

Search Engine, Stage 1

Elisa Breeze
Nicolás Trujillo
Luna Hernández
Carlos Suárez
Ivan Anikin

October 2024

GitHub link

Abstract

In this document we will explain the first stage of creating a search engine with an inverted index. We decided to compare two different data structures: a JSON-based index, and a dictionary-based index. The idea is to evaluate how each structure works to see which one is more efficient when finding the words that want to be looked for. We conducted the performance analysis using Benchmark tools to measure execution time and memory usage of the two query engines. The results indicate that both structures are viable, each with its advantages and disadvantages, but in the end we concluded that the best data structure for our project is the JSON-based one, meaning we will continue using this one for the rest of the project.

1 Introduction

The idea of this work is to create a search engine with an inverted index, that stores words and their positions in documents, allowing efficient searches. The index we created, does not only store which book contains a specific words, but also the location where it appears, allowing a more precise outcome.

By using an inverted index, we can implement fast and effective word searching when a term is queried. In this project we implemented two types of indexers, tested them and chose the best one to continue working with in the next Stages of the Project. The two implementations are: dictionaries and JSON files. We tested them using Benchmark tools, to measure the execution times of the search engines of each version, to see which implementation is the most efficient.

In Stage 1 of the project, a part from the Indexer we also created a Crawler, to download the books periodically from gutenberg.org and save them in a Data-mart.

And lastly, so that we could test our search engine, we created a minimal query engine, one for each data structure, that takes the word or words to be searched for, looks for it and gives back information about the book the word is in, the books name, author, URL, the occurrences and the paragraph it's in.

2 Methodology

We started by researching how to do each part of the project, and how we were going to implement them, where to put the focus and how to divide the work. In the end, we worked together on all the code, by creating different versions and then deciding which implementation worked the best. With this, we managed to write the code to all the parts of the project successfully, and here you can see the result of each part of the project, with a short explanation of the methodology used for the module.

2.1 Crawler

We used a simple approach to develop the crawler: it is designed to periodically download book from the Gutenberg web page. We used the "request" library to send HTTP requests to retrieve the content of the books, and the library "BeautifulSoup" to extract the metadata from the webpages, like for example the book title.

Code Structure:

`get_title(url)`: this function sends a GET request to the Gutenberg page of a specific book and uses BeautifulSoup to parse the HTML and extract the title of the book by locating the first `<h1>` tag in the HTML structure. If it was successful, it returns the title, and if not, it returns a message showing there was an error.

`download_book(book_id, download_route)`: This function forms the URL of the book's text file, by using the book's ID. Then it checks if the download directory exists, and if not, it creates it. Once that has been done, the content of the book and also the title with the `get_title()` function are retrieved, and saved under a name we composed as the title and the book ID.

The crawler goes through a range of book IDs, specified in the code, and downloads and saves each book in the designated directory.

The crawler is important to collect the data for the rest of the project and by periodically downloading new books, it is ensured that the Book Datamart remains up to date and continues to grow.

```
1 import os
2 import requests
3 from bs4 import BeautifulSoup
4
5
6 def get_title(url):
7     """
8     Fetch the title of a webpage by retrieving the
9     text from the first <h1> tag.
10
11     :param url: The URL of the webpage to fetch the
12                 title from.
13     :return: The text of the first <h1> tag if found;
14             otherwise, an error message or indication of
15             failure.
16     """
17     try:
18         # GET request to URL
19         answer = requests.get(url)
20
21         # See if request was successful
22         if answer.status_code == 200:
23             # Analyze HTML
24             soup = BeautifulSoup(answer.text, 'html.
25                                 parser')
26
27             # find the first <h1>
28             h1 = soup.find('h1')
29
30             # If found, give back the text
31             if h1:
32                 return h1.get_text().strip()
33             else:
34                 return "There was no <h1> found."
35         else:
36             return f"Error accessing the page. Status
37                     code: {answer.status_code}"
38     except requests.RequestException as e:
39         return f"Error making the request: {e}"
40
```

```

36
37 def download_book(book_id, download_route):
38     """
39     Download a book from Project Gutenberg using its
        book ID and save it to a specified directory.
40
41     :param book_id: The ID of the book to be
        downloaded from Project Gutenberg.
42     :param download_route: The directory where the
        downloaded book will be saved.
43     :return: None
44     """
45     url = f'https://www.gutenberg.org/files/{book_id}
        {book_id}-0.txt'
46     if not os.path.exists(download_route):
47         os.makedirs(download_route)
48     answer = requests.get(url)
49
50     if answer.status_code == 200:
51         # Extract content as text
52         content = answer.text
53
54         # Call assign_title to find the file name
55         file_name = get_title(f'https://www.gutenberg.
        org/ebooks/{book_id}')
56         file_name = os.path.join(download_route, f'{
        file_name}_{book_id}.txt')
57
58         # Save content in file
59         with open(file_name, 'w', encoding='utf-8') as
        file:
60             file.write(content)
61
62         print(f"Book_{book_id}_downloaded_correctly_
        under_{file_name}")
63     elif answer.status_code == 404:
64         print(f"Book_{book_id}_not_found.")
65     else:
66         print(f"Error_when_downloading_Book_{book_id}:
        _{answer.status_code}")
67     return answer.status_code

```

2.2 Indexers

As explained before, the two indexers implemented use on one hand Dictionaries, and on the other hand, JSON files. The general idea is to save each word with an ID for the book it belongs to, along with its position. Here you can see the explanation of each indexer:

2.2.1 JSON-based Indexer

This indexer stores each word in its own JSON file, which contains a list of all the occurrences of the word across different books, and the basic information (book name, author and ID).

For each word found, a JSON file is created or updated in the case that it already exists. There are 2 versions, a normal one and one that uses parallelism, which is the default one we are going to use in the rest of the project, as it performs much faster. It uses threading, and allows multiple words to be indexed at the same time, which make the process speed it considerably.

It ensures that for every book, all the words and positions are saved in separate JSON files, which creates a modular and distributed system for word searches.

```
1 def indexer5_parallel(datamart_txt_path,
2   datamart_json_path):
3     os.makedirs(datamart_json_path, exist_ok=True)
4     txt_files = [f for f in os.listdir(
5         datamart_txt_path) if re.match(r'^.+?by_+?\d
6         +\.txt$', f)]
7
8     for txt_file in txt_files:
9         match = re.match(r'^(.+?)by_(.+?)_(\d+)\.txt$
10            ', txt_file)
11         if not match:
12             continue
13
14         book_name = match.group(1)
15         author = match.group(2)
16         index = match.group(3)
17         dictionary_key = f"{book_name}_by_{author}_{
18             index}"
19         txt_file_path = os.path.join(datamart_txt_path
20             , txt_file)
21
22         with open(txt_file_path, 'r', encoding='utf-8'
23             ) as file:
24             content = file.read()
```

```

18 words = re.findall(r'\b[a-zA-
19 Z                               ]+\b', content
    .lower())
20 words = [word for word in words if
21 word not in ['in', 'on', 'at', 'by',
    'for', 'with', 'about', 'against',
22 'between',
    'into', 'through', '
    during', 'before', '
    after', 'above', '
    below',
23 'to', 'from', 'up', '
    down', 'of', 'off', '
    over', 'under', '
    again',
24 'further', 'once', 'i',
    'me', 'my', 'myself',
    'we', 'our', 'ours',
    'ourselves',
25 'you', 'your', 'yours',
26 'yourself', 'yourselves'
    , 'he', 'him', 'his',
    'himself', 'she', '
    her', 'hers',
27 'herself',
28 'it', 'its', 'itself', '
    they', 'them', 'their
    ', 'theirs', '
    themselves', 'what',
29 'which',
30 'who', 'whom', 'this', '
    that', 'these', '
    those', 'am', 'is', '
    are', 'was', 'were',
31 'be',
32 'been', 'being', 'have',
    'has', 'had', '
    having', 'do', 'does'
    , 'did', 'doing', 'a'
    ,
33 'an',
34 'the', 'and', 'but', 'if
    ', 'or', 'because', '
    as', 'until', 'while'
    , 'of', 'at',

```

```

35         'by',
36         'for', 'with', 'about',
            'against', 'between',
            'into', 'through', '
            during', 'before',
37         'after',
38         'above', 'below', 'to',
            'from', 'up', 'down',
            'in', 'out', 'on', '
            off', 'over',
39         'under',
40         'again', 'further', '
            then', 'once', 'here'
            , 'there', 'when', '
            where', 'why', 'how',
41         'all',
42         'any', 'both', 'each', '
            few', 'more', 'most',
            'other', 'some', '
            such', 'no', 'nor',
43         'not',
44         'only', 'own', 'same', '
            so', 'than', 'too', '
            very', 's', 't', 'can'
            , 'will', 'just',
45         'don',
46         'should', 'now']]
47
48     with ThreadPoolExecutor() as executor:
49         for position, word in enumerate(words,
50             start=1):
51             executor.submit(process_word, word,
52                 position, dictionary_key,
                    datamart_json_path)
53
54     print("Indexation Completed.")

```

2.2.2 Dictionary-based Indexer

This indexer words by storing the words as keys, and for each word, it maps the book's ID along with the positions where it appears.

`add_words_to_dict`: this function is responsible for inserting the words into the dictionary, and associated each word with the book ID and the positions it has in the text.

The indexer goes through all the files in the Book Datamart, reads the contents of each book, tokenizes it into words and then adds the words to the dictionary.

The dictionary structure is simple: it is divided into smaller partial index files, which are then saves a separate JSON file for more efficient retrieval later on.

```
1 import json
2 import os
3 import re
4
5 from indexer.book_reader import read_words,
    save_metadata_to_json
6 from indexer.path_reader import
    extract_files_from_directory
7
8
9 def add_words_to_dict(words, id_book, dictionary):
10     """
11     Add words and their indexes to a dictionary.
12
13     :param words: A list of words to be added.
14     :param id_book: The ID of the book from which the
15                     words are extracted.
16     :param dictionary: The dictionary where words and
17                       their indexes will be stored.
18     :return: The updated dictionary with words and
19             their corresponding indexes.
20     """
21     for idx, word in enumerate(words):
22         if word not in dictionary:
23             dictionary[word] = {id_book: [idx]}
24         else:
25             if id_book not in dictionary[word]:
26                 dictionary[word][id_book] = [idx]
27             else:
28                 dictionary[word][id_book].append(idx)
29
30     return dictionary
31
32 def id_search(filepath):
33     """
34     Extract the first numeric ID from the given file
```



```

        path.
33
34 :param filepath: The path of the file from which
        to extract the ID.
35 :return: The first numeric ID found in the file
        path, or an empty string if no ID is found.
36 """
37 pattern = r'\d+'
38 coindidencies = re.findall(pattern, filepath)
39
40 if coindidencies:
41     result = coindidencies[0]
42 else:
43     result = ''
44 return result
45
46
47 def save_partial_indexers(indexer, output_directory):
48     """
49     Divide the indexer into smaller parts and save
        them as separate JSON files.
50
51     :param indexer: A dictionary containing the
        indexer data to be divided.
52     :param output_directory: The directory where the
        partial indexer files will be saved.
53     :return: None
54     """
55     if not os.path.exists(output_directory):
56         os.makedirs(output_directory)
57
58     partial_indexers = {}
59
60     for word, data in indexer.items():
61         first_letter = word[0].lower()
62         if first_letter not in partial_indexers:
63             partial_indexers[first_letter] = {}
64         partial_indexers[first_letter][word] = data
65
66     for letter, partial_indexer in partial_indexers.
        items():
67         output_file = os.path.join(output_directory, f
            'indexer_{letter}.json')
68         with open(output_file, 'w', encoding='utf-8')
            as file:
69             json.dump(partial_indexer, file,

```

```

70         ensure_ascii=False, indent=4)
71
72 def indexer_dict(books_datamart, words_datamart,
73                 output_directory_metadata, stopwords_filepath):
74     """
75     Create an indexer from book files, filtering out
76     stopwords, and save metadata and partial
77     indexers.
78
79     :param books_datamart: The path to the directory
80     containing the book files.
81     :param words_datamart: The output directory for
82     saving partial indexers.
83     :param stopwords_filepath: The directory
84     containing the stopwords TXT file.
85     :param output_directory_metadata: The directory
86     containing the metadata JSON file.
87     :return: None
88     """
89     indexer = {}
90     directory_path = books_datamart
91     output_directory = words_datamart
92     filepaths = extract_files_from_directory(
93         directory_path)
94
95     for filepath in filepaths:
96         try:
97             words = read_words(filepath,
98                               stopwords_filepath)
99             if words:
100                 indexer = add_words_to_dict(words,
101                                             id_search(str(filepath)), indexer)
102             else:
103                 print(f"Warning: No words found in {
104                       filepath}")
105
106                 save_metadata_to_json(str(filepath),
107                                       output_directory_metadata)
108
109         except Exception as e:
110             print(f"Error processing {filepath}: {e}")
111
112     save_partial_indexers(indexer, output_directory)

```

2.3 Minimal Query Engine

2.3.1 Minimal Query Engine for JSON

This query engine is responsible for searching books that contain the keywords provided by the user.

`find_book(book_id, book_folder)`: This method looks for a specific book file. "book_id" identifies the book and "book_folder" is the folder where the files are located.

`query_engine(input, book_folder, index_folder, max_occurrences)`: This function searches for books based on the keywords entered, and loads indices from JSON files, searching for matches in the book contents.

What it does exactly is convert the input to lowercase and split it into single words if more than one word was entered. Then, it loads JSON files from the index and builds a dictionary of words, checking if all the words are present in the dictionary. It then identifies books that contain all the keywords, extracting relevant paragraphs from each book, returning a list of dictionaries, each containing all the information about the book and the paragraphs, so that we can print out the information looked for to the user.

```
1 import glob
2 import json
3 import os
4 import re
5
6
7 # uncomment if using memory usage test
8 # from memory_profiler import profile
9
10
11 def find_book(book_id, book_folder):
12     """
13     Searches for a book file in the specified folder
14     by its ID.
15
16     This function looks for a file that ends with '_{
17         book_id}.txt' in the provided
18     book folder and returns the full path if found.
19
20     :param book_id: The ID of the book to search for.
21     :param book_folder: The folder where the book
22         files are stored.
23     :return: The full path to the book file if found,
24         otherwise None.
```

```

21     """
22     for filename in os.listdir(book_folder):
23         if filename.endswith(f"_{book_id}.txt"): #
24             find file that ends with _{book_id}.txt
25             return os.path.join(book_folder, filename)
26     return None
27
28 # uncomment if using memory usage test
29 # @profile
30 def query_engine(input, book_folder, index_folder,
31                 max_occurrences=3):
32     input = input.lower()
33     words = input.split()
34     results = []
35     loaded_words = {}
36
37     # save the JSON objects
38     for filepath in glob.glob(f"{index_folder}/*.json"
39                             ):
40         with open(filepath, "r") as file:
41             data = json.load(file)
42
43             if "id_name" in data and "dictionary" in
44                 data:
45                 word_key = data["id_name"]
46                 dictionary_info = data["dictionary"]
47                 loaded_words[word_key] = {"dictionary"
48                                         : dictionary_info}
49             else:
50                 print(f"Invalid structure in file: {
51                     filepath}")
52
53     # Check if all the words are there
54     words_looked_for = all(word in loaded_words for
55                             word in words)
56
57     if words_looked_for:
58         books_in_common = None
59         for word in words:
60             word_info = loaded_words[word]["dictionary
61                                     "]
62             if books_in_common is None:
63                 books_in_common = set(word_info.keys()
64                                     )
65             else:

```

```

58         books_in_common &= set(word_info.keys
59                                ())
60
61     if books_in_common:
62         for book_key in books_in_common:
63
64             book_info = book_key.split("by")
65             book_name = book_info[0].strip()
66             author_and_id = book_info[1].split("-")
67             author_name = author_and_id[0].strip()
68             book_id = author_and_id[1].strip()
69
70             book_filename = find_book(book_id,
71                                       book_folder)
72
73             if book_filename:
74                 try:
75                     with open(book_filename, "r",
76                               encoding="utf-8") as file:
77                         # we have to specify the
78                         encoding
79                         text = file.read()
80
81                     paragraphs = text.split('\n\n')
82                     relevant_paragraphs = []
83                     occurrences = 0
84
85                     word_pattern = re.compile(rf"\
86                                             b{input}\b", re.IGNORECASE)
87
88                     for paragraph in paragraphs:
89                         if word_pattern.search(
90                             paragraph):
91                             occurrences += len(
92                                 word_pattern.
93                                 findall(paragraph))
94
95                     highlighted_paragraph
96                     = word_pattern.sub(
97                         rf"\033[94m{input}
98                         }\033[0m",
99                         paragraph)
100                     relevant_paragraphs.
101                     append(

```

```

highlighted_paragraph
.strip())
88
89     if relevant_paragraphs:
90         results.append({
91             "book_name": book_name
92             ,
93             "author_name":
94                 author_name ,
95             "URL": f'https://www.
96                 gutenber.org/files
97                 /{book_id}/{book_id
98                 }-0.txt' ,
99             "paragraphs":
100                 relevant_paragraphs
101                 [:max_occurrences] ,
102             "total_occurrences":
103                 occurrences
104         })
105
106     except FileNotFoundError:
107         print(f"Error: The Book {
108             book_filename} was not
109             found.")
110
111     return results
112
113 # uncomment if doing memory usage tests
114 # book_datamart_folder = "../Books_Datamart"
115 # indexer_folder = "../Words_Datamart"
116 # query_engine("wife", book_datamart_folder ,
117               indexer_folder)

```

2.3.2 Minimal Query Engine for Dictionary

The dictionary based query engine looks for books, using a different index structure, based on dictionaries.

`load_json_index(word index_folder)`: This function loads the JSON index file for a word based on its first letter and returns a dictionary with the index data, or an empty dictionary if the file is not found.

`load_metadata(book_id, metadata_folder)`: This function loads the metadata of the book, based on its ID, returning a dictionary with the book's metadata or None if not found.

`extract_paragraphs(book_filename, occurrences_dict, search_words:` This function extracts relevant paragraphs from a book based on occurrences dictionary and the searched words, and returns a list of relevant paragraphs containing the words searched for.

`query_engine(input, index_folder, metadata_folder, book_folder, max_occurrences):` This function looks for the books that contain all the words in the input, and returns the information and paragraphs. It converts the input to lowercase and splits it into words, then loads the word indices for each one in the query, finds common books that contain all the words. Finally it loads the metadata and content of each book, extracting the relevant paragraphs using its function, and returns a list of dictionaries, each containing information about the book and relevant paragraphs, associated to the searched word or words.

```
1 import json
2 import os
3 import re
4
5
6 # uncomment if using memory usage test
7 # from memory_profiler import profile
8
9
10 def load_json_index(word, index_folder):
11     """
12     Loads the JSON file for the word's index based on
13     its first letter.
14
15     :param word: The word whose index is to be loaded.
16     :param index_folder: The directory containing the
17         index JSON files.
18     :return: A dictionary with the index data, or an
19         empty dictionary if the file is not found.
20     """
21     first_letter = word[0].lower()
22     json_path = os.path.join(index_folder, f'indexer_{first_letter}.json')
23
24     if os.path.exists(json_path):
25         with open(json_path, 'r', encoding='utf-8') as file:
26             return json.load(file)
27     else:
28         print(f"Index file for letter '{first_letter}' not found.")
```

```

26         return {}
27
28
29 def load_metadata(book_id, metadata_folder):
30     """
31     Loads the metadata of a book based on its book ID.
32
33     :param book_id: The ID of the book whose metadata
34                     is to be loaded.
35     :param metadata_folder: The directory containing
36                             the metadata JSON files.
37     :return: A dictionary with the book's metadata, or
38             None if not found.
39     """
40
41     hundred_range = (int(book_id) // 100) * 100
42     json_filename = f"books_metadata_{hundred_range}-{
43                     hundred_range+99}.json"
44     json_filepath = os.path.join(metadata_folder,
45                                  json_filename)
46
47     if os.path.exists(json_filepath):
48         with open(json_filepath, 'r', encoding='utf-8'
49                 ) as file:
50             books_data = json.load(file)
51             for book in books_data:
52                 if book['id_book'] == book_id:
53                     return book
54
55     return None
56
57 def extract_paragraphs(book_filename, search_words):
58     """
59     Extract paragraphs from the book based on
60     occurrences dictionary and search words.
61
62     :param book_filename: The path to the book file
63                           from which to extract paragraphs.
64     :param search_words: A list of words to search for
65                          in the book paragraphs.
66     :return: A list of relevant paragraphs containing
67             the search words.
68     """
69
70     try:
71         with open(book_filename, "r", encoding="utf-8"
72                 ) as file:
73             text = file.read()

```



```

61     paragraphs = text.split('\n\n')
62     relevant_paragraphs = []
63     occurrences = 0
64     word_patterns = {word: re.compile(rf"\b{word}\b", re.IGNORECASE) for word in search_words}
65
66
67     for paragraph in paragraphs:
68         for word, pattern in word_patterns.items():
69             if pattern.search(paragraph):
70                 occurrences += len(pattern.findall(paragraph))
71
72                 highlighted_paragraph = pattern.sub(f"\033[94m{word}\033[0m", paragraph)
73                 relevant_paragraphs.append(highlighted_paragraph.strip())
74                 break
75     return relevant_paragraphs, occurrences
76
77 except FileNotFoundError:
78     print(f"Error: Book file not found: {book_filename}")
79     return [], 0
80
81
82 # uncomment if using memory usage test
83 # @profile
84 def query_engine(input_query, index_folder, metadata_folder, book_folder, max_occurrences=3):
85     """
86     Searches for books containing the words in the
87     input query and returns relevant paragraphs.
88
89     :param input_query: The search query (a string of words).
90     :param index_folder: Directory where the word index files are stored.
91     :param metadata_folder: Directory where the book metadata JSON files are stored.
92     :param book_folder: Directory where the book files are stored.
93     :param max_occurrences: Maximum number of

```

```

    paragraphs to return for each book.
93 :return: List of dictionaries with book
    information and paragraphs containing the
    search words.
94 """
95 input_query = input_query.lower()
96 words = input_query.split()
97 results = []
98
99 # Dictionary to store word occurrences across
    books
100 word_occurrences = {}
101
102 # Step 1: Load word indices for all search words
103 for word in words:
104     index_data = load_json_index(word,
        index_folder)
105     if word in index_data:
106         word_occurrences[word] = index_data[word]
107     else:
108         print(f"Word '{word}' not found in any
            index.")
109         return [] # Exit early if any word is
            missing
110
111 # Step 2: Find common books that contain all the
    search words
112 common_books = None
113 for word, occurrences in word_occurrences.items():
114     if common_books is None:
115         common_books = set(occurrences.keys())
116     else:
117         common_books &= set(occurrences.keys())
118 # Step 3: Process each book that contains all the
    words
119 if common_books:
120     for book_id in common_books:
121         # Step 4: Load book metadata
122         metadata = load_metadata(book_id,
            metadata_folder)
123         if not metadata:
124             print(f"Metadata for book ID '{book_id}'
                not found.")
125             continue
126
127         book_name = metadata["book_name"]

```

```

128         author_name = metadata["author"]
129         url = metadata["URL"]
130
131         # Step 5: Load the book content
132         book_path = os.path.join(book_folder, f"{
            book_name}_by_{author_name}_{book_id}.
            txt")
133
134         # Step 6: Extract relevant paragraphs
135         paragraphs, occurrences =
            extract_paragraphs(book_path, words)
136         if paragraphs:
137             results.append({
138                 "book_name": book_name,
139                 "author_name": author_name,
140                 "URL": url,
141                 "paragraphs": paragraphs[:
                    max_occurrences],
142                 "total_occurrences": occurrences
143             })
144
145         return results
146
147     # uncomment if doing memory usage tests
148     # indexer_folder = "../Words_Datamart_Dict"
149     # metadata_datamart_folder = "../Books_Metadata_Dict"
150     # book_datamart_folder = "../Books_Datamart"
151     # query_engine("wife", indexer_folder,
        metadata_datamart_folder, book_datamart_folder)

```

2.3.3 Rest of the code

To make the project work and ensure that all functions work well together, we included several controllers: a crawler controller, an indexer controller for the JSON structure, an indexer controller for the dictionary structure, and a query engine controller. Each of these controllers is executed from the main function within their respective files, allowing the entire search engine to operate correctly.

3 Experiments and Results

To analyze the performance of the two indexers, we executed benchmarks to analyze the execution times, so that we could see which option has a better performance.

To do so, we included a benchmark function in a separate file, to do testing on the two query engines, and measure the execution time and memory usage of each one to find out which one is the best.

In the Benchmark file, there is code for the benchmark test. To execute the tests, use the command:

```
pytest benchmark.py --benchmark-group-by=func
```

To save the results, you can run:

```
pytest benchmark.py --benchmark-save indexer_benchmarks
```

Additionally, to test memory usage, in the actual queryEngine files, there is a `@profile` command and a function call at the bottom so that `memory_profiler` can do the test. To do this, just un-comment those two lines of code.

Here you can see the code for the benchmark:

```
1 from queryEngine import query_engine,
2   query_engine_dict
3
4 def test_query_engine(benchmark):
5     benchmark.extra_info['unit'] = 'ms'
6     book_datamart_folder = "../Books_Datamart"
7     indexer_folder = "../Words_Datamart"
8     benchmark.pedantic(
9         target=query_engine.query_engine,
10        args=("wife", book_datamart_folder,
11            indexer_folder,),
12        iterations=100,
13        rounds=10,
14        warmup_rounds=5
15    )
16
17 def test_query_engine_dict(benchmark):
18     benchmark.extra_info['unit'] = 'ms'
19
20     indexer_folder = "../Words_Datamart_Dict"
21     metadata_datamart_folder = "../Books_Metadata_Dict"
22     book_datamart_folder = "../Books_Datamart"
23
24     benchmark.pedantic(
25         target=query_engine_dict.query_engine,
```

```

26         args=("wife", indexer_folder,
27               metadata_datamart_folder,
28               book_datamart_folder,),
29         iterations=100,
30         rounds=10,
31         warmup_rounds=5
32     )
33
34 # to execute => pytest benchmark.py --benchmark-group-
35 # by=func
36 # to save results => pytest benchmark.py --benchmark-
37 # save indexer_benchmarks
38 # to execute memory usage => in query engine module, I
39 # added @profile, then execute in there, first
40 # uncomment the profile and function call at the enc:
41 # python -m memory_profiler query_engine.py
42 # all of these test, from inside the package

```

And here you can see the results we obtained for each query engine:

Index Type	Execution Time (ms)	Memory Usage (MiB)
Dictionary-based	2.81013e-05	44.1
JSON-based	1.79022e-05	43.8

Table 1: Benchmark test results for the indexers

4 Discussion

It's paradoxical that, although the dictionary type is the fastest during indexing, when it comes to the query engine, it's the opposite.

The performance of the benchmark and memory usage tests show that they are quite similar, but that the JSON query engine is slightly faster and has a slightly better memory usage, with 1.79e-05ms compared to 2.81e-05ms; and 43.8MB compared to 44.1MB.

It's important to note, however, that these figures reflect the query engine performance, that is, the part the user interacts with. When it comes to indexing, the JSON-based indexer actually takes longer to complete the indexing process.

5 Conclusion

After finishing the code and comparing the two indexing approaches for our search engine, we have concluded, that even though they are very similar, for our approach to the whole project, it is better to use JSON-based indexers. As

we observed, the query engine that uses the files generated by the JSON-based indexer provides a slight improvement in speed. All this, together with the flexibility afforded by the modular structure, we have decided to continue with the JSON-based version for the next stages of the project.

On the other hand, the dictionary-based indexer, while also a viable option, did not perform as efficiently, but its simpler structure and fewer file accesses still offer benefits if we see that in the future stages of the project handling a large number of JSON files becomes a problem.