

In this assignment you will work *individually* to build a simple web crawler and search engine. Your main tasks are to crawl a portion of the web, to build an index that allows you to quickly access portions of this web, and to respond to various types of queries—or web searches—much like Google queries. For example, you should be able to type “hoojiedoober” and get a list of pages that contain the word “hoojiedoober.”

This assignment asks you to choose appropriate data structures to support various query operations. We ask you to consider three different types of queries, and we allow you to use different data structures to support each type of query.

1 Words of Wisdom

This is a large assignment with what may seem like a distant deadline. In other words, there is plenty of rope with which to hang yourself. We strongly encourage you to (1) start early, (2) think carefully about your design, (3) think about testing early and often, and (4) bite off and debug small pieces of the assignment before writing more code.

This assignment represents the culmination of your semester, as we’re confident that you can now handle an open-ended assignment of this scope. If you need more incentive to start early, realize that this assignment will be weighted more than any of the others.

2 Web Queries and Search Engines

A typical search engine supports some form of a query language. For example, a simple query might be a string, for which the search engine would return a list of URL’s of web pages that contain that string. More advanced features might include a logical OR operator, the ability to search for synonyms, or the ability to restrict a search to a particular Internet domain. To satisfy such queries, a special HTTP request is sent to the server, which parses the query and then returns the result in the form of web page that contains a set of URL’s that satisfy the query.

Given the enormous size of the web, how do search engines return their results so quickly? These search engines crawl the web and build indices of subsets of the web. Good engines have clever ways of identifying what to index and how to index them.

3 Your Assignment

For this assignment, you will crawl a small self-contained portion of the web. The code that we provide **should not be used to crawl the real web** because our code does not follow established guidelines for crawling. Moreover, if you try to crawl the real web, you will end up building enormous indices. (If you’re curious to understand the guidelines for web crawling, the following URL provides a nice set of guidelines: <http://www.robotstxt.org/guidelines.html>. You should also understand the repercussions of not following such guidelines: <http://xxx.lanl.gov/RobotsBeware.html>.)

More specifically, your assignment is to implement three components. **WebCrawler** is a stand-alone application that crawls the web, starting at a URL specified on the command line. The information gathered by the crawler is stored in a **WebIndex** object, which is saved to disk when WebCrawler terminates. Once a WebIndex has been built, **WebQueryEngine**, which is driven by a GUI that we provide, loads a previously created WebIndex object and is then able to perform queries on the information stored in the WebIndex.

This assignment is structured as a series of increasingly complex types of queries. Before describing these types of queries, we first explain the three components that you will implement. To simplify your task, we provide pieces of each class, along with other support classes that interface with the web server, act as the GUI for your query engine, etc. In addition to the provided code, you are free to use any of the data structures in the Java Collections Framework.

3.1 The WebCrawler class

The **WebCrawler** can be invoked from the command line on a URL. See the Java documentation on the `URL` class for the various formats that work. For example, the TA can point the crawler at the copy of the class website on his hard drive with the command `java WebCrawler file://localhost/home/sharadb/cs314h/www/index.html`. The procedure for specifying files on your own machine may differ depending on how your operating system handles hostnames and paths.

You will need to modify the following methods:

- `public List parse(String url)`
This method is called to parse the specified web page. The actual parsing is performed by calling `super.parse(url)`, but this `parse` method exists so that computations can be performed immediately before and after parsing; for example, you will probably want your crawler to not parse a page that has already been visited. The actual parsing of the page by `super.parse(url)` will invoke the four methods listed below depending on the particular HTML tags that are encountered.
- `public void handleStartTag(HTML.Tag t, MutableAttributeSet a, int pos)`
This method is called when the first part of a two-part tag is encountered, such as `<a>` or ``. The `<a>` tag is especially interesting, as it represents a hyperlink.
- `public void handleEndTag(HTML.Tag t, int pos)`
This method is called when the second part of a two-part tag is encountered, such as `` or ``.
- `public void handleSimpleTag(HTML.Tag t, MutableAttributeSet a, int pos)`
This method is called when a one-part tag is encountered, such as `
`.
- `public void handleText(char[] data, int pos)`
This method is called when content text is encountered.

To show you how the parser will call these methods, the partially implemented `WebCrawler` class that we provide includes code that parses web pages and prints the sequence of callback method calls and encountered HTML elements. You will want to change these methods to actually build your index.

3.2 The WebIndex Class

We provide very little of the `WebIndex` class. The methods in this class will only be called by the code that you write in `WebCrawler` and `WebQueryEngine`, so you can build whatever indexing structures you wish. However, this class should implement the `Serializable` interface, so any data members you use must be `Serializable`. `Serializable` objects can be easily saved to and restored from disk.

You have tremendous freedom to design the index in any manner you like. Therefore, documentation is extremely important. Be sure to include in your documentation a detailed description of your design and *why* you chose your design, including such factors as runtime and space considerations. Analyze and discuss the performance of your design. Also, your code should be well-commented.

3.3 The WebQueryEngine Class

For the `WebQueryEngine` class, you should implement a public constructor along with the following methods:

- `public void useWebIndex(WebIndex index)`
This method sets the `WebIndex` object that will be used for subsequent queries.
- `public Collection query(String query)`
This method takes a query expression as an argument, parses the query, and returns a list of URL's to pages that match the query; this list is returned as a `Collection` of `Strings`. Additional details about parsing can be found in the next section.

We also provide code to use your WebQueryEngine class. We provide an HTML page (index.html) and supporting materials that you can use to interact with your WebQueryEngine via a web browser. This page should be located with your .class files and WebCrawler generated index.db file.

You may find it useful to write an alternative query system for your own testing and debugging. For example, you may find it easier and faster to use a command line tool instead of a web browser for testing. You can adapt the existing code for your own purposes.

You may (and probably will) add additional methods, but you must support the above two methods and not alter their semantics. We will test these two methods using automated techniques, so be sure that there are no hidden assumptions that would cause such a program to fail.

4 The Query Language

We are now ready to explain the various types of queries that your search engine should support. For each type of query, you have two tasks: (1) represent these queries and (2) efficiently perform these queries. Before you start your implementation, think carefully about both of these aspects of the problem. You may implement different types of queries using different indexing structures and different strategies. In your report, be sure to explain why you chose your various data structures.

While the parser for the query language is not the most important part of the assignment, it does require a fair amount of explanation. When designing your index, you should also keep in mind the kinds of queries you want to support, as they will have a great impact on your design decisions.

4.1 Basic Queries

The first part of your assignment is to support simple queries, which consist of individual words, the logical AND (&) operator, the logical OR (|) operator, and parentheses. To simplify the parsing, your language will have to fully parenthesize each query, ie, any use of AND or OR requires a set of parentheses. If you would like to relax these constraints, feel free to do so. Here are some examples of basic queries:

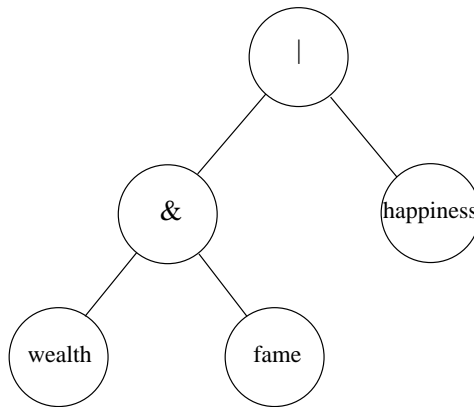
- `snuffleupagus`
Find pages that contain the word “snuffleupagus.”
- `(rosencrantz & guildenstern)`
Find pages that contain both “rosencrantz” and “guildenstern.”
- `(naughty | bear)`
Find pages that contain either “naughty” or “bear.”
- `((wealth & fame) | happiness)`
Find pages that contain both “wealth” and “fame” or pages that contain “happiness.”

4.2 Parsing

You might find it useful for your query engine to parse the String that represents the query into some internal representation before you perform your search. For example, you might represent the above query as the following tree: You do not need an explicit representation of the parse tree, but having one will likely make it *much* easier for you to optimize your search strategies. For example, to satisfy the above query, you could independently search for all pages that contain the word “wealth,” find all pages that contain the word “fame,” and then take the intersection of these two sets. However, it’d probably be much faster to find all pages that contain “wealth” and of these pages search for those that also contain the word “fame.” Having a parse tree can help you do this.

You will learn much more about parsing if you take a compiler course, but for now just use the following two-level approach. First, identify *tokens*, which in our case will be either words or one of the operators. To be more precise, we will treat the left parenthesis as a separate token from the right parenthesis. Second, parse these tokens into a tree.

We have not formally defined what is allowed as a word in our search engine. At a minimum, it should be any combination of letters and numbers. If you want to define words more liberally, you may do so. Be sure to provide



your definition of a word in your code and report. Since most search engines perform case-insensitive searches, you should also do so.

Given our definition of a word, we can identify tokens using the following pseudocode:

```

Token GetToken(String stream)
{
    c = first character in stream;
    if (c == "&")
        return (AndToken);
    else if (c == "|")
        return (OrToken);
    else if (c == "(")
        return (LeftParenToken);
    else if (c == ")")
        return (RightParenToken);
    else
    {
        read until blank or operator;
        rewind the stream one character if it was an operator
        return a Token that contains a reference to the word;
    }
}

```

Notice that you won't always know when a token ends until you have already read the first character of the next token, so you'll have to be careful to make sure you do not accidentally lose characters. You should choose an internal representation for tokens that is easy to work with and easy to understand.

This step is known as *lexical analysis* (or just *lexing* for short) and allows you to iterate through the tokens in a string of input in order.

Once you can identify tokens, you can create the parse tree using the following pseudocode (this type of parser is known as a *recursive descent parser*).

```

Tree parseQuery()
{
    t = getToken();
    if (isLeftParens(t))
    {
        left = parseQuery();           // recursively build the left subtree
        op = getOperator();           // get the binary operator: AND or OR
        right = parseQuery();         // recursively build the right subtree
        op = getToken();              // read the remaining Right Parens
        return makeBinaryNode(op, left, right);
    }
}

```

```

else if (isWord(t))
{
    return makeWord(t);           // return a simple word as a query
}
else
{
    // a parse error has occurred; do something
}
}

```

To understand this parser, it might help you to realize that this parser corresponds to the following set of rules. Starting with *Query*, this set of rules, collectively known as a grammar, can generate the set of all legal queries.

```

Query  → ( Query & Query )
Query  → ( Query | Query )
Query  → word

```

The left hand side of each arrow is called a *non-terminal* and the right hand side is called a *production*. For example, the first line can be read as “A *Query* is an open paren, a *Query*, an ampersand, a *Query*, and a close paren.”

It is easy for an undisciplined programmer to produce a parser that almost works but contains difficult to fix bugs. To avoid this, you should completely understand the grammar and sketch out your design *before* you start coding. There should be a very clear correspondence in your parser between the rules of the grammar and the structure of your code.

One way to structure your code is to have a method for each nonterminal that contains conditions that correspond to each production. This will make your code correspond very cleanly to the grammar and will assist in clarity and debugging.

4.3 Negative Queries

The second part of your assignment will add the ability to make negative queries. The NOT operator (!) matches pages that do not contain the specified word.

Our grammar for this second part includes the following rules.

```

Query  → ( Query & Query )
Query  → ( Query | Query )
Query  → word
Query  → !word

```

The pseudocode from above would be extended with another clause that might look like this:

```

// earlier clauses
if (isNegation(t))
{
    word = getToken();
    return makeNegation(makeWord(word));
}
// later clauses

```

For extra karma, you can support negation of arbitrary queries, not just words. This changes the last grammar rule to *Query* → ! *Query*, with corresponding changes in your code.

4.4 Phrase Queries

The third type of query is a phrase query, which searches for a contiguous sequence of words. The phrase is indicated by surrounding a sequence of words in double quotation marks, for example, "john paul george". The new grammar is shown below:

```

Query  → ( Query & Query )
Query  → ( Query | Query )
Query  → word
Query  → ! word
Query  → " Words "
Words  → word Words
Words  → word

```

The two productions for *Words* define a list of one or more words as a word or as a word followed by more words.

4.5 Implicit AND Queries

Most search engines support implicit logical AND operators: If a query consists of consecutive words (not in quotation marks), the engine searches for pages that contain both words. Modify the parser to support implicit ANDs, matching the following grammar.

```

Query  → word Query
Query  → ( Query & Query )
Query  → ( Query | Query )
Query  → word
Query  → ! word
Query  → " Words "
Words  → word Words
Words  → word

```

Notice that you do not have to worry about precedence in cases like `foo bar | baz`. This isn't a valid query, since the OR is not parenthesized.

5 Testing

We may provide some sample webs for you to play with. You can also download other websites using some tool and crawl the copy on your hard drive. These webs will let you see how your search engine works, but they are no substitute for rigorous testing.

One way to help test your code is to write a program that generates large random graphs representing a “web” with randomly chosen words and randomly chosen links. The program can write this web out to a collection of HTML files for use by your engine. Although the pages generated in this manner will be gibberish, you can in mere seconds create huge webs to test against.

Another problem is knowing whether you are returning the correct URLs in a search. On Unix-based operating systems (like Linux and Mac OS X) there are command line tools like `grep` that allow you to search files for certain patterns. This will let you easily check which pages in your web contain certain words. Other operating systems may have similar tools; ask around.

Lastly, you should keep an eye on efficiency. Obviously, you won't be able to index the billions of pages that Google does, but you should be able to handle tens of thousands of pages and more reasonably. Comment on any scalability issues in your report.

6 Karma

There are numerous things you can do to make your project more fun and entertaining. You can do the following or come up with your own ideas.

- Modify the parser so that it excludes common words such as “where,” “what,” “how,” “and,” “or,” “a,” “an,” “of,” and “I”. If a common word is essential to a search, the plus sign (+) can be used to explicitly include it in the search. For example, `Star Wars Episode +I` will include the “I” in the search. When common words are excluded, Google provides an explanation, so you can determine what Google's common words are

through trial and error. More information about Google's query language can be found at the following URL: <http://www.google.com/support/websearch/>. This page also describes Google's Advanced Search options, which provide additional ideas.

Note that this requires changes to the parser and lexer. Be sure that you do not break anything in the process.

- Support full negative queries. There are a few ways to do this. You can choose to support them directly in your search engine; this will have some impact on your design choices. You can also remember DeMorgan's Law and other logic principles and convert full negative queries into something you can handle. Think about how you want to handle precedence of negation if you do this.
- Remove the restriction on parentheses. You can modify the query language so that parens are no longer required, and the operators have the appropriate precedence. This requires more challenging modifications to the parser. Again, be sure that you do not break support for the required query language in the process.
- Rank your results in some meaningful manner. Many search engines rank pages by frequency of access, but since you will not have such information, you could try to order them by categories (like Yahoo) or by connectedness (like Google). Alternatively, you could provide some other metric of goodness, perhaps by looking at various HTML tags to get clues about the contents of the page; for example, a hit in the title of a web page might be ranked higher than a hit in the body of a web page, or any combination of these methods.
- The existing search engine page and applet isn't very good. In fact, it's pretty awful. If you are so inclined, you may modify it to be more functional and standards-compliant. Google is a good example of an attractive and functional interface. The old MSN Search site is a good example of a bad interface, although it has been recently changed to look like Google.

If you implement highlighted excerpts in the results page or a highlighted "Google cache" you will need to add your own methods to support such functionality. Feel free to do so, but be sure not to break the required `useWebIndex` and `query` methods.

- Learn a bit about database query optimization. While we are doing full text searches and not relational database searches, you might find some inspiration for techniques to make your searches more efficient here.
- Explore the notion of *data mining* to produce a list of related pages even when those related pages do not match the search string. You can use text data mining, also known as *text analytics*, to support such a feature. Be warned that data mining is an open research topic that can consume arbitrary amounts of time, but the basics should be approachable.

The easiest way to support this *related pages* query is to use *cosine similarity*, which is a common distance metric used in classic information retrieval. The idea is that a page can be represented as a high dimension vector with each *meaningful* word representing a different dimension and the value of the dimension being the frequency of that word. To find related pages, simply look for the distances among all the vectors and choose the k nearest vectors using a distance metric; cosine similarity is often used as the distance metric, but any distance metric can work with varying results.

When determining similarity, however, we are relying on the rarity of certain words to provide differentiating abilities. Therefore, it is critical to magnify these differences to provide meaningful separation. Popular techniques include stop-lists, which remove common words such as *the*, *and*, *a*, *I*, etc., and frequency analysis, which removes words that are common in the corpus of documents being considered. You may also want to stress certain types of words, such as verbs and nouns, or you may wish to explore more advanced natural language processing (NLP) or information theoretic ideas.

In selecting a method to enhance the performance of the queries, be careful to account for the efficiency of these queries, which is usually the most difficult part. It's often too expensive to find the *best* match, so consider investigating fast approximate schemes and the use of acceleration structures.

If you do attempt to do text mining, you may also find it helpful to build at least a basic visualization tool to map the vectors from R^n to R^2 or R^3 space. These tools can be invaluable for providing insight into the structure of the documents.

There is also the closely related, but more advanced notion of *concept mining*, which would allow your system to accept a query and find pages related by the concept of the query, rather than its text or boolean operation. The idea is that the web pages that you have crawled are full of information—they consist of far more than the corpus of their words—but the difficulty is in extracting this information in a reasonable way.

For example, a query for American Mustang, should be able to provide at least two groups of pages: those related to the horse and those related to the car. In the former case, the system would ideally find all pages surrounding the concept of Mustang horses—issues such as land management, populations, history, etc. In the latter case, the system would ideally find information about the Ford Motor Company, the car’s history, current sales trends, industry reviews, etc. Identifying and differentiating between such concepts is an active research topic, though a few basic pioneering commercial applications exist.

The general approach is to first annotate the page according to word families (sets of words with similar meanings, as would be found in a thesaurus); this step usually works by clustering word groups together using a tool such as Princeton’s WordNet. This step reduces the number of words in the language to a smaller (but still large) number of clusters. You can then attempt to use techniques analogous to text data mining over these clusters. For better accuracy, *Bayesian models* are often utilized.

7 What to turn in

You will do this assignment individually. Submit your report and program in the usual manner.

Source Code. Turn in `WebCrawler.java`, `WebIndex.java`, `WebQueryEngine.java`, and all other source files that you’ve created or modified that are required to support your implementation. You are allowed to modify any and all files as long as you have the required `useWebIndex` and `query` methods.

Report. Your report will be very important. Since we give you tremendous freedom in design, you should explain your design and analyze your design decisions, in addition to all the usual expectations.

Acknowledgments. This assignment was inspired by a similar assignment in Cornell’s equivalent course, CS312. Many thanks to Walter Chang and Andrew Dreher for improvements to this assignment.