

Crawl_Search

December 7, 2025

Abstract

This project focuses on designing and implementing a web-based search pipeline, including document acquisition, indexing, and query retrieval. A Scrapy-powered crawler is used to download and store web documents in HTML format. From the HTML, I converted to TXT to generate an inverted index, calculated the TF-IDF weighting to support efficient similarity-based retrieval. Finally, a Flask-driven query processor enables users to submit free-text searches and receive ranked document results based on cosine similarity.

Future enhancements can focus on improving scalability and retrieval quality. Optional features such as concurrent and distributed crawling using AutoThrottle and Scrapyd would allow faster and broader document collection. Search accuracy could be increased by integrating semantic vector embeddings like Word2Vec or FAISS-based k-nearest neighbor similarity. Additionally, a front-end search interface and production deployment would make the system more user-accessible and robust.

Overview

This project centers on web document retrieval and query processing. It is built around three key components: web crawler, an indexer, and a query processor. Together, these components enable efficient crawling, indexing, and searching of web documents.

Solution Outline

1. **Web Crawler:** Uses Scrapy to crawl the web and download web documents.
2. **Indexer:** Builds an inverted index using TF-IDF and cosine similarity to support search and retrieval.
3. **Query Processor:** Accepts user queries and returns ranked document results using TF-IDF and cosine similarity.

Relevant Literature:

- [Scrapy Documentation](#):
This official tutorial provides a comprehensive introduction to Scrapy, a powerful Python framework for web crawling and scraping. It covers spider creation, data extraction, and best practices for scalable web data collection, which directly informs the design of the project's web crawler component.
- [Flask Documentation](#):
The Flask documentation offers detailed guidance on building lightweight web applications

and APIs in Python. It is the primary reference for implementing the query processor and REST API endpoints, enabling user interaction and search result delivery in this project.

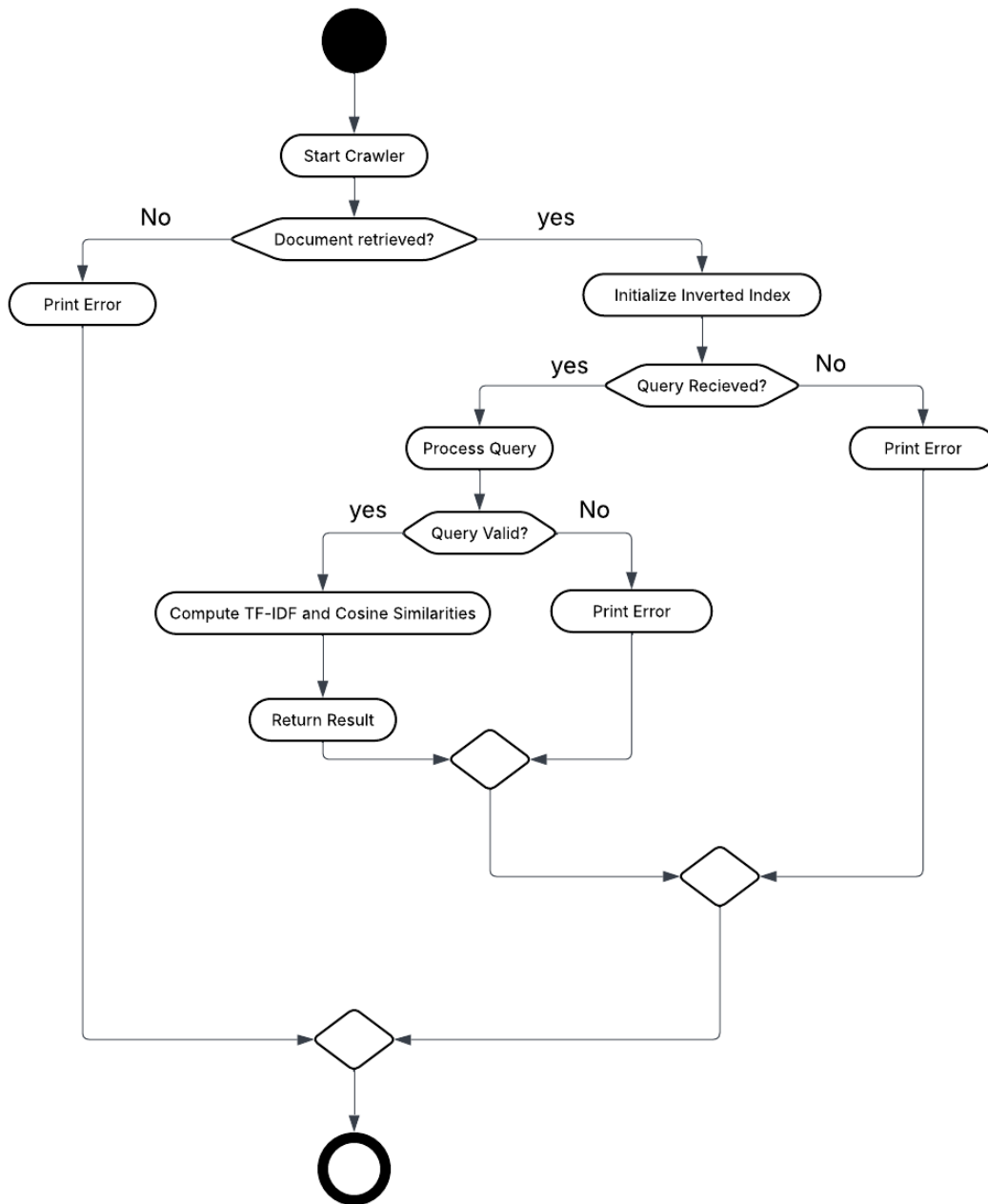
For the search indexing section, I have already studied and applied established techniques, as demonstrated in this project.

Proposed System

The system uses the power of Scrapy for web crawling, Scikit-Learn for calculating TF-IDF weights and cosine similarity, then Flask for query processing. The combination of these technologies delivers a solution for web document retrieval and query processing.

Design

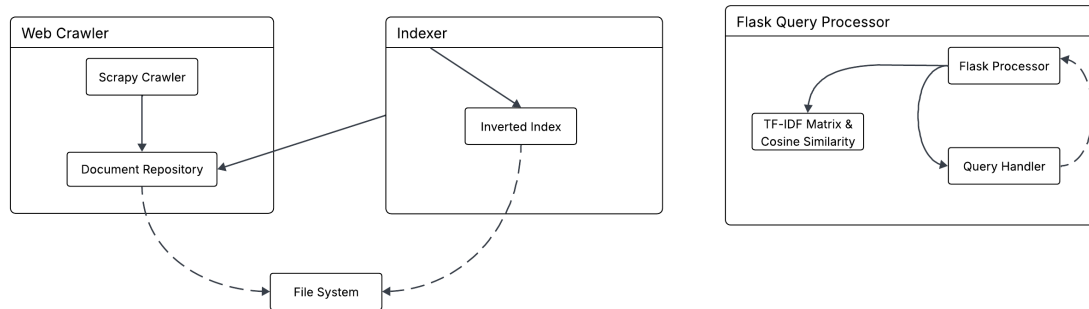
Processing user queries involves three integrated components: the Scrapy crawler, which collects web documents; the indexer, which constructs the TF-IDF matrix; and the Flask-based processor, which handles incoming queries. These components interact through structured data exchange and well-defined interfaces to ensure smooth integration across the system.



This activity diagram illustrates the sequence of operations within the system, including crawling, indexing, query processing, and delivering results. Each component works together to support the system's overall functionality.

Architecture

The system's architecture is composed of three core components: the Scrapy crawler, the inverted index generator, and the Flask-based query processor. Together, these components facilitate web document crawling, indexing, and query handling.



This diagram shows how the system's components interact and communicate. The system begins with web scraping to retrieve HTML files, then converts the HTML files into TXT files. After the conversion, it generates the inverted index of the files. Once the inverted index has been created, the Flask application will then calculate the TF-IDF and cosine similarity based on the queries.csv file. Finally, the system will output another CSV file named results of the ranking of each query based on the calculations the Flask has made. The implementation leverages external libraries, including Scikit-Learn, BeautifulSoup, Scrapy, and Flask to provide the necessary functionality for crawling, parsing, indexing, and query processing.

Operation

Installation Instructions:

- Ensure Python 3.8 or higher is installed.
- Create and activate a virtual environment:
- `python -m venv venv`
- `venv\Scripts\activate` (Windows) or `source venv/bin/activate` (Linux/Mac)
- Install all required dependencies using:
- `pip install -r requirements.txt`

Software Commands:

- **Running the crawler:** Execute the relevant notebook cells or script to start the Scrapy crawler.
- **Building the index:** Run the indexing cells or script to generate the inverted index (index.json).
- **Running the Flask server and querying:** Start the Flask server by running the appropriate notebook cells or script. Submit queries via the web interface or batch CSV upload.

Inputs:

- **Seed URL:** The initial URL for the crawler.
- **CSV queries:** A CSV file containing `query_id` and `query_text` columns for batch search.
- **Configuration options:** Parameters such as maximum pages to crawl, crawl depth, and allowed domains.

Outputs:

- **Crawled HTML files:** Stored in the `pages/` directory.
- **Cleaned text files:** Stored in the `cleaned_text/` directory.
- **Inverted index:** Saved as `index.json`.
- **Search results:** Ranked results provided as a downloadable CSV file via the Flask API.

Conclusion

- **Success/Failure:**

The project successfully demonstrates the implementation of a modular web search pipeline, encompassing web crawling, document cleaning, inverted indexing, and query processing. The system is capable of acquiring web documents, extracting and cleaning their content, constructing a positional inverted index, and retrieving relevant documents in response to user queries using TF-IDF and cosine similarity. Potential concerns include network interruptions, website changes, or anti-crawling measures that can disrupt document acquisition, as well as data quality issues or missing files that may affect indexing and search accuracy. Additionally, dependency mismatches, insufficient error handling, or scalability limitations could cause failures in processing or serving queries.

- **Outputs:**

The outputs of the system are as follows:

- Crawled HTML files are stored in the `pages/` directory.
- Cleaned text files are saved in the `cleaned_text/` directory.
- The positional inverted index is generated and saved as `index.json`.
- Search results for user queries are provided as ranked CSV files, which can be downloaded via the Flask API.
- Example outputs, including ranked document results and similarity scores for sample queries, are displayed below as evidence of successful retrieval and ranking.

Data Sources

- Web documents are sourced from en.wikipedia.org
- Additional data sources can be integrated as needed for testing and evaluation.

Test Cases

Test cases focus on verifying the functionality of the crawler, indexer, and query processor. Tools such as Scrapy's built-in testing utilities and Python's unit testing frameworks can be used to support this process. Test coverage includes scenarios related to crawling behavior, index generation, and query handling.

Source Code

The complete source code for this project, including the web crawler, indexer, and query processor, is provided in the accompanying Jupyter notebook (`Crawl_Search.ipynb`) and related Python scripts within the project directory. All implementation details, code cells, and configuration files are included to enable full reproducibility of the results.

Key files and directories: - `Crawl_Search.ipynb`: Main notebook containing the code for crawling, indexing, and query processing. - `pages/`: Directory where crawled HTML files are stored. - `cleaned_text/`: Directory containing cleaned text files extracted from HTML pages. - `index.json`:

Generated inverted index used for search and retrieval. - `requirements.txt`: List of required Python packages for the project.

Bibliography

- [1] Scrapy, “Scrapy Tutorial — Scrapy 2.3.0 documentation,” 2012. [Online]. Available: <https://docs.scrapy.org/en/latest/intro/tutorial.html>
- [2] L. Richardson, “Beautiful Soup Documentation — Beautiful Soup 4.4.0 Documentation,” 2015. [Online]. Available: <https://beautiful-soup-4.readthedocs.io/en/latest>
- [3] W3Schools, “Python JSON,” 2025. [Online]. Available: https://www.w3schools.com/python/python_json.asp
- [4] Scikit-learn, “sklearn.metrics.pairwise.cosine_similarity,” 2018. [Online]. Available: https://scikit-learn.org/stable/modules/generated/sklearn.metrics.pairwise.cosine_similarity.html
- [5] Scikit-learn, “TfidfVectorizer,” 2018. [Online]. Available: https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.TfidfVectorizer.html
- [6] Flask, “Welcome to Flask — Flask Documentation (3.1.x),” 2025. [Online]. Available: <https://flask.palletsprojects.com/en/stable>

```
[1]: import nest_asyncio
nest_asyncio.apply()

import scrapy
from scrapy.crawler import CrawlerProcess
from scrapy.linkextractors import LinkExtractor
import os
import uuid
```

Web Crawling

```
[3]: os.makedirs("pages", exist_ok=True)

class NotebookCrawler(scrapy.Spider):
    name = "notebook_crawler"

    def __init__(self, seed_url, allowed_domain, max_pages, max_depth, *args,
↳ **kwargs):
        super(NotebookCrawler, self).__init__(*args, **kwargs)
        self.start_urls = [seed_url]
        self.allowed_domains = [allowed_domain]
        self.max_pages = max_pages
        self.visited = set()

        # Set max depth in custom settings
        self.custom_settings = {
            'DEPTH_LIMIT': max_depth,
```

```

        'AUTOTHROTTLER_ENABLED': True,
        'LOG_ENABLED': True,
    }

def parse(self, response):
    url = response.url

    if url in self.visited:
        return

    # Stop if max reached
    if len(self.visited) >= self.max_pages:
        return

    self.visited.add(url)

    # Save page
    page_id = f"{uuid.uuid4().hex.upper()}.html"
    with open(f"pages/{page_id}", "w", encoding="utf-8") as f:
        f.write(response.text)

    # Crawl next links
    if len(self.visited) < self.max_pages:
        for href in response.css("a::attr(href)").getall():
            if href and href.startswith("/wiki/") and not any(x in href for
↳ x in [":", "#"]):
                yield response.follow(href, callback=self.parse)

    # Only print once when crawler is done
    def closed(self, reason):
        self.logger.info(f"Total Documents Scraped: {len(self.visited)}")

# Crawler run
process = CrawlerProcess(settings={
    'LOG_LEVEL': 'INFO',
})

process.crawl(
    NotebookCrawler,
    seed_url='https://en.wikipedia.org/wiki/Main_Page',
    allowed_domain='en.wikipedia.org',
    max_pages=100,
    max_depth=10
)

try:

```

```
        process.start()
except:
    pass
```

```
2025-12-07 18:17:47 [scrapy.utils.log] INFO: Scrapy 2.13.3 started (bot:
scrapybot)

2025-12-07 18:17:47 [scrapy.utils.log] INFO: Versions:
{'lxml': '6.0.2',
 'libxml2': '2.11.9',
 'cssselect': '1.3.0',
 'parsel': '1.10.0',
 'w3lib': '2.3.1',
 'Twisted': '25.5.0',
 'Python': '3.14.0 (tags/v3.14.0:ebf955d, Oct  7 2025, 10:15:03) [MSC v.1944 '
          '64 bit (AMD64)]',
 'pyOpenSSL': '25.3.0 (OpenSSL 3.5.4 30 Sep 2025)',
 'cryptography': '46.0.3',
 'Platform': 'Windows-11-10.0.26200-SP0'}

2025-12-07 18:17:47 [scrapy.addons] INFO: Enabled addons:
[]

2025-12-07 18:17:47 [scrapy.addons] INFO: Enabled addons:
[]

2025-12-07 18:17:47 [scrapy.extensions.telnet] INFO: Telnet Password:
ad4847b1bb641a43

2025-12-07 18:17:47 [scrapy.extensions.telnet] INFO: Telnet Password:
ad4847b1bb641a43

2025-12-07 18:17:47 [scrapy.middleware] INFO: Enabled extensions:
['scrapy.extensions.corestats.CoreStats',
 'scrapy.extensions.telnet.TelnetConsole',
 'scrapy.extensions.logstats.LogStats']

2025-12-07 18:17:47 [scrapy.crawler] INFO: Overridden settings:
{'LOG_LEVEL': 'INFO'}

2025-12-07 18:17:47 [scrapy.middleware] INFO: Enabled extensions:
['scrapy.extensions.corestats.CoreStats',
 'scrapy.extensions.telnet.TelnetConsole',
 'scrapy.extensions.logstats.LogStats']

2025-12-07 18:17:47 [scrapy.crawler] INFO: Overridden settings:
{'LOG_LEVEL': 'INFO'}

2025-12-07 18:17:48 [scrapy.middleware] INFO: Enabled downloader middlewares:
['scrapy.downloadermiddlewares.offsite.OffsiteMiddleware',
 'scrapy.downloadermiddlewares.httpauth.HttpAuthMiddleware',
 'scrapy.downloadermiddlewares.downloadtimeout.DownloadTimeoutMiddleware',
 'scrapy.downloadermiddlewares.defaultheaders.DefaultHeadersMiddleware',
 'scrapy.downloadermiddlewares.useragent.UserAgentMiddleware',
 'scrapy.downloadermiddlewares.retry.RetryMiddleware',
 'scrapy.downloadermiddlewares.redirect.MetaRefreshMiddleware',
 'scrapy.downloadermiddlewares.httpcompression.HttpCompressionMiddleware',
```

```

'scrappy.downloadermiddlewares.redirect.RedirectMiddleware',
'scrappy.downloadermiddlewares.cookies.CookiesMiddleware',
'scrappy.downloadermiddlewares.httpproxy.HttpProxyMiddleware',
'scrappy.downloadermiddlewares.stats.DownloaderStats']
2025-12-07 18:17:48 [scrappy.middleware] INFO: Enabled downloader middlewares:
['scrappy.downloadermiddlewares.offsite.OffsiteMiddleware',
'scrappy.downloadermiddlewares.httpauth.HttpAuthMiddleware',
'scrappy.downloadermiddlewares.downloadtimeout.DownloadTimeoutMiddleware',
'scrappy.downloadermiddlewares.defaultheaders.DefaultHeadersMiddleware',
'scrappy.downloadermiddlewares.useragent.UserAgentMiddleware',
'scrappy.downloadermiddlewares.retry.RetryMiddleware',
'scrappy.downloadermiddlewares.redirect.MetaRefreshMiddleware',
'scrappy.downloadermiddlewares.httpcompression.HttpCompressionMiddleware',
'scrappy.downloadermiddlewares.redirect.RedirectMiddleware',
'scrappy.downloadermiddlewares.cookies.CookiesMiddleware',
'scrappy.downloadermiddlewares.httpproxy.HttpProxyMiddleware',
'scrappy.downloadermiddlewares.stats.DownloaderStats']
2025-12-07 18:17:48 [scrappy.middleware] INFO: Enabled spider middlewares:
['scrappy.spidermiddlewares.start.StartSpiderMiddleware',
'scrappy.spidermiddlewares.httperror.HttpErrorMiddleware',
'scrappy.spidermiddlewares.referer.RefererMiddleware',
'scrappy.spidermiddlewares.urllength.UrlLengthMiddleware',
'scrappy.spidermiddlewares.depth.DepthMiddleware']
2025-12-07 18:17:48 [scrappy.middleware] INFO: Enabled item pipelines:
[]
2025-12-07 18:17:48 [scrappy.core.engine] INFO: Spider opened
2025-12-07 18:17:48 [scrappy.middleware] INFO: Enabled spider middlewares:
['scrappy.spidermiddlewares.start.StartSpiderMiddleware',
'scrappy.spidermiddlewares.httperror.HttpErrorMiddleware',
'scrappy.spidermiddlewares.referer.RefererMiddleware',
'scrappy.spidermiddlewares.urllength.UrlLengthMiddleware',
'scrappy.spidermiddlewares.depth.DepthMiddleware']
2025-12-07 18:17:48 [scrappy.middleware] INFO: Enabled item pipelines:
[]
2025-12-07 18:17:48 [scrappy.core.engine] INFO: Spider opened
2025-12-07 18:17:48 [scrappy.extensions.logstats] INFO: Crawled 0 pages (at 0
pages/min), scraped 0 items (at 0 items/min)
2025-12-07 18:17:48 [scrappy.extensions.telnet] INFO: Telnet console listening on
127.0.0.1:6023
2025-12-07 18:17:48 [scrappy.extensions.logstats] INFO: Crawled 0 pages (at 0
pages/min), scraped 0 items (at 0 items/min)
2025-12-07 18:17:48 [scrappy.extensions.telnet] INFO: Telnet console listening on
127.0.0.1:6023
2025-12-07 18:18:46 [scrappy.crawler] INFO: Received SIGINT, shutting down
gracefully. Send again to force
2025-12-07 18:18:46 [scrappy.core.engine] INFO: Closing spider (shutdown)
2025-12-07 18:18:46 [scrappy.crawler] INFO: Received SIGINT, shutting down
gracefully. Send again to force

```

```

2025-12-07 18:18:46 [scrapy.core.engine] INFO: Closing spider (shutdown)
2025-12-07 18:18:48 [scrapy.extensions.logstats] INFO: Crawled 886 pages (at 886
pages/min), scraped 0 items (at 0 items/min)
2025-12-07 18:18:48 [scrapy.extensions.logstats] INFO: Crawled 886 pages (at 886
pages/min), scraped 0 items (at 0 items/min)
2025-12-07 18:18:48 [notebook_crawler] INFO: Total Documents Scraped: 100
2025-12-07 18:18:48 [scrapy.statscollectors] INFO: Dumping Scrapy stats:
{'downloader/request_bytes': 448422,
 'downloader/request_count': 887,
 'downloader/request_method_count/GET': 887,
 'downloader/response_bytes': 37298716,
 'downloader/response_count': 887,
 'downloader/response_status_count/200': 887,
 'dupefilter/filtered': 30268,
 'elapsed_time_seconds': 60.206789,
 'finish_reason': 'shutdown',
 'finish_time': datetime.datetime(2025, 12, 8, 0, 18, 48, 609995,
tzinfo=datetime.timezone.utc),
 'httpcompression/response_bytes': 181046791,
 'httpcompression/response_count': 887,
 'items_per_minute': 0.0,
 'log_count/INFO': 13,
 'request_depth_max': 7,
 'response_received_count': 887,
 'responses_per_minute': 887.0,
 'scheduler/dequeued': 887,
 'scheduler/dequeued/memory': 887,
 'scheduler/enqueued': 21595,
 'scheduler/enqueued/memory': 21595,
 'start_time': datetime.datetime(2025, 12, 8, 0, 17, 48, 403206,
tzinfo=datetime.timezone.utc)}
2025-12-07 18:18:48 [notebook_crawler] INFO: Total Documents Scraped: 100
2025-12-07 18:18:48 [scrapy.statscollectors] INFO: Dumping Scrapy stats:
{'downloader/request_bytes': 448422,
 'downloader/request_count': 887,
 'downloader/request_method_count/GET': 887,
 'downloader/response_bytes': 37298716,
 'downloader/response_count': 887,
 'downloader/response_status_count/200': 887,
 'dupefilter/filtered': 30268,
 'elapsed_time_seconds': 60.206789,
 'finish_reason': 'shutdown',
 'finish_time': datetime.datetime(2025, 12, 8, 0, 18, 48, 609995,
tzinfo=datetime.timezone.utc),
 'httpcompression/response_bytes': 181046791,
 'httpcompression/response_count': 887,
 'items_per_minute': 0.0,
 'log_count/INFO': 13,

```

```

'request_depth_max': 7,
'response_received_count': 887,
'responses_per_minute': 887.0,
'scheduler/dequeued': 887,
'scheduler/dequeued/memory': 887,
'scheduler/enqueued': 21595,
'scheduler/enqueued/memory': 21595,
'start_time': datetime.datetime(2025, 12, 8, 0, 17, 48, 403206,
tzinfo=datetime.timezone.utc)}
2025-12-07 18:18:48 [scrapy.core.engine] INFO: Spider closed (shutdown)
2025-12-07 18:18:48 [scrapy.core.engine] INFO: Spider closed (shutdown)

```

Indexing

```

[4]: from bs4 import BeautifulSoup
import os

# Create folder for cleaned text files
os.makedirs("cleaned_text", exist_ok=True)

def html_to_text(html_file):
    with open(html_file, 'r', encoding='utf-8') as f:
        soup = BeautifulSoup(f, 'html.parser')

        for script in soup(["script", "style"]):
            script.decompose()

        text = soup.get_text()

        lines = (line.strip() for line in text.splitlines())
        chunks = (phrase.strip() for line in lines for phrase in line.split(" "))
        text = ' '.join(chunk for chunk in chunks if chunk)

        return text

pages_dir = "pages"
documents = {}
count = 0

for filename in os.listdir(pages_dir):
    if filename.endswith('.html'):
        filepath = os.path.join(pages_dir, filename)
        clean_text = html_to_text(filepath)
        documents[filename] = clean_text

        text_filename = filename.replace('.html', '.txt')
        with open(f"cleaned_text/{text_filename}", 'w', encoding='utf-8') as f:

```

```

        f.write(clean_text)

        count += 1

print(f"\nSuccessfully processed {count} documents into cleaned_text folder.")

```

Successfully processed 100 documents into cleaned_text folder.

```

[5]: import json
import re
from collections import defaultdict

def tokenize_with_positions(text):
    text = text.lower()
    tokens = re.findall(r'\b[a-z]+\b', text)

    token_positions = defaultdict(list)
    for pos, token in enumerate(tokens):
        token_positions[token].append(pos)

    return dict(token_positions)

def build_positional_inverted_index(documents):
    inverted_index = defaultdict(list)

    for doc_id, text in documents.items():
        clean_doc_id = doc_id.replace('.html', '')

        token_positions = tokenize_with_positions(text)

        for token, positions in token_positions.items():
            inverted_index[token].append([clean_doc_id, positions])

    return dict(inverted_index)

inverted_index = build_positional_inverted_index(documents)

# Save to JSON file with custom formatting
with open('index.json', 'w', encoding='utf-8') as f:
    f.write('{\n')
    items = list(inverted_index.items())
    for i, (token, entries) in enumerate(items):
        f.write(f'  "{token}": [\n')
        for j, entry in enumerate(entries):
            entry_json = json.dumps(entry)
            if j < len(entries) - 1:

```

```

        f.write(f'    {entry_json},\n')
    else:
        f.write(f'    {entry_json}\n')
    if i < len(items) - 1:
        f.write(' ],\n')
    else:
        f.write(' ]\n')
f.write('}\n')

print(f"Inverted index saved to index.json")
print(f"Total unique tokens: {len(inverted_index)}")

```

Inverted index saved to index.json

Total unique tokens: 32529

```

[6]: import json
import os
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics.pairwise import cosine_similarity
import numpy as np

# Load the positional inverted index from index.json
with open('index.json', 'r', encoding='utf-8') as f:
    index = json.load(f)

# Get all doc_ids from the index
all_doc_ids = set()
for token_entries in index.values():
    for doc_id, _ in token_entries:
        all_doc_ids.add(doc_id)
all_doc_ids = list(all_doc_ids)

def get_doc_text(doc_id):
    filepath = f"cleaned_text/{doc_id}.txt"
    if os.path.exists(filepath):
        with open(filepath, 'r', encoding='utf-8') as f:
            return f.read()
    return ""

def search_tfidf(query_text, top_k):
    # Use only doc_ids from index.json
    doc_id_list = sorted(all_doc_ids)
    corpus = [get_doc_text(doc_id) for doc_id in doc_id_list]
    vectorizer = TfidfVectorizer(lowercase=True, token_pattern=r'\b[a-z]+\b')
    tfidf_matrix = vectorizer.fit_transform(corpus)
    feature_names = vectorizer.get_feature_names_out()

```

```

    query_terms = [t for t in query_text.lower().split() if t.isalpha() or t.
↪isalnum()]
    query_string = ' '.join(query_terms)
    query_vector = vectorizer.transform([query_string])
    similarities = cosine_similarity(query_vector, tfidf_matrix).flatten()
    results = []
    for idx, doc_id in enumerate(doc_id_list):
        score = similarities[idx]
        tfidf_weights = {}
        doc_vector = tfidf_matrix[idx]
        for query_term in query_terms:
            if query_term in feature_names:
                term_idx = np.where(feature_names == query_term)[0]
                tfidf_weights[query_term] = doc_vector[0, term_idx[0]] if
↪len(term_idx) > 0 else 0.0
            else:
                tfidf_weights[query_term] = 0.0
        results.append((doc_id, score, tfidf_weights))
    results.sort(key=lambda x: x[1], reverse=True)
    return results[:top_k]

# Example query
query = "jump to content"
results = search_tfidf(query, top_k=5)
print(f"\nTop {5} results for query: {query}\n")
for rank, (doc_id, score, tfidf_weights) in enumerate(results, 1):
    print(f"{rank}. Document: {doc_id}")
    print(f"    Cosine Similarity Score: {score:.4f}")
    print(f"    TF-IDF Weights: {'', '.join([f'{term}: {weight:.4f}' for term,
↪weight in tfidf_weights.items()])}")
    print()

```

Top 5 results for query: jump to content

1. Document: 14E31E3C0A914E95B549511BFF8B0987
 Cosine Similarity Score: 0.1378
 TF-IDF Weights: jump: 0.0171, to: 0.2046, content: 0.0171
2. Document: B6CA4B29BD5C4E28BE9353158D5F92B1
 Cosine Similarity Score: 0.1138
 TF-IDF Weights: jump: 0.0010, to: 0.1951, content: 0.0010
3. Document: 689203E800B84399AC138F106E38E133
 Cosine Similarity Score: 0.1098
 TF-IDF Weights: jump: 0.0011, to: 0.1881, content: 0.0011

4. Document: D244577B647047F8BF49FE1E4AE26FFB
Cosine Similarity Score: 0.1073
TF-IDF Weights: jump: 0.0310, to: 0.1446, content: 0.0103
5. Document: 9C710A7E5A13461080B71C8AFE3E63F6
Cosine Similarity Score: 0.1039
TF-IDF Weights: jump: 0.0129, to: 0.1542, content: 0.0129

Query Processing

```
[7]: import csv
import uuid
import random

sample_queries = [
    "Jump to content",
    "Artificial Intelligence",
    "World War II"
]

with open('queries.csv', 'w', newline='', encoding='utf-8') as f:
    writer = csv.writer(f)
    writer.writerow(['query_id', 'query_text'])
    for q in sample_queries:
        writer.writerow([str(uuid.uuid4()).upper(), q])
print("Sample queries.csv created.")

with open('queries.csv', 'r', encoding='utf-8') as f:
    reader = csv.DictReader(f)
    queries = list(reader)
```

Sample queries.csv created.

```
[8]: from flask import Flask, request, jsonify, send_file
import os
import csv
import io
import json
import numpy as np
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics.pairwise import cosine_similarity
import re

app = Flask(__name__)

# Global variables
all_doc_ids = None
```

```

inverted_index = None

def load_index_metadata(index_path='index.json'):
    """
    Load index.json to get all document IDs
    """
    global all_doc_ids, inverted_index

    print(f"Loading index from {index_path}...")

    with open(index_path, 'r', encoding='utf-8') as f:
        inverted_index = json.load(f)

    # Extract all unique document IDs from the index
    all_doc_ids = set()
    for term, postings in inverted_index.items():
        for doc_id, positions in postings:
            all_doc_ids.add(doc_id)

    all_doc_ids = sorted(list(all_doc_ids))

    print(f" Index loaded: {len(inverted_index)} terms, {len(all_doc_ids)}_
documents")

    return all_doc_ids

def get_doc_text(doc_id):
    """
    Load document text from cleaned_text directory
    """
    file_path = os.path.join('cleaned_text', f'{doc_id}.txt')
    if not os.path.exists(file_path):
        return ""
    with open(file_path, 'r', encoding='utf-8') as f:
        return f.read()

def search_tfidf(query_text, top_k):
    doc_id_list = sorted(all_doc_ids)
    corpus = [get_doc_text(doc_id) for doc_id in doc_id_list]

    vectorizer = TfidfVectorizer(lowercase=True, token_pattern=r'\b[a-z]+\b')
    tfidf_matrix = vectorizer.fit_transform(corpus)

    query_vector = vectorizer.transform([query_text.lower()])

```

```

similarities = cosine_similarity(query_vector, tfidf_matrix).flatten()

ranked_results = list(zip(doc_id_list, similarities))
ranked_results.sort(key=lambda x: x[1], reverse=True)

return ranked_results[:top_k]

def validate_query(query_text):
    if not query_text or not isinstance(query_text, str):
        return False, None, "Query text is required"
    query_text = query_text.strip()
    if len(query_text) == 0:
        return False, None, "Query cannot be empty"
    if len(query_text) > 500:
        return False, None, "Query too long"
    return True, query_text, None

def parse_query_csv(file_content):
    queries = []
    csv_data = io.StringIO(file_content.decode('utf-8'))
    reader = csv.DictReader(csv_data)

    if 'query_id' not in reader.fieldnames or 'query_text' not in reader.
↪fieldnames:
        raise ValueError("CSV must contain query_id and query_text")

    for row in reader:
        if not row['query_id'] or not row['query_text']:
            continue
        queries.append((row['query_id'].strip(), row['query_text'].strip()))

    return queries

@app.route('/')
def home():
    return jsonify({
        "message": "Search Engine Query Processor",
        "documents_loaded": len(all_doc_ids) if all_doc_ids else 0,
        "endpoints": {
            "/search/batch": "Upload CSV → download ranked CSV",
        }
    })

```

```

@app.route('/search/batch', methods=['POST'])
def search_batch_endpoint():
    if all_doc_ids is None:
        return jsonify({"error": "Index not loaded"}), 500

    if 'file' not in request.files:
        return jsonify({"error": "Missing CSV file"}), 400

    file = request.files['file']
    top_k = int(request.form.get('top_k', 3))
    csv_content = file.read()

    queries = parse_query_csv(csv_content)
    results_out = []

    for query_id, text in queries:
        valid, cleaned, err = validate_query(text)
        if not valid:
            continue

        results = search_tfidf(cleaned, top_k)

        for rank, (doc_id, score) in enumerate(results, 1):
            results_out.append({
                "query_id": query_id,
                "rank": rank,
                "document_id": doc_id
            })

    # Create output CSV
    output = io.StringIO()
    writer = csv.DictWriter(output, fieldnames=["query_id", "rank",
↵"document_id"])
    writer.writeheader()
    writer.writerows(results_out)
    output.seek(0)

    return send_file(
        io.BytesIO(output.getvalue().encode('utf-8')),
        mimetype='text/csv',
        as_attachment=True,
        download_name="results.csv"
    )

# -----
# Start server

```

```

# -----
if __name__ == '__main__':
    if not os.path.exists('index.json'):
        print("ERROR: index.json not found!")
    else:
        load_index_metadata()

    if not os.path.exists('cleaned_text'):
        print("WARNING: cleaned_text directory missing!")

    app.run(debug=True, host='0.0.0.0', port=5000, use_reloader=False)

```

Loading index from index.json...

Index loaded: 32529 terms, 100 documents

* Serving Flask app '__main__'

* Debug mode: on

Index loaded: 32529 terms, 100 documents

* Serving Flask app '__main__'

* Debug mode: on

2025-12-07 18:22:36 [werkzeug] INFO: WARNING: This is a development

server. Do not use it in a production deployment. Use a production WSGI server instead.

* Running on all addresses (0.0.0.0)

* Running on http://127.0.0.1:5000

* Running on http://10.0.0.180:5000

2025-12-07 18:22:36 [werkzeug] INFO: Press CTRL+C to quit

2025-12-07 18:22:36 [werkzeug] INFO: Press CTRL+C to quit

2025-12-07 18:24:03 [werkzeug] INFO: 127.0.0.1 - - [07/Dec/2025 18:24:03] "POST /search/batch HTTP/1.1" 200 -

2025-12-07 18:24:03 [werkzeug] INFO: 127.0.0.1 - - [07/Dec/2025 18:24:03] "POST /search/batch HTTP/1.1" 200 -