

# Cuteserver: A Web Server Written in C

Project for the Operating Systems Lecture

Rahel Kempf

Ephraim Siegfried

June 11, 2024

# Introduction

In this project we aspire to understand the inner workings of web servers. It is of interest for us, as we both recently set up our own home servers and came in contact with different tools such as nginx. Building our own web server helped us be more competent in using those tools. It also helped us deepen our knowledge of OS topics such as socket programming, thread/process creation & management and inter-process communication.

This project resulted in Cuteserver, a simple web server written in C. It can handle various HTTP/1.1 requests from multiple clients concurrently. It is configurable and supports hosting content for multiple domains and is thus similar to a reverse proxy. It can be hosted easily with a prepared docker image. In the following chapters, we will discuss the steps taken and the challenges we faced in this project.

## Background

In this section we explain some concepts we had to understand implementing our project.

### HTTP Protocol

The Hypertext Transfer Protocol (HTTP) is an application layer protocol which defines how data can be exchanged over the Internet. It is generally sent over TCP, as it relies on a reliable transport protocol. As a client-server protocol, requests are initiated by the recipient (browser) and responses are served by the provider (web server). The structure of these requests and responses is defined by HTTP.<sup>1</sup> For our project, we mainly focused on version HTTP 1.1.

The general structure of a request is the following: The first line denotes the **Method, Path and Version of the Protocol**. The following lines are called **Headers** and contain information for the servers. An empty line signals the end of the headers, for some request types (such as POST) the **Request Body** follows. A header to point out would be the **Connection** header, which was introduced in version 1.1 to provide persistent connections by reusing a TCP connection for multiple request/responses instead of opening a new connection for each response. The response follows the same structure, just that the first line contains the Version of the Protocol First, then the Status Code and Status Message<sup>2</sup>

There are several Request Methods which signal what kind of action should be performed.<sup>3</sup> For example: **GET** retrieve state representation of target resource. Parameters sent through URL

**HEAD** retrieve only metadata.

**POST** request that target resource processes the representation in the request body.

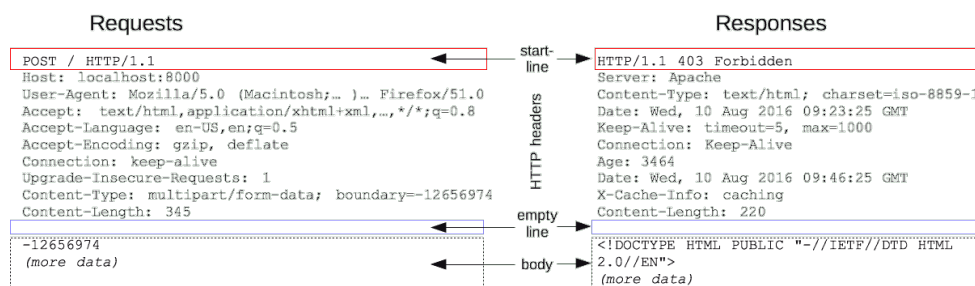


Figure 1: HTTP Requests and Responses

Source: <https://developer.mozilla.org/en-US/docs/Web/HTTP/Messages>

<sup>1</sup>HTTP - Hypertext Transfer Protocol. <https://developer.mozilla.org/en-US/docs/Web/HTTP/Overview>

<sup>2</sup>HTTP Response Status Codes. <https://developer.mozilla.org/en-US/docs/Web/HTTP/Status>

<sup>3</sup>RFC 1945, HTTP Specification. <https://datatracker.ietf.org/doc/html/rfc1945#section-8>

## Common Gateway Interface

The Common Gateway Interface (CGI) <sup>4</sup> is a technology we encountered during our research on handling HTTP POST requests. CGI is an interface that defines how the web server interacts with external programs. This enables the creation of dynamic content, such as generating user-specific webpages or processing form submissions, and facilitates communication between the web server and application backend services.

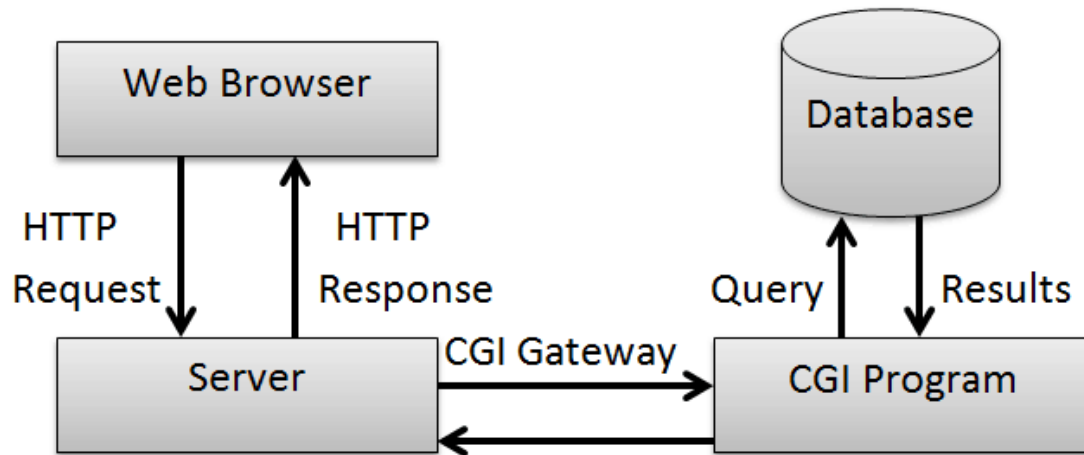


Figure 2: The Common Gateway Interface

Source: <https://www.educba.com/what-is-common-gateway-interface/>

When a request with a URL denoting a CGI script, like `example.com/backend.cgi`, enters the web server, the server locates the executable file called `backend.cgi`. It then spawns a process to execute the CGI script and communicates with this process via piped standard input/output and environment variables. After the CGI script sends its response to the web server, the server may verify the response and forward it to the client.

Although CGI is now quite outdated and less frequently used due to the inefficiency of spawning a new process for each request, enhanced versions such as FastCGI or SCGI are widely used. These versions improve performance by spawning CGI scripts only once. In our project we implemented the original CGI specification from 1997.

## Implementation

### Development Environment

The development of Cuteserver was primarily accomplished through pair programming. We held regular meetings where we shared our terminals using a shell sharing tool called **sshx** <sup>5</sup>. This tool was particularly effective for collaboration since we both used **Neovim** as our code editor.

Initially, we used **Makefiles** to build the source code, but this approach was inefficient due to frequent adjustments. As the project grew more complex, we decided to use **CMake** for more efficient project building and file generation.

Our implementation strategy for the web server was iterative, starting with small, manageable tasks and progressively tackling more complex ones. We regularly tested the code by running it with different inputs to ensure functionality.

<sup>4</sup>RFC 3875. CGI Specification. <https://datatracker.ietf.org/doc/html/rfc3875>

<sup>5</sup>SSHx. <https://sshx.io/>

## Project Structure

We tackled the implementation of our web server by structuring the code into 6 core modules. In the main module the server socket is set up and incoming TCP connections are handled through a threadpool. The parser module deserializes the request and puts the headers into a hashmap to simplify further processing. The request module contains the logic for handling static and CGI requests and set appropriate response headers. In the CGI handler module CGI scripts are spawned and their output is read. The response module serializes the response and sends it. Additionally, the config module handles TOML configurations set by the user.

## Static File Requests

When a static request enters the server, the requested file is located and read. We assemble metadata such as content length and type to form the response headers. If it's a GET-Request, we send the headers and the file, for HEAD-Requests only the headers are sent. If the file cannot be found or another error is encountered, we send an error response with the help of `http.cat` <sup>6</sup>.

## CGI

Dynamic Request are handled with CGI. This means the incoming request is parsed, then passed to the CGI handler which creates a new process running the corresponding cgi script. The CGI Protocol defines "meta-variables"<sup>7</sup> which have to be passed from server to CGI script using environment variables.

The CGI process is started using the **execve** function, allowing us to pass the environment variables as an argument. Additional data has to be passed to the child process through standard in/output.<sup>8</sup> Communication between cgi process and server (parent process) is enabled through pipes. A pair of pipe descriptors is created before the child process is forked. Then, the child's descriptors are redirected to standard input and standard output<sup>9</sup> (See Figure 3).

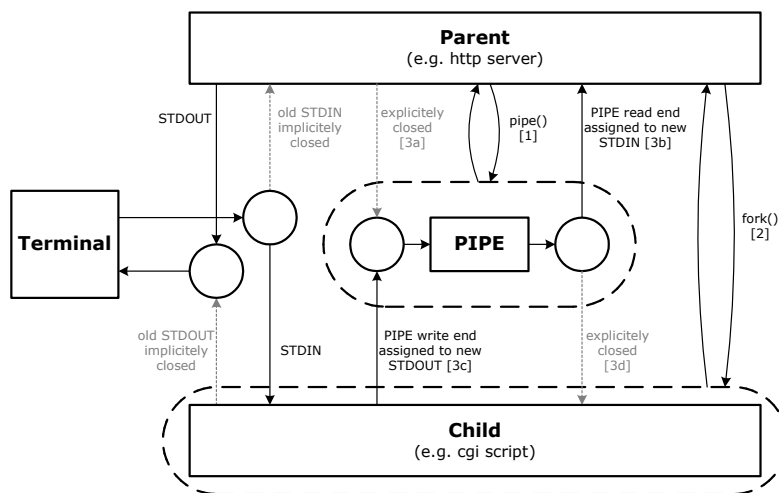


Figure 3: Pipe Redirection

Source: [http://www.fmc-modeling.org/category/projects/apache/amp/A\\_4\\_Pipes.html](http://www.fmc-modeling.org/category/projects/apache/amp/A_4_Pipes.html)

## Server Configuration

From the beginning, we aimed for our web server to support hosting multiple resources connected to different domain names. To achieve this multidomain support, we needed to make the web server configurable. To ensure user-friendly configuration, we used **TOML** files. The web server parses a configu-

<sup>6</sup><https://http.cat>

<sup>7</sup>RFC 3875. CGI Specification. Section 4.1 Request Meta-Variables. <https://datatracker.ietf.org/doc/html/rfc3875#section-4.1>

<sup>8</sup>Mapping UNIX pipe descriptors to stdin and stdout. <http://www.unixwiz.net/techtips/remap-pipe-fds.html>

<sup>9</sup>Pipes. [http://www.fmc-modeling.org/category/projects/apache/amp/A\\_4\\_Pipes.html](http://www.fmc-modeling.org/category/projects/apache/amp/A_4_Pipes.html)

ration file listing all resources and their locations for the server to serve. Additionally, users can specify resource-independent settings, such as the number of active worker threads and log storage locations. We also provide a command line interface where users can specify the TCP address and the path to the configuration file.

## Example Application

We had several web applications to test our web server. Initially, we used simple HTML files and gradually added more complex ones. Eventually, we developed a more sophisticated application using the React framework for the front end. This application is a simple chat app that makes POST requests when users submit chat messages. These POST requests are handled by our backend service, which stores the chat messages in a file and sends the content to clients. We initially wrote this service using the Flask framework in Python to test if interpreted languages could work as CGI scripts. While this worked, we later translated the backend Python script to a C script to achieve faster response times by running compiled executables.

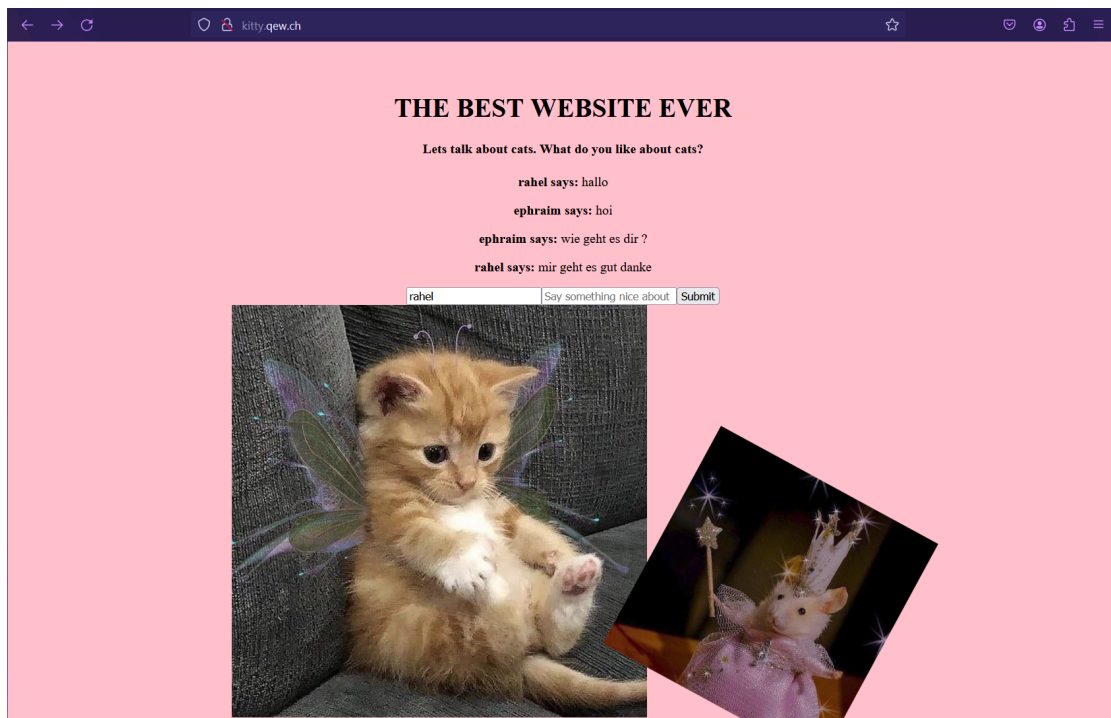


Figure 4: Our Example Application

## Containerization

We decided to containerize our application for two main reasons. Initially, we used chroot jails to sandbox our application, i.e., a mechanism that isolates a process and its children from the rest of the system by changing their apparent root directory to a specified path. This approach worked until we implemented CGI support, which required the web server to start new processes with the necessary libraries for the CGI scripts. Manually copying those dependencies into the chroot jail would have been too laborious. Docker containers offered an easier solution by securely sandboxing execution and including dependencies. Additionally, Dockerfiles enabled us to automate the building of the web server and the example application, making it easy to use and run. We used multistage builds, first creating containers for building and then copying only the necessary files to the final containers. Multistage build reduced the container size from 1.2 GB to approximately 200 MB. This process enhanced our understanding of containerization.

## Results

In Figure 4, we showcase our first successful keep-alive connection. Here, two threads each handle a series of requests over a single connection, rather than multiple separate connections. In Figure 5, we present our first running CGI script, which outputs dynamic data, such as the user's IP address, that cannot be generated by a static file.

```
cyberpumucki@pumucklaptop:~/Documents/code/webServer$ sudo ./cuteserver 1999
19:56:17 INFO src/main.c:99: Server listening on 127.0.0.1:1999
19:56:21 DEBUG src/main.c:116: Created sock: 4
19:56:21 INFO src/main.c:121: New connection accepted from 127.0.0.1:40928
19:56:21 DEBUG src/main.c:28: TID: 41716 sock: 4
19:56:21 DEBUG src/main.c:29:
19:56:21 INFO src/parser.c:33: GET / HTTP/1.1
19:56:21 DEBUG src/main.c:56: Client sent: keep-alive
19:56:21 INFO src/response.c:59: Sending /index.html over socket: 4
19:56:22 DEBUG src/main.c:116: Created sock: 5
19:56:22 INFO src/main.c:121: New connection accepted from 127.0.0.1:50792
19:56:22 DEBUG src/main.c:28: TID: 41717 sock: 5
19:56:22 DEBUG src/main.c:29:
19:56:22 INFO src/parser.c:33: GET /script.js HTTP/1.1
19:56:22 DEBUG src/main.c:56: Client sent: keep-alive
19:56:22 INFO src/response.c:59: Sending /script.js over socket: 5
19:56:22 DEBUG src/main.c:28: TID: 41716 sock: 4
19:56:22 DEBUG src/main.c:29:
19:56:22 INFO src/parser.c:33: GET /style.css HTTP/1.1
19:56:22 DEBUG src/main.c:56: Client sent: keep-alive
19:56:22 INFO src/response.c:59: Sending /style.css over socket: 4
19:56:23 DEBUG src/main.c:28: TID: 41717 sock: 5
19:56:23 DEBUG src/main.c:29:
19:56:23 INFO src/parser.c:33: GET /chaetzli.jpg HTTP/1.1
19:56:23 DEBUG src/main.c:56: Client sent: keep-alive
19:56:23 INFO src/response.c:59: Sending /chaetzli.jpg over socket: 5
19:56:23 DEBUG src/main.c:28: TID: 41716 sock: 4
19:56:23 DEBUG src/main.c:29:
19:56:24 DEBUG src/main.c:28: TID: 41717 sock: 5
19:56:24 DEBUG src/main.c:29:
19:56:24 INFO src/parser.c:33: GET /favicon.ico HTTP/1.1
19:56:24 DEBUG src/main.c:56: Client sent: keep-alive
```

Figure 5: Keep-Alive: Multiple Requests handled over one Connection

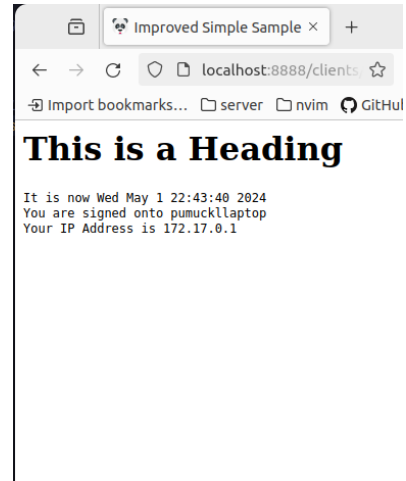


Figure 6: First simple CGI Script with Environment Variables

## Performance Testing

We wanted to test the influence of the number of worker threads on the performance of our application. For this we used **Locust**<sup>10</sup> to tests POST and GET requests to our application, as well as static file requests.

### Static File Tests

We conducted these tests using the following Locust Configuration: 500 virtual users, ramping up at 50 users/second, for a runtime of 60 seconds. Requests were made to a static file hosted on a home server with 4 cores, and multithreading was managed by OpenMP. The graphs for 2 versus 3 worker threads (Figures 7 and 8) exhibit similar trends, so the application runs relatively stable. The graph shows an initial peak in response time, followed by a gradual decline and stabilization at lower values. We can only explain this result by assuming that the responses are being cached somewhere. The number of requests is similar, but the average response time with 2 threads is 239.07ms, while with 3 threads, it is 154.19ms—a reduction of about 85ms. This indicates that multithreading reduces the response time per request.

### CGI Tests

We also tested the response times of requests made to CGI scripts using the same Locust configuration and hosting conditions. Under heavy load with two worker threads, the average response time was 1506.15 ms. With four threads, the average response time improved to 882.46 ms, showing a speedup of 1.7. Additionally, the number of failures decreased from 222 to 11 with four threads (See Figures 9 and 10).

<sup>11</sup> This demonstrates a significant overall improvement when using more threads. Comparing requests to the static file with two threads, requests to the CGI script were 1267.08 ms slower. This is due to the spawning of new processes for each request in the CGI case and the additional code required for handling CGI scripts. These results would likely improve if we had implemented FastCGI.

<sup>10</sup><https://locust.io/>

<sup>11</sup>Detailed results can be found in \testing\results relative to the project root directory

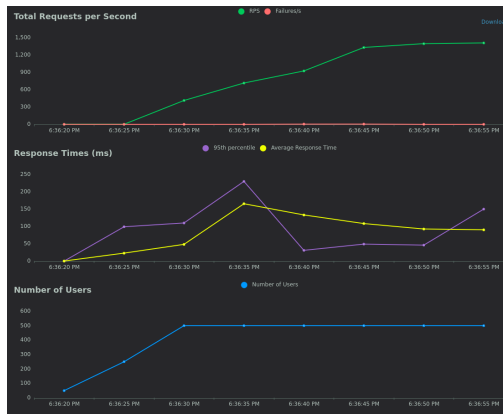


Figure 7: Requests to a static file, 2 Threads

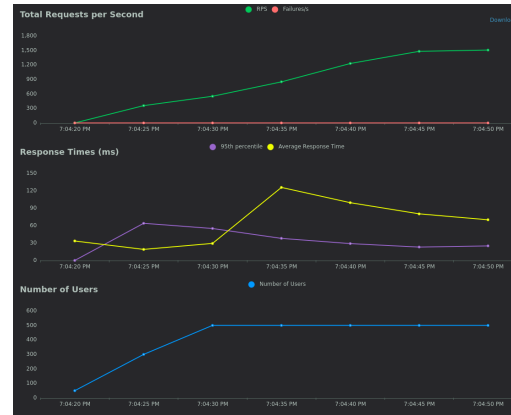


Figure 8: Requests to a static file, 3 Threads

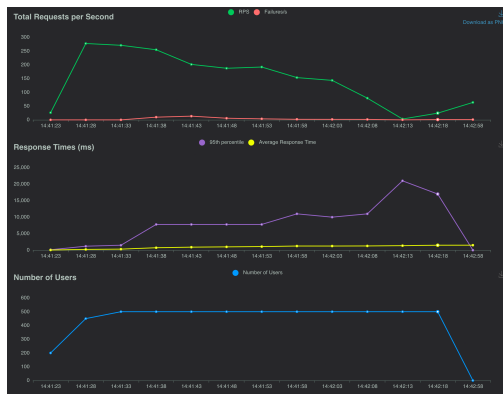


Figure 9: Requests to a CGI script, 2 Threads

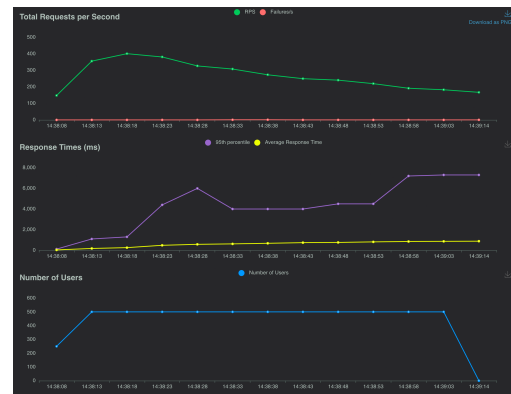


Figure 10: Requests to a CGI script, 4 Threads

## Conclusion

Our initial goal was to gain a deeper understanding of the internal workings of web servers and their configurations. We successfully achieved this goal through this project.

In Comparison to our initial project plan, our web server fulfilled the following requirements: process HTTP requests, allow multiple clients to connect, handle concurrent access to resources, handle errors, log web traffic, handle multiple domains and resources. We implemented everything but websockets. Overall, we are satisfied with our result.

## Lessons learned

**Deepen OS Knowledge** This project allowed us to apply a wide range of OS concepts covered in the OS course. We used socket communication to serve clients, implemented concurrency with OpenMP, created processes using `fork()`, managed inter-process communication with pipes, and implemented containerization with Docker.

**C Programming** We gained a deeper understanding of C programming principles, especially in handling strings. Writing HTTP parsers required extensive work with strings, leading to many bugs related to string handling in C. As a result, we have learned a lot about working with char arrays and continue to improve in this area.

**Project Plan** We rarely looked at the JIRA Project Plan we created in the beginning. Maybe next time we can find a more appropriate planning tool for our workflow. Where we spent more time as initially planned was the POST-Request handling, as we didn't know about CGI (Common Gateway Interface) before, and had to spend more time than anticipated understanding and implementing the protocol.

**Reproducibility** We encountered difficulties in ensuring our application ran consistently across different machines. Factors such as the OS, hardware, and network configurations on our servers caused errors. We

resolved some of these issues by containerizing our application, learning a lot in the process. However, despite containerizing everything, there were still differences in application behavior.

## **Future Outlook**

Below are some features we wanted to implement but didn't have the time for:

**Fast CGI** Implementing Fast CGI to enhance performance and keep up with industry standards.

**Security** Testing and enhancing the security of our web server.

**Web sockets** Implementing web sockets to eliminate the need for client polling for application state changes, thereby improving overall performance.



## References

**Source Code** <https://github.com/EphraimSiegfried/cuteserver>

## Libraries

**Hashmap** <https://github.com/tezc/sc>

**Logging** <https://github.com/rxi/log.c>

**Threadpool** (used in "main" branch) <https://github.com/Pithikos/C-Thread-Pool>

**OpenMP** (used in "openmp" branch) <https://www.openmp.org>

**TOML Parser** <https://github.com/cktan/tomlc99>

**JSON Parser** <https://github.com/json-c/json-c>

## Inspirations

We used the following web server implementations as inspirations for our own web server. Initially, we copied some code snippets from these sources, but over time, we refactored and significantly modified these code parts.

**Nweb** <https://nmon.sourceforge.io/pmwiki.php?n=Site.Nweb>

**Simple web server in C** <https://github.com/bloominstituteoftechnology/C-Web-Server>

## Use of AI

We primarily used ChatGPT for the example application part of our project. It wrote some functions and translated python Flask code to C code. However, we always had to modify much of the generated output.

## Declaration of Independent Authorship

We attest with our individual signatures that we have written this report independently and without outside help. We also attest that the information concerning the sources used in this work is true and complete in every respect. All sources that have been quoted or paraphrased have been marked accordingly.

Additionally, we affirm that any text passages written with the help of AI-supported technology are marked as such, including a reference to the AI-supported program used.

This report may be checked for plagiarism and use of AI-supported technology using the appropriate software. We understand that unethical conduct may lead to a grade of 1 or "fail" or expulsion from the study program.