

AI Text Agent: Entity Extraction & Summarizer Assistant

1. Introduction

This document provides a detailed explanation of the project 'AI Text Agent: Entity Extraction & Summarizer Assistant'. The aim of this assignment is to build a robust agent that can extract key information (entities) and generate summaries from unstructured text documents. We utilize state-of-the-art NLP models, combine classic corpora (Reuters, Wikipedia, CSV data), and implement an orchestrated agent pipeline using modern open-source frameworks. The end goal is a flexible system that can answer complex queries using both extraction and summarization techniques.

2. Environment Setup and Library Imports

The project leverages Python and a suite of libraries for NLP, data handling, and orchestration. Key libraries used include:

- NLTK: Natural Language Toolkit for classic corpora and preprocessing
- Transformers: State-of-the-art models for NER and summarization
- pandas: Tabular data manipulation
- wikipedia-api: Scraping Wikipedia content
- matplotlib, WordCloud: Visualization
- langchain: Building agentic systems and tool orchestration
- Google Generative AI (for advanced LLM capabilities)

The code below sets up all necessary imports, ensures required NLTK datasets are downloaded, and initializes some key variables (e.g., stopwords, lemmatizer).

3. Data Loading and Exploration

The agent works with multiple data sources, demonstrating its flexibility in handling real-world, unstructured data. This includes:

- Reuters news corpus (classic benchmark for NLP)
- Live Wikipedia article scraping
- CSV files (tabular text data)

3.1 Loading Reuters

Reuters is a classic dataset in NLP. The code extracts all 'train' documents from the Reuters corpus, resulting in a collection of raw news articles for further analysis.

3.2 Loading Articles from Wikipedia

The Wikipedia API is used to fetch articles on specified topics. This enables dynamic, real-world data ingestion for testing the agent's adaptability.

3.3 Loading Data from CSV

To simulate tabular text input (e.g., support tickets, reviews), the code reads and processes a CSV file, extracting all document text.

4. Text Preprocessing

Effective text preprocessing is critical in NLP pipelines. The notebook applies the following steps:

- Lowercasing
- Removal of punctuation and special characters
- Tokenization
- Stopword removal (using NLTK stopwords)
- Lemmatization (WordNetLemmatizer) for root form normalization

These steps clean and standardize raw text, improving downstream model performance.

5. Entity Extraction with DistilBERT

Entity extraction (NER) identifies important terms like names, organizations, locations, dates, etc. The notebook uses Hugging Face's DistilBERT-based NER model via the `pipeline` interface. Grouped entities are returned for clarity. Sample output might be:

- PER: John Smith
- ORG: Google
- LOC: New York

6. Summarization with BART

Automatic summarization condenses long articles or tickets into concise, meaningful summaries. Here, we use Facebook's BART model via Hugging Face's summarization pipeline. This is well-suited for abstractive summaries. Also any LLM can also be used to enhance accuracy and handling large documents.

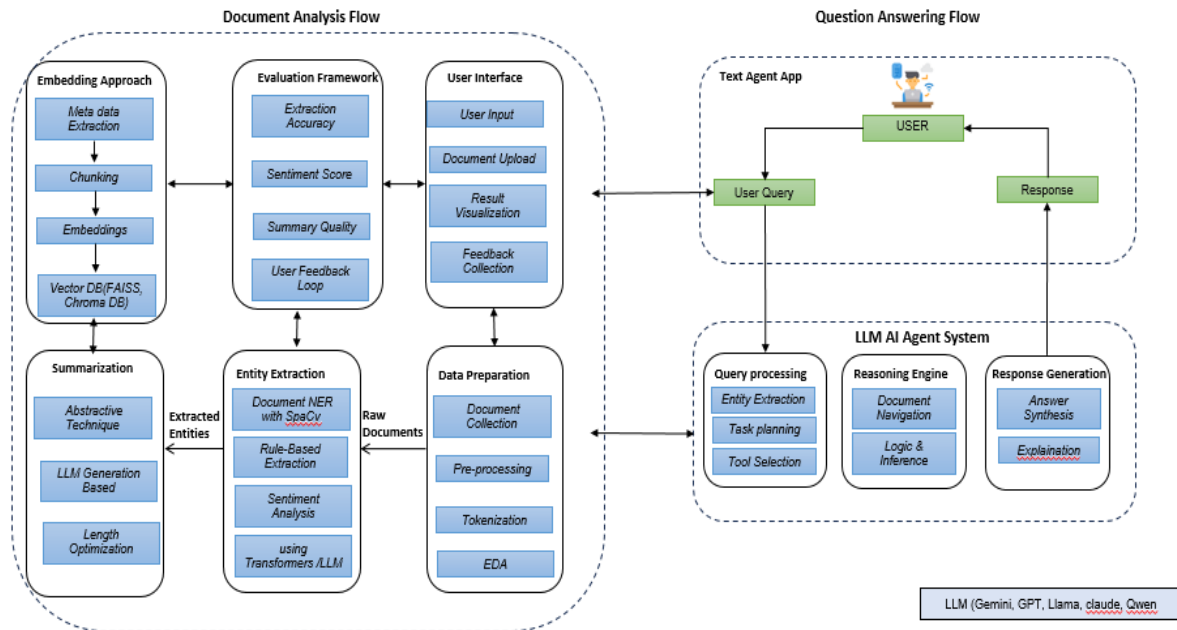
7. User Interaction flow

7.1. Document Analysis Flow

User -> Upload Documents -> System Processing -> Extracted Information & Summaries -> User

7.2. Question Answering Flow

User -> Ask Question -> LLM AI Agent System -> Document Retrieval -> Information Extraction -> Answer Composition -> User



8. Agent Construction and Tool Orchestration

The core innovation is orchestrating extraction and summarization as 'tools' using an agentic framework. Each tool is registered with a clear function (entity extraction or summarization). The agent (via LangChain) uses a large language model (GoogleGenerativeAI) to interpret user queries and select the correct tool.

9. Demonstration of Query Handling

The agent is now ready to process arbitrary queries. For example:

- 'Extract all names and locations from this document.'
- 'Summarize this support ticket.'

The agent uses the LLM to interpret the query, select the right tool, and run the appropriate extraction or summarization logic on the supplied text.

10. Conclusion and Extensions

This assignment demonstrates a flexible agent that bridges entity extraction and summarization using state-of-the-art models. By separating these capabilities as tools and orchestrating them via an agent, the system is easily extensible—other models/tools can be added (sentiment analysis).

Possible extensions include integrating RAG (retrieval augmented generation), supporting

more data formats, or deploying as a web service for real-world use.