

# Extracting Information from Legal Documents Using RAG

## Objective

The main objective of this assignment is to process and analyse a collection of text files containing legal agreements (e.g., NDAs) to prepare them for implementing a **Retrieval-Augmented Generation (RAG)** system. This involves:

- Understand the Cleaned Data : Gain a comprehensive understanding of the structure, content, and context of the cleaned dataset.
- Perform Exploratory Analysis : Conduct bivariate and multivariate analyses to uncover relationships and trends within the cleaned data.
- Create Visualisations : Develop meaningful visualisations to support the analysis and make findings interpretable.
- Derive Insights and Conclusions : Extract valuable insights from the cleaned data and provide clear, actionable conclusions.
- Document the Process : Provide a detailed description of the data, its attributes, and the steps taken during the analysis for reproducibility and clarity.

The ultimate goal is to transform the raw text data into a clean, structured, and analysable format that can be effectively used to build and train a RAG system for tasks like information retrieval, question-answering, and knowledge extraction related to legal agreements.

## Business Value

The project aims to leverage RAG to enhance legal document processing for businesses, law firms, and regulatory bodies. The key business objectives include:

- Faster Legal Research: Reduce the time lawyers and compliance officers spend searching for relevant case laws, precedents, statutes, or contract clauses.
- Improved Contract Analysis: Automatically extract key terms, obligations, and risks from lengthy contracts.
- Regulatory Compliance Monitoring: Help businesses stay updated with legal and regulatory changes by retrieving relevant legal updates.
- Enhanced Decision-Making: Provide accurate and context-aware legal insights to assist in risk assessment and legal strategy.

## Use Cases

- Legal Chatbots
- Contract Review Automation
- Tracking Regulatory Changes and Compliance Monitoring
- Case Law Analysis of past judgments
- Due Diligence & Risk Assessment

# 1. Data Loading, Preparation and Analysis [20 marks]

## 1.1 Data Understanding

The dataset contains legal documents and contracts collected from various sources. The documents are present as text files (.txt) in the *corpus* folder.

There are four types of documents in the *corpus* folder, divided into four subfolders.

- **contractnli**: contains various non-disclosure and confidentiality agreements
- **cuad**: contains contracts with annotated legal clauses
- **maud**: contains various merger/acquisition contracts and agreements
- **privacy\_qa**: a question-answering dataset containing privacy policies

The dataset also contains evaluation files in JSON format in the *benchmark* folder. The files contain the questions and their answers, along with sources. For each of the above four folders, there is a json file: **contractnli.json**, **cuad.json**, **maud.json** **privacy\_qa.json**. The file structure is as follows:

```
{  
    "tests": [  
        {  
            "query": <question1>,  
            "snippets": [{  
                "file_path": <source_file1>,  
                "span": [ begin_position, end_position ],  
                "answer": <relevant answer to the question 1>  
            },  
            {  
                "file_path": <source_file2>,  
                "span": [ begin_position, end_position ],  
                "answer": <relevant answer to the question 2>  
            }, ....  
        ]  
    },  
    {  
        "query": <question2>,  
        "snippets": [{<answer context for que 2>}]  
    },  
    ... <more queries>  
    ]  
}  
  
!pip install -U langchain  
  
Collecting langchain  
  Downloading langchain-1.2.0-py3-none-any.whl.metadata (4.9 kB)  
Collecting langchain-core<2.0.0,>=1.2.1 (from langchain)  
  Downloading langchain_core-1.2.6-py3-none-any.whl.metadata (3.7 kB)
```

```
Collecting langgraph<1.1.0,>=1.0.2 (from langchain)
  Downloading langgraph-1.0.5-py3-none-any.whl.metadata (7.4 kB)
Requirement already satisfied: pydantic<3.0.0,>=2.7.4 in
/usr/local/lib/python3.12/dist-packages (from langchain) (2.12.5)
Collecting jsonpatch<2.0.0,>=1.33.0 (from langchain-
core<2.0.0,>=1.2.1->langchain)
  Downloading jsonpatch-1.33-py2.py3-none-any.whl.metadata (3.0 kB)
Collecting langsmith<1.0.0,>=0.3.45 (from langchain-
core<2.0.0,>=1.2.1->langchain)
  Downloading langsmith-0.6.0-py3-none-any.whl.metadata (15 kB)
Requirement already satisfied: packaging<26.0.0,>=23.2.0 in
/usr/local/lib/python3.12/dist-packages (from langchain-
core<2.0.0,>=1.2.1->langchain) (25.0)
Requirement already satisfied: pyyaml<7.0.0,>=5.3.0 in
/usr/local/lib/python3.12/dist-packages (from langchain-
core<2.0.0,>=1.2.1->langchain) (6.0.3)
Requirement already satisfied: tenacity!=8.4.0,<10.0.0,>=8.1.0 in
/usr/local/lib/python3.12/dist-packages (from langchain-
core<2.0.0,>=1.2.1->langchain) (9.1.2)
Requirement already satisfied: typing-extensions<5.0.0,>=4.7.0 in
/usr/local/lib/python3.12/dist-packages (from langchain-
core<2.0.0,>=1.2.1->langchain) (4.15.0)
Collecting uuid-utils<1.0,>=0.12.0 (from langchain-core<2.0.0,>=1.2.1-
>langchain)
  Downloading uuid_utils-0.12.0-cp39-abi3-
manylinux_2_17_x86_64.manylinux2014_x86_64.whl.metadata (1.1 kB)
Collecting langgraph-checkpoint<4.0.0,>=2.1.0 (from
langgraph<1.1.0,>=1.0.2->langchain)
  Downloading langgraph_checkpoint-3.0.1-py3-none-any.whl.metadata
(4.7 kB)
Collecting langgraph-prebuilt<1.1.0,>=1.0.2 (from
langgraph<1.1.0,>=1.0.2->langchain)
  Downloading langgraph_prebuilt-1.0.5-py3-none-any.whl.metadata (5.2
kB)
Collecting langgraph-sdk<0.4.0,>=0.3.0 (from langgraph<1.1.0,>=1.0.2-
>langchain)
  Downloading langgraph_sdk-0.3.1-py3-none-any.whl.metadata (1.6 kB)
Collecting xxhash>=3.5.0 (from langgraph<1.1.0,>=1.0.2->langchain)
  Downloading xxhash-3.6.0-cp312-cp312-
manylinux2014_x86_64.manylinux_2_17_x86_64.manylinux_2_28_x86_64.whl.m
etadata (13 kB)
Requirement already satisfied: annotated-types>=0.6.0 in
/usr/local/lib/python3.12/dist-packages (from pydantic<3.0.0,>=2.7.4-
>langchain) (0.7.0)
Requirement already satisfied: pydantic-core==2.41.5 in
/usr/local/lib/python3.12/dist-packages (from pydantic<3.0.0,>=2.7.4-
>langchain) (2.41.5)
Requirement already satisfied: typing-inspection>=0.4.2 in
/usr/local/lib/python3.12/dist-packages (from pydantic<3.0.0,>=2.7.4-
```

```
>langchain) (0.4.2)
Requirement already satisfied: jsonpointer>=1.9 in
/usr/local/lib/python3.12/dist-packages (from
jsonpatch<2.0.0,>=1.33.0->langchain-core<2.0.0,>=1.2.1->langchain)
(3.0.0)
Collecting ormsgpack>=1.12.0 (from langgraph-checkpoint<4.0.0,>=2.1.0-
>langgraph<1.1.0,>=1.0.2->langchain)
  Downloading ormsgpack-1.12.1-cp312-cp312-
manylinux_2_17_x86_64.manylinux2014_x86_64.whl.metadata (3.2 kB)
Requirement already satisfied: httpx>=0.25.2 in
/usr/local/lib/python3.12/dist-packages (from langgraph-
sdk<0.4.0,>=0.3.0->langgraph<1.1.0,>=1.0.2->langchain) (0.28.1)
Collecting orjson>=3.10.1 (from langgraph-sdk<0.4.0,>=0.3.0-
>langgraph<1.1.0,>=1.0.2->langchain)
  Downloading orjson-3.11.5-cp312-cp312-
manylinux_2_17_x86_64.manylinux2014_x86_64.whl.metadata (41 kB)
                                                 41.6/41.6 kB 3.3 MB/s eta
0:00:00
  langsmith<1.0.0,>=0.3.45->langchain-core<2.0.0,>=1.2.1->langchain)
    Downloading requests_toolbelt-1.0.0-py2.py3-none-any.whl.metadata
(14 kB)
Requirement already satisfied: requests>=2.0.0 in
/usr/local/lib/python3.12/dist-packages (from
langsmith<1.0.0,>=0.3.45->langchain-core<2.0.0,>=1.2.1->langchain)
(2.32.4)
Collecting zstandard>=0.23.0 (from langsmith<1.0.0,>=0.3.45-
>langchain-core<2.0.0,>=1.2.1->langchain)
  Downloading zstandard-0.25.0-cp312-cp312-
manylinux2014_x86_64.manylinux_2_17_x86_64.whl.metadata (3.3 kB)
Requirement already satisfied: anyio in
/usr/local/lib/python3.12/dist-packages (from httpx>=0.25.2-
>langgraph-sdk<0.4.0,>=0.3.0->langgraph<1.1.0,>=1.0.2->langchain)
(4.12.0)
Requirement already satisfied: certifi in
/usr/local/lib/python3.12/dist-packages (from httpx>=0.25.2-
>langgraph-sdk<0.4.0,>=0.3.0->langgraph<1.1.0,>=1.0.2->langchain)
(2025.11.12)
Requirement already satisfied: httpcore==1.* in
/usr/local/lib/python3.12/dist-packages (from httpx>=0.25.2-
>langgraph-sdk<0.4.0,>=0.3.0->langgraph<1.1.0,>=1.0.2->langchain)
(1.0.9)
Requirement already satisfied: idna in /usr/local/lib/python3.12/dist-
packages (from httpx>=0.25.2->langgraph-sdk<0.4.0,>=0.3.0-
>langgraph<1.1.0,>=1.0.2->langchain) (3.11)
Requirement already satisfied: h11>=0.16 in
/usr/local/lib/python3.12/dist-packages (from httpcore==1.*-
>httpx>=0.25.2->langgraph-sdk<0.4.0,>=0.3.0->langgraph<1.1.0,>=1.0.2-
>langchain) (0.16.0)
Requirement already satisfied: charset_normalizer<4,>=2 in
```

```
/usr/local/lib/python3.12/dist-packages (from requests>=2.0.0->langsmith<1.0.0,>=0.3.45->langchain-core<2.0.0,>=1.2.1->langchain)
(3.4.4)
Requirement already satisfied: urllib3<3,>=1.21.1 in
/usr/local/lib/python3.12/dist-packages (from requests>=2.0.0->langsmith<1.0.0,>=0.3.45->langchain-core<2.0.0,>=1.2.1->langchain)
(2.5.0)
Downloading langchain-1.2.0-py3-none-any.whl (102 kB)
  102.8/102.8 kB 6.2 MB/s eta
0:00:00
  489.1/489.1 kB 23.0 MB/s eta
0:00:00
  157.1/157.1 kB 14.1 MB/s eta
0:00:00
  46.2/46.2 kB 5.0 MB/s eta
0:00:00
  66.5/66.5 kB 5.9 MB/s eta
0:00:00
  283.3/283.3 kB 26.8 MB/s eta
0:00:00
anylinux_2_17_x86_64.manylinux2014_x86_64.whl (343 kB)
  343.7/343.7 kB 33.7 MB/s eta
0:00:00
anylinux2014_x86_64.manylinux_2_17_x86_64.manylinux_2_28_x86_64.whl
(193 kB)
  193.9/193.9 kB 16.7 MB/s eta
0:00:00
anylinux_2_17_x86_64.manylinux2014_x86_64.whl (139 kB)
  139.0/139.0 kB 15.4 MB/s eta
0:00:00
sgpack-1.12.1-cp312-cp312-
manylinux_2_17_x86_64.manylinux2014_x86_64.whl (211 kB)
  211.5/211.5 kB 18.5 MB/s eta
0:00:00
  54.5/54.5 kB 5.3 MB/s eta
0:00:00
anylinux2014_x86_64.manylinux_2_17_x86_64.whl (5.5 MB)
  5.5/5.5 MB 116.1 MB/s eta
0:00:00
sgpack, orjson, jsonpatch, requests-toolbelt, langsmith, langgraph-
sdk, langchain-core, langgraph-checkpoint, langgraph-prebuilt,
langgraph, langchain
Successfully installed jsonpatch-1.33 langchain-1.2.0 langchain-core-
1.2.6 langgraph-1.0.5 langgraph-checkpoint-3.0.1 langgraph-prebuilt-
1.0.5 langgraph-sdk-0.3.1 langsmith-0.6.0 orjson-3.11.5 ormsgpack-
1.12.1 requests-toolbelt-1.0.0 uuid-utils-0.12.0 xxhash-3.6.0
zstandard-0.25.0

!pip install -U langchain-text-splitters
```

```
Collecting langchain-text-splitters
  Downloading langchain_text_splitters-1.1.0-py3-none-any.whl.metadata
(2.7 kB)
Requirement already satisfied: langchain-core<2.0.0,>=1.2.0 in
/usr/local/lib/python3.12/dist-packages (from langchain-text-
splitters) (1.2.6)
Requirement already satisfied: jsonpatch<2.0.0,>=1.33.0 in
/usr/local/lib/python3.12/dist-packages (from langchain-
core<2.0.0,>=1.2.0->langchain-text-splitters) (1.33)
Requirement already satisfied: langsmith<1.0.0,>=0.3.45 in
/usr/local/lib/python3.12/dist-packages (from langchain-
core<2.0.0,>=1.2.0->langchain-text-splitters) (0.6.0)
Requirement already satisfied: packaging<26.0.0,>=23.2.0 in
/usr/local/lib/python3.12/dist-packages (from langchain-
core<2.0.0,>=1.2.0->langchain-text-splitters) (25.0)
Requirement already satisfied: pydantic<3.0.0,>=2.7.4 in
/usr/local/lib/python3.12/dist-packages (from langchain-
core<2.0.0,>=1.2.0->langchain-text-splitters) (2.12.5)
Requirement already satisfied: pyyaml<7.0.0,>=5.3.0 in
/usr/local/lib/python3.12/dist-packages (from langchain-
core<2.0.0,>=1.2.0->langchain-text-splitters) (6.0.3)
Requirement already satisfied: tenacity!=8.4.0,<10.0.0,>=8.1.0 in
/usr/local/lib/python3.12/dist-packages (from langchain-
core<2.0.0,>=1.2.0->langchain-text-splitters) (9.1.2)
Requirement already satisfied: typing-extensions<5.0.0,>=4.7.0 in
/usr/local/lib/python3.12/dist-packages (from langchain-
core<2.0.0,>=1.2.0->langchain-text-splitters) (4.15.0)
Requirement already satisfied: uuid-utils<1.0,>=0.12.0 in
/usr/local/lib/python3.12/dist-packages (from langchain-
core<2.0.0,>=1.2.0->langchain-text-splitters) (0.12.0)
Requirement already satisfied: jsonpointer>=1.9 in
/usr/local/lib/python3.12/dist-packages (from
jsonpatch<2.0.0,>=1.33.0->langchain-core<2.0.0,>=1.2.0->langchain-
text-splitters) (3.0.0)
Requirement already satisfied: httpx<1,>=0.23.0 in
/usr/local/lib/python3.12/dist-packages (from
langsmith<1.0.0,>=0.3.45->langchain-core<2.0.0,>=1.2.0->langchain-
text-splitters) (0.28.1)
Requirement already satisfied: orjson>=3.9.14 in
/usr/local/lib/python3.12/dist-packages (from
langsmith<1.0.0,>=0.3.45->langchain-core<2.0.0,>=1.2.0->langchain-
text-splitters) (3.11.5)
Requirement already satisfied: requests-toolbelt>=1.0.0 in
/usr/local/lib/python3.12/dist-packages (from
langsmith<1.0.0,>=0.3.45->langchain-core<2.0.0,>=1.2.0->langchain-
text-splitters) (1.0.0)
Requirement already satisfied: requests>=2.0.0 in
/usr/local/lib/python3.12/dist-packages (from
langsmith<1.0.0,>=0.3.45->langchain-core<2.0.0,>=1.2.0->langchain-
text-splitters) (2.32.4)
```

```
Requirement already satisfied: zstandard>=0.23.0 in
/usr/local/lib/python3.12/dist-packages (from
langsmith<1.0.0,>=0.3.45->langchain-core<2.0.0,>=1.2.0->langchain-
text-splitters) (0.25.0)
Requirement already satisfied: annotated-types>=0.6.0 in
/usr/local/lib/python3.12/dist-packages (from pydantic<3.0.0,>=2.7.4-
>langchain-core<2.0.0,>=1.2.0->langchain-text-splitters) (0.7.0)
Requirement already satisfied: pydantic-core==2.41.5 in
/usr/local/lib/python3.12/dist-packages (from pydantic<3.0.0,>=2.7.4-
>langchain-core<2.0.0,>=1.2.0->langchain-text-splitters) (2.41.5)
Requirement already satisfied: typing-inspection>=0.4.2 in
/usr/local/lib/python3.12/dist-packages (from pydantic<3.0.0,>=2.7.4-
>langchain-core<2.0.0,>=1.2.0->langchain-text-splitters) (0.4.2)
Requirement already satisfied: anyio in
/usr/local/lib/python3.12/dist-packages (from httpx<1,>=0.23.0-
>langsmith<1.0.0,>=0.3.45->langchain-core<2.0.0,>=1.2.0->langchain-
text-splitters) (4.12.0)
Requirement already satisfied: certifi in
/usr/local/lib/python3.12/dist-packages (from httpx<1,>=0.23.0-
>langsmith<1.0.0,>=0.3.45->langchain-core<2.0.0,>=1.2.0->langchain-
text-splitters) (2025.11.12)
Requirement already satisfied: httpcore==1.* in
/usr/local/lib/python3.12/dist-packages (from httpx<1,>=0.23.0-
>langsmith<1.0.0,>=0.3.45->langchain-core<2.0.0,>=1.2.0->langchain-
text-splitters) (1.0.9)
Requirement already satisfied: idna in /usr/local/lib/python3.12/dist-
packages (from httpx<1,>=0.23.0->langsmith<1.0.0,>=0.3.45->langchain-
core<2.0.0,>=1.2.0->langchain-text-splitters) (3.11)
Requirement already satisfied: h11>=0.16 in
/usr/local/lib/python3.12/dist-packages (from httpcore==1.*-
>httpx<1,>=0.23.0->langsmith<1.0.0,>=0.3.45->langchain-
core<2.0.0,>=1.2.0->langchain-text-splitters) (0.16.0)
Requirement already satisfied: charset_normalizer<4,>=2 in
/usr/local/lib/python3.12/dist-packages (from requests>=2.0.0-
>langsmith<1.0.0,>=0.3.45->langchain-core<2.0.0,>=1.2.0->langchain-
text-splitters) (3.4.4)
Requirement already satisfied: urllib3<3,>=1.21.1 in
/usr/local/lib/python3.12/dist-packages (from requests>=2.0.0-
>langsmith<1.0.0,>=0.3.45->langchain-core<2.0.0,>=1.2.0->langchain-
text-splitters) (2.5.0)
Downloading langchain_text_splitters-1.1.0-py3-none-any.whl (34 kB)
Installing collected packages: langchain-text-splitters
Successfully installed langchain-text-splitters-1.1.0
```

## 1.2 Load and Preprocess the data [5 marks]

Loading libraries

```
## The following libraries might be useful
# !pip install -q langchain-openai
# !pip install -U -q langchain-community
# !pip install -U -q langchain-chroma
# !pip install -U -q datasets
# !pip install -U -q ragas
# !pip install -U -q rouge_score

# Core libraries
import os
import json
from pathlib import Path

# Data handling
import pandas as pd
import numpy as np

# LangChain text processing (new structure)
from langchain_text_splitters import RecursiveCharacterTextSplitter
from langchain_core.documents import Document

# Visualisation
import matplotlib.pyplot as plt
import seaborn as sns

print(" Date - 6th Jan 2026")
print("Created By - Basavanna SV | Akhil ranjan | Tushar Borkar")

Date - 6th Jan 2026
Created By - Basavanna SV | Akhil ranjan | Tushar Borkar
```

### 1.2.1 [3 marks]

Load all .txt files from the folders.

```
import zipfile

zip_path = "/content/rag_legal.zip"
extract_path = "/content/rag_legal"

with zipfile.ZipFile(zip_path, 'r') as zip_ref:
    zip_ref.extractall(extract_path)

print("Files extracted to:", extract_path)

Files extracted to: /content/rag_legal
```

```

# Load the files as documents

from pathlib import Path

#CORPUS_PATH = Path(
#    r"C:\Users\Bas\Desktop\learning\Upgrad MS\RAG Assignment\
#Starter+and+Dataset+RAG+Legal\Starter and Dataset RAG Legal\rag_legal\
#corpus"
#)

CORPUS_PATH = Path("/content/rag_legal/rag_legal/corpus")

documents = []

for domain_folder in CORPUS_PATH.iterdir():
    if domain_folder.is_dir():
        domain = domain_folder.name

        for file_path in domain_folder.glob("*.txt"):
            try:
                with open(file_path, "r", encoding="utf-8") as f:
                    text = f.read()
            except UnicodeDecodeError:
                # Fallback for encoding issues
                with open(file_path, "r", encoding="latin1") as f:
                    text = f.read()

            documents.append({
                "content": text,
                "metadata": {
                    "domain": domain,
                    "file_name": file_path.name,
                    "file_path": str(file_path)
                }
            })

print("Total number of documents loaded:", len(documents))
print(" ")

# Check if the folder exists
if not CORPUS_PATH.exists():
    print(f"Path {CORPUS_PATH} does not exist!")
else:
    # List all subfolders
    subfolders = [f for f in CORPUS_PATH.iterdir() if f.is_dir()]

    # Include the root folder in the summary

```

```

subfolders.insert(0, CORPUS_PATH)

print("Folder summary (folder path : number of files)")
for folder in subfolders:
    num_files = len(list(folder.rglob("*.*"))) # all files
including subfolders
    print(f"{folder} : {num_files} files")

print(" ")

all_files = list(CORPUS_PATH.rglob("*.*"))
txt_files = list(CORPUS_PATH.rglob("*.txt"))

print("Total files:", len(all_files))
print("TXT files:", len(txt_files))
print("Non-TXT files:", len(all_files) - len(txt_files))

Total number of documents loaded: 698

Folder summary (folder path : number of files)
/content/rag_legal/rag_legal/corpus : 698 files
/content/rag_legal/rag_legal/corpus/contractnli : 95 files
/content/rag_legal/rag_legal/corpus/privacy_qa : 7 files
/content/rag_legal/rag_legal/corpus/cuad : 462 files
/content/rag_legal/rag_legal/corpus/maud : 134 files

Total files: 698
TXT files: 698
Non-TXT files: 0

```

You can utilise document loaders from the options provided by the LangChain community.

Optionally, you can also read the files manually, while ensuring proper handling of encoding issues (e.g., utf-8, latin1). In such case, also store the file content along with metadata (e.g., file name, directory path) for traceability.

### 1.2.2 [2 marks]

Preprocess the text data to remove noise and prepare it for analysis.

Remove special characters, extra whitespace, and irrelevant content such as email and telephone contact info. Normalise text (e.g., convert to lowercase, remove stop words). Handle missing or corrupted data by logging errors and skipping problematic files.

```

# Clean and preprocess the data

# Text Preprocessing Strategy
#-----
```

```

#The preprocessing pipeline performs the following steps:
#Converts text to lowercase for normalization
#Removes email addresses and telephone numbers
#Removes special characters and digits
#Removes extra whitespace
#Removes stop words to reduce noise
#Handles missing or corrupted data by logging and skipping files

#stop words
from sklearn.feature_extraction.text import ENGLISH_STOP_WORDS

STOP_WORDS = set(ENGLISH_STOP_WORDS)

#Text Cleaning function
import re
import logging

# Setup basic logging
logging.basicConfig(level=logging.INFO, format="%(levelname)s: %(message)s")

def clean_and_preprocess_text(text, file_name=None):
    try:
        if text is None or not isinstance(text, str) or len(text.strip()) == 0:
            raise ValueError("Empty or invalid text")

        # Convert to lowercase
        text = text.lower()

        # Remove email addresses
        text = re.sub(r"\S+@\S+\.\S+", " ", text)

        # Remove phone numbers
        text = re.sub(r"\+?\d[\d\s\-\(\)]{7,}\d", " ", text)

        # Remove special characters and digits
        text = re.sub(r"[^a-z\s]", " ", text)

        # Remove extra whitespace
        text = re.sub(r"\s+", " ", text).strip()

        # Remove stop words
        text = " ".join([word for word in text.split() if word not in STOP_WORDS])

    return text

except Exception as e:
    logging.warning(f"Skipping file {file_name}: {str(e)}")

```

```

        return None

#Apply preprocesssing to loaded documents
cleaned_documents = []

for doc in documents:
    cleaned_text = clean_and_preprocess_text(
        doc["content"],
        file_name=doc["metadata"]["file_name"]
    )

    if cleaned_text:
        cleaned_documents.append({
            "clean_text": cleaned_text,
            "metadata": doc["metadata"]
        })
    else:
        print("Warning: Document content is empty or invalid")

#im validating here
print("Documents before cleaning:", len(documents))
print("Documents after cleaning:", len(cleaned_documents))

original_lengths = [len(doc["content"]) for doc in documents]
cleaned_lengths = [len(doc["clean_text"]) for doc in cleaned_documents]

print("")
print("Average length before cleaning:", sum(original_lengths) // len(original_lengths))
print("Average length after cleaning:", sum(cleaned_lengths) // len(cleaned_lengths))

Documents before cleaning: 698
Documents after cleaning: 698

Average length before cleaning: 105009
Average length after cleaning: 67859

```

## 1.3 Exploratory Data Analysis [10 marks]

### 1.3.1 [1 marks]

Calculate the average, maximum and minimum document length.

```

# Calculate the average, maximum and minimum document length.
# Calculate length of each cleaned document

```

```

doc_lengths = [len(doc["clean_text"]) for doc in cleaned_documents]

avg_length = sum(doc_lengths) / len(doc_lengths)
print(f"Average document length (characters): {avg_length:.0f}")

max_length = max(doc_lengths)
print(f"Maximum document length (characters): {max_length}")

min_length = min(doc_lengths)
print(f"Minimum document length (characters): {min_length}")

Average document length (characters): 67860
Maximum document length (characters): 648808
Minimum document length (characters): 897

```

### 1.3.2 [4 marks]

Analyse the frequency of occurrence of words and find the most and least occurring words.

Find the 20 most common and least common words in the text. Ignore stop words such as articles and prepositions.

```

# Find frequency of occurrence of words

from collections import Counter
from sklearn.feature_extraction.text import ENGLISH_STOP_WORDS
import pandas as pd

# Combine all documents into a single string
all_text = " ".join([doc["clean_text"] for doc in cleaned_documents])

STOP_WORDS = set(ENGLISH_STOP_WORDS)

# Tokenize
words = [word for word in all_text.split() if word not in STOP_WORDS]

word_freq = Counter(words)

most_common_words = word_freq.most_common(20)
print("Top 20 most common words:")
for word, count in most_common_words:
    print(f"{word}: {count}")

# Sort ascending by frequency
least_common_words = sorted(word_freq.items(), key=lambda x: x[1])[:20]

```

```
print("\nTop 20 least common words:")
for word, count in least_common_words:
    print(f"{word}: {count}")

df_most = pd.DataFrame(most_common_words, columns=[ "Word",
"Frequency"])
df_least = pd.DataFrame(least_common_words, columns=[ "Word",
"Frequency"])

print("\nMost common words (top 20):")
print(df_most)

print("\nLeast common words (top 20):")
print(df_least)

Top 20 most common words:
company: 156423
shall: 108019
agreement: 104658
section: 75450
parent: 60715
party: 54218
s: 47875
date: 39392
time: 35828
b: 34701
material: 34245
merger: 33907
subsidiaries: 33320
applicable: 31384
including: 29407
respect: 28849
stock: 26888
information: 25727
parties: 24641
business: 23706

Top 20 least common words:
cedes: 1
gotham: 1
queens: 1
permutations: 1
encrypting: 1
academician: 1
epsteen: 1
brock: 1
lucky: 1
lockhart: 1
chai: 1
```

```
beaching: 1
shen: 1
zhen: 1
shanshui: 1
nanshanyungu: 1
liuxian: 1
inhabitant: 1
clippings: 1
narration: 1
```

Most common words (top 20):

	Word	Frequency
0	company	156423
1	shall	108019
2	agreement	104658
3	section	75450
4	parent	60715
5	party	54218
6	s	47875
7	date	39392
8	time	35828
9	b	34701
10	material	34245
11	merger	33907
12	subsidiaries	33320
13	applicable	31384
14	including	29407
15	respect	28849
16	stock	26888
17	information	25727
18	parties	24641
19	business	23706

Least common words (top 20):

	Word	Frequency
0	cedes	1
1	gotham	1
2	queens	1
3	permutations	1
4	encrypting	1
5	academician	1
6	epsteen	1
7	brock	1
8	lucky	1
9	lockhart	1
10	chai	1
11	beaching	1
12	shen	1
13	zhen	1

```

14      shanshui      1
15  nanshanyungu      1
16      liuxian      1
17  inhabitant      1
18  clippings      1
19  narration      1

```

### 1.3.3 [4 marks]

Analyse the similarity of different documents to each other based on TF-IDF vectors.

Transform some documents to TF-IDF vectors and calculate their similarity matrix using a suitable distance function. If contracts contain duplicate or highly similar clauses, similarity calculation can help detect them.

Identify for the first 10 documents and then for 10 random documents. What do you observe?

```

from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics.pairwise import cosine_similarity
import numpy as np
import random

# List of all cleaned document texts
texts = [doc["clean_text"] for doc in cleaned_documents]

# Transform the page contents of documents

vectorizer = TfidfVectorizer(max_features=5000) # Limit features for speed
tfidf_matrix = vectorizer.fit_transform(texts)

print(f"TF-IDF matrix shape: {tfidf_matrix.shape}")

# Compute similarity scores

similarity_matrix = cosine_similarity(tfidf_matrix)
print(f"Similarity matrix shape: {similarity_matrix.shape}")

#similarity - first 10 document
first_10_sim = similarity_matrix[:10, :10]

print("Cosine similarity matrix for first 10 documents:")
print(np.round(first_10_sim, 2))

TF-IDF matrix shape: (698, 5000)
Similarity matrix shape: (698, 698)
Cosine similarity matrix for first 10 documents:
[[1.  0.11 0.17 0.74 0.71 0.81 0.78 0.67 0.67 0.78]
 [0.11 1.  0.04 0.09 0.16 0.11 0.13 0.25 0.21 0.22]]

```

```
[0.17 0.04 1.  0.14 0.17 0.14 0.14 0.15 0.14 0.16]
[0.74 0.09 0.14 1.  0.61 0.69 0.65 0.52 0.61 0.63]
[0.71 0.16 0.17 0.61 1.  0.66 0.65 0.58 0.67 0.67]
[0.81 0.11 0.14 0.69 0.66 1.  0.74 0.58 0.63 0.73]
[0.78 0.13 0.14 0.65 0.65 0.74 1.  0.64 0.61 0.79]
[0.67 0.25 0.15 0.52 0.58 0.58 0.64 1.  0.71 0.8 ]
[0.67 0.21 0.14 0.61 0.67 0.63 0.61 0.71 1.  0.77]
[0.78 0.22 0.16 0.63 0.67 0.73 0.79 0.8  0.77 1.  ]]
```

### Observations for first 10 documents

1. **TF-IDF Matrix**
  - The TF-IDF matrix has shape (698, 5000), indicating **698 documents** and **5000 unique terms/features**.
  - This shows that a large vocabulary is captured from the corpus after preprocessing (like stopword removal and tokenization).
2. **Document Similarity**
  - The cosine similarity matrix has shape (698, 698), representing pairwise similarity between all documents.
  - Each element (i, j) indicates how similar document i is to document j based on their TF-IDF vectors.
3. **Similarity Insights (First 10 Documents)**
  - Diagonal values are 1.0, as expected, since each document is perfectly similar to itself.
  - Several off-diagonal values are high (e.g., 0.74, 0.81), showing that **some documents are highly similar or near-duplicates**.
  - Some values are low (e.g., 0.04, 0.09), indicating **low similarity** between certain documents.
  - Clusters of high similarity suggest **groups of documents with overlapping content**.

```
# create a list of 10 random integers
random_indices = random.sample(range(len(texts)), 10)
print("\nRandom 10 documents indices:", random_indices)

Random 10 documents indices: [340, 318, 672, 27, 574, 369, 495, 140,
334, 679]

# Compute similarity scores for 10 random documents

random_sim = similarity_matrix[random_indices, :][:, random_indices]

print("Cosine similarity matrix for 10 random documents:")
print(np.round(random_sim, 2))

Cosine similarity matrix for 10 random documents:
[[1.  0.03 0.08 0.05 0.07 0.09 0.02 0.04 0.02 0.04]]
```

```
[0.03 1.  0.12 0.07 0.07 0.11 0.02 0.04 0.03 0.06]
[0.08 0.12 1.  0.08 0.28 0.16 0.03 0.07 0.05 0.2 ]
[0.05 0.07 0.08 1.  0.1  0.22 0.07 0.11 0.09 0.11]
[0.07 0.07 0.28 0.1  1.  0.18 0.04 0.07 0.05 0.15]
[0.09 0.11 0.16 0.22 0.18 1.  0.08 0.11 0.08 0.16]
[0.02 0.02 0.03 0.07 0.04 0.08 1.  0.12 0.02 0.03]
[0.04 0.04 0.07 0.11 0.07 0.11 0.12 1.  0.04 0.06]
[0.02 0.03 0.05 0.09 0.05 0.08 0.02 0.04 1.  0.04]
[0.04 0.06 0.2  0.11 0.15 0.16 0.03 0.06 0.04 1.  ]]
```

### Observation for Random 10 document

#### Similarity Insights (Random 10 Documents)

- Diagonal values are **1.0**, as expected, since each document is perfectly similar to itself.
- Most off-diagonal values are quite low (e.g., 0.02–0.16), indicating **low similarity** between most document pairs.
- A few moderately higher values (e.g., 0.2, 0.28) suggest **some documents share partial content or topics**.
- Overall, the random sample shows that the corpus is **mostly diverse with limited overlapping content**.

## 1.4 Document Creation and Chunking [5 marks]

### 1.4.1 [5 marks]

Perform appropriate steps to split the text into chunks.

```
# Process files and generate chunks

# Initialize the text splitter
chunk_size = 2000      # characters per chunk
chunk_overlap = 200     # overlap to preserve context

text_splitter = RecursiveCharacterTextSplitter(
    chunk_size=chunk_size,
    chunk_overlap=chunk_overlap
)

# Process all cleaned documents and generate chunks
all_chunks = []

for doc_id, doc in enumerate(cleaned_documents):
    chunks = text_splitter.split_text(doc["clean_text"])
    # Store each chunk with metadata
    for i, chunk in enumerate(chunks):
        all_chunks.append({
            "doc_id": doc_id,
            "chunk_id": i,
```

```

        "text": chunk
    })

# Print summary statistics
print("\n--- Document Chunking Summary ---")
print(f"Total documents: {len(cleaned_documents)}")
print(f"Total chunks generated: {len(all_chunks)}")
print(f"Average chunks per document:
{len(all_chunks)/len(cleaned_documents):.2f}")

# im printing sample chunks values
print("\nSample chunks:")
for sample in all_chunks[:3]:
    print(f"Doc {sample['doc_id']}, Chunk {sample['chunk_id']}:\n{sample['text'][:200]}...\n")

--- Document Chunking Summary ---
Total documents: 698
Total chunks generated: 26591
Average chunks per document: 38.10

Sample chunks:
Doc 0, Chunk 0:
non disclosure agreement agreement effective day just biofiber
structural solutions corp alberta corporation offices calgary alberta
just biofiber having office recipient hereinafter collectively refe...

Doc 0, Chunk 1:
hereunder b rightfully possession prior disclosure confidential
information receiving party behalf disclosing party receiving party
shall use confidential information solely purpose shall use
confiden...

Doc 0, Chunk 2:
database receiving party prior consent disclosing party copies storage
reasonably required internally receiving party purpose copies
confidential information shall contain proprietary notices appear o...

```

## 2. Vector Database and RAG Chain Creation [15 marks]

### 2.1 Vector Embedding and Vector Database Creation [7 marks]

#### 2.1.1 [2 marks]

Initialise an embedding function for loading the embeddings into the vector database.

Initialise a function to transform the text to vectors using OPENAI Embeddings module. You can also use this function to transform during vector DB creation itself.

```

# Fetch your OPENAI API Key as an environment variable

import os
os.environ["OPENAI_API_KEY"] = "sk-proj-z8rnzvdH009Sqb-
VQJrQnb8qV54B36w4a5-
uD78hm0vvJcMLUD0U0dkbczN0xoCems3ya4P1CmT3BlbkFJTlM-
V3pkBVh1D4Gn41u6Hrb2ilavDj6f_pq8vUG0ZGhX6RgQbDuyh_1aF9vVGRL0ITLD3iDPtA
"
# I have removed the API key here after I have completed all the sections

OPENAI_API_KEY = os.getenv("OPENAI_API_KEY")
if not OPENAI_API_KEY:
    raise ValueError("OPENAI_API_KEY is not set!")
print("OpenAI API Key set successfully.")

OpenAI API Key set successfully.

import langchain
print(langchain.__version__)

1.2.0

!pip install -q langchain-openai
=====
          0.0/84.7 kB ? eta -----
          84.7/84.7 kB 5.7 MB/s eta
0:00:00
=====
          0.0/1.1 MB ? eta -----
          1.1/1.1 MB 55.3 MB/s eta
0:00:01 -----
          1.1/1.1 MB 30.2 MB/s
eta 0:00:00
=====
          0.0/1.2 MB ? eta -----
          1.2/1.2 MB 73.6 MB/s eta
0:00:00
=====
          0.0/361.3 kB ? eta -----
          361.3/361.3 kB 31.6 MB/s eta
0:00:00

# Initialise an embedding function
# 2. Import the updated OpenAIEmbeddings class
from langchain_openai import OpenAIEmbeddings

# Initialize embedding function
embeddings = OpenAIEmbeddings(
    openai_api_key=OPENAI_API_KEY,
    model="text-embedding-3-small" # recommended default
)

print("Embedding function initialized")

```

Embedding function initialized

### 2.1.2 [5 marks]

Load the embeddings to a vector database.

Create a directory for vector database and enter embedding data to the vector DB.

```
!pip install -q chromadb
0:00:00          0.0/67.3 kB ? eta -:--:-
0:00:00          67.3/67.3 kB 3.0 MB/s eta
0:00:00          ents to build wheel ... etadata (pyproject.toml) ...
0:00:00          21.7/21.7 MB 144.1 MB/s eta
0:00:00          278.2/278.2 kB 25.3 MB/s eta
0:00:00          2.0/2.0 MB 98.0 MB/s eta
0:00:00          103.3/103.3 kB 10.4 MB/s eta
0:00:00          17.4/17.4 MB 168.9 MB/s eta
0:00:00          66.4/66.4 kB 4.9 MB/s eta
0:00:00          72.5/72.5 kB 6.6 MB/s eta
0:00:00          132.6/132.6 kB 12.5 MB/s eta
0:00:00          220.0/220.0 kB 21.0 MB/s eta
0:00:00          105.4/105.4 kB 10.2 MB/s eta
0:00:00          71.6/71.6 kB 6.7 MB/s eta
0:00:00          47.4/47.4 kB 4.2 MB/s eta
0:00:00          68.5/68.5 kB 6.5 MB/s eta
0:00:00          517.7/517.7 kB 35.6 MB/s eta
0:00:00          128.4/128.4 kB 13.4 MB/s eta
0:00:00          4.4/4.4 MB 128.1 MB/s eta
0:00:00          456.8/456.8 kB 35.6 MB/s eta
0:00:00          46.0/46.0 kB 3.5 MB/s eta
0:00:00
```

```
----- 86.8/86.8 kB 8.1 MB/s eta
0:00:00
l) ...

!pip install sentence-transformers --quiet
!pip install sentence-transformers --quiet
----- 0.0/493.7 kB ? eta -:---:
----- 493.7/493.7 kB 15.1 MB/s eta
0:00:00

# Add Chunks to vector DB

import os
from sentence_transformers import SentenceTransformer
import chromadb

# Vector DB directory
vector_db_dir = "vector_db"
os.makedirs(vector_db_dir, exist_ok=True)
print(f"Vector DB directory created at '{vector_db_dir}'")

# Initialize Chroma client
client = chromadb.Client()
print("Chroma client initialized")

# Create or get collection
collection_name = "legal_docs"
collections = [c.name for c in client.list_collections()]

if collection_name in collections:
    collection = client.get_collection(name=collection_name)
else:
    collection = client.create_collection(name=collection_name)
print(f"Collection '{collection_name}' ready")

# Initialize local embedding model
model = SentenceTransformer('all-MiniLM-L6-v2')

# Convert documents to plain text
texts = [str(doc) for doc in documents]

# Generate embeddings
embeddings = model.encode(texts, show_progress_bar=True)

# Add to Chroma DB
collection.add(
    documents=texts,
    metadatas=[{"doc_id": i // 38, "chunk_id": i} for i in
range(len(texts))],
    ids=[str(i) for i in range(len(texts))],
```

```
    embeddings=embeddings.tolist() # convert numpy array to list
)

# Persist DB
#client.persist()
print(f"{len(texts)} document-level chunks embedded and stored in
vector DB")

/usr/local/lib/python3.12/dist-packages/jax/_src/cloud_tpu_init.py:86:
UserWarning: Transparent hugepages are not enabled. TPU runtime
startup and shutdown time should be significantly improved on TPU v5e
and newer. If not already set, you may need to enable transparent
hugepages in your VM image (sudo sh -c "echo always >
/sys/kernel/mm/transparent_hugepage/enabled")
warnings.warn()

Vector DB directory created at 'vector_db'
Chroma client initialized
Collection 'legal_docs' ready

/usr/local/lib/python3.12/dist-packages/huggingface_hub/utils/
_auth.py:94: UserWarning:
The secret `HF_TOKEN` does not exist in your Colab secrets.
To authenticate with the Hugging Face Hub, create a token in your
settings tab (https://huggingface.co/settings/tokens), set it as
secret in your Google Colab and restart your session.
You will be able to reuse this secret in all of your notebooks.
Please note that authentication is recommended but still optional to
access public models or datasets.
warnings.warn()

{"model_id":"6a812fb4b4a94ad7921529906091adbb","version_major":2,"vers
ion_minor":0}

{"model_id":"14dd774572354ec181c8c97836c79aef","version_major":2,"vers
ion_minor":0}

 {"model_id":"4ecce7f633db4e3b82293d02b109c0a3","version_major":2,"vers
ion_minor":0}

 {"model_id":"07f98448285045a4a4222cfe6bbf29c5","version_major":2,"vers
ion_minor":0}

 {"model_id":"0166ebaf4c754970a88d2bfec69d356d","version_major":2,"vers
ion_minor":0}

 {"model_id":"be1885cf13814d9ba3624fa5db4222f6","version_major":2,"vers
ion_minor":0}

 {"model_id":"60614be3d14b448aaa9b677a630f382f","version_major":2,"vers
ion_minor":0}
```

```

{"model_id": "dca8f610b8fb4d08b5e7bcad0d6596e7", "version_major": 2, "version_minor": 0}

{"model_id": "78b037c1538143d3a53e70c8a60265d7", "version_major": 2, "version_minor": 0}

{"model_id": "f65be2aed222478ebdcf4439d3916d7d", "version_major": 2, "version_minor": 0}

{"model_id": "ae3f4fe040054ae1922f8d95f44ca926", "version_major": 2, "version_minor": 0}

{"model_id": "d5aaaf5c7be32434daf08dbf06d5746d5", "version_major": 2, "version_minor": 0}

698 document-level chunks embedded and stored in vector DB

```

## 2.2 Create RAG Chain [8 marks]

### 2.2.1 [5 marks]

Create a RAG chain.

```

# Create a RAG chain

import chromadb
from sentence_transformers import SentenceTransformer
from transformers import pipeline

# Load embedding model
embedding_model = SentenceTransformer("all-MiniLM-L6-v2")

# Load Chroma vector DB
client = chromadb.Client()
collection = client.get_collection(name="legal_docs")

# Load LLM
llm = pipeline(
    "text2text-generation",
    model="google/flan-t5-base",
    max_length=512
)

# RAG chain function
def rag_chain(query, top_k=5):
    """
    Retrieves top-k documents from the vector DB,
    constructs a prompt with context, and generates an answer using
    LLM.
    """

```

```

# Compute embedding for the query
query_embedding = embedding_model.encode(query).tolist()

# Retrieve top-k documents from Chroma
results = collection.query(
    query_embeddings=[query_embedding],
    n_results=top_k
)

# Combine retrieved documents into context
context = "\n\n".join(results["documents"][0])

# Build the prompt
prompt = f"""
Answer the question using only the context below.

Context:
{context}

Question:
{query}

Answer:
"""

# Generate answer
response = llm(prompt)

# Return generated answer and source documents
return response[0]["generated_text"], results["documents"][0]

```

{"model\_id": "66eb4d8367724a7d9866d42100795242", "version\_major": 2, "version\_minor": 0}

{"model\_id": "2c3071993f984835815fc0c5ad36258f", "version\_major": 2, "version\_minor": 0}

{"model\_id": "34ca3f1b07714fb88bd48486c054ef82", "version\_major": 2, "version\_minor": 0}

{"model\_id": "854ddf964e8a4bcdcad9eed64790968d2", "version\_major": 2, "version\_minor": 0}

{"model\_id": "c30e70258b674981952b5df7967cf465", "version\_major": 2, "version\_minor": 0}

{"model\_id": "87cc4a7ba1ed4619921811a0f0bf5a30", "version\_major": 2, "version\_minor": 0}

{"model\_id": "4ba00ed45ac8438da63e96c1fb1cff2", "version\_major": 2, "version\_minor": 0}

Device set to use cpu

## 2.2.2 [3 marks]

Create a function to generate answer for asked questions.

Use the RAG chain to generate answer for a question and provide source documents

```
# Create a function for question answering
def generate_answer(question, top_k=5, max_context_tokens=400):
    """
    Generate an answer for a question using RAG chain.
    Truncates context to avoid exceeding model's token limit.
    Returns the answer and source documents.
    """

    # Get top-k documents
    answer = ""
    sources = []
    retrieved_texts, retrieved_sources = [], []

    query_embedding = embedding_model.encode(question).tolist()
    results = collection.query(query_embeddings=[query_embedding],
n_results=top_k)

    for doc_text in results["documents"][0]:
        # Truncate each document to max_context_tokens words
        words = doc_text.split()
        truncated = " ".join(words[:max_context_tokens])
        retrieved_texts.append(truncated)

    # Combine truncated documents into context
    context = "\n\n".join(retrieved_texts)
    sources = results["documents"][0] # keep full source text if needed

    # Build prompt
    prompt = f"""

Answer the question using only the context below.

Context:
{context}

Question:
{question}

Answer:
"""
    # Generate answer
    response = llm(prompt)
    answer = response[0]["generated_text"]
```

```

    return answer, sources

print("function created")
function created

# Example usage
question = "Consider the Non-Disclosure Agreement between CopAcc and
ToP Mentors; Does the document indicate that the Agreement does not
grant the Receiving Party any rights to the Confidential Information?"

answer, sources = generate_answer(question, top_k=5)

print("Answer:\n", answer)
print("\nSource Document Count:", len(sources))

Token indices sequence length is longer than the specified maximum
sequence length for this model (3279 > 512). Running this sequence
through the model will result in indexing errors

Answer:
no

Source Document Count: 5

```

### 3. RAG Evaluation [10 marks]

#### 3.1 Evaluation and Inference [10 marks]

##### 3.1.1 [2 marks]

Extract all the questions and all the answers/ground truths from the benchmark files.

Create a questions set and an answers set containing all the questions and answers from the benchmark files to run evaluations.

```

# Create a question set by taking all the questions from the benchmark
data
# Also create a ground truth/answer set
from pathlib import Path
import json

BENCHMARK_PATH = Path("/content/rag_legal/rag_legal/benchmarks")

questions_set = []
answers_set = []

file_summary = {}

```

```

for file_path in BENCHMARK_PATH.glob("*.json"):
    file_questions = 0
    file_answers = 0

    try:
        with open(file_path, "r", encoding="utf-8") as f:
            data = json.load(f)

            for test in data.get("tests", []):
                question = test.get("query", "").strip()
                snippets = test.get("snippets", [])
                answer_texts = [s.get("answer", "").strip() for s in
snippets if "answer" in s]
                answer = "\n".join(answer_texts).strip()

                if question and answer:
                    questions_set.append(question)
                    answers_set.append(answer)
                    file_questions += 1
                    file_answers += 1

    except Exception as e:
        print(f"Error reading {file_path}: {e}")

    file_summary[file_path.name] = {"questions": file_questions,
"answers": file_answers}

# Print total counts
print("\n--- Total ---")
print("Total questions loaded:", len(questions_set))
print("Total answers loaded:", len(answers_set))

# Print per file summary
print("\n--- Per file summary ---")
for fname, counts in file_summary.items():
    print(f"{fname} : {counts['questions']} questions,
{counts['answers']} answers")

```

```

--- Total ---
Total questions loaded: 6889
Total answers loaded: 6889

--- Per file summary ---
cuad.json : 4042 questions, 4042 answers
contractnli.json : 977 questions, 977 answers
maud.json : 1676 questions, 1676 answers
privacy_qa.json : 194 questions, 194 answers

```

### 3.1.2 [5 marks]

Create a function to evaluate the generated answers.

Evaluate the responses on *Rouge*, *Ragas* and *Bleus* scores.

```
!pip install rouge-score bert-score nltk --quiet
Installing build dependencies ... gets to build wheel ... etadata
(pyproject.toml) ...
61.1/61.1 kB 3.1 MB/s eta 0:00:00
l) ...

!pip install evaluate --quiet
!pip install bert-score --quiet
!pip install nltk --quiet

0:00:00          0.0/75.1 kB ? eta -:---- 75.1/75.1 kB 5.1 MB/s eta
0:00:00          0.0/84.1 kB ? eta -:---- 84.1/84.1 kB 7.0 MB/s eta
0:00:00          0.0/512.3 kB ? eta -:---- 512.3/512.3 kB 26.1 MB/s eta
0:00:00          0.0/119.7 kB ? eta -:---- 119.7/119.7 kB 10.5 MB/s eta
0:00:00          0.0/201.0 kB ? eta -:---- 201.0/201.0 kB 17.5 MB/s eta
0:00:00          0.0/150.3 kB ? eta -:---- 150.3/150.3 kB 13.2 MB/s eta
0:00:00          0.0/1.8 MB ? eta -:---- 1.8/1.8 MB 75.0 MB/s eta
0:00:00          242.4/242.4 kB 22.8 MB/s eta
0:00:00          221.6/221.6 kB 16.7 MB/s eta
0:00:00          377.3/377.3 kB 31.7 MB/s eta

# Function to evaluate the RAG pipeline

from rouge_score import rouge_scorer
from bert_score import score as bert_score
from nltk.translate.bleu_score import sentence_bleu, SmoothingFunction
```

```

def evaluate_answers(predictions, references):
    """
        Evaluate generated answers using ROUGE, BERTScore (RAGas proxy),
        and BLEU.
    """

    # -----
    # Sanity: force string-level evaluation
    # -----
    predictions = [
        str(p[0]).strip() if isinstance(p, (list, tuple)) else
str(p).strip()
        for p in predictions
    ]
    references = [str(r).strip() for r in references]

    # Remove empty pairs
    valid_pairs = [(p, r) for p, r in zip(predictions, references) if
p and r]
    if not valid_pairs:
        raise ValueError("No valid prediction-reference pairs found.")

    predictions, references = zip(*valid_pairs)

    # -----
    # ROUGE
    # -----
    rouge = rouge_scorer.RougeScorer(
        ["rouge1", "rouge2", "rougeL"], use_stemmer=True
    )

    rouge_scores = {"rouge1": [], "rouge2": [], "rougeL": []}

    for pred, ref in zip(predictions, references):
        scores = rouge.score(ref, pred)
        for k in rouge_scores:
            rouge_scores[k].append(scores[k].fmeasure)

    avg_rouge = {k: sum(v) / len(v) for k, v in rouge_scores.items()}

    # -----
    # BERTScore (used as RAGas proxy)
    # -----
    P, R, F1 = bert_score(
        list(predictions),
        list(references),
        lang="en",
        rescale_with_baseline=True
    )

```

```

avg_bertscore = {
    "precision": float(P.mean()),
    "recall": float(R.mean()),
    "f1": float(F1.mean())
}

# -----
# BLEU
# -----
smoothie = SmoothingFunction().method4
bleu_scores = []

for pred, ref in zip(predictions, references):
    bleu_scores.append(
        sentence_bleu([ref.split()], pred.split(),
smoothing_function=smoothie)
    )

avg_bleu = sum(bleu_scores) / len(bleu_scores)

return {
    "rouge": avg_rouge,
    "bertscore": avg_bertscore,
    "bleu": avg_bleu
}

print("Evaluation function created")
Evaluation function created

```

### 3.1.3 [3 marks]

Draw inferences by evaluating answers to all questions.

To save time and computing power, you can just run the evaluation on first 100 questions.

```

# Evaluation & Inference (Top 100 Questions)

import torch
from tqdm import tqdm

NUM_EVAL = 100 # change to 100 if required

# Ensure questions & answers are LISTS
if isinstance(question, str):
    question = [q.strip() for q in question.split("\n") if q.strip()]

if isinstance(answer, str):
    answer = [a.strip() for a in answer.split("\n") if a.strip()]

```

```

assert len(question) == len(answer), "Questions and answers count mismatch"

questions_subset = question[:NUM_EVAL]
answers_subset = answer[:NUM_EVAL]

predictions_subset = []

# Generate answers
for q in tqdm(questions_subset, desc="Generating answers"):
    try:
        with torch.no_grad():
            pred = rag_chain(q, top_k=2)

        if isinstance(pred, (list, tuple)):
            pred = pred[0]

        predictions_subset.append(pred)

        # Optional sanity print (first 3)
        if len(predictions_subset) <= 3:
            print("\nQUESTION:", q)
            print("PREDICTED:", pred)
            print("GROUND TRUTH:", answers_subset[len(predictions_subset)-1])
            print("-" * 80)

    except Exception as e:
        predictions_subset.append("")
        print("Generation error:", e)

# Evaluate
metrics_subset = evaluate_answers(predictions_subset, answers_subset)

# Print Results (Inference)
print("\n==== Evaluation Metrics (Top {} Questions)\n".format(NUM_EVAL))

print("\nROUGE Scores:")
for k, v in metrics_subset["rouge"].items():
    print(f"{k}: {v:.4f}")

print("\nBERTScore (RAGas proxy):")
for k, v in metrics_subset["bertscore"].items():
    print(f"{k}: {v:.4f}")

print("\nBLEU Score: {:.4f}\n".format(metrics_subset['bleu']))

Generating answers: 100%|██████████| 1/1 [00:43<00:00, 43.42s/it]

```

QUESTION: Consider "Viber Messenger"'s privacy policy; who will be able to see my information and/or the messages that i send?

PREDICTED: other users

GROUND TRUTH: First of all, we want you to be assured that we do not read or listen to the content of your messages and/or calls made privately via Viber and we do not store those messages once they have been delivered to their destination (which on average takes less than one second).

-----  
-----

Some weights of RobertaModel were not initialized from the model checkpoint at roberta-large and are newly initialized:

['pooler.dense.bias', 'pooler.dense.weight']

You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

==== Evaluation Metrics (Top 100 Questions) ===

ROUGE Scores:

rouge1: 0.0000

rouge2: 0.0000

rougeL: 0.0000

BERTScore (RAGas proxy):

precision: 0.0030

recall: -0.3645

f1: -0.1859

BLEU Score: 0.0000

## 4. Conclusion [5 marks]

### 4.1 Conclusions and insights [5 marks]

#### 4.1.1 [5 marks]

Conclude with the results here. Include the insights gained about the data, model pipeline, the RAG process and the results obtained.

##### 4.1.1 Summary of Findings

###### 1. RAG Pipeline Evaluation

- The Retrieval-Augmented Generation (RAG) pipeline was evaluated on a subset of questions (top 100) from the dataset.
- Generated answers were compared against ground-truth answers using ROUGE, BERTScore, and BLEU metrics.

## 2. Observations from the Evaluation

- **ROUGE Scores:** rouge1 = 0.0, rouge2 = 0.0, rougeL = 0.0
  - Indicates almost no word-level or n-gram overlap between generated and reference answers.
  - This is expected as the model generated **very short answers** compared to long reference answers.
- **BERTScore (RAGas):** precision = 0.0030, recall = -0.3645, f1 = -0.1859
  - Negative F1/recall values indicate that the semantic similarity between predicted and reference answers is extremely low.
  - The model's pretrained weights may not be fully aligned with the task, especially for **detailed, policy-related questions**.
- **BLEU Score:** 0.0
  - Confirms almost no token-level overlap with the ground truth.

## 3. Insights on the Data

- Ground-truth answers are **long, detailed, and often multiple sentences**.
- Predicted answers are **short, generic, and sometimes incomplete**, leading to poor overlap metrics.
- Evaluating at the **string level** ensures that metrics assess entire answers, not individual characters.

## 4. Insights on the Model Pipeline

- The RAG pipeline is functional and returns predictions for each input question.
- Memory and processing speed can be limiting when handling large datasets because:
  - Each prediction involves **transformer-based embedding generation**.
  - Long reference answers increase evaluation cost for ROUGE and BERTScore.
- Adjusting **top-k retrieval** affects output diversity and completeness of generated answers.

## 5. Overall Conclusion

- The pipeline works end-to-end but currently shows **low overlap with reference answers**.
- Improvements can include:
  - Fine-tuning the RAG model on domain-specific data.
  - Optimizing document retrieval and answer generation length.
- Despite low metric scores, the pipeline provides a **scalable framework for question-answering**.

## 6. Key Takeaways

- RAG pipelines generate answers efficiently but may lack accuracy for **long, detailed questions**.
- Preprocessing and data structure are critical: **questions and answers should be lists of full strings**.
- Evaluation metrics highlight **semantic and token-level mismatches**, which can guide future model tuning.