

Let's Search

Scenario

MustangWiki's search engine is broken! MustangWiki is a collaborative, open place where students, faculty, staff and the community can create online, freely-accessible content. Please help MustangWiki get its search feature back by implementing a fast search engine for all of the pages in the collection. For this project, you'll be using a dump of the WikiBooks (<http://en.wikibooks.org>) as the basis for your search engine. Wiki Books contains more than 200,000 entries stored in XML.

Search Engine Architecture

Search engines are designed to allow users to quickly locate the data they want or need. Input to a search engine is a set of documents commonly referred to as the **corpus**. Typically, the user will enter a search query, and any documents that satisfy that query are returned to the user. Another task of a search engine is to rank the results based upon relevancy.

The four major components¹ of a typical search engine are the following:

1. Document parser/processor,
2. Query processor,
3. Search processor, and
4. Ranking processor.

Figure 1 provides a general overview of a potential system architecture:

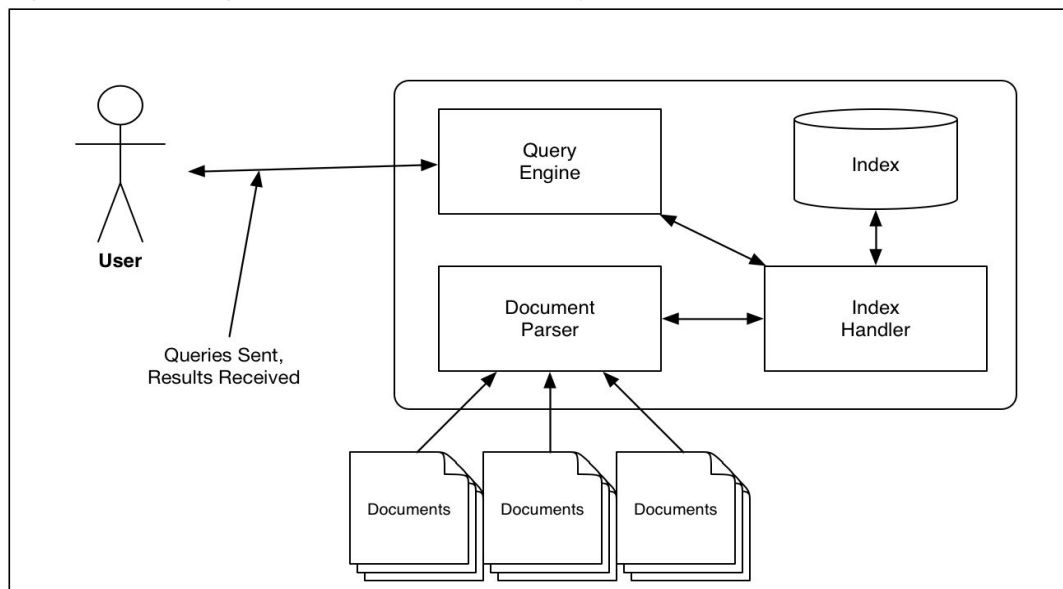


Figure 1 - Sample Search Engine System Architecture

The basic document for this project is a page from WikiBooks. The dataset you will be using contains over 200,000 documents. The entire corpus is stored in a single XML file (<http://en.wikipedia.org/wiki/XML>). XML

¹<http://www.infoday.com/searcher/may01/liddy.htm>

allows for marking up of a document using a hierarchical set of “tags”. You'll be considering `<page> . . . </page>` tags as one document as you generate your index.

Below is an overview of the major tasks/responsibilities of each of the components of the search engine.

The **index handler**, the workhorse of the search engine, is responsible for the following:

- *Reading from and writing to the main index.* You'll be creating an **inverted file index** which stores references from each element to be indexed to the corresponding document(s) in which those elements exist. This is essentially the same data structure you created for Program 2.
- *Searching the inverted file index based on a request from the query processor.*
- *Storing other data with each word.*

The **document parser/processor** is responsible for the following tasks:

- *Processing each document in the corpus.* While there is only one XML file, that file contains more than 200K individual pages. The fact that the documents are marked up in XML allows you to use an XML parser to help parse the documents. However, you can brute-force parse the XML docs if you'd like. Note that the XML document contains other interesting information besides just the text of the document that you may want to include in your index to use for ranking.
- *Removing stopwords from the documents.* Stopwords are common words that appear in text but that provide little discriminatory power with respect to the value of a document relative to a query because of the commonality of the words. Example stop words include “a”, “the”, and “if”. One possible list of stop words to use for this project can be found at <http://www.webconfs.com/stop-words.php>. You may use other stop word lists you find online.
- *Stemming words.* Stemming² refers to removing certain endings of words in the English language. For instance, the stemmed version of “running” may be “run”. For this project, you may make use of any previously implemented stemming algorithm that you can find online. One such algorithm is the Porter Stemming algorithm. More information as well as implementations can be found at <http://tartarus.org/~martin/PorterStemmer/>. Another option is <http://www.oleandersolutions.com/stemming/stemming.html>. You may use others.
- C++ implementation of Porter 2: https://bitbucket.org/smassung/porter2_stemmer/src
- *Computing/maintaining term frequencies.* Term frequency is how often a word appears in a particular document as well as in the overall corpus. This will be used in relevance ranking of documents returned from a search.

The **query processor** is responsible for:

- *Parsing of queries entered by the user of the search engine.* For this project, you'll implement functionality to handle **simple** prefix Boolean queries entered by the user. The Boolean expression will be prefixed with a Boolean operator of either AND or OR if there more than one word is of interest. Trailing search terms may be preceded with NOT to indicate documents including that term should be removed from the result set. For simple one-term searches, an operator is not required. Here are some examples:
 - **Boston**
 - This query should return all documents that contain the word Boston.
 - **AND programming computer java**
 - This query should return all documents that contain the words programming **and** computer **and** java
 - **OR Boston Seattle**
 - This query should return all documents that contain either Boston **OR** Seattle **OR** both.
 - **AND book Boston NOT Seattle**
 - This query should return all documents that contain book and Boston, but not Seattle.

²See <https://en.wikipedia.org/wiki/Stemming> for more information.

○ Boston NOT Seattle

■ This query should return all document that contain Boston, but not Seattle.

- *Ranking the Results.* **Relevancy ranking** refers to organizing the results of a query so that “more relevant” documents are higher in the result set than less relevant documents. The difficulty here is determining what the concept of “more relevant” means. One way of handling relevancy is by using a basic **term frequency – inverse document frequency** (tf/idf) statistic³. tf/idf is used to determine how important a particular word is to a document from the corpus. If a word appears frequently in document d_t but infrequently in other documents, then document d_t would be ranked higher than another document d_s in which a query term appears frequently, but it also appears frequently in other documents as well.

The **user interface** is responsible for:

- Receiving queries from the user
- Communicating with the Search Engine
- Formatting and displaying results in an organized, logical fashion

More info on the UI later.

The Index

The **inverted file index**⁴ is a data structure that relates each unique word to the document(s) in which it appears. It allows for efficient execution of a query to quickly determine in which documents a particular query term appears. For instance, let's assume we have the following documents with ascribed contents:

- d_1 = Computer network security
- d_2 = network cryptography
- d_3 = database security

The inverted file index for these documents would contain, at a very minimum, the following:

- computer = d_1
- network = d_1, d_2
- security = d_1, d_3
- cryptography = d_2
- database = d_3

The query “AND computer security” would find the intersection of the documents that contained *computer* and the documents that contained *security*.

- set of documents containing computer = d_1
- set of documents containing security = d_1, d_3
- the intersection of the set of documents containing computer AND security = d_1

³<http://en.wikipedia.org/wiki/Tf-idf> or <http://nlp.stanford.edu/IR-book/html/htmledition/tf-idf-weighting-1.html> for more information

⁴See http://en.wikipedia.org/wiki/Inverted_index for more information.

Inverted File Index Implementation Details

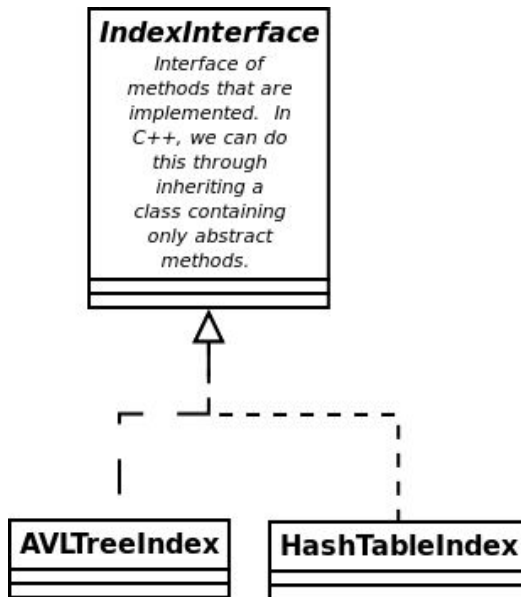


Figure 2: Index Interface Class Diagram

The heart of this project is the **inverted file index**. Notice from the example above that what is being stored is essentially words and a list of documents in which each appears (accompanied by some bookkeeping information). Once the index is created, retrieving the set of documents for various words will be the central task of the index. Therefore, we should use a data structure that will allow for efficient searching.

You will implement at least two different data structures to store the index.

- AVL Tree
- Hash Table with collisions handled by separate chaining

You may also include a list implementation of the index. This will provide a good starting place.

In your implementation, you should strive to abstract the index data structure idea from the underlying storage data structure implementation. What this means is that the AVL tree and hash table implementation should share a common interface. In Figure 2, IndexInterface may contain abstract methods such as `addWord()`, `getDocsForWord(word:string)`, etc.

Each class that inherits from IndexInterface would be required to implement those methods. So that would mean that we could do something like:

```

IndexInterface* if = new AVLTreeIndex;
    or
IndexInterface* if = new HashTableIndex;
  
```

Index Persistence

The index must also be persistent once it is created. This means that the contents of the index should be written to disk when the program ends and read when the program starts. The user should have the option of clearing the index and starting over.

User Interface

The user interface of the application should provide the following options:

- The program should have **two modes** (controlled by menu or command line parameters)
 - maintenance mode –
 - allows the user to add documents to the index by supplying the path to a new file containing properly marked-up documents
 - allows the user to clear the index completely
 - interactive mode –
 - allow the user to indicate if they want the index loaded into an AVL structure or a hash table structure (if a persisted index exists).
 - allow the user to enter a properly formatted Boolean query (as described above).
 - The results should display the article/page title, contributor's username, date associated with the article, and the TF/IDF value (or other relevancy ranking statistic). The result set shown to the user need not contain any more than 15 articles. You may paginate all postings if you wish.
 - The user should be allowed to choose one of the pages from the result set and have its contents displayed.
 - Note that the query terms should have stop words removed and stemmed before querying the index.
 - Upon request, print basic statistics of the search engine including:
 - Total number of pages indexed
 - Total number of words indexed
 - Top 50 most frequent words
 - Any other options you deem appropriate and useful.

Document Data Set

We will be using a complete export of the WikiBooks website. The original data set can be found here: <http://dumps.wikimedia.org/enwikibooks>. Each of the directories on that page are dates. Just choose the most recent date. The main file that you'll be indexing is named like:

enwikibooks-<date>-pages-meta-current.xml.bz2

The file contains more than 200,000 pages. There are other export files on the webpage above that you may find useful in your project. As you begin processing the file and working with the data, you'll notice that there are a number of articles that contain no **text** section or non-english characters. You may skip (meaning, skip indexing) those pages.

Mechanics of Implementation

Some things to note:

- This project may be done individually, in teams of two students, or in teams of three students.
 - Individually: Finish all work on your own.
 - Team of 2 students:
 - Each team member must contribute to both the design AND implementation of the project.
 - Each class in the design must have an "owner". The owner is a group member that is principally responsible for its design, implementation and integration into the overall project.
 - Team of 3 students:
 - Complete all work for this project and, additionally, ONE of the two following features:
 - Implement 2-word phrase searching. A 2-word phrase search will be

indicated by square brackets (e.g. []) around the 2-word phrase. Reject any query that contains more than 2 words in the brackets. Multiple 2-word phrases may be found in one query (e.g. AND [color printer] coffee [digital camera]).

- Implement nested boolean queries such as **AND coffee (OR dark medium) roast**. Levels of nesting will not exceed 2 (the example AND coffee (OR dark medium) roast would be 1 level of nesting).
- This project must be implemented using an object-oriented design methodology.
- You are free to use as much of the C++ standard library as you would like. In fact, I encourage you to make generous use of it. You may use other libraries as well except for the caveat below.
- You must implement your own version of an AVL tree and Hash Table (the storage data structures for the index). You may, of course, refer to other implementations for guidance, but you **MAY NOT** incorporate the total implementation from another source.
- All of your code must be properly documented and formatted
- Each class should be separated into interface and implementation (.h and .cpp) files unless templated.
- Each file should have appropriate header comments to include owner of the class and a history of updates/modifications to the class

Submission Schedule

You must submit the following:

- **Teams:** Due Monday April 4 @ 5pm submitted to Canvas
- **Design Documents:** Due Monday Apr 11, 2016 at 8am (Submit to Canvas)
 - Class diagram indicating
 - Interface for each class
 - Role and responsibilities of each class
 - Owner for each class (if done in groups)
- Implementation Milestone 1: Due In Lab – Week of April 11.
 - Demonstrate to your TA that you've made substantial progress on the document processor and one of the index data structures.
- Implementation Milestone 2: Due In Lab – Week of April 18.
 - Demonstrate that you've made adequate progress on query parser and have almost completed the indexer for one data structure.
- Sanity Check/Parsing Speed Check with Prof. Fontenot – Apr 22 (Friday). Sign up sheet will be made available soon.
- Lab time during week of Apr 25 for final polishing and tweaks.
- **Final Project: Due Monday May 2, 2016 @ 6:00am (no extensions!)**
 - Complete project with full user interface
 - Your goal for parsing the entire wikibooks export file is 4 minutes. You should be able to implement parsing that can rip through all 200,000 pages in 4 minutes or under.
 - Users Manual to include information on management piece as well as regular user piece
 - Documentation
 - updated UML diagrams
 - documentation about each class in your project. Consider using Doxygen⁵ for this.
 - Report that compares the underlying functionality of the AVL implementation vs the Hash Table implementation.
 - Which one is better?
 - How do you know?
 - Can you quantify this?
 - Is one better for small data sets compared to large data sets?
- Demonstration of functionality to Professor Fontenot and TAs on Monday May 2, 2016 (sign up sheet to be distributed).

⁵See <http://www.stack.nl/~dimitri/doxygen/> for more information.

Thoughts and Suggestions

- If you wait even 1 week to start this project, you will very likely not finish.
- A significant portion of your grade will come from your demonstration of the project to Prof. Fontenot and the TAs. Be ready for this.
- Take an hour to read about the various parts of the C++ STL, particularly the container classes. They can help you immensely in the project.
- As mentioned previously, beware of code that you find on the Internet. It isn't always as good as it may seem upon initial inspection. Make sure that any code you use in the project is cited/referenced in the header comments of the project.
- Don't take the Thanksgiving break off from the project. Work on it every day. Make use of a variety of options to communicate with your team. You can have virtual meetings through Google Hangout or use GroupMe to keep in touch with text messages.
- Take the large Wikibooks export file and examine it. Data is rarely beautiful and nicely formatted. Use the file to extract some sample test files of various sizes (10 pages, 100 pages, 1000 pages) to use during testing. Don't start out trying to index the whole thing.

Grading:

This project is worth 25% of your final grade in this course (all other implementation projects are worth 35% percent of your final grade).

	<i>Points possible</i>	<i>Points awarded</i>
Early Design Documents	10	
Milestone 1	10	
Milestone 2	10	
<u>Completed Project</u>		
Completed Project	50	
Documentation	25	
Demonstration	25	

Sample “page” from the wikibooks xml file:

```
<page>
  <title>Talk:Scratch</title>
  <ns>1</ns>
  <id>123374</id>
  <revision>
    <id>877966</id>
    <parentid>877964</parentid>
    <timestamp>2007-05-27T15:22:13Z</timestamp>
    <contributor>
      <username>Robert Horning</username>
      <id>1227</id>
    </contributor>
    <minor />
    <comment>/* Central Discussion */</comment>
    <text xml:space="preserve">== Keep it simple ==
```

Please keep this page clean and simple. This is the “front page” that has links to other parts of this Wikibook, and should only have major sections listed. Individual chapters/modules and other parts of this book ought to be listed in the table of contents and as needed in other modules.

If there are ideas on how to improve this page, please feel free to add comments and suggestions below. --[[User:Robert Horning|Rob Horning]] 23:45, 26 May 2007 (UTC)

== Central Discussion ==

If you are making a comment about the general organization of this Wikibook, I would encourage you to use [[Scratch/Planning]] instead of this page for those discussion. This talk page should be limited to the appearance of this “title page” and any general organization thoughts to help improve this first impression people may have about this book. --[[User:Robert Horning|Rob Horning]] 15:21, 27 May 2007 (UTC)</text>

```
  <sha1>m7955bdj4askeelqqoriyeuv6sgw9lh</sha1>
  <model>wikitext</model>
  <format>text/x-wiki</format>
</revision>
</page>
```