

COMP 221: Introduction to Computer Systems

Project 6: Web Proxy
Assigned: 4/8/14, Due: 4/24/14

Important: This project may be done individually or in pairs. Be sure to carefully read the course policies for assignments (including the honor code policy) on the assignments page of the course web site:

<http://www.clear.rice.edu/comp221/html/assignments.html>

Overview

In this project, you will write a concurrent Web proxy that logs requests. A Web proxy is a program that acts as a middleman between a Web browser and an *end server*. Instead of contacting the end server directly to get a Web page, the browser contacts the proxy, which forwards the request on to the end server. When the end server replies to the proxy, the proxy sends the reply on to the browser.

Proxies are used for many purposes. Sometimes proxies are used in firewalls, such that the proxy is the only way for a browser inside the firewall to contact an end server outside. The proxy may do translation on the page, for instance, to make it viewable on a Web-enabled cell phone. Proxies are also used as *anonymizers*. By stripping a request of all identifying information, a proxy can make the browser anonymous to the end server. Proxies can even be used to cache Web objects, by storing a copy of, say, an image when a request for it is first made, and then serving that image in response to future requests rather than going to the end server.

In the first part of the project, you will write a simple sequential proxy that repeatedly waits for a request, forwards the request to the end server, and returns the result back to the browser, keeping a log of such requests in a disk file. This part will help you understand the basics about network programming and the HTTP protocol.

In the second part of the project, you will upgrade your proxy so that it uses threads to deal with multiple clients concurrently. This part will give you some experience with concurrency and synchronization, which are crucial computer systems concepts.

Part I: Implementing a Sequential Web Proxy

In this part, you will implement a sequential logging web proxy. Your proxy should open a socket and listen for a connection request. When it receives a connection request, it should accept the connection, read the HTTP request, and parse the request to determine the name of the end server. Your proxy should then open a connection to that end server, forward the request to it, receive the reply, and forward the reply to the browser.

Your proxy need only support the HTTP GET method. It is not required to support HEAD, POST, or any other HTTP method, aside from GET.

Since your proxy is a middleman between client and end server, it will have elements of both. It will act as a server to the web browser, and as a client to the end server. Thus you will get experience with both client and server programming.

Logging

Your proxy should keep track of all requests in a log file named `proxy.log`. Each log file entry should be a line of the form:

```
Date: browserIP URL size
```

where `browserIP` is the IP address of the client, `URL` is the requested URL, `size` is the total size in bytes of the response that was returned, including the first line and headers. For instance:

```
Wed 30 Mar 2011 02:51:02 CDT: 128.2.111.38 http://www.rice.edu/ 34314
```

In effect, if your proxy does not support persistent connections, then `size` is the number of bytes that your proxy receives from the end server between the opening and closing of the connection.

Only requests that are met by a response from an end server should be logged. We have provided the function `format_log_entry` in `proxy.c` to create a log entry in the required format. Note that `format_log_entry` does not put a trailing newline on the returned string.

Port Numbers

Your proxy should listen for its connection requests on the port number passed in on the command line:

```
unix> ./proxy 18181
```

You may use any port number p , where $18000 \leq p \leq 18200$, and where p is not currently being used by any other user services (including other students' proxies).

Note that the port number restrictions are imposed by CLEAR. In general, you would be able to run a proxy on a wider range of ports.

Part II: Dealing with multiple requests concurrently

Real proxies do not process requests sequentially. They deal with multiple requests concurrently. Once you have a working sequential logging proxy, you should alter it to handle multiple requests concurrently. To do this, you should use a thread pool to handle each new connection request (CS:APP 12.5.5).

With this approach, it is possible for multiple peer threads to write to the log file concurrently. Thus, you will need to ensure that the log file entries are written atomically. Otherwise, the log file might become corrupted. For instance, one line in the file might begin in the middle of another.

Provided Files

The provided files are all available in

```
/clear/www/htdocs/comp221/assignments/proxy/
```

To begin working on the project, do the following:

- Copy all of the files in this directory to the directory in which you plan to do your work.
- Type your team member names and Rice Net IDs in the header comment at the top of `proxy.c`.

The `proxy.c` file contains a few support routines to help you parse URIs (`parse_uri`) and format log entries (`format_log_entry`).

The provided `Makefile` is set up to compile the `csapp.c` file along with your `proxy.c` file and to link in all of the appropriate libraries.

The echo client and server developed in CS:APP Chapters 11 and 12 is also available in the `echo` subdirectory as an example of a client/server application.

A compiled proxy (`proxyref`) is provided as a reference. The provided proxy runs concurrently, and outputs to log file `proxyref.log`. Your final proxy should behave as the reference proxy (remember that the required command line syntax for your proxy is `proxy <portnum>`). You are not required to match the debugging output of the reference proxy (which is also provided to help you develop your own proxy).

Hints

- A good way to get going on your proxy is to start with the basic echo server (CS:APP 11.4.9) and then gradually add functionality that turns the server into a proxy.
- Initially, you should debug your proxy using telnet as the client (CS:APP 11.5.3).
- Later, test your proxy with a real browser. Explore the browser settings until you find “proxies”, then enter the host and port where you’re running yours. You should only set the HTTP proxy, as that is all your code is going to be able to handle. Here are the instructions for a few different browsers:

In FireFox, choose “Firefox”, then “Options”, then “Advanced”, then “Network”, then “Settings”. Use the manual proxy configuration to set the HTTP proxy.

In Internet Explorer, choose the gear in the upper right, then “Internet Options”, then “Connections”, then “LAN Settings”. Check ‘Use a proxy server,’ and click “Advanced”. Just set the HTTP proxy.

In Chrome, choose the menu in the upper right, then “Show advanced settings...”, then “Change proxy settings...”, then follow the instructions for Internet Explorer starting with “LAN Settings”.

- When you test your proxy with a real browser, you should start with very simple web pages. We have provided a simple, text-only web page at:

`http://www.clear.rice.edu/comp221/test-cases/proxy/proxytest.html`

Once you get that working, we suggest you try the course web page and the links therein:

`http://www.clear.rice.edu/comp221/html/`

Then you can try more complicated sites with multiple files and images, such as `www.rice.edu` and `www.yahoo.com`. If you have trouble with a particular site, you may always use `proxyref` to see what features may be required to access that site.

- Be careful about memory leaks. When the processing for an HTTP request fails for any reason, the thread must close all open socket descriptors and free all memory resources before terminating.
- You will find it very useful to assign each thread a small unique integer ID (such as the current request number) and then pass this ID as one of the arguments to the thread routine. If you display this ID in each of your debugging output statements, then you can accurately track the activity of each thread.

- To avoid a potentially fatal memory leak, your threads should run as detached, not joinable (CS:APP 12.3.6).
- Since the log file is being written to by multiple threads, in order for the output from different threads to not be scrambled, all of the output from a single thread needs to be printed out atomically. Fortunately, the POSIX standard requires that individual stream operations, such as `fprintf(3)` and `fwrite(3)`, be atomic. However, without the use of a mutex to serialize access to the log file, even back-to-back stream operations by one thread may be interleaved with stream operations by another thread.
- To successfully implement a proxy using a thread pool, you must buffer the requests as they come in. In implementing this buffer, you should use mutex and condition variables, which are talked about in Lab 12.
- Be very careful about calling thread-unsafe functions inside a thread. In particular, the `open_clientfd` function in `csapp.c` is thread-unsafe because it calls `gethostbyname`, a Class-3 thread unsafe function (CS:APP 12.7.1). You will need to write a thread-safe version of `open_clientfd`, called `open_clientfdts` using the thread-safe functions `getaddrinfo` and `getnameinfo`.
- Use the RIO (Robust I/O) package (CS:APP 10.4) for all I/O on sockets. Do not use standard I/O on sockets. You will quickly run into problems if you do. However, standard I/O calls, such as `fopen` and `fwrite`, are fine for I/O on the log file.
- The `Rio_readn`, `Rio_readlineb`, and `Rio_writen` error checking wrappers in `csapp.c` are not appropriate for a realistic proxy because they terminate the process when they encounter an error. Instead, you should write new wrappers (`Rio_readnw`, `Rio_readlinebw`, and `Rio_writenw`) that simply return after printing a warning message when I/O fails (you may want to use `strerror` to print a useful message based on the value of `errno`). When either of the read wrappers detects an error, it should return 0, as though it encountered EOF on the socket.
- Reads and writes can fail for a variety of reasons. The most common read failure is an `errno == ECONNRESET` error caused by reading from a connection that has already been closed by the peer on the other end, typically an overloaded end server. The most common write failure is an `errno == EPIPE` error caused by writing to a connection that has been closed by its peer on the other end. This can occur, for example, when a user hits their browser's "Stop" button during a long transfer.
- Writing to a connection that has been closed by the peer elicits an error with `errno` set to `EPIPE` the first time. Writing to such a connection a second time elicits a `SIGPIPE` signal whose default action is to terminate the process. One relatively easy way to keep your proxy from crashing is to use the `SIG_IGN` argument to the `Signal` function (CS:APP 8.5.3) to explicitly ignore these `SIGPIPE` signals (or you can catch them and print an appropriate warning message).
- Modern web browsers and servers support persistent connections, which allows multiple requests to use the same connection. Your proxy will not do so. However, your browser is likely to set the headers `Proxy-Connection` or `Connection` to indicate that it would like to use persistent connections. If you pass these headers on to the end server, it will assume that you can support them. If you do not support persistent connections, then subsequent requests on that connection will fail, so some or all of the web page will not load in your browser. Therefore, you should strip the `Connection` and `Proxy-Connection` headers out of all requests, if they are present. Furthermore, HTTP/1.1

requires a `Connection: close` header be sent if you want the connection to close. Note that you must leave the other headers intact as many browsers and servers make use of them and will not work correctly without them.

- Your proxy should be able to accept both HTTP/1.0 and HTTP/1.1 requests. The server will make certain assumptions about the client if the browser uses HTTP/1.1. You should ensure that your implementation is compatible with these assumptions.

Evaluation

To turnin your code, you should use the turnin process on CLEAR (<https://docs.rice.edu/confluence/x/qAiUAQ>). First, create a proxy directory within your comp221 directory. Then, be sure that your code is entirely contained in a file named `proxy.c` inside this proxy directory. (Note that we will compile only your `proxy.c` file along with the provided `csapp.c`.) Finally, submit your solution by running the following command from inside your comp221 directory:

```
UNIX% turnin comp221-S14:proxy
```

After which, you should receive an e-mail confirming the submission of your solution.

Instead of a writeup, each individual/group will be evaluated on the basis of a demonstration to either the instructor or a TA. During the demo, you will demonstrate that your proxy behaves correctly and you will answer basic questions about your proxy. These demos will occur during the last day of classes and will be scheduled near the end of the term.

The project will be graded as follows:

- *Basic proxy functionality (50 points).*

Your proxy should correctly accept connections, forward the requests to the end server, and pass the response back to the browser, making a log entry for each request. As a starting point, your program should be able to proxy browser requests to the following Web sites and correctly log the requests:

- `http://www.clear.rice.edu/comp221/test-cases/proxy/proxytest.html`
- `http://www.rice.edu`
- `http://www.yahoo.com`
- `http://www.aol.com`
- `http://www.nfl.com`

Additional web sites will be tested at the demo.

- *Handling concurrent requests (20 points).*

Your proxy should be able to handle multiple concurrent connections. The following test is one way to determine this: (1) Open a connection to your proxy using `telnet`, and then leave it open without typing in any data. (2) Use a Web browser (pointed at your proxy) to request content from some end server.

- *Style (20 points).*

This includes general program style, the thoroughness of your comments, and whether you check the return code from all system calls. Furthermore, your threads should run detached, your code should not have any memory leaks, and your proxy should be thread-safe, protecting all updates of the log file and protecting calls to any thread unsafe functions, such as `gethostbyaddr`.

- *Demonstration (10 points).*

At the demonstration you will be asked basic questions about your proxy that would normally have been in your writeup on previous projects (i.e. descriptions of your design and testing strategy). This is your opportunity to explain any interesting design features and show off your proxy.