

# Inverted Index

Jorge Lorenzo Lorenzo, Owen Urdaneta Morales, Luis Perera Pérez,  
Juan López de Hierro, Carlos Jose Moreno Vega, Toffalori Christian

October 20, 2024

## Abstract

The methods used to store information, such as the structure or the order in which data is stored, are crucial when it comes to data access and maintaining coherence between different datasets. While these two aspects do not always align, the ideal approach is to strike a balance between them. In this project, multiple books are downloaded from the Project Gutenberg website using a Crawler, with the aim of processing the words and indexing them in two distinct structures to compare their performance: one based on a dictionary structure, and the other utilizing a local network NoSQL database, MongoDB.

## 1 Introduction

Storing books can be a highly resource-intensive task, especially when the primary goal is to process the words for a user. While real-time reading is possible, if the aim is to create a service-oriented project, the wait time for the user can be substantial. This issue becomes more pronounced as the number of books and their sizes increase. Therefore, it is optimal to preprocess the information to minimize response times for potential clients.

This project, developed in python, applies the same concept of preprocessing words to ensure fast access. This is achieved through two indexing methods: one based on a matrix structure, and the other leveraging a NoSQL database, MongoDB. The process begins by downloading the books from the web, and then, through the execution of these two indexers, we store the words along with relevant information to locate them, such as which books they appear in and their position within those books. The ultimate goal of the program is to provide a search engine for the user.

## 2 Method

### 2.1 Crawler

The Crawler is the part of the code responsible for downloading the data to be used in the project, as well as preprocessing it. In this case, that data consists of books from the Gutenberg Project.

#### Full Code

```
1 import requests
2 import os
3 import csv
4 from bs4 import BeautifulSoup
5 import re
6
7 # Create directories to store books and metadata if they do not exist
8 os.makedirs('datalake/books', exist_ok=True)
9 os.makedirs('datalake', exist_ok=True)
10
11 BOOKS_TO_DOWNLOAD = 5
12 METADATA_CSV_FILE = 'datalake/metadata.csv'
13
14 def execute():
15     """Main function to download eBooks and their metadata."""
```

```

16 id_limit_book_downloaded = check_limit_id_books()
17 id_book = id_limit_book_downloaded + 1
18 count_downloaded_books = 0
19
20 # Create the CSV file and write the header if it doesn't exist
21 if not os.path.exists(METADATA_CSV_FILE):
22     with open(METADATA_CSV_FILE, 'w', newline='', encoding='utf-8') as csvfile:
23         writer = csv.writer(csvfile)
24         writer.writerow(["ID", "Title", "Author", "Release Date", "Most Recently
25                             Updated", "Language"])
26
27 while count_downloaded_books < BOOKS_TO_DOWNLOAD:
28     if not check_book_exist(id_book):
29         title = download_ebook(id_book)
30         if title: # Proceed only if the download was successful
31             metadata = obtain_metadata(id_book, title)
32             if metadata: # Only process if metadata was retrieved successfully
33                 append_metadata_to_csv(metadata)
34                 processing_book(id_book, title)
35                 os.remove(f"datalake/books/{id_book} {title}.txt")
36                 count_downloaded_books += 1
37             id_book += 1
38
39 # Save the ID and overwrite only the first line of the lastBookId.txt file
40 with open("resources/lastBookId.txt", "r+") as file:
41     lines = file.readlines() # Read all lines from the file
42     if lines:
43         lines[0] = str(id_book - 1) + "\n" # Update only the first line
44     else:
45         lines.append(str(id_book - 1) + "\n") # If empty, add the first line
46
47 # Go back to the start of the file and overwrite lines
48 file.seek(0)
49 file.writelines(lines)
50 file.truncate() # Remove any leftover content if lines were shortened
51
52 def download_ebook(id_book):
53     """Download the eBook by its ID and return its title."""
54     try:
55         book_url = f"https://www.gutenberg.org/ebooks/{id_book}"
56         response = requests.get(book_url) # Make a request to the book's URL (ID)
57         if response.status_code == 200:
58             soup = BeautifulSoup(response.text, 'html.parser')
59             title_tag = soup.find('h1') # Find the h1 tag for the title
60             title = title_tag.get_text() if title_tag else f"Unknown Title id {id_book}"
61             text_link = soup.find('a', href=True, string="Plain Text UTF-8") # Find
62                 the link to the plain text file
63             if text_link:
64                 link_txt = 'https://www.gutenberg.org' + text_link['href']
65                 response_txt = requests.get(link_txt)
66                 if response_txt.status_code == 200:
67                     title = re.sub(r'[\W\s]', '', title) # Clean title
68                     with open(f"datalake/books/{id_book} {title}.txt", 'w', encoding='
69                         utf-8') as file:
70                         file.write(response_txt.text) # Write the book content to the
71                         text file
72                     print(f"Book with ID {id_book} has been downloaded.")
73                     return title
74                 else:
75                     print(f"Error accessing the text file of the book with ID {id_book}
76                         : {response_txt.status_code}")
77             else:
78                 print(f"The book with ID {id_book} does not have a text file available
79                     .")
80             else:
81                 print(f"Error accessing the book with ID {id_book}: {response.status_code}
82                     ")
83     except Exception as e:
84         print(f"Error accessing the book with ID {id_book}: {e}")

```

```

79 def obtain_metadata(id_book, title):
80     """Obtain metadata from the downloaded eBook."""
81     try:
82         with open(f"datalake/books/{id_book} {title}.txt", 'r', encoding='utf-8') as
            file:
83             text = file.read()
84             metadata = {
85                 "ID": id_book,
86                 "Title": re.search(r"Title: (.+)", text).group(1) if re.search(r"Title
                    : (.+)", text) else "Unknown Title",
87                 "Author": re.search(r"Author: (.+)", text).group(1) if re.search(r"
                    Author: (.+)", text) else "Unknown Author",
88                 "Release Date": re.search(r"Release Date: (.+)", text).group(1) if re.
                    search(r"Release Date: (.+)", text) else "Unknown",
89                 "Most Recently Updated": re.search(r"Most recently updated: (.+)",
                    text).group(1) if re.search(r"Most recently updated: (.+)", text)
                    else "Unknown",
90                 "Language": re.search(r"Language: (.+)", text).group(1) if re.search(r
                    "Language: (.+)", text) else "Unknown",
91             }
92             return metadata # Return the metadata for further processing
93     except Exception as e:
94         print(f"Error accessing the book with ID {id_book}: {e}")
95         return None # Return None if an error occurred
96
97 def append_metadata_to_csv(metadata):
98     """Append the metadata of the book to the CSV file."""
99     with open(METADATA_CSV_FILE, 'a', newline='', encoding='utf-8') as csvfile:
100         writer = csv.writer(csvfile)
101         writer.writerow([metadata["ID"], metadata["Title"], metadata["Author"],
            metadata["Release Date"],
102             metadata["Most Recently Updated"], metadata["Language"]])
103
104 def processing_book(id_book, title):
105     """Process the downloaded eBook to extract the main content."""
106     try:
107         with open(f"datalake/books/{id_book} {title}.txt", 'r', encoding='utf-8') as
            file:
108             text = file.read()
109             content_start = re.search(r"\*\*\* START OF THE PROJECT GUTENBERG EBOOK .+
                \*\*\*", text)
110             content_end = re.search(r"\*\*\* END OF THE PROJECT GUTENBERG EBOOK .+
                \*\*\*", text)
111
112             if content_start and content_end:
113                 raw_content = text[content_start.end():content_end.start()].strip()
114                 lines = raw_content.splitlines()
115
116                 # Process the lines into paragraphs
117                 paragraphs = []
118                 current_paragraph = []
119                 empty_line_count = 0
120
121                 for line in lines:
122                     stripped_line = line.rstrip()
123                     if stripped_line:
124                         if empty_line_count > 2:
125                             paragraphs.append("") # Add a line break for paragraph
                                separation
126                             current_paragraph.append(stripped_line)
127                             empty_line_count = 0
128                         else:
129                             if current_paragraph:
130                                 paragraphs.append(" ".join(current_paragraph))
131                                 current_paragraph = []
132                                 empty_line_count += 1
133
134                     if current_paragraph:
135                         paragraphs.append(" ".join(current_paragraph))
136
137             final_content = "\n".join(paragraphs)

```

```

138         # Save processed content to a new file
139         with open(f"datalake/books/{id_book}.txt", 'w', encoding='utf-8') as
140             file:
141             file.write(final_content)
142
143     except Exception as e:
144         print(f"Error accessing the book with ID {id_book}: {e}")
145
146 def check_book_exist(id_book):
147     """Check if the book with the given ID already exists in the download directory.
148     """
149     for filename in os.listdir('datalake/books/'):
150         if filename.startswith(f"{id_book}"):
151             print(f"The book with ID {id_book} has already been downloaded.")
152             return True
153     return False
154
155 def check_limit_id_books():
156     """Check the last downloaded book ID from the lastBookId.txt file."""
157     if os.path.exists("resources/lastBookId.txt"):
158         with open("resources/lastBookId.txt", "r") as file:
159             first_line = file.readline().strip() # Read the first line and strip
160             whitespace
161             return int(first_line) if first_line.isdigit() else 0 # Check if it's a
162             number
163     else:
164         with open("resources/lastBookId.txt", "w") as file:
165             file.write("0\n")
166         return 0
167
168 if __name__ == "__main__":
169     execute()

```

Listing 1: Crawler for downloading books from Gutenberg

### 2.1.1 Code Explanation

**Directory Setup** The crawler first ensures that the necessary directories for storing books and metadata exist. This is achieved by using the `os.makedirs()` function, which creates the directories `datalake/books` for storing downloaded books, and `datalake` for storing the corresponding metadata if they do not already exist.

It is important to note that for the crawler to function properly, a directory named `resources` must exist, containing a file called `lastBookId.txt`. This file is essential as it keeps track of the last downloaded book's ID, allowing the crawler to resume downloading from the next available book. Without this file, the system would not be able to determine which book to start downloading next.

**Main Execution Flow** The core functionality is handled by the `execute()` function, which controls the downloading and processing of books. It first checks the ID of the last downloaded book by calling the `check_limit_id_books()` function, which reads from the `lastBookId.txt` file. Starting from the next book ID, the crawler attempts to download and process up to 5 books, as defined by the `BOOKS_TO_DOWNLOAD` constant.

If a book has already been downloaded, this is verified by the `check_book_exist()` function, which looks for the book in the `datalake/books` directory. If the book has not yet been downloaded, it proceeds to the `download_ebook()` function, which attempts to download the book and save it locally.

Once the book is successfully downloaded, the `obtain_metadata()` function retrieves relevant metadata, such as the title, author, release date, and language, by parsing the book's content. This metadata is then saved in a CSV file using the `append_metadata_to_csv()` function.

Afterward, the book content is processed by the `processing_book()` function to extract only the main text, removing any headers, footers, and irrelevant content. Finally, the processed book file is deleted to free up space, and the loop continues until 5 books have been successfully downloaded and processed.

**Downloading Books** The `download_ebook()` function is responsible for downloading a book from the Gutenberg Project by its ID. It constructs the URL for the book and retrieves the HTML content of the page. Using BeautifulSoup, the crawler identifies the title and the link to the plain text version of the book. If the plain text version is available, the book's content is downloaded and saved locally. The title is also cleaned to remove any special characters to ensure a valid file name.

If the book does not have a plain text version, or if any error occurs during the download process, appropriate error messages are logged.

**Extracting Metadata** Once a book is downloaded, the `obtain_metadata()` function extracts key metadata from the text, such as the title, author, release date, and language. This is done using regular expressions, which search for specific metadata tags within the text. If the metadata is found, it is stored in a dictionary and returned to be saved in the metadata CSV file.

If any error occurs while accessing or extracting the metadata, the function logs an error message and returns `None` to avoid interrupting the workflow.

**Processing Book Content** After the book is downloaded and the metadata is extracted, the `processing_book()` function processes the content to isolate the main text of the book. This involves stripping out the header and footer content added by the Gutenberg Project, which is enclosed between specific tags ("START OF THE PROJECT GUTENBERG EBOOK" and "END OF THE PROJECT GUTENBERG EBOOK").

The main content is then processed into paragraphs, and excessive line breaks are removed to ensure a cleaner format. The processed content is saved into a new file, which only contains the body of the book.

**Metadata Storage** All the collected metadata, such as the book's ID, title, author, release date, and language, is stored in a CSV file named `metadata.csv`. This allows easy reference and provides a structured way to access the information later. The function `append_metadata_to_csv()` ensures that each downloaded book's metadata is appended to the CSV file.

**Book ID Management** To prevent downloading the same books multiple times, the crawler keeps track of the last downloaded book's ID in the file `lastBookId.txt`. The `check_limit_id_books()` function reads this file to retrieve the ID of the last downloaded book, while `execute()` updates the file each time new books are downloaded. This process ensures that downloads can resume seamlessly in subsequent executions.

**Error Handling and Logging** The crawler includes comprehensive error handling throughout its functions to ensure robustness. If any error occurs while downloading a book, extracting metadata, or processing the content, the error is logged, and the crawler continues with the next book. This ensures that one failed download does not interrupt the entire process.

## 2.2 Inverted Indexer

In this work, an inverse indexer has been implemented to process a collection of books and store each word in individual JSON files. These files are organized into directories labeled from A to Z, where each folder contains words that begin with that letter. The indexer processes books that have not been previously indexed, ensuring no duplicate processing, and updates the corresponding file with the number of occurrences of the word, along with the line and position within the line.

The indexing process ensures that each word from a book is stored in a dedicated JSON file. For example, if the word "book" appears in multiple books, all occurrences are recorded in a single `book.json` file, located in the directory corresponding to the letter B. Inside this file, the exact position of each appearance of the word is stored, indicating the line number and the word's position in that line.

To prevent repeated processing of books, the system uses a tracking file, `resources/lastBookId.txt`, where the ID of the last indexed book is stored. This way, every time the program is executed, the indexer starts with the next book, ensuring that only new books are processed. At the end of the process, the tracking file is updated with the ID of the most recently indexed book.

A key aspect of this indexer is the lemmatization of words. Using the `simplemma` library, words are normalized so that morphological variations like plurals or verb conjugations are stored under the same root. This means that words like "books" and "book" are grouped into the same JSON file, facilitating more efficient searches. Additionally, the use of `unidecode` normalizes words by removing special characters and accents, ensuring greater consistency in the stored data.

The indexer also tracks the exact line number and the position within each line where a word appears. This functionality provides precise context about the location of each word in the indexed books, which is essential for conducting detailed reverse searches.

## Full Code

```

1 import os
2 import json
3 import re
4 from unidecode import unidecode
5 from simplemma import lemmatize
6
7 # Create necessary directories
8 OUTPUT_DIR = 'datamart/reverse_indexes/'
9 GLOBAL_INDEX_FILE = 'datamart/global_index.txt'
10
11 def create_directory(path):
12     os.makedirs(path, exist_ok=True)
13
14 def read_lines(file_path):
15     try:
16         with open(file_path, 'r') as file:
17             return file.readlines()
18     except FileNotFoundError:
19         print(f"The file '{file_path}' does not exist.")
20         return []
21
22 def write_lines(file_path, lines):
23     with open(file_path, 'w') as file:
24         file.writelines(lines)
25
26 def read_json(file_path):
27     if os.path.exists(file_path):
28         with open(file_path, 'r', encoding='utf-8') as json_file:
29             return json.load(json_file)
30     return None
31
32 def write_json(file_path, data):
33     with open(file_path, 'w', encoding='utf-8') as json_file:
34         json.dump(data, json_file, ensure_ascii=False, indent=4)
35
36 def read_global_index():
37     global_index_dict = {}
38     if os.path.exists(GLOBAL_INDEX_FILE):
39         lines = read_lines(GLOBAL_INDEX_FILE)
40         for line in lines:
41             line = line.strip()
42             if line:
43                 parts = line.split(" : ")
44                 if len(parts) == 2:
45                     word, books = parts
46                     global_index_dict[word] = books.split(", ")
47     return global_index_dict
48
49 def write_global_index(global_index_dict):
50     lines = [f"{word} : {' '.join(books)}\n" for word, books in global_index_dict.items()]
51     write_lines(GLOBAL_INDEX_FILE, lines)
52
53 def clean_word(word):
54     word = re.sub(r'^_.*|_.*$|[\d]+|[\^\\w\s]', '', word)
55     return unidecode(word).strip('_')
56

```

```

57 def lemm_add(dictionary, word, lang='en'):
58     lemm = lemmatize(word, lang)
59     if lemm != word:
60         if lemm not in dictionary:
61             dictionary[lemm] = {'allocations': {}, 'total': 0}
62         if "allocations" in dictionary[word]:
63             pop = dictionary.pop(word)
64             for book_key, new_info in pop['allocations'].items():
65                 if book_key in dictionary[lemm]['allocations']:
66                     existing_positions = set(tuple(pos) for pos in dictionary[lemm]['
67                         allocations'][book_key]['position'])
68                     existing_positions.update(tuple(pos) for pos in new_info['position
69                         '])
70                     dictionary[lemm]['allocations'][book_key]['position'] = list(
71                         existing_positions)
72                     dictionary[lemm]['allocations'][book_key]['times'] = len(
73                         existing_positions)
74                 else:
75                     dictionary[lemm]['allocations'][book_key] = new_info
76                     dictionary[lemm]['total'] += pop['total']
77
78 def index_book(book_id):
79     dictionary = {}
80     book_file_path = f"datalake/books/{book_id}.txt"
81
82     with open(book_file_path, encoding='utf-8') as fp:
83         for line_number, line in enumerate(fp, start=1):
84             words = [clean_word(word) for word in line.lower().split()]
85             for position, word in enumerate(words):
86                 if word:
87                     if word not in dictionary:
88                         dictionary[word] = {'allocations': {}, 'total': 0}
89                         book_key = f"BookID_{book_id}"
90
91                     if book_key not in dictionary[word]['allocations']:
92                         dictionary[word]['allocations'][book_key] = {'times': 0, '
93                             position': []}
94
95                     pos_tuple = (line_number, position + 1)
96                     if pos_tuple not in dictionary[word]['allocations'][book_key]['
97                         position']:
98                         dictionary[word]['allocations'][book_key]['position'].append(
99                             pos_tuple)
100                     dictionary[word]['allocations'][book_key]['times'] += 1
101
102                     dictionary[word]['total'] += 1
103
104     for word in list(dictionary.keys()):
105         lemm_add(dictionary, word)
106
107     save_index(dictionary, book_id)
108
109 def save_index(dictionary, book_id):
110     reserved_words = ['con', 'prn', 'aux', 'nul', 'com1', 'com2', 'com3',
111         'com4', 'com5', 'com6', 'com7', 'com8', 'com9',
112         'lpt1', 'lpt2', 'lpt3', 'lpt4', 'lpt5', 'lpt6',
113         'lpt7', 'lpt8', 'lpt9']
114
115     for word, data in dictionary.items():
116         first_letter = word[0].lower()
117         if word.lower() in reserved_words:
118             print(f"Word '{word}' skipped because it is a reserved name.")
119             continue
120
121         create_directory(f"datamart/reverse_indexes/{first_letter}")
122         json_file_path = f"datamart/reverse_indexes/{first_letter}/{word}.json"
123
124         existing_data = read_json(json_file_path) or {
125             'word': word,
126             'allocations': data['allocations'],
127             'total': 0

```

```

121     }
122
123     if existing_data:
124         existing_allocations = existing_data['allocations']
125         for book_key, new_info in data['allocations'].items():
126             if book_key in existing_allocations:
127                 existing_positions = set(tuple(pos) for pos in
128                                         existing_allocations[book_key]['position'])
129                 existing_positions.update(tuple(pos) for pos in new_info['position'])
130                 existing_allocations[book_key]['position'] = list(
131                     existing_positions)
132                 existing_allocations[book_key]['times'] = len(existing_positions)
133             else:
134                 existing_allocations[book_key] = new_info
135
136         existing_data['total'] = sum(info['times'] for info in
137                                     existing_allocations.values())
138
139     write_json(json_file_path, existing_data)
140
141     update_global_index(dictionary, book_id)
142
143 def update_global_index(dictionary, book_id):
144     global_index_dict = read_global_index()
145
146     for word in dictionary.keys():
147         if word in global_index_dict:
148             if f"BookID_{book_id}" not in global_index_dict[word]:
149                 global_index_dict[word].append(f"BookID_{book_id}")
150         else:
151             global_index_dict[word] = [f"BookID_{book_id}"]
152
153     write_global_index(global_index_dict)
154
155 def read_numbers_from_file(file_path):
156     try:
157         lines = read_lines(file_path)
158         first_number = int(lines[0].strip())
159         second_number = int(lines[1].strip()) if len(lines) > 1 else 1
160
161         write_lines(file_path, [f"{first_number}\n", f"{second_number}\n"])
162         return first_number, second_number
163     except ValueError:
164         print("The file content is not a valid integer.")
165
166 def execute():
167     n, m = read_numbers_from_file("resources/lastBookId.txt")
168     for i in range(m, n + 1):
169         index_book(i)
170         print(f"Book with ID {i} has been indexed.")
171
172 if __name__ == "__main__":
173     execute()

```

Listing 2: Inverted indexer Code

### 2.2.1 Code Explanation

**Directory and File Handling** The code begins by defining helper functions for file operations like reading and writing text files (`read_lines()` and `write_lines()`), and working with JSON files (`read_json()` and `write_json()`). These are critical for interacting with book content and saving the indexed word information.

**Word Normalization** The `clean_word()` function cleans individual words by removing special characters, numbers, and diacritics. This ensures that words are stored in a uniform format. The `lemm_add()` function is responsible for adding lemmatized words to the dictionary, grouping together different word forms (e.g., "running" and "run") under a single entry.



**Book Indexing** The core functionality of indexing each book is handled by `index_book()`. It processes the book file line by line, identifies the words in each line, and tracks their position within the book. For each word, it creates an "allocation" object that stores the positions where the word appears in the book.

**Saving the Index and Reserved Words** After processing, the `save_index()` function saves the word allocations and positions in a JSON file. Words that are reserved in the system (like "con", "aux", etc.) are skipped to avoid system conflicts. The file path for each word is based on its first letter, which helps organize the data for fast access.

**Global Index Update** The `update_global_index()` function updates a global index file that keeps track of which books contain each word. This ensures that the system can quickly identify the presence of words across all books in the system.

**Tracking the Last Book** The `read_numbers_from_file()` function reads the last processed book's ID from `lastBookId.txt`, and the `execute()` function ensures the indexer resumes from the correct point, allowing new books to be indexed without reprocessing old ones.

## 2.3 MongoDB Indexer

In order to use MongoDB, it is necessary to carry out some preliminary installations. First, in the Python environment where the project is being developed, execute the following installation command in the terminal:

```
1 pip install pymongo
```

This is necessary to perform MongoDB operations through Python code. Afterward, you need to install the MongoDB services by downloading them from their official website:

<https://www.mongodb.com/try/download/community>

Since this project runs a local database, it is required to install the Community Server version of MongoDB. It is recommended to download and install the '.msi' file, as it provides a guided setup. Perform a full installation and accept the installation of MongoDB Compass, which is the GUI used to visually inspect the database content.

At this point, it is important to highlight that these steps have been executed on a Windows operating system, as it involves interacting with directories and system variables. After the installation is complete, navigate to your system drive and create a new folder called 'data' with another folder inside named 'db'. This directory is where MongoDB will store the data from the connections.

Once these steps are completed, go to the Program Files directory on your system and find the MongoDB directory. Then, proceed to enter the 'Server', '8.0' (or whichever version you choose to install; in this project, version 8.0 was used), and 'bin' directories. Copy the directory path. Afterward, access the system's environment variables and locate the 'Path' under the user variables. Select it and add a new entry with the copied directory path.

After doing this, you should be able to launch MongoDB by running the command 'mongod' from the system terminal, which will start a local network database that will continue to listen as long as it remains running.

Finally, using MongoDB Compass, you can connect to the database by adding the connection string 'mongodb://localhost:27017' to visually manage the databases in the connection.

### Full Code

```
1 import os
2 from pymongo import MongoClient
3 import re
4 from unicode import unicode
5 from simplemma import lemmatize
6
7 # Configure the MongoDB connection
8 client = MongoClient('mongodb://localhost:27017/')
9 db = client['SEARCH_ENGINE'] # Database name
```

```

10 words_collection = db['WORDS'] # Collection name for words
11 processed_books_collection = db['PROCESSED_BOOKS'] # Collection name for processed
    books
12
13 # Predefined pre-path
14 PRE_PATH = "datalake/books/"
15
16 # Clean and lemmatize words
17 def clean_word(word):
18     word = re.sub(r'^_*|_.*_|_.*$|[\d]+|[\^w\s]', '', word)
19     word = unicode(word).strip('_').strip()
20     return word if word else None
21
22 def lemm_add(dictionary, word, lang='en'):
23     if word:
24         lemm = lemmatize(word, lang)
25         if lemm != word:
26             if lemm not in dictionary:
27                 dictionary[lemm] = {'allocations': {}, 'total': 0}
28             if "allocations" in dictionary[word]:
29                 pop = dictionary.pop(word)
30                 for book_key, new_info in pop['allocations'].items():
31                     if book_key in dictionary[lemm]['allocations']:
32                         existing_positions = set(tuple(pos) for pos in dictionary[lemm]
33                                                 ['allocations'][book_key]['position'])
34                         existing_positions.update(tuple(pos) for pos in new_info['position'])
35                         dictionary[lemm]['allocations'][book_key]['position'] = list(
36                             existing_positions)
37                         dictionary[lemm]['allocations'][book_key]['times'] = len(
38                             existing_positions)
39                     else:
40                         dictionary[lemm]['allocations'][book_key] = new_info
41                         dictionary[lemm]['total'] += pop['total']
42
43 def process_book(file_name):
44     file_path = os.path.join(PRE_PATH, file_name)
45     book_name = file_name
46
47     if processed_books_collection.find_one({'book_name': book_name}):
48         print(f"The book '{book_name}' has already been processed.")
49         return
50
51     with open(file_path, 'r', encoding='utf-8') as file:
52         text = file.read()
53
54     word_occurrences = {}
55     paragraphs = text.split("\n\n")
56
57     for i, paragraph in enumerate(paragraphs):
58         words = [clean_word(word) for word in paragraph.lower().split()]
59         words = [word for word in words if word]
60
61         for word in set(words):
62             word_count = words.count(word)
63
64             if word not in word_occurrences:
65                 word_occurrences[word] = {'count': 0, 'books': {}}
66                 word_occurrences[word]['count'] += word_count
67
68             if book_name not in word_occurrences[word]['books']:
69                 word_occurrences[word]['books'][book_name] = {}
70
71                 word_occurrences[word]['books'][book_name][f"paragraph {i}"] = word_count
72
73     for word in list(word_occurrences.keys()):
74         lemm_add(word_occurrences, word)
75
76     existing_data = words_collection.find_one({'word': word})
77
78     if existing_data:

```

```

76         new_count = existing_data['count'] + word_occurrences[word]['count']
77
78     for book, paragraphs in word_occurrences[word]['books'].items():
79         if book in existing_data['books']:
80             for paragraph, count in paragraphs.items():
81                 if paragraph not in existing_data['books'][book]:
82                     existing_data['books'][book][paragraph] = count
83                 else:
84                     existing_data['books'][book][paragraph] += count
85             else:
86                 existing_data['books'][book] = paragraphs
87
88     words_collection.update_one(
89         {'word': word},
90         {'$set': {'count': new_count, 'books': existing_data['books']}}
91     )
92     print(f"Updated {word}: {new_count} occurrences across books.")
93 else:
94     words_collection.insert_one({
95         'word': word,
96         'count': word_occurrences[word]['count'],
97         'books': word_occurrences[word]['books']
98     })
99     print(f"Inserted {word}: {word_occurrences[word]['count']} occurrences.")
100
101     processed_books_collection.insert_one({'book_name': book_name})
102     print(f"The book '{book_name}' has been processed and recorded.")
103
104 def execute():
105     choice = input("Do you want to process a (1) single file or (2) all files in the
106         folder? (1/2): ")
107     if choice == '1':
108         file_name = input("Enter the name of the file to process (without the path): ")
109         process_book(file_name)
110     elif choice == '2':
111         folder_path = PRE_PATH
112         for filename in os.listdir(folder_path):
113             if filename.endswith('.txt'):
114                 process_book(filename)
115                 print(f"Processed file: '{filename}'")
116     else:
117         print("Invalid choice. Please try again.")
118
119     print("Data saved in the MongoDB database.")
120
121 if __name__ == "__main__":
122     execute()

```

Listing 3: MongoDB Indexer Code

### 2.3.1 Code Explanation

The MongoDB indexer script serves to process the text files of books stored locally and to store the relevant word data in a MongoDB database. The process includes cleaning, lemmatizing, and organizing words to optimize their searchability in a future retrieval system. Below is an explanation of the major components of the script:

**MongoDB Connection Setup:** The script begins by establishing a connection to a local MongoDB instance using the `MongoClient` from the `pymongo` library. It connects to the `SEARCH_ENGINE` database and creates two collections: `WORDS` and `PROCESSED_BOOKS`. Importantly, MongoDB does not require the database or collections to be created in advance. They are automatically generated when the first write operation occurs.

**Cleaning and Lemmatizing Words:** The `clean_word()` function is used to clean the text by removing unwanted characters, digits, and punctuation. The `lemm.add()` function then lemmatizes

the words using the `simpllemma` library. This process normalizes the words, reducing them to their base forms, and ensures that similar words are grouped together.

**Processing a Book:** The `process_book()` function handles the main logic for processing the text of each book. It reads the text file, splits it into paragraphs, and extracts words from each paragraph. Each word's occurrence is counted and stored, along with its paragraph and the total count across the book. The processed word occurrences are then either inserted into the MongoDB collection or updated if they already exist.

**Lemmatization and Insertion in MongoDB:** After the words are cleaned and lemmatized, they are checked against the existing entries in the MongoDB `WORDS` collection. If the word already exists, its count and paragraph occurrences are updated. If not, it is inserted as a new entry. MongoDB automatically manages the insertion and updating of documents in the collections.

**Storage Format:** The word occurrences are stored in MongoDB using a nested document structure. Each word is a document in the `WORDS` collection, and the associated metadata for each word is stored in a dictionary format. The structure of each word document is as follows:

```
{
  "word": "example", # The word itself
  "count": 15, # Total occurrences across all books
  "books": {
    "book_name_1": {
      "paragraph 1": 3, # Occurrences in paragraph 1
      "paragraph 4": 2 # Occurrences in paragraph 4
    },
    "book_name_2": {
      "paragraph 2": 1
    }
  }
}
```

This structure allows for efficient querying of specific words, as well as tracking in which books and paragraphs each word appears. MongoDB's flexible schema makes it ideal for this kind of dynamic storage, where the number of books or paragraphs can vary from word to word. The `PROCESSED_BOOKS` collection keeps track of books that have already been processed to avoid redundant operations.

**Execution Flow:** The `execute()` function allows the user to either process a single file or all files in the folder. Once processed, the data is stored in MongoDB, allowing for future queries and retrievals to be executed efficiently.

## 2.4 Inverted Indexer API

This work implements a terminal-based API that allows users to search for words within a collection of books that have been previously indexed. The user enters a word into the terminal, and the system looks up that word in the corresponding JSON files, organized by the first letter of the word. If the word is found in the index, the system returns detailed information about its occurrences, including the book title, the exact line where the word appears, and its position within the line. Furthermore, the system considers lemmatization, so variations of the word, such as plurals or conjugations, are treated as the same base word. It also handles case-insensitive searches, ensuring that both uppercase and lowercase versions of the word are correctly matched. Below is a detailed explanation of the process and the key parts of the code that make this functionality possible.

### Full Code

```

1 import csv
2 import json
3 from unicode import unicode
4 from simplemma import lemmatize
5
6 # Function to obtain the book title from the metadata.csv file
7 def get_book_title(book_id):
8     with open('datalake/metadata.csv', mode='r', encoding='utf-8') as file:
9         reader = csv.DictReader(file)
10        for row in reader:
11            if row['ID'] == str(book_id):
12                return row['Title']
13        return f"BookID_{book_id}" # If the title is not found, return the BookID.
14
15 # Function to get word occurrences from the corresponding JSON file
16 def get_word_occurrences_in_books(word):
17     try:
18         with open(f"datamart/reverse_indexes/{word[0].lower()}/{word}.json", 'r',
19                 encoding='utf-8') as json_file:
20             existing_data = json.load(json_file)
21             return existing_data['allocations'] # Returns a dictionary with books and
22                                                # positions
23     except FileNotFoundError:
24         raise ValueError(f"The word '{word}' has not been indexed yet.")
25
26 # Function to read paragraphs and highlight the word where it appears
27 def read_paragraphs(book_id, occurrences, search_word):
28     with open(f"datalake/books/{book_id}.txt", encoding='utf-8') as fp:
29         result = []
30         current_occurrence = 0
31         occurrences.sort(key=lambda x: x[0]) # Ensure positions are in order
32         for i, line in enumerate(fp, start=1):
33             # If the current line is where the word appears
34             while current_occurrence < len(occurrences) and i == occurrences[
35                 current_occurrence][0]:
36                 # Highlight the word in purple using ANSI escape codes
37                 highlighted_line = line.replace(search_word, f"\033[35m{search_word}
38                 \033[0m")
39                 result.append(f"Line {i}: {highlighted_line.strip()} (Position {
40                     occurrences[current_occurrence][1]})")
41                 current_occurrence += 1
42                 if current_occurrence >= len(occurrences):
43                     break
44         return result
45
46 # Main function to search for a word across all books
47 def search_word_across_books(word, lang='en'):
48     word = unicode(word.lower()).strip() # Clean the word
49     word = lemmatize(word, lang) # Lemmatize the word
50
51     try:
52         occurrences_by_book = get_word_occurrences_in_books(word) # Retrieve books
53         # and positions
54         if occurrences_by_book:
55             print(f"The word '{word}' appears in the following books:")
56             for book_id, info in occurrences_by_book.items():
57                 title = get_book_title(book_id.split('_')[-1]) # Get the book title
58                 # from metadata.csv
59                 valid_positions = [pos for pos in info['position'] if pos]
60                 if valid_positions:
61                     print(f"\nBook Title: {title}")
62                     paragraphs = read_paragraphs(book_id.split('_')[-1],
63                                                 valid_positions, word)
64                     for paragraph in paragraphs:
65                         print(paragraph)
66             else:
67                 print(f"The word '{word}' was not found in any books.")
68     except ValueError as e:
69         print(e)

```

```

63 if __name__ == "__main__":
64     word = input("Enter the word to search: ")
65     search_word_across_books(word)

```

Listing 4: Inverted Indexer Word Search API

## Code Explanation

**Retrieving Book Titles** The `get_book_title()` function is responsible for obtaining the title of a book based on its ID. This function reads the `metadata.csv` file where book metadata is stored, such as the book ID and title. If the title is found, it is returned. If not, the book's ID is returned as a fallback. This ensures that even if the title isn't available, the system has a way to identify the book.

**Looking Up Word Occurrences** The `get_word_occurrences_in_books()` function retrieves the occurrences of a word from the corresponding JSON file. The file is organized in folders based on the first letter of the word. If the word is not found (i.e., if the file does not exist), the function raises an exception, informing the user that the word has not yet been indexed.

**Displaying Word Occurrences** The function `read_paragraphs()` is responsible for displaying the occurrences of the word within each book. It reads through the text file of the book and highlights the occurrences of the word using ANSI escape codes. These escape codes change the color of the word to purple when displayed in the terminal. The function also shows the line and position of each occurrence of the word, which helps the user locate the word precisely within the book.

**Main Search Function** The main function, `search_word_across_books()`, ties everything together. It first normalizes the word by converting it to lowercase and applying lemmatization to handle variations (such as singular/plural or verb forms). Then, it retrieves the books where the word occurs using `get_word_occurrences_in_books()`. For each book, it fetches the title and displays the lines where the word appears using the `read_paragraphs()` function. If the word is not found in any book, a message is displayed indicating that no occurrences were found.

## 2.5 MongoDB Indexer API

The API performs a search on a MongoDB collection where the words are stored. Each word is associated with a list of books where it appears, and in each book, the paragraphs and the number of times the word occurs in those paragraphs are saved. When the user enters a word in the terminal, the API queries MongoDB, searches for the corresponding books, and returns the exact positions within the paragraphs where the word was found.

In addition to returning the results, the API uses the `lemmatize` function to normalize the words, so that variations such as "book" and "books" are considered the same word. It also handles case-insensitivity to ensure that all possible occurrences are captured.

Once the information is obtained, the API prints the book titles, which are extracted from a `metadata.csv` file, and then shows the lines and positions where the word appears, highlighting it in magenta to enhance visualization.

## Full Code

```

1 import csv
2 import re
3 from pymongo import MongoClient
4 from unidecode import unidecode
5 from simplemma import lemmatize
6 from termcolor import colored
7
8 # MongoDB connection configuration
9 client = MongoClient('mongodb://localhost:27017/')
10 db = client['SEARCH_ENGINE']
11 words_collection = db['WORDS']
12

```

```

13 # Load metadata from the CSV file to get the book title
14 def get_book_title_from_csv(book_id):
15     with open('datalake/metadata.csv', mode='r', encoding='utf-8') as file:
16         csv_reader = csv.DictReader(file)
17         for row in csv_reader:
18             if row['ID'] == book_id:
19                 return row['Title']
20     return f"BookID_{book_id}" # Default title if not found
21
22 # Read the lines of a book and search for those containing the search word (case-
    insensitive)
23 def read_lines_with_word(book_file, paragraphs, search_word):
24     try:
25         with open(f"datalake/books/{book_file}", encoding='utf-8') as file:
26             lines = file.readlines()
27
28             result = []
29             paragraph_count = 0
30             current_paragraph_lines = []
31
32             search_word_lower = search_word.lower() # Convert the word to lowercase for
                case-insensitive search
33
34             for i, line in enumerate(lines):
35                 # If the line is empty, the paragraph has ended
36                 if line.strip() == "":
37                     if current_paragraph_lines:
38                         # Only if the current paragraph is in the selected ones
39                         if paragraph_count in paragraphs:
40                             for current_line in current_paragraph_lines:
41                                 # Case-insensitive search
42                                 if search_word_lower in current_line.lower():
43                                     # Highlight the word in magenta while maintaining the
                                        original case
44                                     highlighted_line = re.sub(f"(?i)({re.escape(
                                        search_word)})", colored(r"\1", 'magenta'),
                                        current_line)
45                                     result.append(f"Line {i + 1}: {highlighted_line.strip(
                                        ())}")
46                                     current_paragraph_lines = [] # Reset the paragraph
47                                     paragraph_count += 1 # Increase the paragraph counter
48                     else:
49                         # Add the current line to the paragraph
50                         current_paragraph_lines.append(line)
51
52             return result
53
54     except FileNotFoundError:
55         print(f"The file '{book_file}' was not found.")
56         return []
57
58 # Extract the paragraph number from MongoDB keys
59 def extract_paragraph_numbers(paragraph_keys):
60     paragraph_numbers = []
61     for key in paragraph_keys:
62         match = re.search(r'\d+$', key) # Look for numbers at the end of the string
63         if match:
64             paragraph_numbers.append(int(match.group()))
65         else:
66             print(f"Skipping invalid paragraph key: {key}")
67     return paragraph_numbers
68
69 # Search for a word in MongoDB and print the lines that contain it
70 def search_word(word, lang='en'):
71     word = unidecode(word.lower()).strip() # Normalize and clean the word
72     word = lemmatize(word, lang) # Lemmatize the word
73
74     # Query MongoDB to search for the word
75     result = words_collection.find_one({'word': word})
76
77     if result:

```

```

78         # Print the total occurrences
79         print(f"\nThe word '{word}' appears in the following books:")
80
81         # Iterate over the books where the word appears
82         for book_file, book_data in result['books'].items():
83             book_id = book_file.split('.')[0] # Extract the book ID
84             book_title = get_book_title_from_csv(book_id) # Get the title from the
                        CSV
85             print(f"\nBook Title: {book_title}")
86
87             # Extract the paragraph numbers
88             paragraphs = extract_paragraph_numbers(book_data.keys())
89             lines_with_word = read_lines_with_word(book_file, paragraphs, word)
90             for line in lines_with_word:
91                 print(line)
92         else:
93             print(f"\nThe word '{word}' was not found in the database.\n")
94
95 # Main function
96 def main():
97     word_to_search = input("Enter the word you want to search: ")
98     search_word(word_to_search)
99
100 # Execute the program
101 if __name__ == "__main__":
102     main()

```

Listing 5: MongoDB Word Search API

## Code Explanation

**MongoDB Configuration and Book Title Retrieval** The first part of the code establishes a connection with a local MongoDB instance through the `MongoClient`. The database used for the search system is `SEARCH_ENGINE`, and the collection where the words are stored is `WORDS`. To retrieve the title of a book, the function `get_book_title_from_csv()` is used, which looks up the title in a `metadata.csv` file that stores book metadata (like title and ID). If the title is not found, the function returns a default string based on the book ID.

**Reading Lines Containing the Word** The function `read_lines_with_word()` processes each book to find the lines that contain the search word. It operates on a case-insensitive basis by converting both the line and the word to lowercase. Additionally, it highlights the found word using the `termcolor` package, displaying it in magenta. The paragraphs are processed by splitting the book's content into sections and scanning through each paragraph.

**Extracting Paragraph Information** The function `extract_paragraph_numbers()` takes the keys from MongoDB (which are strings that identify the paragraphs) and extracts the corresponding paragraph numbers. If a key does not contain a valid number, it is skipped, and an appropriate message is printed.

**Searching for a Word in MongoDB** The `search_word()` function handles the main logic of searching for a word in the database. It first cleans the word, converts it to lowercase, and applies lemmatization to normalize it. Afterward, it queries the MongoDB collection for the word. If found, it prints the books where the word appears and calls `read_lines_with_word()` to display the relevant lines. If the word is not found, an appropriate message is displayed.

**Main Program Execution** Finally, the `main()` function prompts the user to input a word they wish to search for and executes the search using the previously defined functions. The word is searched for in the MongoDB database, and the relevant results are printed in the terminal.



## 3 Experiments

### 3.1 Benchmarking the Indexing Function

To evaluate the performance of the `execute` function, the `pytest-benchmark` tool was used. This module in `pytest` allows measuring the time it takes for a function to complete, along with additional data such as standard deviation, mean, and operations per second (OPS).

The following code was used to perform the benchmark:

```
1 import pytest
2 import shutil
3 import os
4 from indexer1 import execute, update_last_book_id
5
6 @pytest.fixture(scope="function")
7 def setup_teardown():
8     # Setup: Remove the "datamart" directory and reset the ID file
9     shutil.rmtree("datamart", ignore_errors=True) # Remove datamart if it exists
10    update_last_book_id("resources/lastBookId.txt", 0, 1) # Reset the ID file
11    yield
12    # Additional cleanup can be done here if needed
13
14 def test_execute_benchmark(benchmark, setup_teardown):
15     # Benchmark the execute function
16     benchmark(execute)
```

Listing 6: Benchmark Code for the `execute` Function (First Test)

This test uses a fixture to set up the environment before each benchmark execution. It deletes the `datamart` directory and resets the `lastBookId.txt` file, ensuring that every test runs under the same initial conditions.

#### 3.1.1 First Benchmark Results

The results of the benchmark provide the following execution times for the `execute` function:

- **Minimum:** 0.4507 ms
- **Maximum:** 30.4305 ms
- **Mean:** 0.6052 ms
- **Standard Deviation:** 0.8659 ms
- **Median:** 0.5446 ms
- **Interquartile Range (IQR):** 0.1190 ms
- **Operations Per Second (OPS):** 1.6523 Kops/s

The benchmark results show a wide range of execution times, with a mean of 0.6052 milliseconds and a standard deviation of 0.8659 milliseconds, indicating variability in execution speed across different runs. The function achieved approximately 1.65 operations per second (Kops/s).

#### 3.1.2 Second Benchmark Results

In the second benchmark, the `execute` function was tested with a different input choice ("2"), which processes all files. The following results were obtained:

```
1 import pytest
2 from indexer2 import execute
3
4 @pytest.mark.benchmark
5 def test_execute_benchmark(benchmark):
6     # Benchmark the 'execute()' function with '2' as the input to process all files
7     benchmark(execute, choice='2')
```

Listing 7: Benchmark Code for the `execute` Function (Second Test)

The results of the second benchmark are as follows:

- **Minimum:** 4.3335 ms
- **Maximum:** 5.0354 ms
- **Mean:** 4.6000 ms
- **Standard Deviation:** 0.3195 ms
- **Median:** 4.4035 ms
- **Interquartile Range (IQR):** 0.5250 ms
- **Operations Per Second (OPS):** 217.3904 ops/s

The benchmark results show a narrower range of execution times compared to the first test, with a mean of 4.6000 milliseconds and a standard deviation of 0.3195 milliseconds. The function achieved approximately 217.39 operations per second (ops/s).

## 4 Conclusion

As we can see on the benchmarks, handling the storage and indexing with ad-hoc functions yields better results than using an external database, especially for smaller books, where the cost of connecting to the database is disproportionate. For larger books however, the solutions we used seem to be less effective than the solutions implemented in the database, so future versions of the program have room for optimization when working with large texts.

### 4.0.1 Future improvements

In terms of possible improvements, we have considered adding multilanguage support, requesting specific works from project Gutenberg, and generating a summary of the work through a LLM.

### 4.0.2 Source Code on GitHub

The source code and the benchmark tests are available in the following GitHub repository:

[\*https : //github.com/LuisPereraPerez/BigData\\_proyect\*](https://github.com/LuisPereraPerez/BigData_proyect)