# PaperIQ Al Powered Research Insight Analyzer

By: Ishan Patil

# Overview

Introduction	01
Objectives	02
The Research Bottleneck	03
Technical Architecture	04
Powerful Feature Set	05
Future Roadmap	06
Conclusion and Impact	07

# Introduction

With the growing volume of research papers, reports, and documents, extracting relevant information has become tedious and time-consuming. Traditional keyword search often misses contextual meaning, making precise retrieval difficult.

The Document Analyzer Chatbot addresses this by processing PDF/DOCX files and answering queries in a context-aware manner. It leverages lightweight NLP techniques (tokenization, lemmatization, stopword removal, topic extraction, TF-IDF) combined with regex-based rule matching to deliver accurate and efficient information retrieval.

# Objectives

## Objective 01

Extract and structure text from PDF and DOCX documents for analysis.

## Objective 02

Understand user queries and identify key topics for retrieval.

## Objective 03

Retrieve relevant content using regex-based pattern matching and TF-IDF cosine similarity.

## Objective 04

Generate natural, context-aware responses tailored to query types (definition, features, how/why).

## Objective 05

Provide document summarization and an interactive chat interface for easy access to information.

## The Research Bottleneck

Time-Intensive Reading

Researchers spend 40% of their time manually reading lengthy academic papers and reports

Manual Extraction Burden

Extracting definitions, key features, processes, and summaries requires painstaking manual effort across multiple documents

No Interactive Q&A

Traditional PDF readers lack conversational capabilities - researchers can't ask questions and get instant answers from their documents

# Technical Architecture



#### **Document Processing**

- Extract text from PDF (PyMuPDF) and DOCX (python-docx) files
- Maintain paragraph and section-level structure for downstream analysis



#### Text preprocessing

- Tokenization, lemmatization, and stopword removal using NLTK
- Clean text and normalize content for accurate retrieval



#### **Query Understanding**

- Identify headings and split text into structured sections
- Break sections into manageable paragraph chunks (about 200 chars)



#### Smart Retrieval

- Regex-based and NLP-assisted extraction of query topics
- Classify queries into types: definition, features, working, literature review
- Display page number for the reference



#### **Response Generation**

- Regex + TF-IDF similarity to find relevant paragraphs
- Rule-based NER
- Context-aware natural responses formatted for readability

# Extract



# **Tech Stack**

# **Frontend**

- HTML
- CSS
- JavaScript

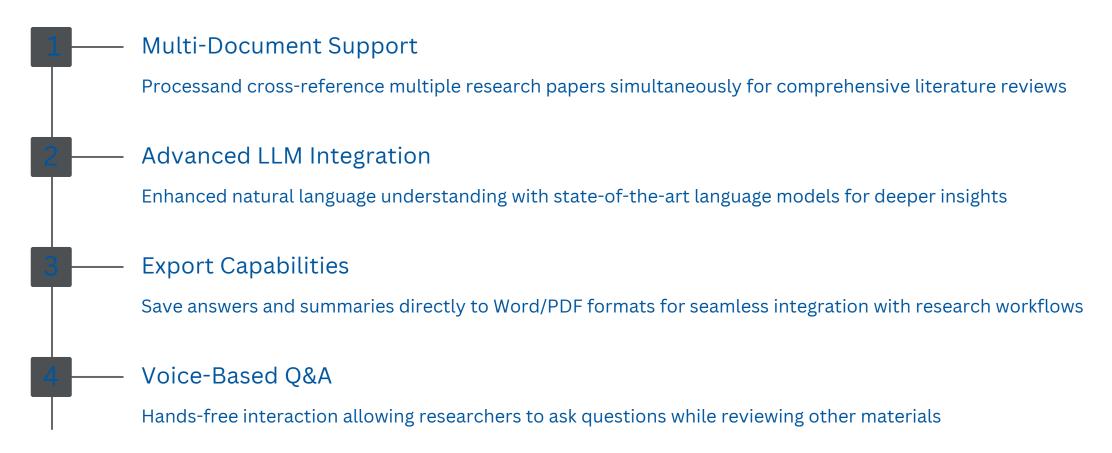
# Other libraries used

- re
- math
- io (BytesIO)

# **Backend**

- Python
- NLTK (Natural Language Toolkit)
- scikit-learn
- PyMuPDF (fitz)
- python-docx
- NumPy

# **Future Roadmap**



# **Conclusion & Impact**

The Document Analyzer Chatbot provides an effective solution for extracting and retrieving relevant information from large PDF/DOCX documents. By combining preprocessing, topic extraction, regex-based rules, and TF-IDF similarity, it delivers accurate, context-aware responses to user queries. The system reduces manual effort, improves efficiency, and demonstrates how lightweight NLP can be used for practical document analysis. With further enhancements, it can evolve into a robust research and enterprise assistant.

#### **Immediate Benefits**

PaperIQ delivers instant insights from long documents, reducing manual effort and accelerating research timelines

### **Productivity Enhancement**

Transforms research review from hoursto minutes, enabling teams to focus on analysis rather than information gathering

# THANK YOU!