

# CSDS 325/425: Computer Networks

## Project #2

Due: October 2, 11:59 PM

The second project of the semester involves writing a simple command line-based web client. The aim of this project is (i) to get your feet wet with writing C/C++, (ii) to write a program that exchanges information with another computer over a network and (iii) to start concretely thinking about protocols.

### Overview

The usage of your program—which will be called *proj2*—is as follows:

```
./proj2 -u URL [-d] [-q] [-r] -o filename
```

Specifically:

- The “-u” option specifies the URL your web client will access. The “-u” option *must* be present on the command line.
- The “-d”, “-q” and “-r” arguments are all optional and any number and combination of these could be given. These options all trigger output that is described below.
- The “-o” option specifies a filename where the downloaded contents of the supplied URL will be written. The “-o” option must be present on the command line.
- The command line arguments may appear in any order.
- Unknown command line arguments must trigger meaningful error messages and termination of the program.

### -u option

The “-u” option is used to supply the web server and page the client will access. The URL format your program will be expected to deal with is:

```
http://hostname[/path/to/file]
```

- The path in brackets is optional and may or may not be in the URLs your program must accept. Note: There will be no brackets in the actual options given to your program.
- Every URL must begin with “http://”. While in the general case alternate protocols can be encoded in URLs—e.g., “https://” or “ftp://”—your program will only support HTTP. Further, the “http” is not case sensitive. I.e., “http”, “Http”, “HTTP”, etc. must all be accepted.
- Following the “http://” will be a hostname. The hostname portion of the URL ends with the first “/”. If a “/” does not appear after the hostname begins, the hostname continues to the end of the URL.
- Anything after the hostname is the path and filename to be sent verbatim to the web server (including the leading “/” character). If the URL does not contain a filename you must use the default filename of “/”.

**Hint:** Strings in C and C++ are tedious. A set of standard string processing routines helps (a little!). Use the manual pages to look for routines such as *strncasecmp()*, *strstr()*, *strncmp()*, *strtok()*, etc.

## **-d option**

The “-d” option will be used to print debugging information about the given command line parameters to standard output (i.e., the screen). When “-d” is given on the command line your program will output the following lines:

```
DBG: host: [hostname]
DBG: web_file: [url_filename]
DBG: output_file: [local_filename]
```

The format for these lines must follow these requirements:

- The “DBG:” must be at the beginning of the line.
- A single space follows “DBG:”.
- The labels that appear after “DBG:<space>” above must be exactly as they appear above (e.g., using all lower case letters).
- After the label, add a colon (“:”) and then a single space before printing the value.
- The “[hostname]” value comes from the URL given on the command line, as described in the “-u” discussion above.
- The “[url\_filename]” value is the filename portion of the URL given on the command line, as described in the “-u” discussion above. If no filename is given on the command line, the default filename of “/” must be printed.
- The “[local\_filename]” value is the name of the file on the local system where the web page at the given URL will be stored. This is the filename given with the “-o” option on the command line.
- The three lines must appear in the order given above.
- Do not print extra lines—including blank lines.
- The “-d” output is to be printed regardless of whether there are errors fetching the web page. The “-d” output will not be printed if there are errors in the command line options given by the user.

The following are several illustrative examples with the “-d” option:

```
./proj2 -d -u http://www.icir.org -o testing.html
DBG: host: www.icir.org
DBG: web_file: /
DBG: output_file: testing.html

./proj2 -o mallman.html -d -u http://www.icir.org/mallman/
DBG: host: www.icir.org
DBG: web_file: /mallman/
DBG: output_file: mallman.html

./proj2 -u http://www.icir.org/mallman/index.html -o /tmp/mallman.html -d
DBG: host: www.icir.org
DBG: web_file: /mallman/index.html
DBG: output_file: /tmp/mallman.html
```

## **-q option**

When the “-q” option is present on the command line, your program must print the HTTP request sent by your web client to the web server to standard output (the screen). The HTTP request you will transmit to the web server will look like this:

```
GET [url_filename] HTTP/1.0\r\n
Host: [hostname]\r\n
User-Agent: CWRU CSDS 325 SimpleClient 1.0\r\n
\r\n
```

Notes:

- The HTTP “GET” method must be used and expressed in all capital letters.
- The “[url\_filename]” and “[hostname]” values are taken from the URL furnished on the command line and explained in the “-u” discussion above.
- “HTTP/1.0” must be the specified version of the HTTP protocol used.
- A single space separates “GET” and the [url\_filename].
- A single space separates the [url\_filename] and “HTTP/1.0”.
- A single space follows “Host:”.
- The “User-Agent” line above must be used verbatim. A single space appears between each word.
- All lines must end with a carriage return (\r) and a newline (\n).
- A line with only a carriage return (\r) and newline (\n) ends the HTTP request (per the HTTP specification).

The output corresponding to the above HTTP request when “-q” is specified will look like this:

```
OUT: GET [url_filename] HTTP/1.0
OUT: Host: [hostname]
OUT: User-Agent: CWRU CSDS 325 SimpleClient 1.0
```

In other words, the output should exactly mirror what was sent to the web server, with two exceptions:

- When printing to the screen, each HTTP request line must begin with “OUT:” followed by a single space.
- The blank line that terminates the HTTP request must be excluded from the “-q” output.

Examples:

```
./proj2 -q -u http://sigcomm.org -o sigcomm.html
OUT: GET / HTTP/1.0
OUT: Host: sigcomm.org
OUT: User-Agent: CWRU CSDS 325 SimpleClient 1.0

./proj2 -u http://www.icir.org/mallman/ -q -o mallman.html
OUT: GET /mallman/ HTTP/1.0
OUT: Host: www.icir.org
OUT: User-Agent: CWRU CSDS 325 SimpleClient 1.0
```

## **-r option**

When the “-r” option is present on the command line, your program must print the HTTP response header received from the web server to standard output (the screen). Each line must be printed exactly as the line was received with the exception that “INC:” followed by a single space must begin each line. Examples:

```
./proj2 -u http://www.icir.org/mallman/ -o mallman.html -r
INC: HTTP/1.1 200 OK
INC: Date: Fri, 08 Sep 2023 18:38:34 GMT
INC: Server: Apache
INC: Accept-Ranges: bytes
INC: Content-Length: 6569
INC: Connection: close
INC: Content-Type: text/html; charset=UTF-8
```

```
./proj2 -r -u http://case.edu/ -o case.html
INC: HTTP/1.1 301 Moved Permanently
INC: Date: Fri, 08 Sep 2023 18:42:06 GMT
INC: Server: Apache
INC: Location: https://case.edu/
INC: Content-Length: 225
INC: Connection: close
INC: Content-Type: text/html; charset=iso-8859-1
```

The output must not include the blank line that terminates the HTTP response header (and, hence, separates the header from the content).

## **Printing Order**

Note: any number, order and combination of “-d”, “-q” and “-r” is allowed. This includes no printing options. When multiple options are given, the output will be printed in a standard order. First the “-d” output will be given, if needed (the “DBG:” lines). Second, the “-q” output will be given, if needed (the “OUT:” lines). Finally, the “-r” output will be given, if needed (the “INC:” lines). This order holds regardless of the order the options are provided on the command line. Examples:

```
./proj2 -u http://www.icir.org/mallman/ -o mallman.html -q -r -d
DBG: host: www.icir.org
DBG: web_file: /mallman/
DBG: output_file: mallman.html
OUT: GET /mallman/ HTTP/1.0
OUT: Host: www.icir.org
OUT: User-Agent: CWRU CSDS 325 SimpleClient 1.0
INC: HTTP/1.1 200 OK
INC: Date: Fri, 08 Sep 2023 18:43:13 GMT
INC: Server: Apache
INC: Accept-Ranges: bytes
INC: Content-Length: 6569
INC: Connection: close
INC: Content-Type: text/html; charset=UTF-8
```

```
./proj2 -q -u http://www.icir.org/mallman/ -o mallman.html -d
DBG: host: www.icir.org
DBG: web_file: /mallman/
DBG: output_file: mallman.html
OUT: GET /mallman/ HTTP/1.0
OUT: Host: www.icir.org
OUT: User-Agent: CWRU CSDS 325 SimpleClient 1.0
```

```
./proj2 -r -q -u http://www.icir.org/mallman/ -o mallman.html
OUT: GET /mallman/ HTTP/1.0
OUT: Host: www.icir.org
OUT: User-Agent: CWRU CSDS 325 SimpleClient 1.0
INC: HTTP/1.1 200 OK
INC: Date: Fri, 08 Sep 2023 18:44:14 GMT
INC: Server: Apache
INC: Accept-Ranges: bytes
INC: Content-Length: 6569
INC: Connection: close
INC: Content-Type: text/html; charset=UTF-8
```

```
./proj2 -u http://www.icir.org/mallman/ -o mallman.html
```

## **-o option**

The “-o” command line argument is used to tell the client where to save the contents of the downloaded URL. The web page content is everything received from the web server that follows the HTTP response header and the blank line that appears after the header. The file must contain exactly the data received from the web server. The client should only create the target file when the server returns an “OK” code of 200 in the HTTP header. When the client receives a non-200 response, it must print a meaningful error message. There are test URLs on the class web page, but you can test with arbitrary web servers, as well. The *wget* tool is available on the class servers and can help in your testing, as follows.

```
wget -O wget_mallman.html http://www.icir.org/mallman/
[...]
./proj2 -o proj2_mallman.html -u http://www.icir.org/mallman/
ls -l *.html
-rw----- 1 mallman staff 6569 Sep 11 16:49 proj2_mallman.html
-rw----- 1 mallman staff 6569 Sep 11 16:49 wget_mallman.html
diff proj2_mallman.html wget_mallman.html
sha1sum proj2_mallman.html wget_mallman.html
56ddd089280c8491d832ddd40de5bfd69ade5b85 proj2_mallman.html
56ddd089280c8491d832ddd40de5bfd69ade5b85 wget_mallman.html

./proj2 -r -o testing.html -u http://www.icir.org/mallman/doesnt-exist
INC: HTTP/1.1 404 Not Found
INC: Date: Sun, 11 Sep 2022 20:50:26 GMT
INC: Server: Apache
INC: Content-Length: 258
INC: Content-Type: text/html; charset=iso-8859-1
ERROR: non-200 response code
ls -l testing.html
ls: testing.html: No such file or directory
```

## Final Bits

1. Submission specifications:
  - (a) All project files must be submitted to Canvas in a gzip-ed tar file called “[CaseID]-proj2.tar.gz” (without the brackets).
  - (b) Your submission must contain all code and a Makefile that by default produces an executable called “proj2” (i.e., when typing “make”).
  - (c) Do not include executables or object files in the tarball.
  - (d) Do not include sample input or output files in your tarball.
  - (e) Do not include multiple versions of your program in your submission.
  - (f) Do not include directories in your tarball.
  - (g) Do not use spaces in file names.
  - (h) Every source file must contain a header comment that includes (i) your name, (ii) your Case network ID, (iii) the filename, (iv) the date created and (v) a brief description of the code contained in the file.
2. Your project must be written using sockets-based code and implement the required parts of HTTP. You may not leverage a third-party library for these tasks. Projects that rely on extra libraries will be returned ungraded.
3. Your submission may include a “notes.txt” file for any information you wish to convey during the grading process. We will review the contents of this file, but not of arbitrary files in your tarball (e.g., “readme.txt”).
4. There will be sample reference output on the class web page by the end of the day on September 11. (Not all sample input will have (all) corresponding reference output.)
5. Print only what is described above. Extra debugging information must not be included. Adding an extra option (e.g., “-v” for verbose mode) to dump debugging information is always fine.
6. Do not make assumptions about the size of the URL contents. Your project should handle arbitrary size downloads.
7. If you encounter or envision a situation not well described in this assignment, please Do Something Reasonable in your code and include an explanation in the “notes.txt” file in your submission. If you’d like to ensure you’re on the right track, please feel free to discuss these situations with me.
8. Hints / tips:
  - (a) Needlessly reserving large amounts of memory in case it may be needed (e.g., for a large content download) is unreasonable.
  - (b) Errors will be thrown at your project.
  - (c) You may leverage the example sockets code we reviewed in class. This code is available from the class web page.
  - (d) A simple C program and Makefile are available on the class web page. The program illustrates the use of *getopt()* to parse the command line arguments. The Makefile can be easily adapted for this project.
9. *WHEN STUCK, ASK QUESTIONS!*

# CSDS 425: Computer Networks

## Project #2 Extensions

Due: October 3, 11:59 PM

Graduate students will be responsible for writing the basic web client described above, as well as two extensions described below. Undergraduates can do some or all of the following for extra credit.

### Following Redirects

When the “-f” option is given on the command line the client will follow redirections from the web server. For instance, consider this sample invocation of the web client:

```
./proj2 -r -o testing.html -u http://www.icir.org/mark/
INC: HTTP/1.1 301 Moved Permanently
INC: Date: Fri, 08 Sep 2023 19:04:17 GMT
INC: Server: Apache
INC: Location: http://www.icir.org/mallman/
INC: Content-Length: 298
INC: Connection: close
INC: Content-Type: text/html; charset=iso-8859-1
ERROR: non-200 response code
```

In this case, the desired web page is not at the URL given on the command line (“http://www.icir.org/mark/”). Rather, the web server uses a response with a code of 301 to redirect the client to a different URL—which is given in the “Location:” line of the response header. When the “-f” option is given, the web client will download the URL given in the redirection message. Redirection can happen more than once. E.g., “www.foo.com” could redirect to “www.bar.com” which could in turn redirect to “bar.com”. When redirects happen and the “-q” or “-r” options are given the client should print every request and/or response header encountered in the order it is encountered. The output file given with “-o” will contain the contents of the ultimate (last) response. An example:

```
./proj2 -q -r -u http://www.icir.org/mark/ -o mark.html -f
OUT: GET /mark/ HTTP/1.0
OUT: Host: www.icir.org
OUT: User-Agent: CWRU CSDS 325 Client 1.0
INC: HTTP/1.1 301 Moved Permanently
INC: Date: Fri, 08 Sep 2023 19:05:05 GMT
INC: Server: Apache
INC: Location: http://www.icir.org/mallman/
INC: Content-Length: 298
INC: Connection: close
INC: Content-Type: text/html; charset=iso-8859-1
OUT: GET /mallman/ HTTP/1.0
OUT: Host: www.icir.org
OUT: User-Agent: CWRU CSDS 325 Client 1.0
INC: HTTP/1.1 200 OK
INC: Date: Fri, 08 Sep 2023 19:05:42 GMT
INC: Server: Apache
INC: Accept-Ranges: bytes
INC: Content-Length: 6569
INC: Connection: close
INC: Content-Type: text/html; charset=UTF-8
```

```
wget -O wget_mark.html http://www.icir.org/mark/
[...]
ls -l mark.html wget_mark.html
-rw----- 1 mallman staff 6569 Sep 11 17:09 mark.html
-rw----- 1 mallman staff 6569 Sep 11 17:09 wget_mark.html
sha1sum mark.html wget_mark.html
56ddd089280c8491d832ddd40de5bfd69ade5b85 mark.html
56ddd089280c8491d832ddd40de5bfd69ade5b85 wget_mark.html
```

## Chunked Encoding

One of the mechanisms that was added in HTTP version 1.1, but not included in HTTP version 1.0 is “chunked encoding”. When the “-C” option is given on the command line the web client will correctly download and store HTTP responses that arrive using this encoding strategy. When using “-C”, the requests will advertise the client as “HTTP/1.1” instead of “HTTP/1.0”. (Note: the client will still use “HTTP/1.0” in the absence of the “-C” option.) Further, the URL’s content will not come as a stream of bytes, but as a series of chunks, which include non-content information about their size. You will need to research the details of chunked encoding before you complete this portion of the project. Example:

```
wget -O wget_sigcomm.html http://www.sigcomm.org/
[...]

./proj2 -u http://www.sigcomm.org/ -o 1.0-sigcomm.html -q
OUT: GET / HTTP/1.0
OUT: Host: www.sigcomm.org
OUT: User-Agent: CWRU CSDS 325 SimpleClient 1.0

./proj2 -u http://www.sigcomm.org/ -o 1.1-sigcomm.html -q -C
OUT: GET / HTTP/1.1
OUT: Host: www.sigcomm.org
OUT: User-Agent: CWRU CSDS 325 SimpleClient 1.0

sha1sum *html
db688fc539883fca69fd7320f4a74c39b8261984 wget_sigcomm.html
db688fc539883fca69fd7320f4a74c39b8261984 1.0-sigcomm.html
db688fc539883fca69fd7320f4a74c39b8261984 1.1-sigcomm.html
```