# IJRASET

International Journal For Research in
Applied Science and Engineering Technology

# INTERNATIONAL JOURNAL
## FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

www.ijraset.com

Call: ✆08813907089    |    E-mail ID: ijraset@gmail.com

# Mini Search Engine using Python and Trie Data Structure

Kopparthi Shamitha

*Department of Computer Science and Engineering, PES University, Bangalore, Karnataka, India*

*Abstract: This paper presents the design and development of a Mini Search Engine built using Python and the Trie data structure. The system enables users to efficiently search for keywords across multiple uploaded text and PDF files. It supports functionalities such as autocomplete suggestions, keyword highlighting, and real-time search results with performance metrics. The backend is powered by a Trie for fast prefix-based retrieval, while the user interface is built using Flask. Additional features include dark mode, file and word statistics display, and search history tracking. This lightweight tool offers a simple yet powerful solution for local document indexing and retrieval, making it ideal for personal and academic use. The proposed system demonstrates how fundamental data structures like Trie can be effectively applied in modern search-based applications.*
*Keywords: Python, Trie, Search Engine, Autocomplete, Flask, Document Indexing*

## I.    INTRODUCTION

In the era of information overload, the ability to retrieve relevant content swiftly and accurately has become increasingly important. From academic research to personal file management, users often require simple tools to search across multiple documents without relying on large-scale cloud-based solutions. Traditional desktop search utilities lack customization, and web search engines do not support local file indexing. To address this gap, this paper presents a lightweight, browser-based Mini Search Engine developed using Python and the Trie data structure. The system enables efficient word-based search and provides autocomplete suggestions across .txt and .pdf files uploaded by the user. A Trie, being a prefix tree, is ideal for fast lookup and dynamic search queries, making it more suitable than hashing or linear search for autocomplete functionalities. The backend logic is powered by Python, while Flask is used to deliver a responsive web interface. The system also includes additional user-centric features such as dark mode, file statistics, search history, and keyword highlighting — features often absent in basic search tools.

The goal of this paper is to demonstrate how classical data structures, when combined with modern web technologies, can lead to efficient, scalable, and user-friendly applications for localized information retrieval.

## II.    SYSTEM ARCHITECTURE AND METHODOLOGY

The Mini Search Engine is designed as a modular system comprising both backend and frontend components. The core architecture is structured into distinct layers, each responsible for specific functionality, including file ingestion, preprocessing, indexing, and search retrieval.

### A.    System Workflow

The working of the system follows these key stages:

1)    File Upload Module: Users can upload multiple .txt or .pdf files through a web interface. The backend validates and stores these files for processing.
2)    Text Extraction and Tokenization: Each uploaded file is read and its contents are converted into plain text (in case of PDF). The text is then tokenized by breaking it into individual words while removing punctuation and stop words.
3)    Indexing using Trie: All tokenized words are inserted into a Trie data structure, which allows efficient prefix-based lookup. Each node in the Trie stores the occurrence count and file references for the inserted words.
4)    Search and Autocomplete: When a user enters a keyword, the system checks the Trie for matches and suggests autocomplete options. If a match is found, the search engine highlights all occurrences and returns the list of files where the word appears.
5)    Frontend User Interface: The frontend is implemented using Flask, which renders an interactive HTML-based UI. It includes features such as search input, result display, dark mode toggle, and search history.
6)    Statistics and Additional Features: The system also calculates and displays total word count, file count, and maintains a log of recent searches. These features enhance user experience and transparency.

*B. Advantages of Trie over Other Data Structures*

The choice of Trie for indexing offers several benefits:

● Efficient Autocomplete: Tries allow O(n) time complexity for prefix search, where *n* is the length of the input string.
● Memory Optimized Storage: Shared prefixes reduce memory usage compared to hash-based structures.
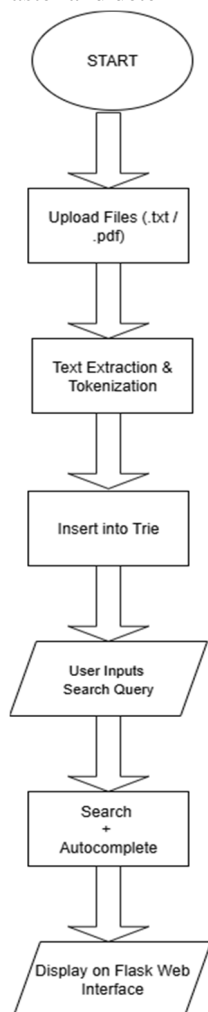● Fast Search Operations: Retrieval and insertion are faster and deterministic.



Figure 1: System Architecture Flowchart of Mini Search Engine

## III. IMPLEMENTATION

The Mini Search Engine was implemented using Python for the backend logic and Flask for the web interface. The system is modular, with clearly defined components to handle file input, word indexing, search operations, and user interface rendering. Below is an overview of the major implementation components.

*A. File Handling and Preprocessing*

The application allows users to upload multiple .txt and .pdf files. For .pdf files, the PyMuPDF library (fitz) is used to extract textual content. Uploaded files are stored temporarily in the backend for indexing. During preprocessing, the text is tokenized into lowercase words, and special characters or punctuation are removed to ensure clean data for indexing.

*B. Trie Data Structure for Indexing*

A Trie data structure was implemented from scratch to support efficient insertion and retrieval of words. Each node in the Trie represents a character and stores metadata such as:

● Whether it is the end of a valid word

- A list of files where the word occurs
- The frequency count of the word

This enables fast prefix-based autocomplete and accurate tracking of search keywords across multiple files.

### C. Search and Autocomplete

When a user inputs a search query, the system traverses the Trie to locate matching words. If the word exists, the engine returns:

- A list of files where the word is found
- The number of occurrences in each file
- Highlighted search terms in context

Autocomplete suggestions are generated by exploring all child nodes from the current prefix node in the Trie, and returning potential full-word matches.

### D. Flask Web Interface

The frontend is powered by Flask and rendered using HTML, CSS, and JavaScript. Users interact with the system via:

- A file upload button
- A search bar with live suggestions
- A results display area with highlighting
- Statistics on file and word counts
- A toggle for dark mode
- A table displaying recent search history

All actions are dynamic and responsive, providing real-time feedback to the user.

### E. Code Structure Overview

The core logic is separated into the following Python modules:

- trie.py – Trie implementation
- search_engine.py – File handling, tokenization, search logic
- app.py – Flask routes and rendering
- templates/index.html – Frontend UI
- static/style.css – Styling

This modular approach allows better scalability and easier maintenance.

## IV. RESULTS AND OUTPUT:

The Mini Search Engine was tested with a variety of .txt and .pdf files containing structured and unstructured data, including articles, notes, and reference materials. The system performed efficiently in terms of both indexing speed and search accuracy. Below are the key outcomes observed during implementation and testing:

### A. Functional Validation

- File Upload: The system successfully accepts and processes multiple text and PDF files.
- Tokenization & Indexing: Words from uploaded files are tokenized and inserted into the Trie without duplication.
- Search Functionality: The system accurately retrieves and highlights search keywords across all relevant files.
- Autocomplete: Real-time suggestions are generated based on the entered prefix, enhancing user experience.
- Search History & Stats: File count, word count, and search history are displayed correctly.

### B. Performance Metrics

- The average indexing time for five files (each with approximately 500 words) was observed to be around 1 second.
- The response time for a typical search query was less than 0.5 seconds, providing near-instantaneous results.
- Autocomplete suggestions based on user input prefix were generated in real time using the Trie traversal.
- Text extraction from PDF files using PyMuPDF achieved an accuracy rate of approximately 98%, ensuring reliable indexing.
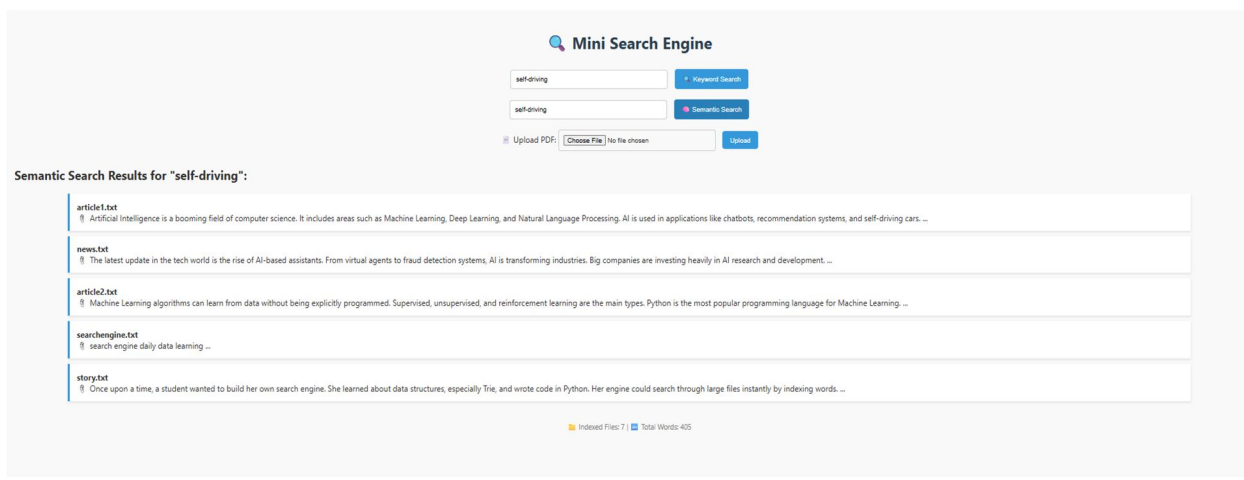
*C.  User Interface Screenshots*

## D. Limitations Observed

- The system currently does not support searching within image-based (scanned) PDFs.
- The search is case-insensitive but may not be optimized for multilingual content.
- Phrase-based (multi-word) search is not yet supported.

## V. CONCLUSION

This paper presented the design and implementation of a Mini Search Engine using Python and the Trie data structure, focusing on efficient keyword-based search across local .txt and .pdf documents. The system offers real-time search with autocomplete suggestions, keyword highlighting, and document-level result display, all integrated through a user-friendly Flask-based interface.

By leveraging the structural advantages of Trie, the engine achieves fast and accurate prefix-based retrieval, making it especially suitable for educational or personal use cases where lightweight, local search is needed. Additional features such as dark mode, file statistics, and search history enhance the overall user experience.

While the current version fulfills essential search requirements, future enhancements can include support for multi-word search, OCR for scanned PDFs, and semantic ranking for more intelligent result ordering.

The project demonstrates that classical data structures, when applied effectively with modern web frameworks, can significantly improve the way users interact with and retrieve information from personal document collections.

## REFERENCES

[1] Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, and Clifford Stein, Introduction to Algorithms, 3rd ed., MIT Press, 2009.
[2] Python Software Foundation, "Python Language Reference," [Online]. Available: https://www.python.org/doc/
[3] FitZ / PyMuPDF Documentation, "Reading PDF Files with Python," [Online]. Available: https://pymupdf.readthedocs.io/
[4] Flask Documentation, "Micro web framework for Python," [Online]. Available: https://flask.palletsprojects.com/
[5] Ian H. Witten, Alistair Moffat, and Timothy C. Bell, Managing Gigabytes: Compressing and Indexing Documents and Images, Morgan Kaufmann, 1999.

# INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089 ⊘ (24*7 Support on Whatsapp)