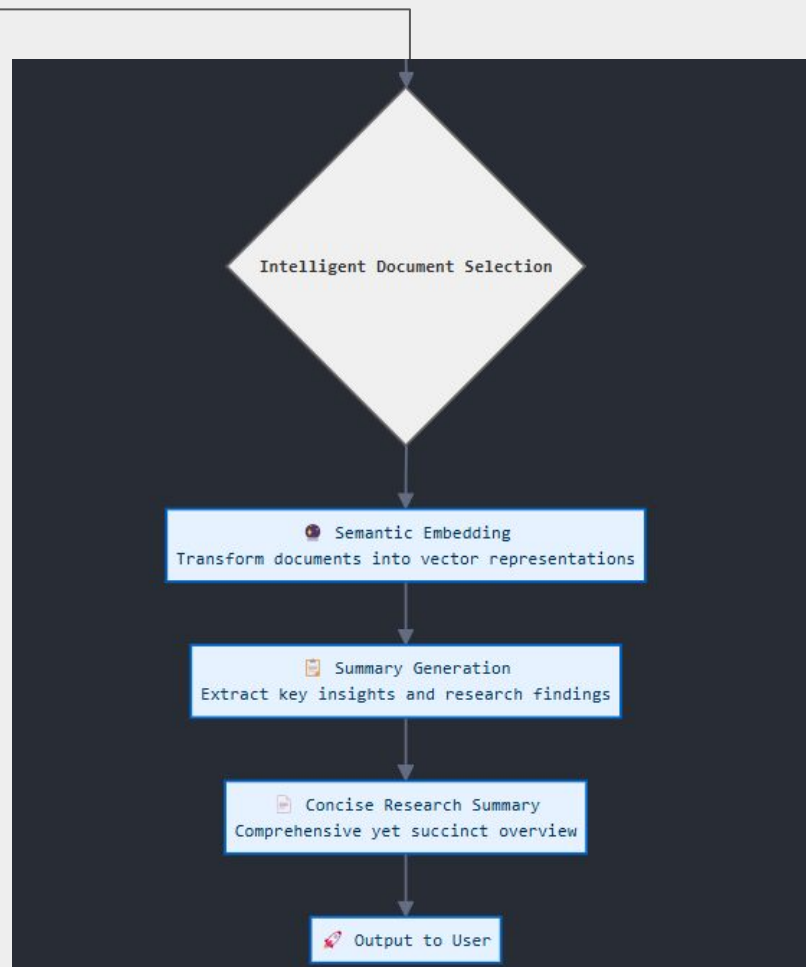
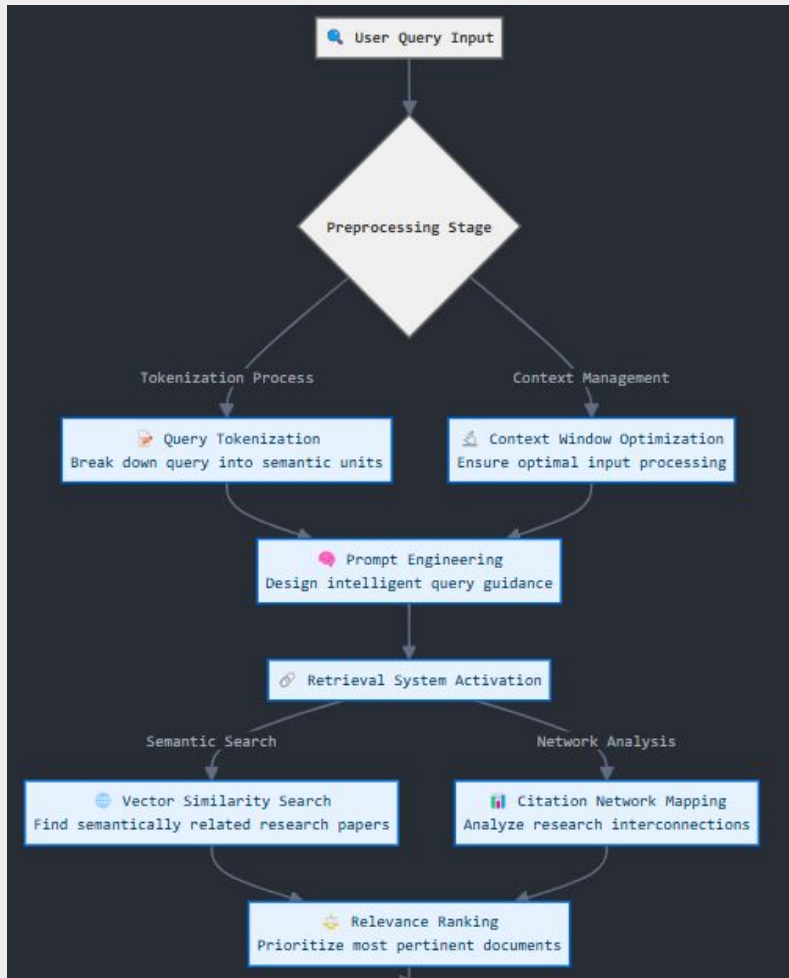
- **Vector Database**

Employ FIASS for storing document embeddings efficiently. The system is optimized to enable fast and accurate retrieval of relevant papers based on input queries.

- **Sentence Transformers**

Leverage sentence transformers for creating high-quality document embeddings. This ensures semantic search capabilities, allowing for precise matching of user queries with research content.

Model Architecture



Model Architecture: Summarization Pipeline and Retrieval System

- **Summarization Pipeline**

- *Input Processing*: Prepares user queries by tokenizing and formatting them for optimal model comprehension.
- *Context Window Management*: Ensures relevant content fits within the model's input limitations for accurate processing.
- *Prompt Engineering*: Designs effective prompts tailored for academic summarization tasks to guide the LLM's responses.
- *Output Generation*: Produces concise and coherent summaries, highlighting key insights from the research papers.

Model Architecture: Summarization Pipeline and Retrieval System

- **Retrieval System**

- *Vector Similarity Search*: Uses embeddings to find papers closely matching the input query based on semantic similarity.
- *Relevance Ranking*: Ranks retrieved documents by their relevance to the query, improving the precision of results.
- *Citation Network Analysis*: Incorporates citation relationships to enhance retrieval by linking related research papers.

Results & Performance

- **Score used for Evaluation**

- We utilized **BERTScore** for evaluating the quality of generated summaries. BERTScore effectively measures semantic similarity between the generated and reference summaries by leveraging embeddings from advanced transformer-based models.

- **Score used for Evaluation**

- **Precision (P)**: Measures how much of the generated summary aligns with the reference summary.
- **Recall (R)**: Measures how much of the reference summary is covered by the generated summary.
- **F1 Score (F1)**: Harmonic mean of Precision and Recall, providing a balanced measure of both.

Results & Performance

- **Summary Evaluation Results**

- **BERTScore Precision (P): 0.5246**
- **BERTScore Recall (R): 0.6752**
- **BERTScore F1 (F1): 0.5905.**

```
Summary Evaluation Results:
-----
BERTScore P: 0.5246
BERTScore R: 0.6752
BERTScore F1: 0.5905
```

- **Strengths**

- High **Recall** indicates the model retrieves most of the relevant content from the research papers.
- Effective tokenization and context management ensure comprehensive summaries.

- **Improvement Areas:**

- **Precision** is lower compared to Recall, suggesting some irrelevant content may be included.
- improving prompt engineering generate more focused outputs.

Future Work & Applications

- Implement advanced fine-tuning of the LLM model using domain-specific datasets to improve the accuracy and relevance of academic summarizations.
- Integrate multi-modal data processing capabilities, such as combining textual data with visual data like charts or graphs from research papers.
- Develop a real-time research assistant API for integration with academic platforms, enhancing accessibility for researchers.
- Expand the retrieval system with citation and co-authorship networks to provide deeper contextual insights and related work suggestions.
- Extend the evaluation framework by incorporating human feedback loops for iterative improvement of summarization and retrieval quality.