

Project #1: Web Crawler

1 Objectives

This project is designed to reinforce a number of learning objectives. First, and most important, you will gain experience designing and implementing a small (but interesting) Java application. This will reinforce your object-oriented design skills, as well as your ability to effectively apply the Java development tools. More specifically, the project will reinforce your understanding of Java interfaces, the Singleton pattern, Java exceptions, and the Java IO library. It will also improve your ability to debug Java-based systems.

Before describing the project requirements, two words of caution are in order. First, you should be prepared to experience some frustration with the level of detail provided in the assignment description. In order to learn object-oriented design skills, you need to be given a chance to do some design — which means that all of the implementation steps can't be laid out for you in advance. This level of flexibility is more in line with what you will encounter when you are employed as a developer.

The second warning is that this is a time-consuming assignment. You are expected to take the time to investigate the Java class libraries. Some of the most relevant classes are identified in this document. Others could significantly simplify your implementation — but you're expected to identify those through independent exploration. You will also be expected to learn basic HTML syntax, if you haven't done so already. There are a number of websites devoted to this topic. Have fun! Explore!

2 Overview

For this project, you will implement a simple command-line *web crawler*. The application will accept three arguments. The first will specify the web address from which the crawl should begin. The second (an integer) will specify the maximum crawl depth. Finally, the third argument will specify the local directory used to store results from the crawl.

When the crawl begins, it will perform a basic parse of the web page specified at the command-line. The goal of the parse is to identify the *(i)* images used within the page, the *(ii)* other files linked from within the page¹, and the *(iii)* web pages linked from within the page. The first two element sets will be saved to disk in the directory specified at the command-line. The third set of elements will be used as input to the next step in the page crawl; this process will be applied recursively up to the maximum specified crawl depth. No web page should ever be crawled more than once, even if it is encountered through independent search paths.

It is worth emphasizing that the term “*parse*” is being used loosely here. You are not required to do any sophisticated HTML processing; use the simplest solution that works. (Hint: The String

¹Excluding web pages and image files

class provides several useful methods for identifying substrings.) If you wish, you may also use the jsoup library to assist in the parsing. This is the only library you may use on this assignment. If you do use the jsoup library, you should include the jar file in your submission.

3 Requirements

Again, much of the design work is left to you. So while you must satisfy the requirements listed here, many of the details are left to you. This is the fun part of programming!

- You must follow good design practices throughout your implementation. In particular: (i) Program to interfaces wherever possible. (ii) Never use public fields. (iii) Follow Java naming conventions. (iv) Simplify your code by eliminating redundancy. (v) Catch and handle all checked exceptions.
- All of your classes and interfaces must be defined in the `cu.cs.cpsc215.project1` package.
- Your main application class must be named `WebCrawler`.
- Your application must include a `WebPage` class used to represent a web page. The class must include the following methods: `crawl()`, `getImages()`, `getFiles()`, and `getWebPages()`. The `crawl()` method is used to parse the web page represented by the target object. The remaining methods are used to retrieve the results of the parse. Hence, the `get` methods must only be invoked after `crawl()` has terminated.
- Your application must include classes for representing the images, files, and web pages linked from within a page. These classes must be named `WebImage`, `WebFile`, and `WebPage`, respectively. (The last class is discussed above.) Instances of these classes will be returned by the `WebPage` *getter* methods discussed above.
- `WebImage`, `WebFile`, and `WebPage` must each implement the `WebElement` interface. This interface must provide methods to enable a client component (i.e., a user of a `WebElement` instance) to save the corresponding element to disk.

Note that while the `WebPage` class must correctly implement `WebElement`, your crawler application must *not* save `WebPage` objects to disk. `WebPage` implements the `WebElement` interface only for the sake of future program enhancements.

- Your application must implement a *singleton* class, `DownloadRepository`. An instance of this class must be suitably initialized at system startup using an appropriate static initializer. This singleton must expose an interface that can be used by the crawler class (`WebPage`) to save a `WebElement` object to the repository. Your class must handle filename collisions via suitable renaming, should they occur.

4 Submission Instructions

This project is due before class on Monday, October 6th. Absolutely no late assignments will be accepted.

When your project is complete, archive your source materials and use `handin.cs.clemson.edu` to submit your work. You must name your archive **project1.zip**, and it must compile and run when unpacked with the following script (note that the arguments to `WebCrawler` are representative):

```
unzip project1.zip
javac -cp ./jsoup*.jar cu/cs/cpsc215/project1/WebCrawler.java
java -cp ./jsoup*.jar cu.cs.cpsc215.project1.WebCrawler http://some-address 3 /path/to/output/folder
```

Note that, if you use `jsoup`, you must include the `jsoup` jar in the root of your zip archive. If you do not use `jsoup`, you do not need to include the jar in your archive. Projects that fail to build and execute correctly with this script will receive a 0.

You must also submit (on handin) a PDF document containing your source code before the beginning of class on the due date. Each source file should begin on a new page and have a header stating the file name. The document should also include a cover sheet containing the title of the assignment, the date, your name, and any other relevant information you wish to provide.

NOTE: Technical problems are not an excuse for late assignments. Prepare your submission materials in advance to avoid last-minute crises.

5 Grading

Your project will be graded based on the quality of your design and implementation, as well as the functionality of the application itself. You are expected to correctly handle all basic web pages, under the assumption that the web page content is correct. If you encounter an error during processing — due to improper HTML formatting, embedded scripts that are not handled by your parser, etc. — your application must handle the error gracefully — this means terminating only the affected branch in your crawl.

This is an intermediate course in software development. Your source materials should be properly documented. Your source must compile. Your application must not crash. A violation of any of these requirements will result in an automatic zero. **Test your application thoroughly.**

6 Collaboration

You may work independently or with one partner. You must *not* discuss the problem or the solution with classmates outside of your group. Keep this in mind before you choose to work independently.

Collaborating in any manner with peers outside of your group will be treated as academic misconduct.