

Document Content Search using Elasticsearch

Students:

Milconu Adina-Maria

Vlășceanu Daniela-Claudia

Contents

Project Description.....	2
List of Use Cases.....	2
UC1 – Upload Document	2
UC2 – List Documents	2
UC3 – View Document Content	2
UC4 – Search in Document Content	2
UC5 – Delete Document.....	2
Elasticsearch Mapping	3
Implementation description	3
Backend.....	3
Text Extraction.....	3
Search.....	4
Frontend.....	4
API Documentation – Swagger.....	4
Postman Testing	5

Project Description

This project implements a web-based document management and search system that allows users to upload documents (PDF and Word), extract their textual content, store it in Elasticsearch, and perform full-text and partial searches with highlighted results.

The system provides a REST API and a simple web interface where users can:

- Upload documents
- View list uploaded documents
- Delete documents
- Search inside document content using single or multiple terms
- View highlighted matches inside documents

Elasticsearch is used as the search engine, while Apache Tika is used for text extraction from documents. The backend is implemented using Spring Boot (Java).

List of Use Cases

UC1 – Upload Document

- Description: The user uploads a PDF or Word document.
- Result: The document content is extracted and indexed in Elasticsearch.

UC2 – List Documents

- Description: The user requests a list of all uploaded documents.
- Result: The system returns document IDs and filenames.

UC3 – View Document Content

- Description: The user selects a document.
- Result: The full extracted content of the document is displayed.

UC4 – Search in Document Content

- Description: The user searches using one or multiple keywords.
- Result: Matching document fragments are returned with highlighted terms.

UC5 – Delete Document

- Description: The user deletes a document.

- Result: The document is removed from Elasticsearch.

Elasticsearch Mapping

Mapping used in the project

Index: documents

PUT documents

```
{  
  "mappings": {  
    "properties": {  
      "filename": { "type": "keyword" },  
      "content": { "type": "text" },  
      "file_type": { "type": "keyword" },  
      "word_count": { "type": "integer" },  
      "upload_date": { "type": "date" }  
    }  
  }  
}
```

Implementation description

Backend

- Implemented using Spring Boot
- REST APIs handle: upload document, get list of all documents, view content of document, search in the content of a document, delete document

Text Extraction

- Apache Tika extracts content from PDF and Word documents
- Extracted text is indexed into Elasticsearch

Search

- Supports: Partial search, Multi-term search
- Elasticsearch highlights matched terms

Frontend

- HTML + JavaScript UI
- Sidebar displays documents
- Clicking search results opens the correct document with highlight

API Documentation – Swagger

APIs used are documented in swagger:

The screenshot shows the API Documentation – Swagger interface. It displays two main sections: 'Upload API' and 'Documents API'. The 'Upload API' section has one endpoint listed: 'POST /upload Upload a document'. The 'Documents API' section has four endpoints listed: 'GET /search Search in document content', 'GET /documents Get all documents', 'GET /documents/{id} Get document content by ID', and 'DELETE /documents/{id} Delete document by ID'. The 'DELETE' endpoint is highlighted with a red background.

Postman Testing

Collection of Postman was created to test the APIs:

The screenshot shows the Postman application interface. On the left, the 'My Workspace' sidebar lists collections, environments, history, and flows. The 'DocumentContentSearch' collection is selected, and its 'Upload document' POST request is currently being tested. The request URL is `http://localhost:8080/upload`. The 'Body' tab is selected, showing a file input field containing 'Document Content Search using Elasticsearch.docx'. The 'Headers' tab shows a Content-Type header set to 'application/json'. The 'Test Results' section at the bottom shows a single entry: 'filename': 'Document Content Search using Elasticsearch.docx', 'message': 'upload successful'. The status bar indicates a 200 OK response with 145 ms latency and 257 B size.