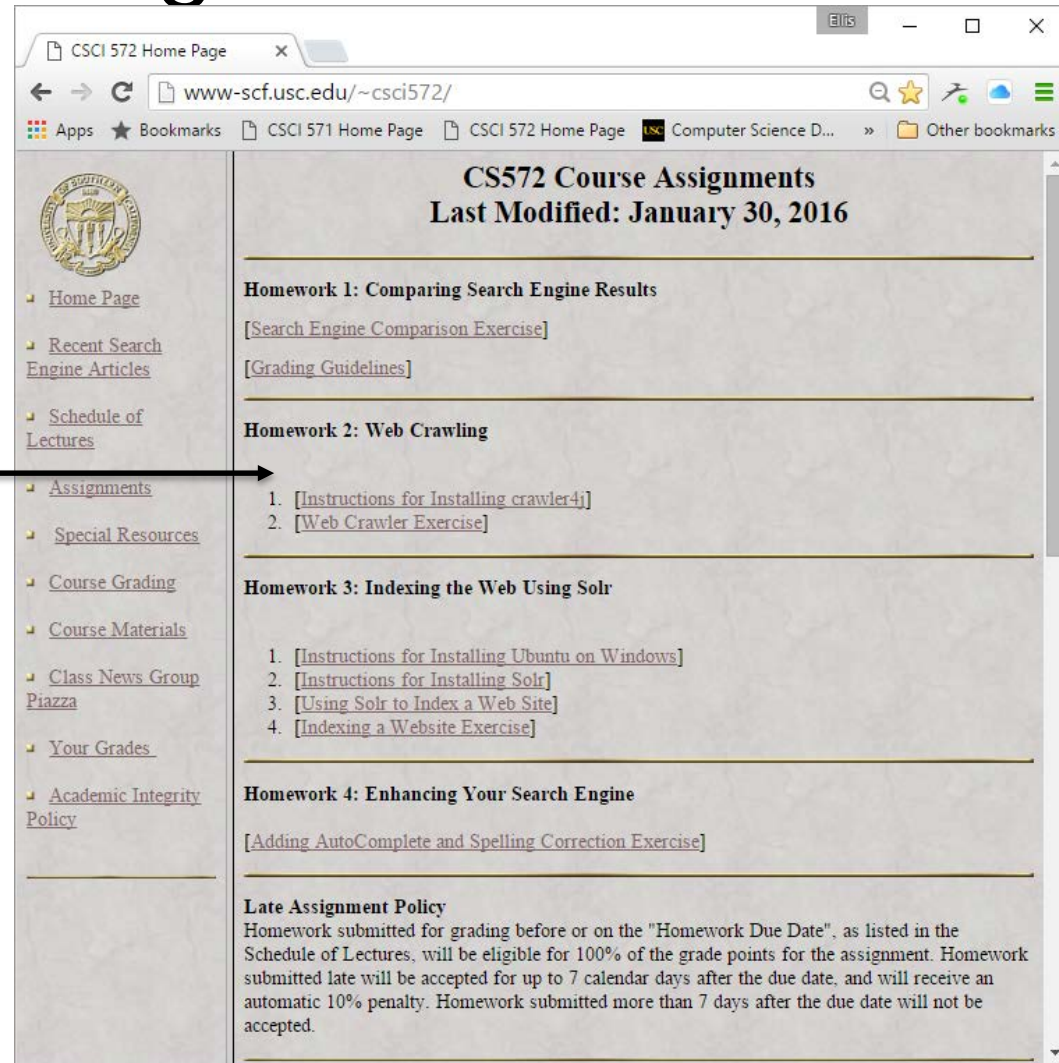# HW2 Assignment

- Involves
  - Java programming
    - I assume all of you know how to program in Java!
  - Eclipse Software Development Environment
  - crawler4j, an open source java web crawler
  - a crawl and analysis of a web site and an analysis of the crawl

# What is Eclipse?

- Eclipse started as a proprietary IBM product (IBM Visual age for Smalltalk/Java)
  - Embracing the open source model IBM opened the product up
- Open Source
  - It is a general purpose open platform that facilitates and encourages the development of third party plug-ins
- Best known as an Integrated Development Environment (IDE)
  - Provides tools for coding, building, running and debugging applications
- Originally designed for Java, now supports many other languages
  - Good support for C, C++
  - Python, PHP, Ruby, etc…

# Prerequisites for Running Eclipse

- Eclipse is written in Java and will thus need an installed JRE (Java Runtime Environment) or JDK (Java Development Kit) in which to execute
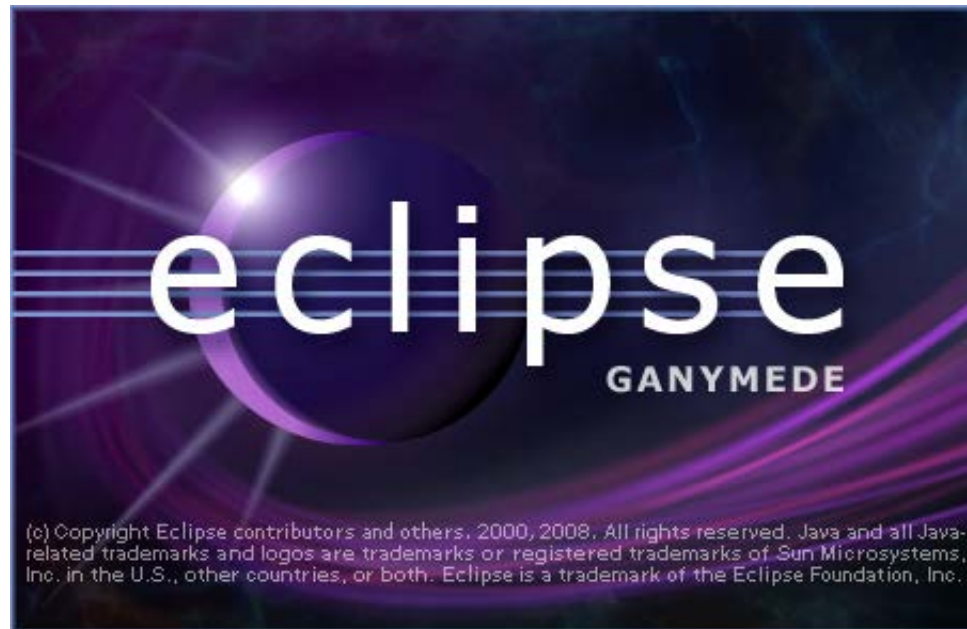  - JDK recommended

# Obtaining Eclipse

- Eclipse can be downloaded from…
  - http://www.eclipse.org/downloads/packages/
- Eclipse comes bundled as a zip file (Windows) or a tarball (all other operating systems)

# Installing Eclipse

- Simply unwrap the zip file to some directory where you want to store the executables
- On windows
  - I typically unwrap the zip file to C:\eclipse\
  - I then typically create a shortcut on my desktop to the eclipse executable
    - C:\eclipse\eclipse.exe
- Under Linux
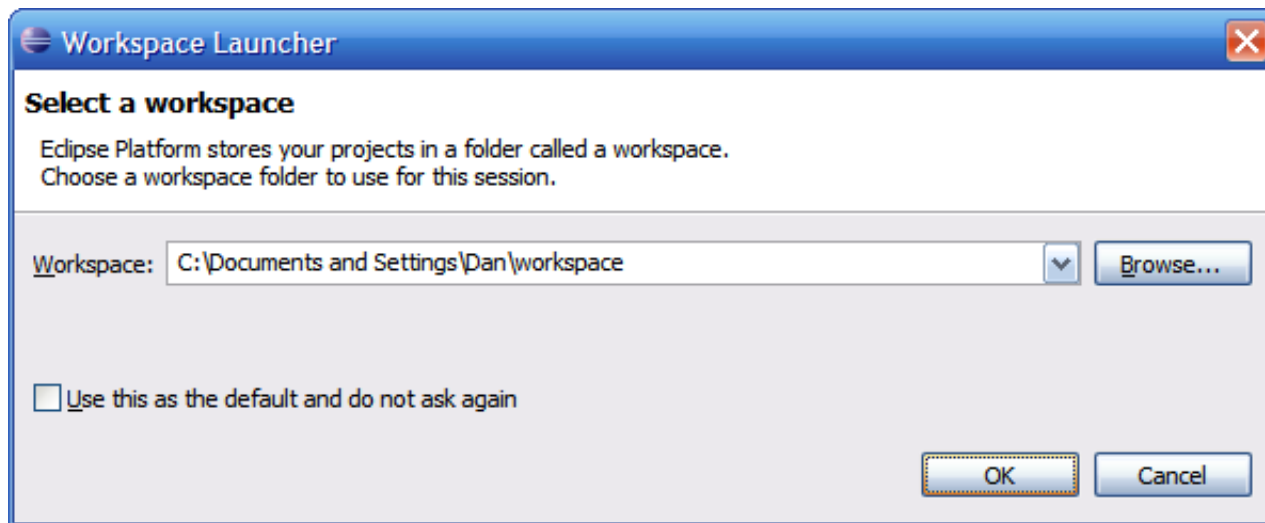  - I typically unwrap to /eclipse/

# Launching Eclipse

- Once you have the environment setup, go ahead and launch eclipse

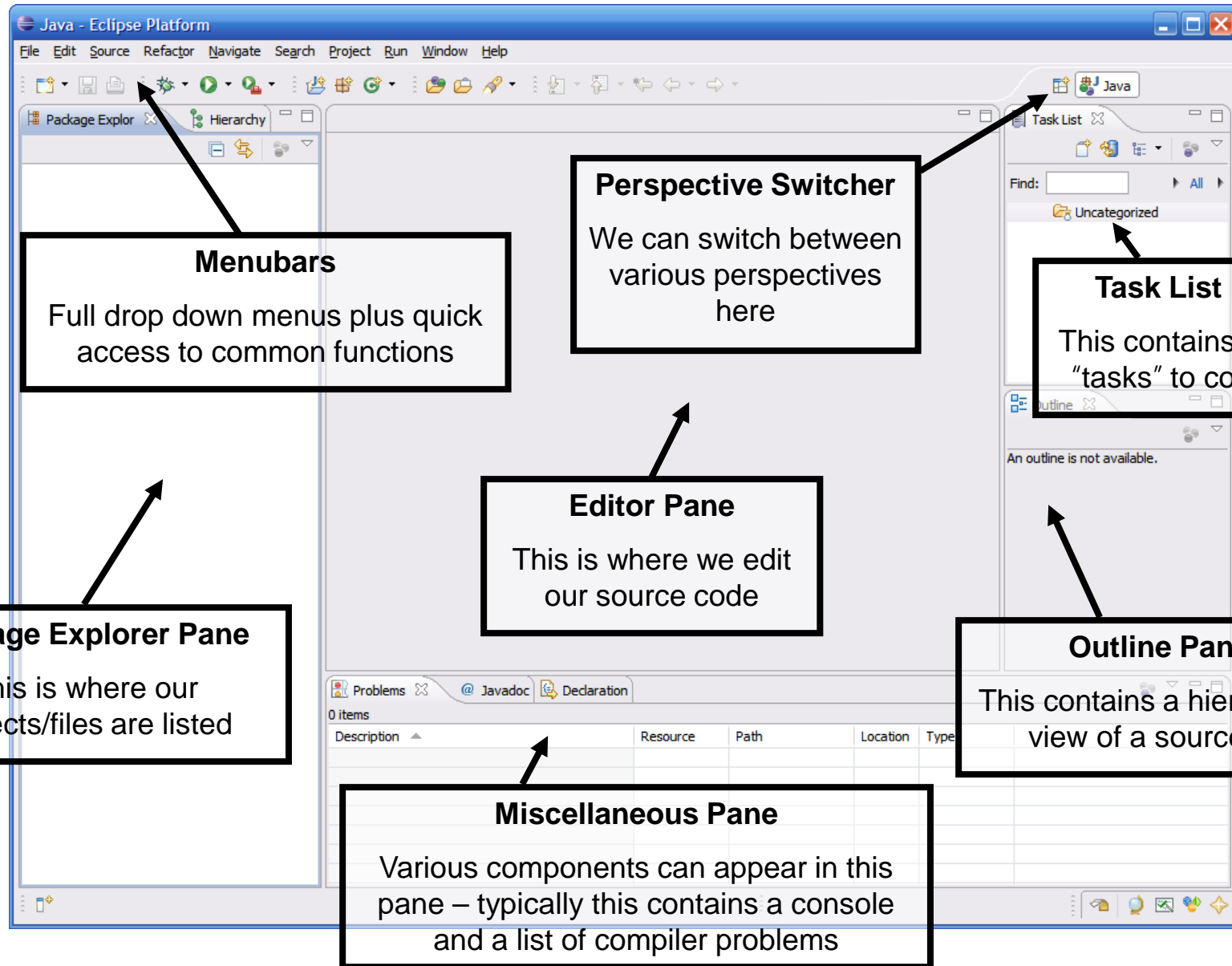- You should see the following splash screen…

# Selecting a Workspace

- In Eclipse, all of your code will live under a workspace
- A workspace is nothing more than a location where we will store the source code and where Eclipse will write out preferences
- Eclipse allows you to have multiple workspaces – each tailored in its own way
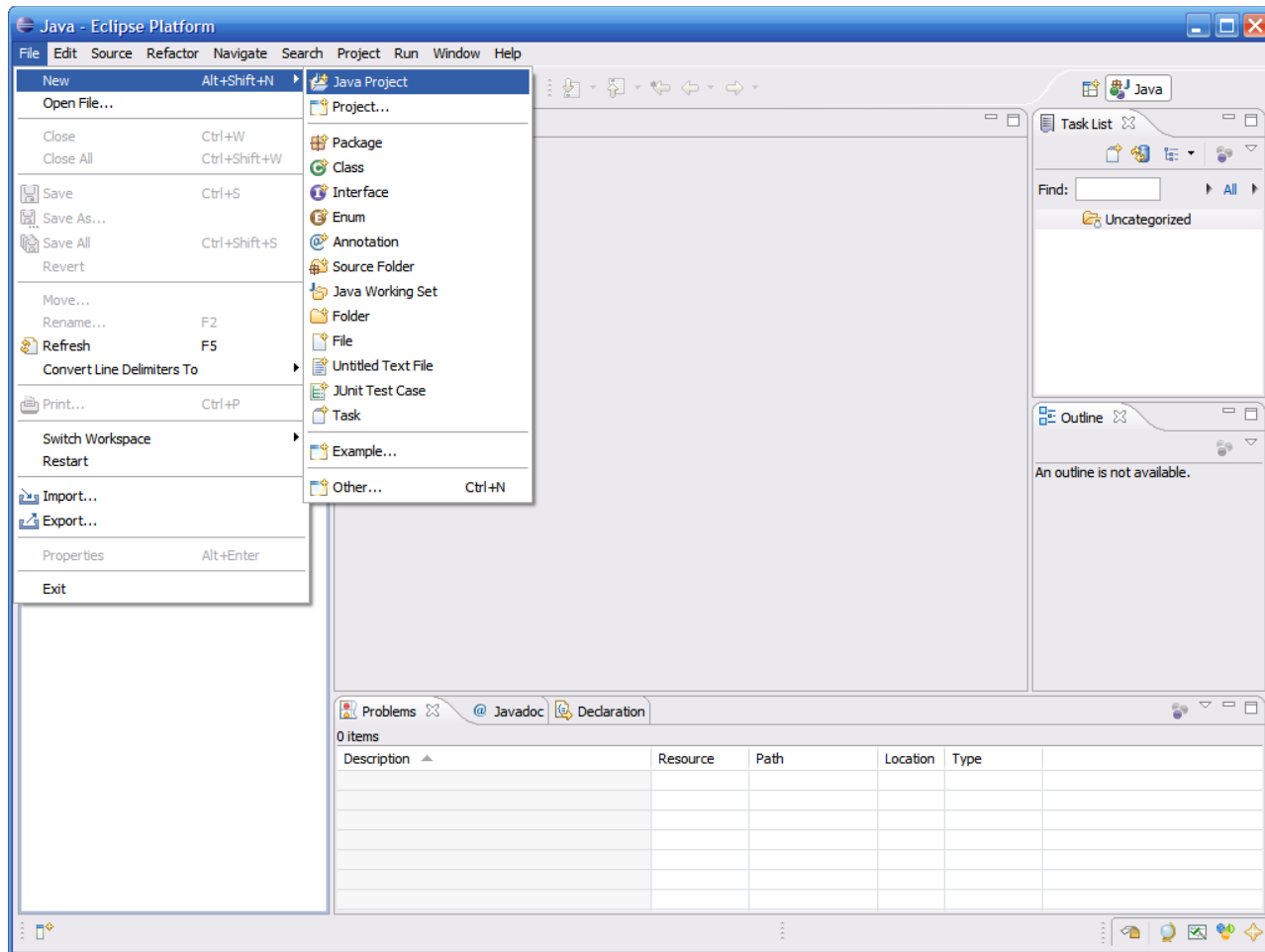- Choose a location where you want to store your files, then click OK
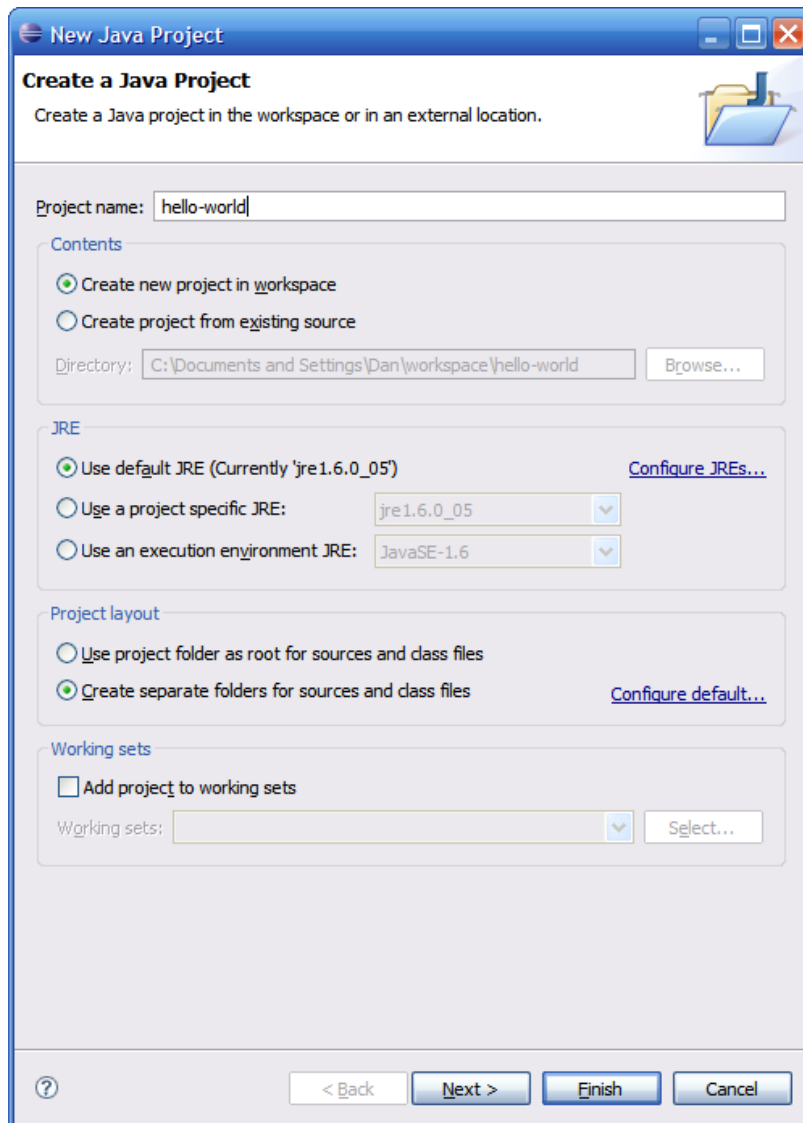
# Eclipse IDE Components



**Menubars**

Full drop down menus plus quick access to common functions

**Perspective Switcher**

We can switch between various perspectives here

**Task List Pane**

This contains a list of "tasks" to complete

**Editor Pane**

This is where we edit our source code

**Package Explorer Pane**

This is where our projects/files are listed

**Outline Pane**

This contains a hierarchical view of a source file

**Miscellaneous Pane**

Various components can appear in this pane – typically this contains a console and a list of compiler problems

# Creating a New Project

- All code in Eclipse needs to live under a project
- To create a project: File → New → Java Project
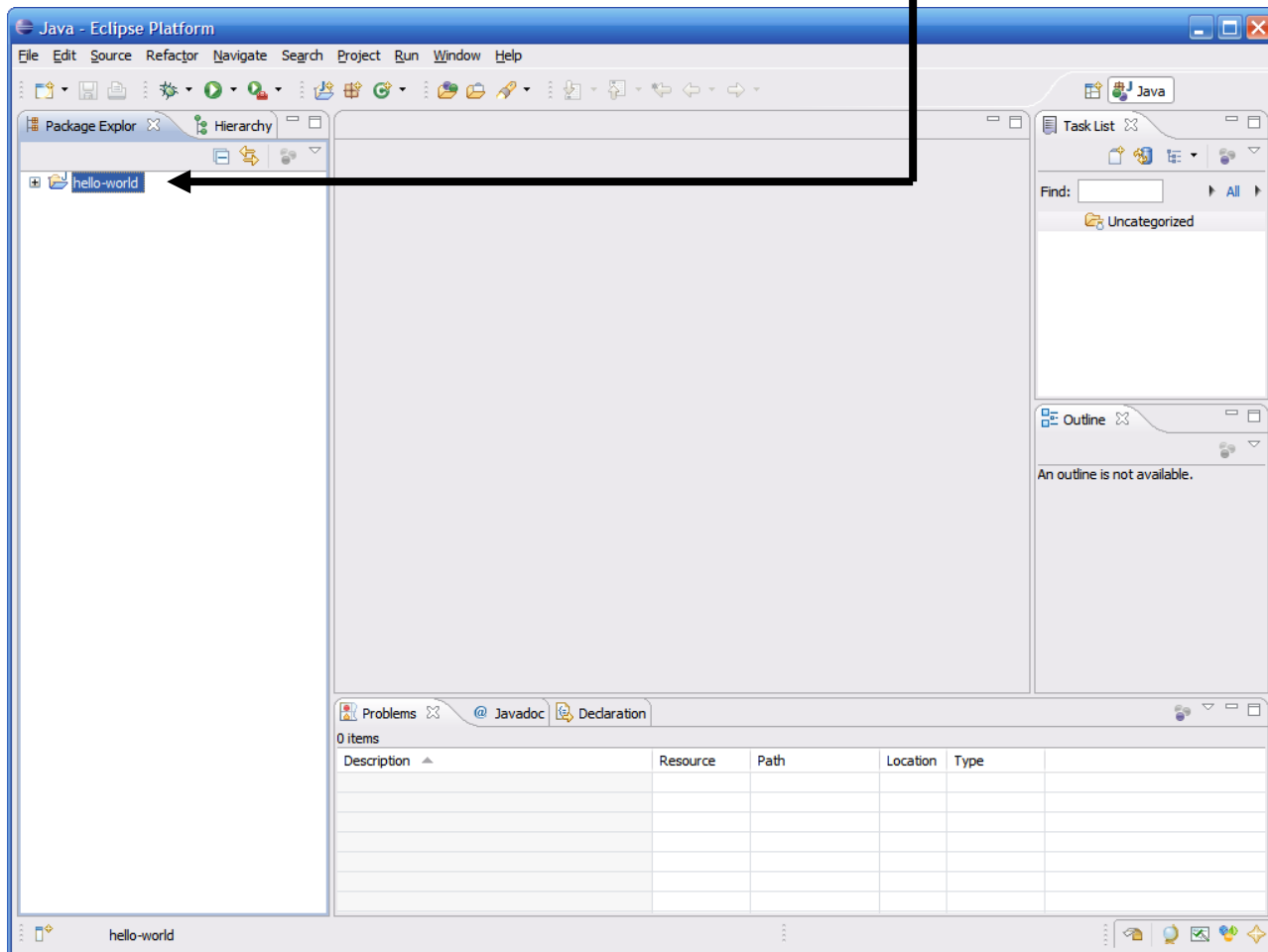
# Creating a New Project (continued)



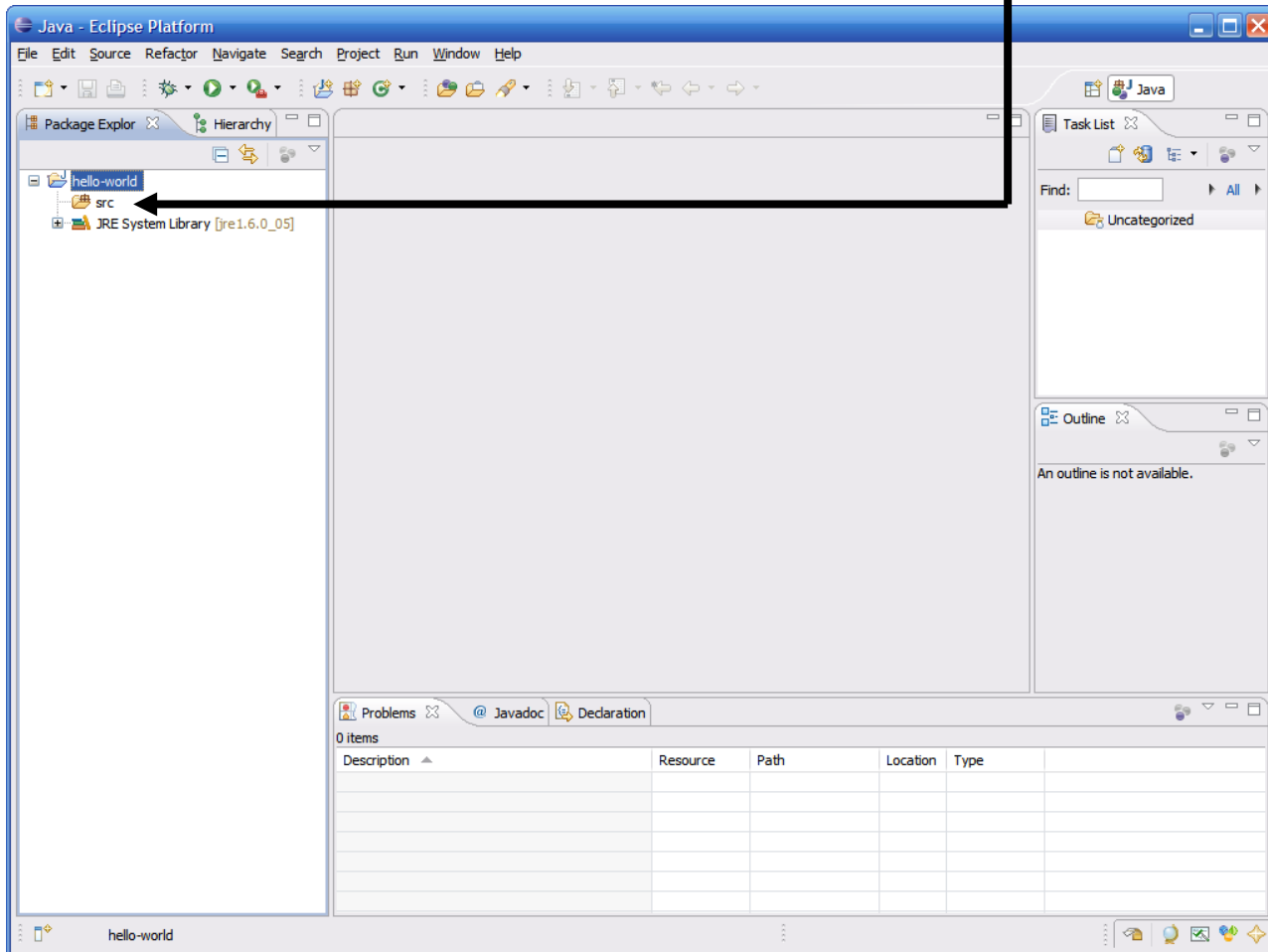- Enter a name for the project, then click Finish

# Creating a New Project (continued)

- The newly created project should then appear under the Package Explorer
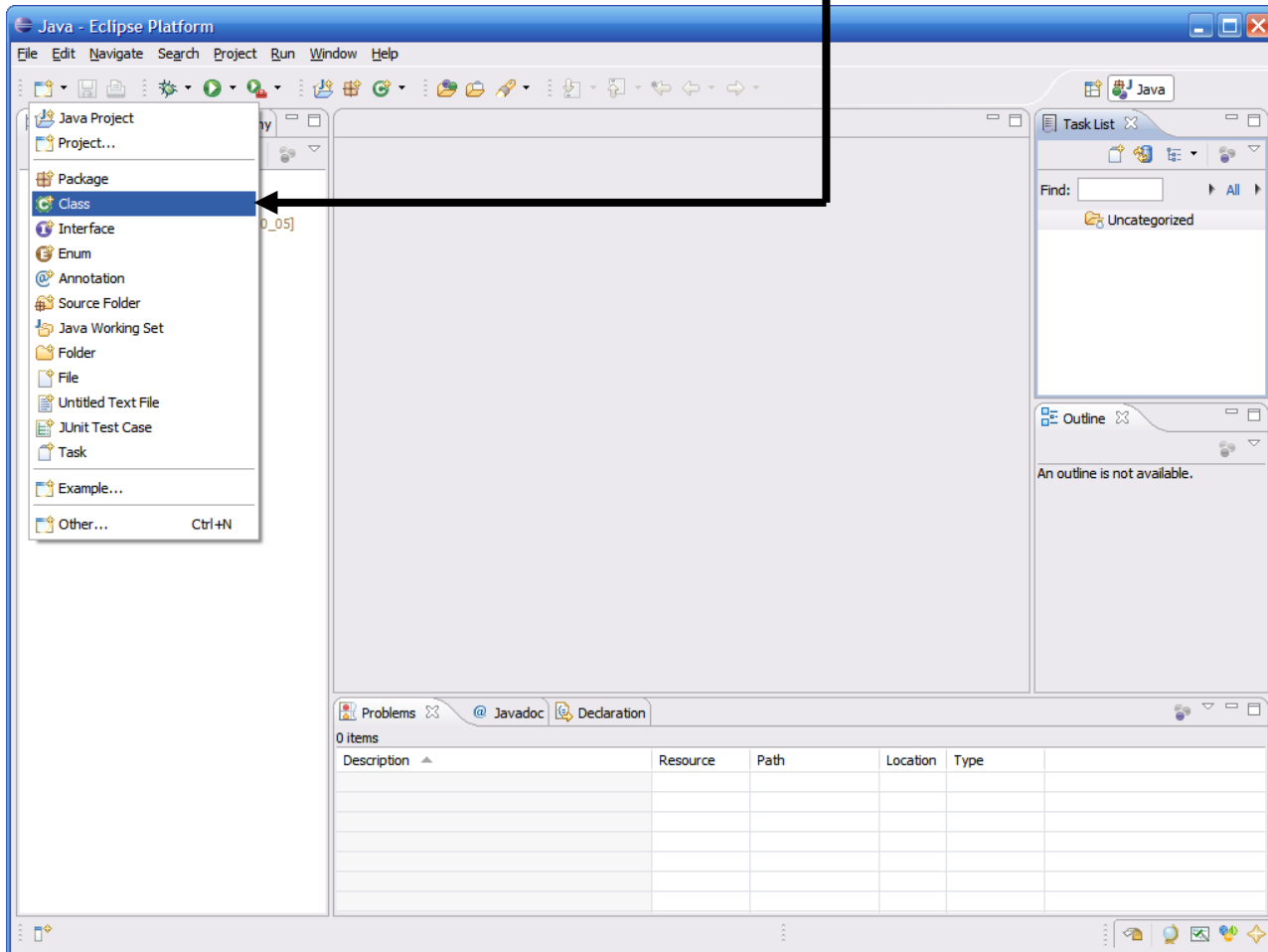
# The src folder

- Eclipse automatically creates a folder to store your source code in called src

# Creating a Class

- To create a class, simply click on the New button, then select Class

# Creating a Class (continued)



- This brings up the new class wizard
- From here you can specify the following...
  - Package
  - Class name
  - Superclass
  - Whether or not to include a main
  - Etc…
- Fill in necessary information then click Finish to continue

# The Created Class

- As you can see a number of things have now happened…



Directory structure for package and actual java file created automatically

Source is loaded into the editor pane, already stubbed out

Source displayed in a hierarchical fashion listing each method name

# Compiling Source Code

- One huge feature of Eclipse is that it automatically compiles your code in the background
  - You no longer need to go to the command prompt and compile code directly
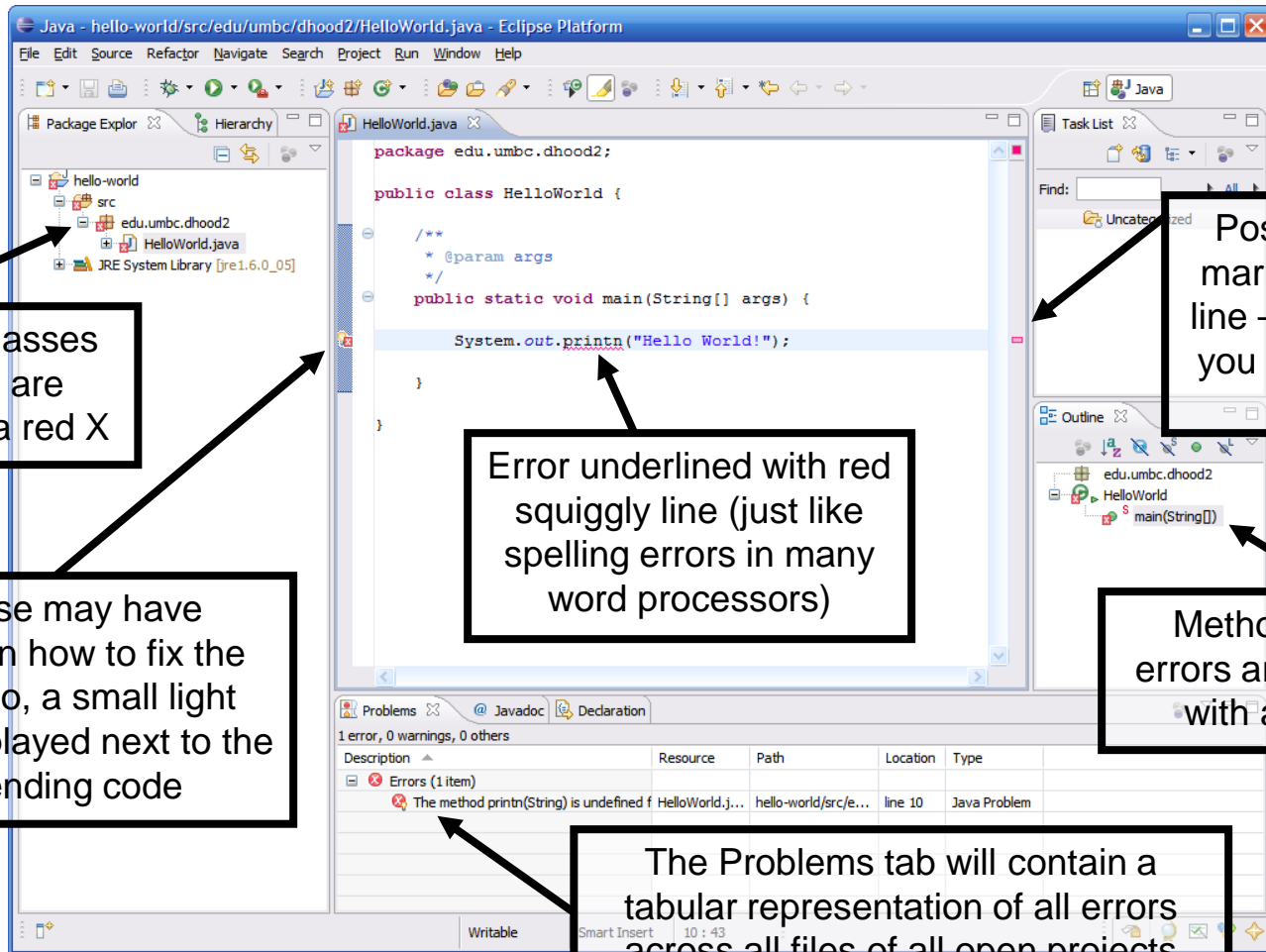- This means that errors can be corrected when made
  - We all know that iterative development is the best approach to developing code, but going to shell to do a compile can interrupt the normal course of development
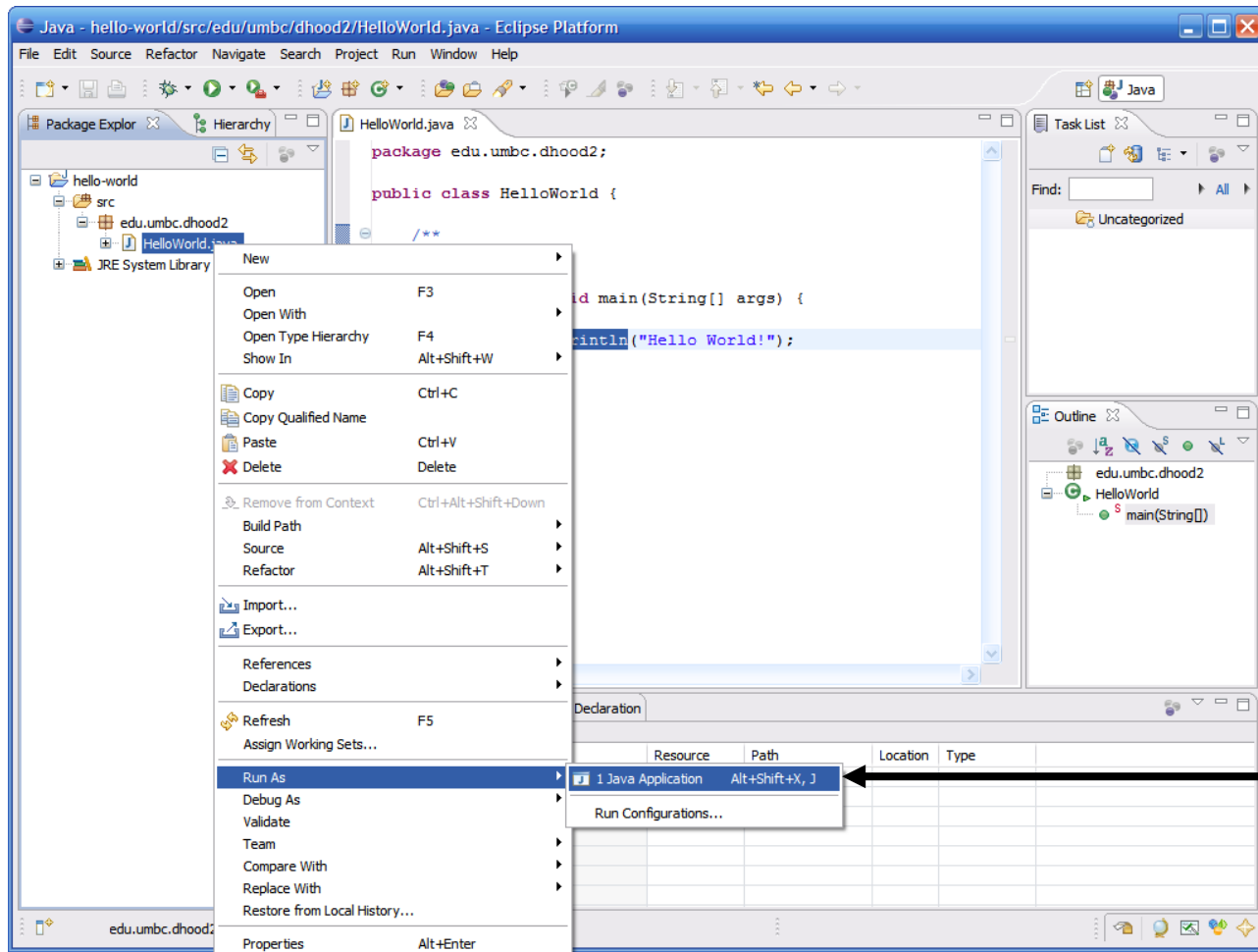
# Example Compilation Error

- This code contains a typo in the println statement…



Packages/Classes with errors are marked with a red X

Position in file is marked with a red line – 1 click allows you to jump to line with error

Error underlined with red squiggly line (just like spelling errors in many word processors)

Methods with errors are marked with a red X

Often Eclipse may have suggestions on how to fix the problem – if so, a small light bulb will be displayed next to the line of offending code

The Problems tab will contain a tabular representation of all errors across all files of all open projects
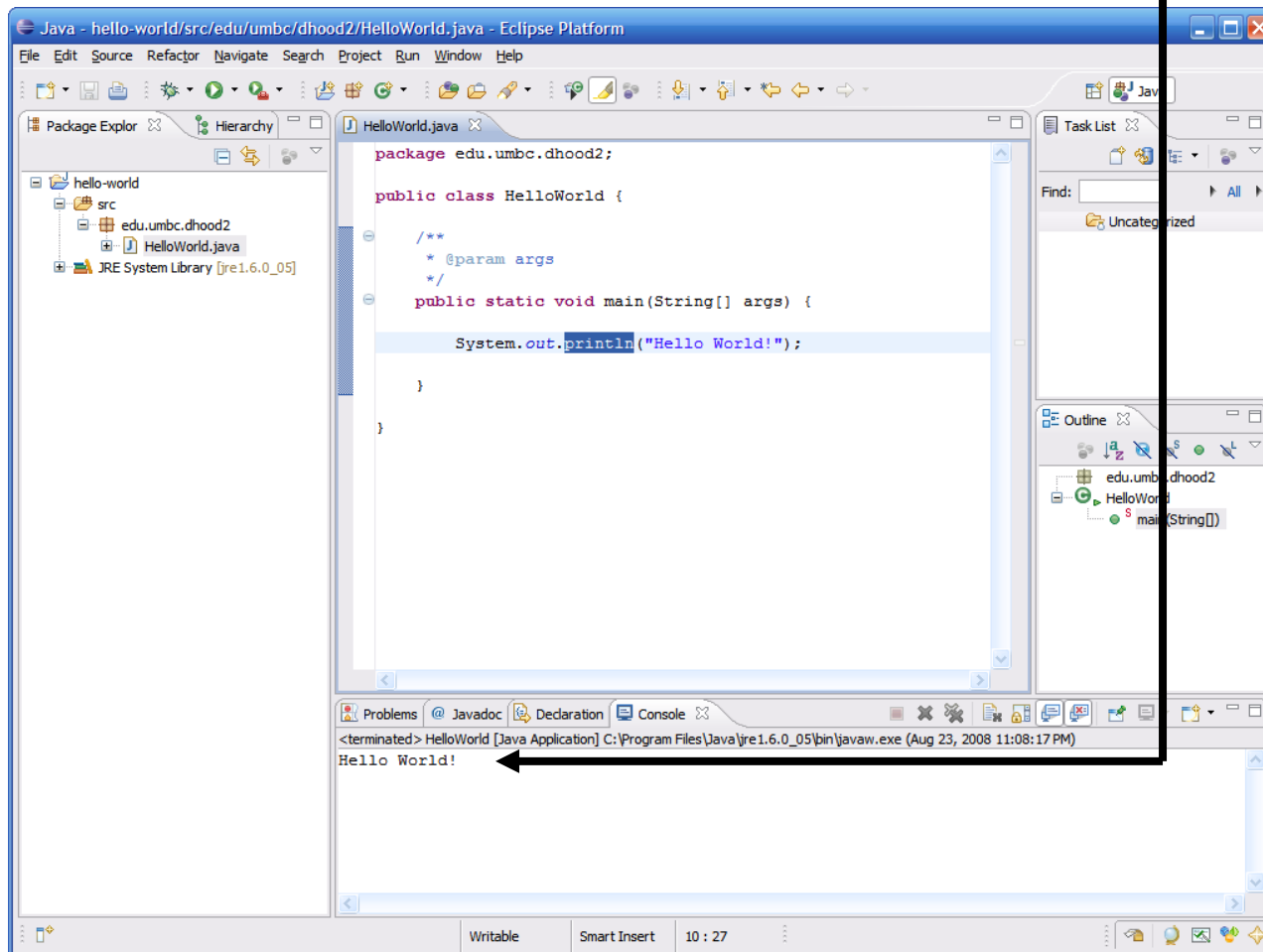
# Running Code

- An easy way to run code is to right click on the class and select Run As → Java Application
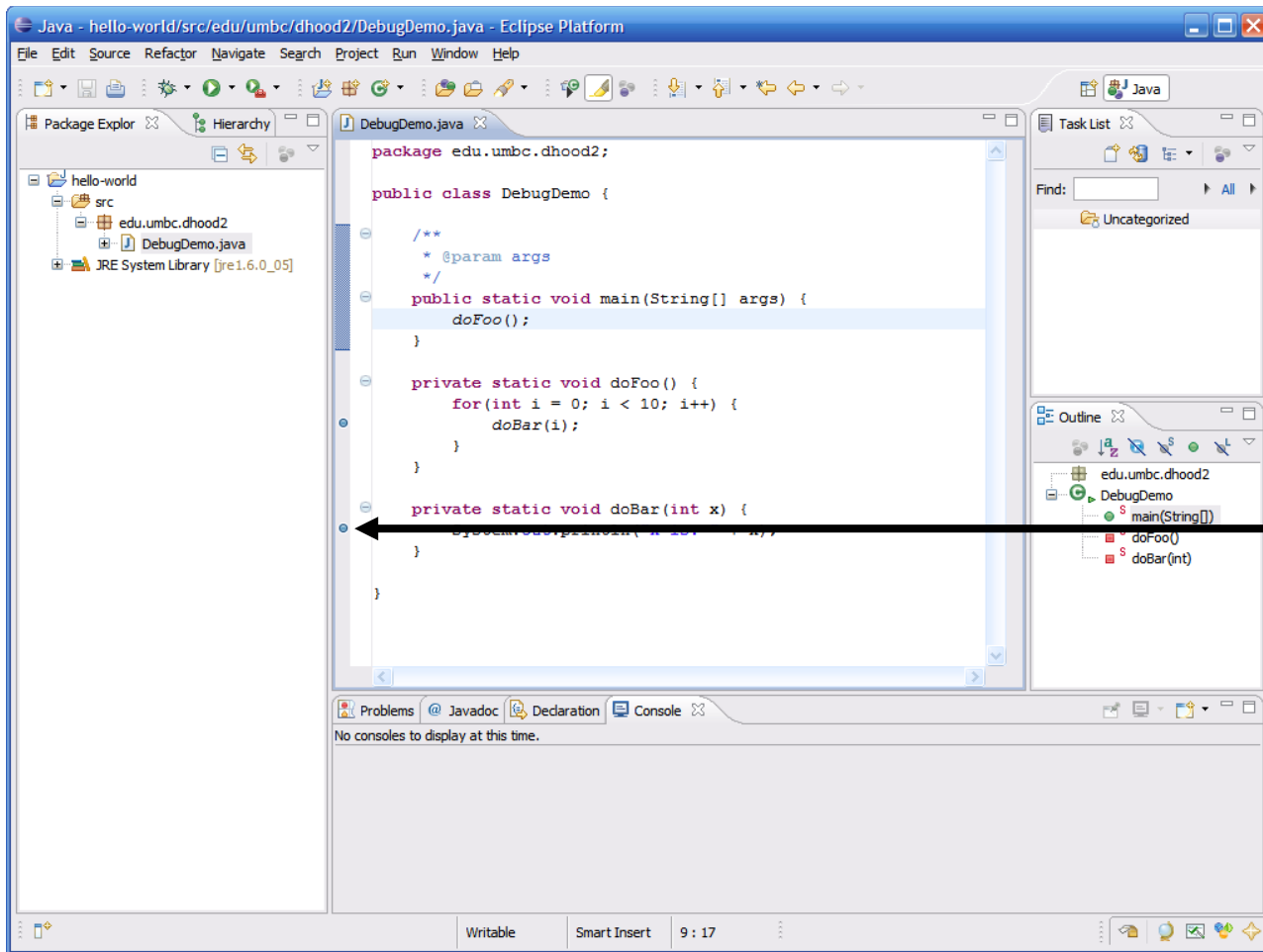
# Running Code (continued)

- The output of running the code can be seen in the Console tab in the bottom pane
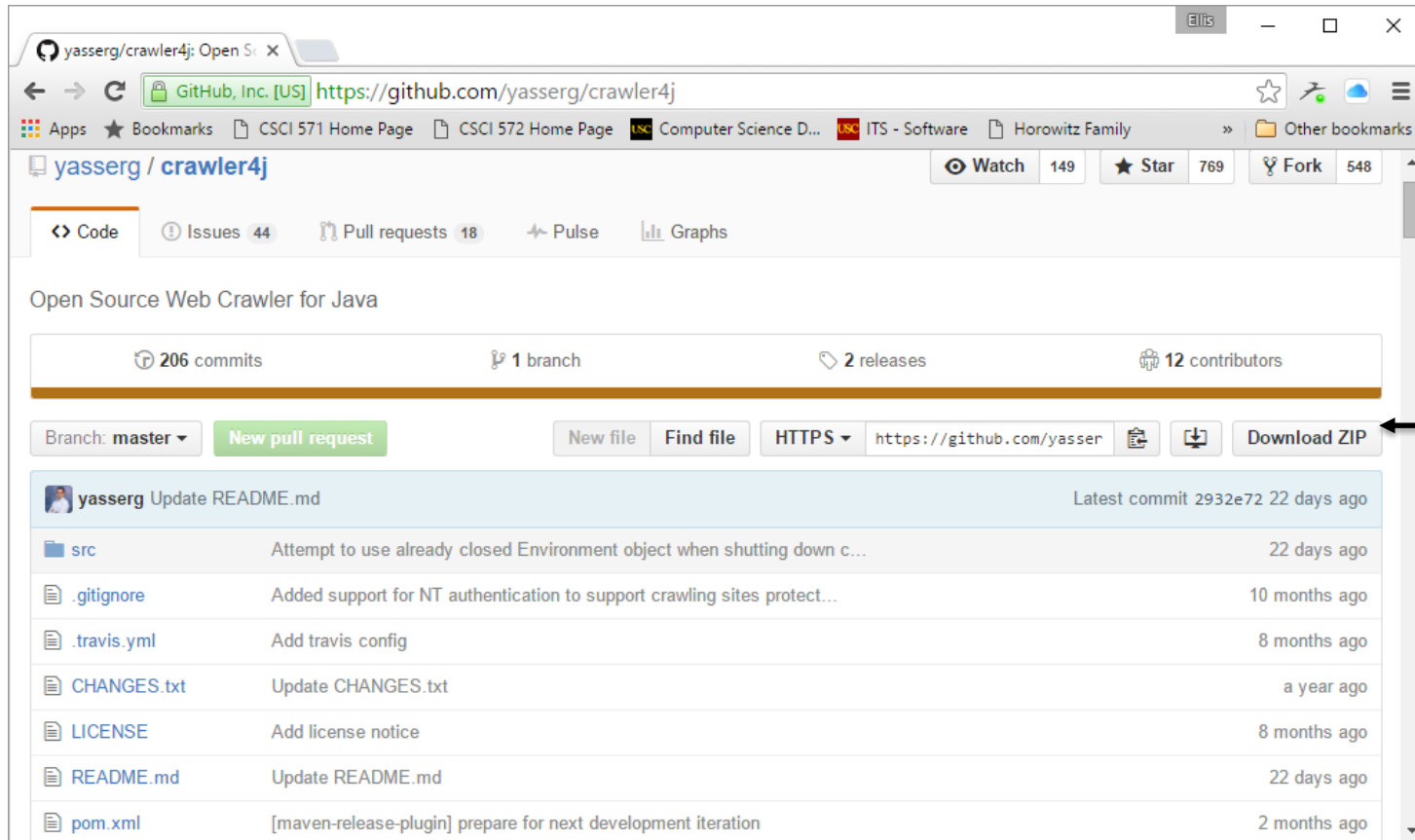
# Debugging Code

- Eclipse comes with a pretty good built-in debugger
- You can set break points in your code by double clicking in the left hand margin – break points are represented by these blue bubbles

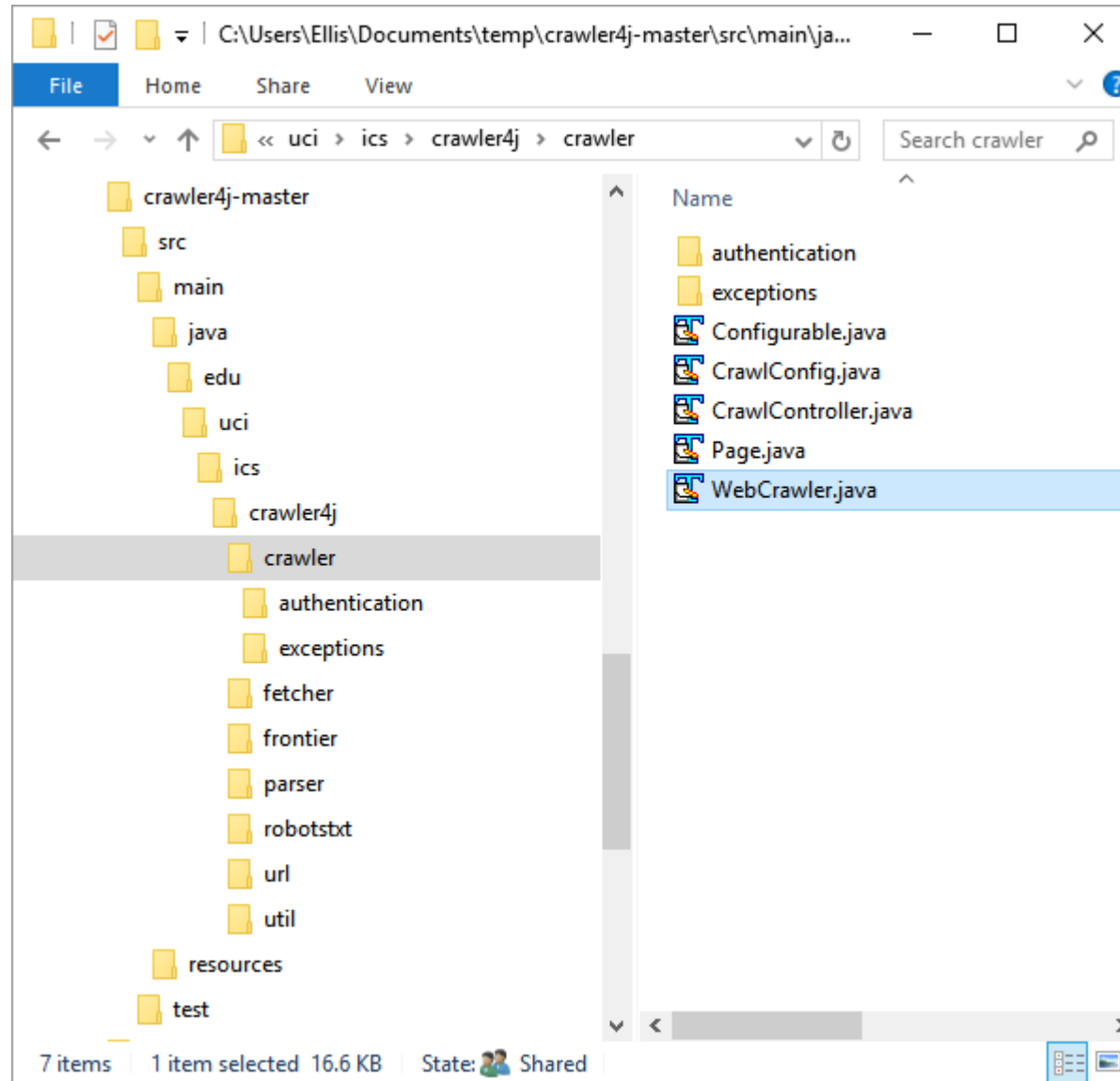# Instructions for Installing Eclipse and Crawler4j

- the document describes
  - installation of Eclipse on either Windows or MacOS
    - you will need to set Windows environment variables pointing to your JDK
- download crawler4j from github
  - **GitHub** is a web-based repository hosting service for software. Originally the Git system offered distributed revision control and source code management (SCM) functionality, but on the command line; GitHub offers a web interface and some additional features.
  - As of 2015, GitHub reports having over 11 million users and over 29.4 million repositories
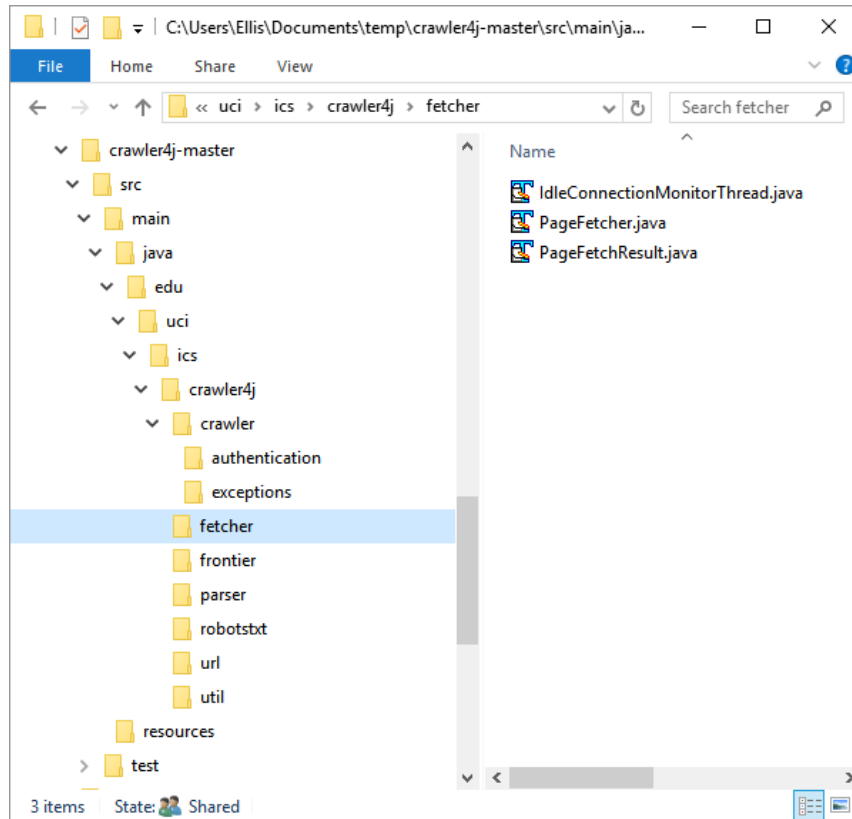
# Downloading Crawler4j from GitHub



See especially the README file page at
https://github.com/yasserg/crawler4j/blob/master/README.md
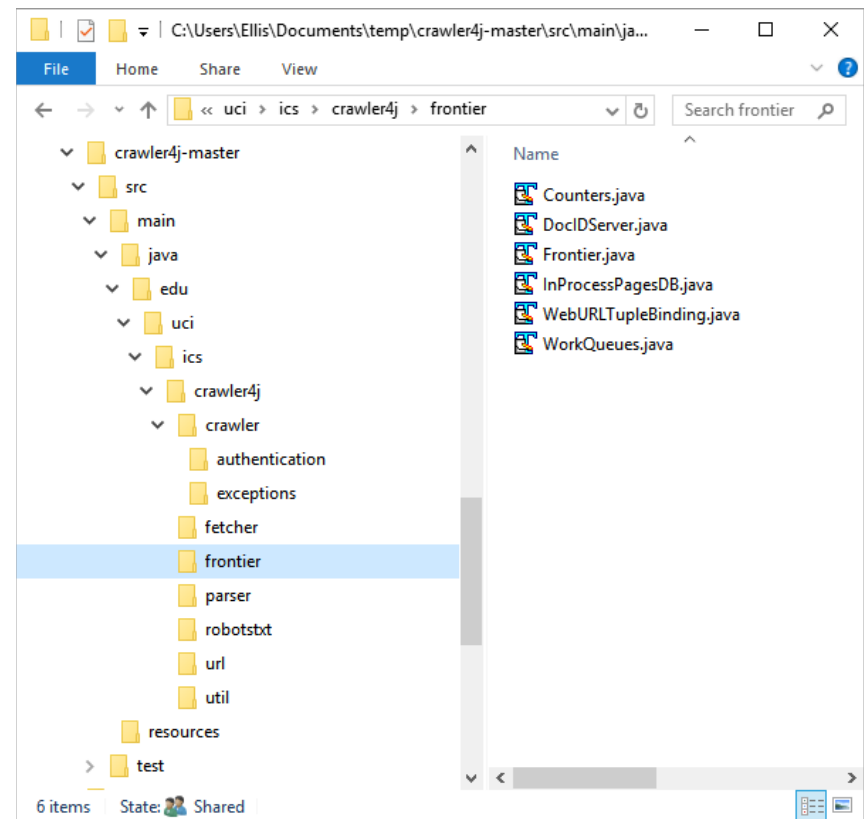
# Crawler4j Source code



Crawler folder, a good place to start; look especially at WebCrawler.java
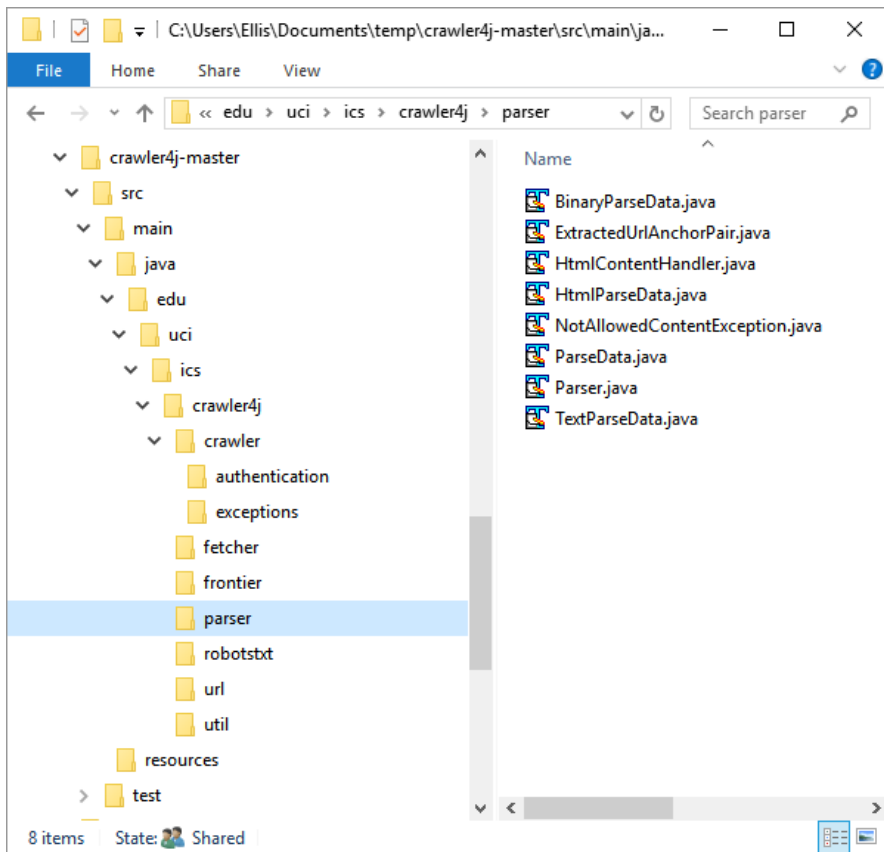
*Fetcher Code handles*:
- schemes: http, https
- politeness delay;
- redirects;
- max-size settings;
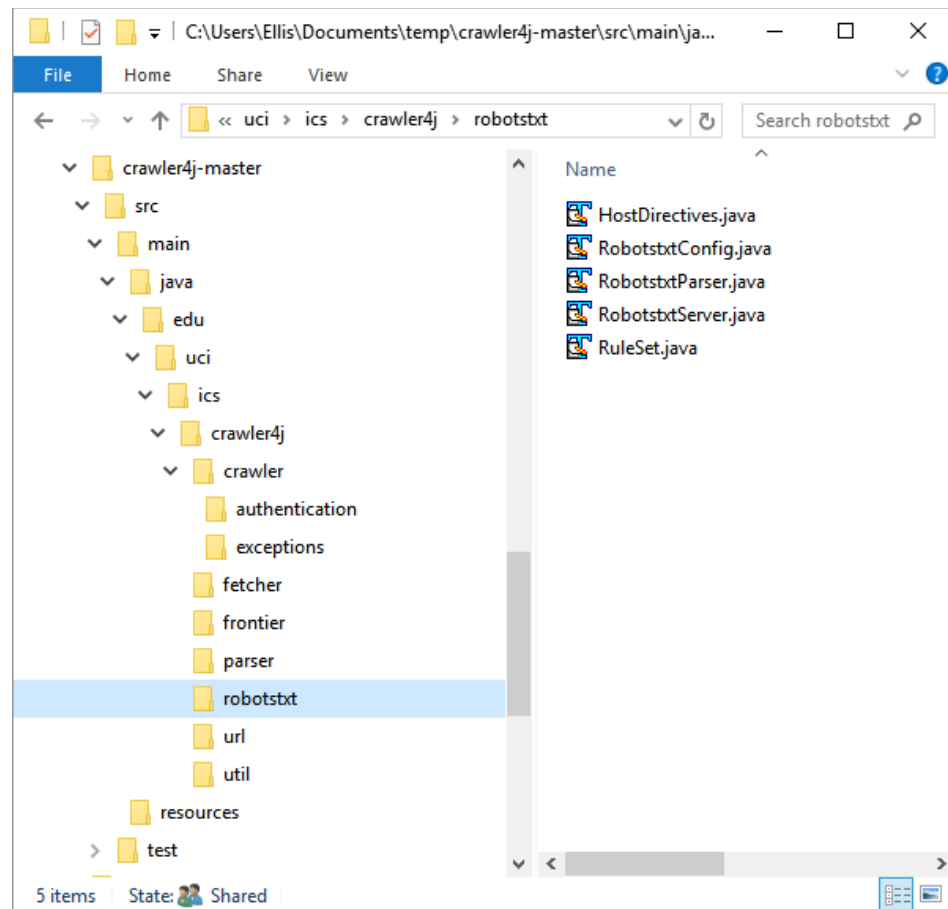- expired connections

*Frontier Code handles*:
- statistics database;
- previously seen URLs
- queue of pending URLs

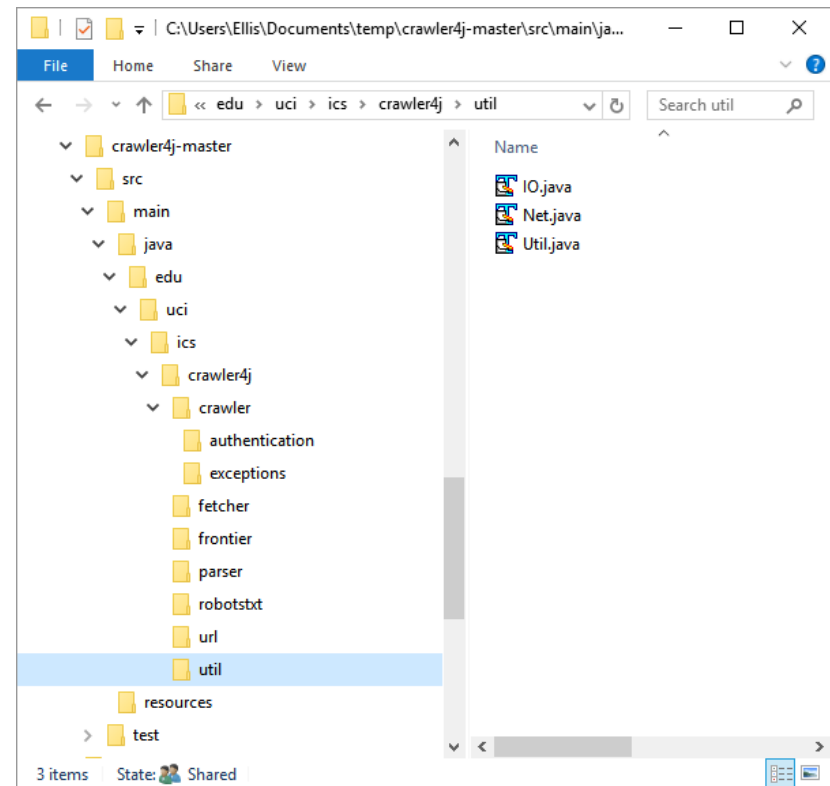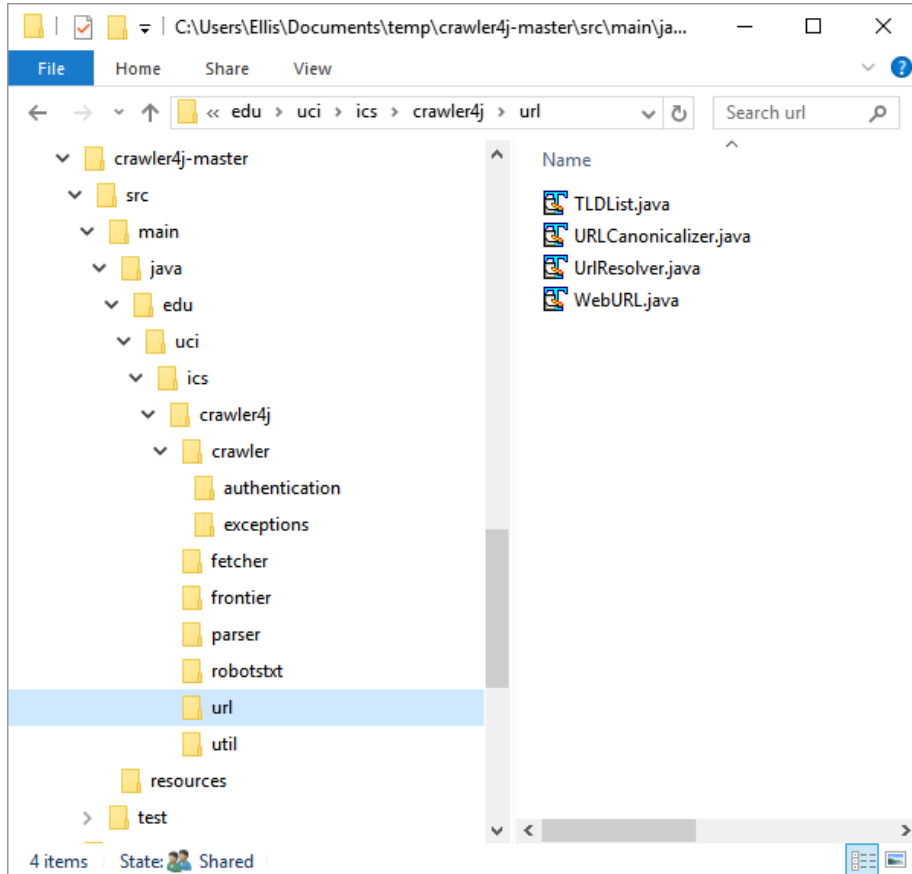Parser Code handles:
- binary data
- html pages
- extracting links

Robots.txt Code handles:
- fetching and re-fetching robots.txt
- caching robots.txt files
- interpreting commands
- working with Page Fetcher

utility routines

URL resolver and canonicalizer handles:
- checking against list of TLDs
- normalizes URL, removes . or .., etc
- alters name/value pairs
- converts #nn values
- evaluates <base>

# Configuring the Crawler and Seeding it

```
public class Controller {
    public static void main(String[] args) throws Exception {
        String crawlStorageFolder = "/data/crawl";
        int numberOfCrawlers = 7;
        CrawlConfig config = new CrawlConfig();
        config.setCrawlStorageFolder(crawlStorageFolder);
        /*  Instantiate the controller for this crawl.*/
        PageFetcher pageFetcher = new PageFetcher(config);
        RobotstxtConfig robotstxtConfig = new RobotstxtConfig();
        RobotstxtServer robotstxtServer = new RobotstxtServer(robotstxtConfig,
pageFetcher);
        CrawlController controller = new CrawlController(config, pageFetcher,
robotstxtServer);
        /*  For each crawl, you need to add some seed urls. These are the first
         * URLs that are fetched and then the crawler starts following links
         * which are found in these pages   */
        controller.addSeed("http://www.viterbi.usc.edu/");
        /*  Start the crawl. This is a blocking operation, meaning that your code
         * will reach the line after this only when crawling is finished.   */
        controller.start(MyCrawler.class, numberOfCrawlers);
    }
}
```

folder to store downloads; #crawlers

set up pagefetcher and robots.txt going

crawling viterbi.usc.edu

# Defining Which Pages to Crawl

```java
public class MyCrawler extends WebCrawler {
    private final static Pattern FILTERS =
Pattern.compile(".*(\\.(css|js|gif|jpg"  + "|png|mp3|mp3|zip|gz))$");    ⟵
   /** This method receives two parameters. The first parameter is the page
     * in which we have discovered this new url and the second parameter is
     * the new url. You should implement this function to specify whether
     * the given url should be crawled or not (based on your crawling logic).
     * In this example, we are instructing the crawler to ignore urls that
     * have css, js, git, … extensions and to only accept urls that start
     * with "http://www.viterbi.usc.edu/". In this case, we didn't need the
     * referring Page parameter to make the decision.    */
    @Override
    public boolean shouldVisit(Page referringPage, WebURL url) {
        String href = url.getURL().toLowerCase();
        return !FILTERS.matcher(href).matches()
              && href.startsWith("http://www.viterbi.usc.edu/");
    }
```

# Matching URLs

- ".*(\\.(css|js|gif|jpg" + "|png|mp3|mp4|zip|gz))$"

- A regular expression, specified as a string, must first be compiled into an instance of this class.
- a Matcher object that can match arbitrary character sequences against the regular expression
- See https://docs.oracle.com/javase/7/docs/api/java/util/regex/Pattern.html

- In the above there are two strings concatenated by plus; consider the simpler form:
- ".*(\\.(css|js|zip|gz))$"

- . matches any character
- * matches zero or more of preceding character
- \. matches a literal dot
- $ anchors the pattern at the end of the string

# Parsing the Downloaded Page

```java
/**   This function is called when a page is fetched and ready
  * to be processed by your program.   */
@Override
public void visit(Page page) {
    String url = page.getWebURL().getURL();
    System.out.println("URL: " + url);
    if (page.getParseData() instanceof HtmlParseData) {
        HtmlParseData htmlParseData = (HtmlParseData) page.getParseData();
        String text = htmlParseData.getText();
        String html = htmlParseData.getHtml();
        Set<WebURL> links = htmlParseData.getOutgoingUrls();
        System.out.println("Text length: " + text.length());
        System.out.println("Html length: " + html.length());
        System.out.println("Number of outgoing links: " + links.size());
    }
}
```

# The Actual Exercise

- *the URLs it attempts to fetch*, **fetch.csv.** The number of rows should be close to 5,000 as that is our pre-set limit.

- *the files it successfully downloads*, **visit.csv**; clearly the number of rows will be less than the number of rows in fetch.csv

- *all of the URLs that were discovered* and processed in some way; **urls.csv**. This file will be much larger than 5,000 rows as it will have numerous repeated URLs

# More things to Save

- you should save all of the downloaded web pages, etc. for processing in the next exercise.

# Even More things to Save

- Fetch statistics:
  - # fetches attempted:
    The total number of URLs that the crawler attempted to fetch. This is usually equal to the MAXPAGES setting if the crawler reached that limit; less if the website is smaller than that.
  - # fetches succeeded:
    The number of URLs that were successfully downloaded in their entirety, i.e. returning a HTTP status code of 2XX.
  - # fetches failed or aborted:
    The number of fetches that failed for whatever reason, including, but not limited to: HTTP redirections (3XX), client errors (4XX), server errors (5XX) and other network-related errors.
-

# Outgoing URLs

- Outgoing URLs: statistics about URLs extracted from visited HTML pages
  - Total URLs extracted:
    The grand total number of URLs extracted from all visited pages
  - # unique URLs extracted:
    The number of *unique* URLs encountered by the crawler
  - # unique URLs within School:
    The number of *unique* URLs encountered that are associated with the school website,
    i.e. the URL begins with the given root URL of the school.
  - # unique USC URLs outside School:
    The number of *unique* usc.edu URLs encountered that were *not* from the school website.
  - # unique URLs outside USC:
    The number of all other *unique* URLs encountered

Sample statistics Output

Sample
of
fetch.csv
file

# Page Rank

- For Assignment #3 you will need to compute the Page Rank of each page you download in this assignment.

- To do that you will need to keep a record of all URLs contained in a given page. Therefore you should also generate a csv file that includes every successfully downloaded html file, in column 1, and the outgoing URLs that were contained in the page, in subsequent columns 2, 3, 4, etc. You need not submit this file with your assignment #2, but you will need it when you get to assignment #3.