

Project 3: Multi-Threaded Web Server

CSCi 4061: Introduction to Operating Systems - Weissman – Fall 2022

To be completed by your group of size 3

Posted	Nov. 7, 05:00 PM
Intermediate Submission Due	Nov. 16, 11:59 PM
Final Submission Due	Nov. 23, 11:59 PM

1. Overview

The purpose of this lab is to construct a multi-threaded web server using POSIX threads (pthreads) in the C language to learn about thread programming and synchronization methods. Your web server should be able to handle these file types: HTML, GIF, JPEG, TXT, and of any arbitrary size. It should handle a limited portion of the HTTP web protocol (namely, the GET command to fetch a web page / file). We will provide pieces of code (some already compiled into object files, and some source code) that will help you complete the lab. You will also gain experience with `malloc/free`.

2. Description

Your server will be composed of two types of threads: **dispatcher threads** and **worker threads**. The purpose of the dispatcher threads is to repeatedly accept an incoming connection, read the client request (i.e. URL) from the connection, and place the request in a queue. We will assume that there will only be one request per incoming connection. The purpose of the worker threads is to monitor the request queue, retrieve requests and serve the request's result back to the client. The request queue is a bounded buffer and will need to be properly synchronized. Unlike Project 2, the URL will be dynamic length. The associated content (when read into memory) is NOT a string, just bytes. So, `memcpy` may be handy.

Since the server will continuously wait for client requests, you will need to use `^C` to terminate. Note: unfortunately, this will be a messy termination without cleanup of memory (valgrind may complain) and cleanup of child threads, but that is ok (for now) since the process will be terminated. Doing so when the server is idle may be best. You will use the following functions which have been precompiled into object files that we have provided (**full documentation of these functions has been provided in util.h**). More will be said below about how to use these functions. You can assume that these are thread safe:

```
void init (int port);
int accept_connection(void);
int get_request(int fd, char *filename);
int return_result(int fd, char *content_type, char *buf, int numbytes);
int return_error(int fd, char *buf);
```

3. Thread pool

Your server should create a fixed pool of dispatcher and worker threads when the server starts. The dispatcher thread pool size should be `num_dispatch` and the worker thread pool size should be `num_workers`.

4. Incoming Requests

An HTTP request has the form: **GET /dir1/subdir1/.../target_file HTTP/1.1** where `/dir1/subdir1/.../` is assumed to be located under your web tree root location. Our `get_request()` automatically parses this for you and gives you the `/dir1/subdir1/.../target` file portion and stores it in the `filename` parameter. Your web tree will be rooted at a specific location, specified by one of the arguments to your server (See section 8 for details). For example, if your web tree is rooted at `/home/user/joe/html`, then a request for `/index.html` would map to `/home/user/joe/html/index.html`. You can `chdir()` into the Web root to retrieve files using relative paths.

5. Returning Results

You will use `return_result()` to send the webpage data back to the web browser from the worker threads provided the file was opened and read correctly, or the data was found in the cache (see section 6). If there was any problem with accessing the file, then you should use `return_error()` instead. Our code will automatically append HTTP headers around the data before sending it back to the browser/user. Part of returning the result is sending back a special parameter to the browser: the content-type of the data. The file extension determines the content-type:-

- Files ending in `.html` or `.htm` are content-type “text/html”
- Files ending in `.jpg` are content-type “image/jpeg”
- Files ending in `.gif` are content-type “image/gif”
- Any other files may be considered, by default, to be of content-type “text/plain”.

6. Caching

To improve runtime performance, you will implement caching which stores web page content in memory for faster access. When a worker serves a request, it will look up the request in the cache first. If the request is in the cache (Cache HIT), it will get the result from the cache and return it to the user. If the request is not in the cache (Cache MISS), it will get the result from disk as usual, put the entry in the cache and then return result to the user. The cache size is set by a command-line argument (# of entries) and you need to log information about the cache (HIT or MISS) (see section 7 for more details). How to implement caching is up to you. You can implement any cache replacement policy like random, FIFO, LRU or LFU policy to choose an entry to evict when the cache is full. You must explain the policy you implemented in the README file. **Your cache contains a fixed number of entries each of which point to a dynamically created chunk of memory large enough to store the current file contents.** Remember to `free` the entry when evicted.

7. Request Logging

From the worker threads, you must carefully log each request (normal or error-related) to a file called “web_server_log” and also to the terminal (stdout) in the format below. The logfile should be created in the same directory where the final executable “web_server” exists. You must also protect the log file from races. The format is:

[threadId][reqNum][fd][Request string][bytes/error][Cache HIT/MISS]

- **threadId** is an integer from 0 to num_workers -1 indicating the thread index of request handling worker. (Note: this is **not** the pthread_t returned by pthread_create).
- **reqNum** is total number of requests a specific worker thread has handled so far, including the current request.
- **fd** is the file descriptor given to you by accept_connection() for this request
- **Request string** is the filename returned by get_request
- **bytes/error** is either the number of bytes returned by a successful request, or -1 if an error occurred
- **Cache HIT/MISS** is either “TRUE” or “FALSE” depending on whether this specific request was found in the cache or not.

The log (in the “web_server_log” file and in the terminal) should look something like the example below. We provide the code for this.

```
[8][1][5][image/jpg/30.jpg][17772][ FALSE]
```

```
[9][1][5][image/jpg/30.jpg][17772][ TRUE]
```

```
[2][1][5][image/jpg/282.jpg][Requested file not found.][ FALSE]
```

8. How to run the server

Your server should be run as:

```
% ./web_server port path num_dispatch num_workers queue_length cache_length
```

The server will be configurable in the following ways:

- **port** number can be specified (you may only use ports 1025 - 65535 by default). When you send wget request (Section 11) to the server, the port number should be the same with this.
- **path** is the path to your web root location from where all the files will be served
- **num_dispatcher** is how many dispatcher threads to start up
- **num_workers** is how many worker threads to start up
- **queue_length** is the fixed, bounded length of the request queue
- **cache_length** is the number of entries available in the cache.

9. Provided Files and How to Use Them

We have provided many functions which you must use in order to complete this assignment. We have handled all of the networking system calls for you. We have also handled the HTTP protocol parsing for you. Some of the library function calls assumes that the program has “chdir”ed to the web tree root directory. You need to make sure that you chdir to the web tree root somewhere in the beginning.

We have provided a makefile you can use to compile your server. Here is a list of the files we have provided.

1. `server.c`: You only need to modify this file to implement a server.
2. `server.h`: You can modify this file if you need to, but we do not suggest it.
3. `makefile`: You can use this to compile your program using our object files, or you can make your own. You can study this to see how it compiles our object code along with your server code to produce the correct binary executables.
4. `util.h`: This contains a very detailed documentation of each function that you must study and understand before using the functions.
5. `util.o`: This is the compiled code of the functions described in `util.h` to be used for the web server. Link `util.o` into your multi-threaded server code and it will produce a fully-functioning web server.
6. `testing.zip`: This file includes images, texts and url files to test your server. See section 10 for more detailed information.
7. `web_server_sol`: This is a solution webserver executable. You can test it by sending `wget` requests to this server. Note that, therefore, the logs might be different from yours because of different algorithms used for caching. Use it after running “`chmod +x web_server_sol`”.

10. Makefile Tips:

- `make clean` Will remove all executables and logs to clean up directory
- `make` Will compile solution into an executable called `web_server`
- `make test` Will request you enter a random port and will start `web_server` **without** a cache or dynamic thread pool
- `make test_full` Will request you enter a random port and will start `web_server` **with** a cache or dynamic thread pool
- `make solution` Will run the solution `web_server` so you can see what the results should look like
- `make force_kill` Having trouble with killing your process after using a signal handler? Use “CTRL Z” and then call `make force_kill` to **try** and kill the process
- `make submission` Will tar up the necessary files into a tar file with your group number

11. How to test your server

We will test your server with “wget”, the non-interactive webpage downloader. After you run the server, open a new terminal to test the server. You can try to download a file using this command (where 9000 is the port specified to the server):

```
-> wget http://127.0.0.1:9000/image/jpg/29.jpg
```

Please note that 127.0.0.1 means localhost, the port number used in wget request should match the port number specified to the web_server.

This command will try to download the file at root/image/jpg/29.jpg. If it failed to download the file for some reason, it will show an error message. You can also test your server with any web browser (Internet Explorer, Chrome, Firefox, and so on). We will provide “testing.zip” to make testing the server easier. Once you extract it, you can find an instruction file “how_to_test” which explains how to use and test your program with these files. The server should be able to serve requests from other machines, so you need to run the server on a machine and run “wget” command on another machine to test your server (use the server's IP address instead of 127.0.0.1, run ‘ip addr’ – look for inet). Note, if you try to connect to the server which runs on CSELabs machine with your own machine (laptop), it may not be able to connect to the server because CSELabs machines are protected by a firewall. If you want to test it using different two machines, you can use two CSELabs machines; one for the server and another one for client.

Testing Concurrency: Bash has a nifty command called xargs, which allows a set of arguments to be piped to a command such that multiple executions will be run concurrently.

The format of the command is `xargs -n num_args -P num_procs cmd`

So, for our purposes, the command “`cat urls | xargs -n 1 -P 8 wget`” will run wget 8 times simultaneously (-P 8) with 1 argument each time (-n 1) from the pipe produced by cat urls.

NOTE: If -P is given an argument that is smaller than the number of args in the pipe, then multiple sets of size P will be run, with each set being concurrent.

Ex. If I had 10 arguments in the pipe, and I set -P 5, 2 sets of 5 concurrent processes will run. This can be used to test that your server really does permit concurrent execution of multiple requests.

12. Simplifying Assumptions

- The maximum number of dispatcher threads will be 100.
- The maximum number of worker threads will be 100.
- The maximum length of the request queue will be 100 requests.
- The maximum size of the cache will be 100 entries.
- The maximum length of filename will be 1024.
- Any HTTP request for a filename containing two consecutive periods or two consecutive slashes (“..” or “//”) will automatically be detected as a bad request by our compiled code for security.

13. Documentation

Within your code you should use one or two sentences to describe each function that you write. You do not need to comment every line of your code. However, you might want to comment portions of your code to increase readability.

14. Deliverables

A tar file containing your code and a Makefile that can be used to build and execute your code, and a README file. The README indicates how the work was partitioned within your group. There could be overlap. We will check this to see that everyone participated. You submit ONLY ONE tar file per group. You will be docked points if you do not follow these instructions.

At the top of your README file, please include the following comment:

```
/* CSCI-4061 Fall 2022 – Project #3
 * Group Member #1: <name> <x500>
 * Group Member #2: <name> <x500>
 * Group Member #3: <name> <x500> */
*/
```

15. Intermediate Submission

Your group will be responsible for implementing enough of the project such that a single request can be handled by a single dispatcher thread (no request queue needed). The request information should be printed to the terminal in the format:

```
fprintf(stderr, "Dispatcher Received Request: fd[%d] request[%s]\n", <insert_fd>, <insert_str>);
```

16. Grading (Tentative)

5% README

10% Documentation within code, coding, and style: indentations, readability of code, use of defined constants rather than numbers

10% Interim evaluation

70% Correctness, error handling, meeting the specification. Broken down as follows:

- 15% for handling one request successfully.
- 5% error checking.
- 20% for handling multiple requests.
- 20% for handling multiple concurrent requests (test with xargs).
- 10% for handling caching.
- 5% for a correct log file.