# SCALABILITY REPORT

# ECE 568 Engineering Robust Server Software

**Author**

Zixiao Ma (zm96)

Shuqi Shen (ss1481)

April 8, 2024

# 1 Experimental Results

In this experiment, we tested different number of tasks on different number of CPU cores. The experimental results are shown below. Our application uses a thread pool with a fixed number of 10 threads.

## 1.1 Visualizing Experimental Data

### 1.1.1 Table of Data for Scalability

| Cores | Tasks | Time / s |
|:-----:|:-----:|:--------:|
| 1 | 30 | 5.383 |
| 1 | 100 | 6.378 |
| 1 | 200 | 7.372 |
| 1 | 500 | 10.632 |
| 1 | 1000 | 14.437 |
| 2 | 30 | 3.149 |
| 2 | 100 | 3.896 |
| 2 | 200 | 5.357 |
| 2 | 500 | 8.336 |
| 2 | 1000 | 12.882 |
| 4 | 30 | 2.874 |
| 4 | 100 | 3.558 |
| 4 | 200 | 5.250 |
| 4 | 500 | 8.376 |
| 4 | 1000 | 12.826 |

Table 1: Table of Data for Scalability
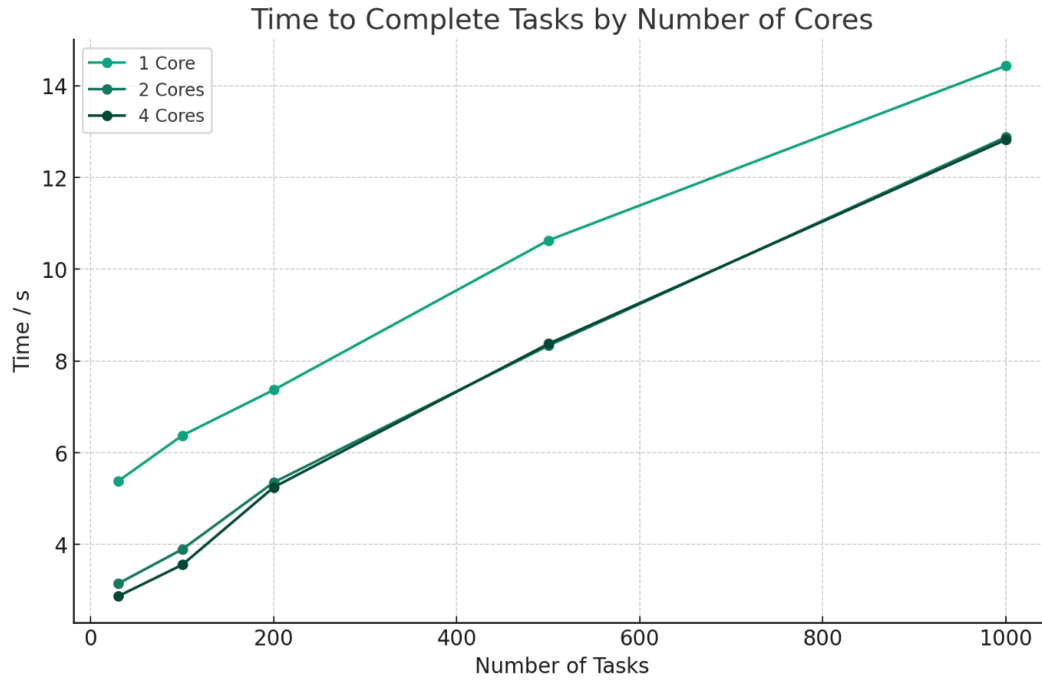
**1.1.2 Line graph of data for scalability**



Figure 1: Line Graph of Experimental Results

## 1.2 Performance Description

Increasing the number of cores significantly reduces the time taken for tasks. This improvement is most noticeable when moving from 1 core to 2 cores, and still significant but less drastic when moving to 4 cores.The diminishing returns effect becomes apparent with higher core counts. The jump from 1 to 2 cores offers a larger performance gain compared to the jump from 2 to 4 cores.

# 2   Scalability Analysis

## 2.1   Analysis of the Performance Divergence on Different CPU Cores

As it is mentioned above, when moving from 1 core to 2 cores, the performance improves a lot, and the performance does not change much when moving from 2 cores to 4 cores. The analysis of this experimental observation suggests that this is due to the CPU utilization drop.

From the provided screenshots and the data on CPU utilization, it is observed that when the server is operating with a single core, the CPU utilization is near maximum. This suggests that the CPU is fully occupied with the workload, leading to a bottleneck where the process cannot be executed any faster because there is no more CPU capacity available. When the server operates with two and then four cores, there appears to be some unused capacity, as indicated by the lower CPU utilization percentages. Therefore, the completion times of tasks do not significantly change as might be expected with the additional cores.
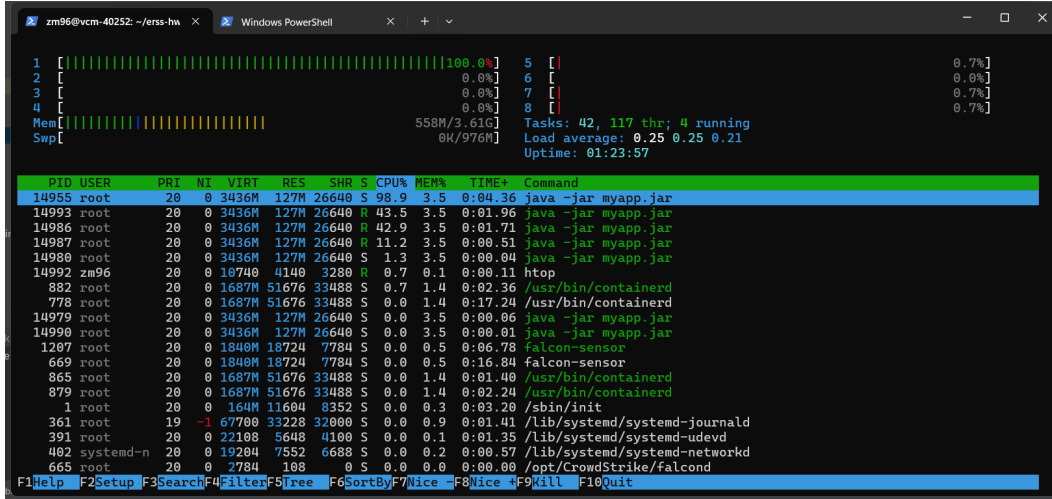


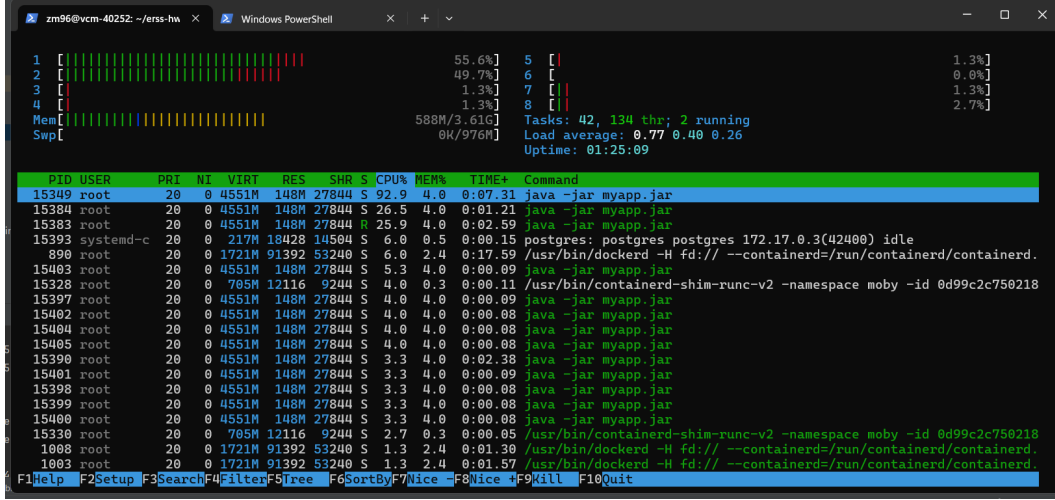Figure 2: CPU utilization for 1 CPU core
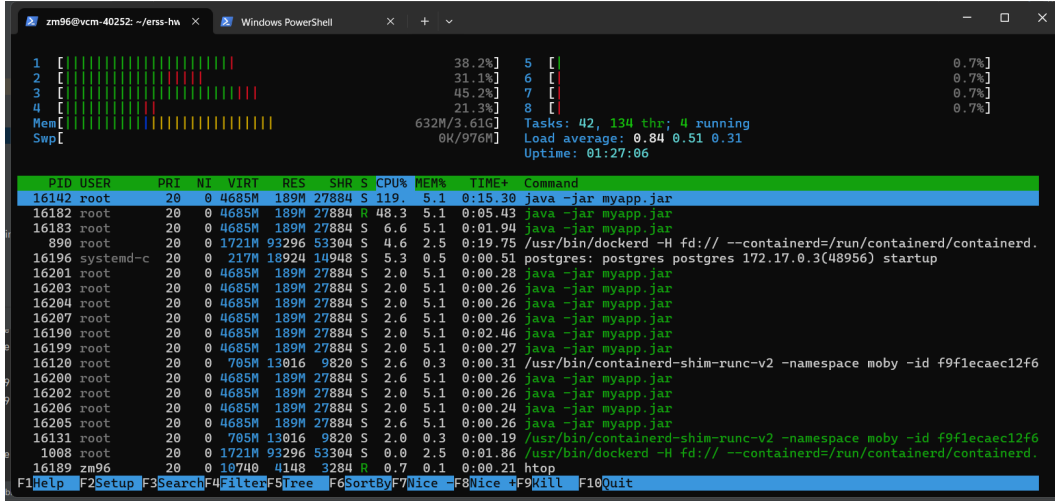
Figure 3: CPU utilization for 2 CPU cores



Figure 4: CPU utilization for 4 CPU cores

## 2.2 Analysis of the Performance Divergence on Different Task Numbers

As it is mentioned above, as the number of tasks increases, the response times increase linearly. The analysis of this experimental observation suggests that this is due to the bottlenecks of IO including Disk I/O, Network I/O and Database I/O.

### 2.2.1 Disk I/O Impact

Increased Read/Write Operations: As the number of tasks increases, if each task involves disk I/O operations, the total number of read and write operations also increases. This can lead to queuing delays as the disk works to process a larger queue of operations.

5

Concurrency and Access Times: Disks have a limit on how many concurrent operations they can handle. As more tasks attempt simultaneous access, the time to complete each operation can increase due to the need for the disk to switch between tasks.

### 2.2.2 Network I/O Impact

Bandwidth Saturation: With more tasks potentially sending or receiving data, the demand for network bandwidth can exceed supply, leading to congestion and increased latency.

Packet Queuing: Network devices may need to queue packets for transmission if the volume exceeds the immediate processing capability, which can result in increased network I/O latency.

### 2.2.3 Database I/O Impact

Connection and Query Load: More tasks often mean more database connections and queries. If a database cannot handle the increased load efficiently, it can become a bottleneck, with tasks waiting in line to access the database. Lock Contention: In a multi-task environment, tasks may compete for locks on database resources. As the number of tasks increases, so does lock contention, which can delay the completion of database operations.

In summary, as the number of tasks increases, the potential for I/O bottlenecks becomes more pronounced due to increased demands on disk, network, and database resources. These bottlenecks can result in tasks taking longer to complete, which may explain why the time does not decrease proportionally with additional CPU cores, as seen in the provided data. The increase in time with the increase in the number of tasks, even with more CPU cores available, underscores the importance of I/O in overall system performance.

# 3 Instructions on Experimental Result Reproduction

## 3.1 Run the Server

To run the server with different number of CPU cores, you should firstly use **sudo docker-compose up** and then **sudo docker-compose down -v** to make sure that all images are created. After that, you should use the first two commands in testing/commands.sh to run the application with specific CPU cores. After each test, you should run the last three commands in commands.sh to clean the stored and cached data for following tests.

## 3.2 Run the Tesing client

To reproduce the experimental result, you should open testing/client/SocketClient.java and go to the main function. If you prefer to run the test in the IDE like IDEA in local machine, make sure the working directory is testing/client, and change SERVER_HOST to the IP address of the server in the main function. If you prefer to run the test in the same server of the application, change SERVER_HOST to the localhost in the main function, and enter the testing/client directory and then compile and run SocketClient.java.

```
sudo docker run --name db --cpuset-cpus="0-3" -e POSTGRES_DB=postgres -e POSTGRES_USER=postgres -e POSTGRES_PASSWORD=postgres

sudo docker run --name app --cpuset-cpus="0-3" -p 12345:12345 --link db:db -d erss-hwk4-ss1481-zm96_app

sudo docker stop $(sudo docker ps -aq)

sudo docker rm $(sudo docker ps -aq)

sudo docker volume rm $(sudo docker volume ls -q)
```

Figure 5: commands.sh

To reproduce the scalability test, uncomment the code below and change 1000 to the number of tasks you want to test. To test the functionalities, we have provided comprehensive tests in the testing/testxml directory, and you can test them in the loop code above the scalability testing code, or read and send contents in whatever test files as you like. After you shut the program down, the total operation time will be printed.

```java
public static void main(String[] args) {
    SERVER_HOST = "localhost";
    SERVER_PORT = 12345;

    int numberOfThreads = 1;
    ExecutorService executor = Executors.newFixedThreadPool(numberOfThreads);

      for (int i = 1; i <= 20; i++) {
          String fileName = "../testxml/test" + i + ".xml";
          String xmlData = readXmlFromResources(fileName);

          if (xmlData == null) {
              System.err.println("Failed to read XML data from file: " + fileName);
              continue;
          }

          // log
          System.out.println("Thread will send data from: " + fileName);
          executor.submit(() -> sendXmlToServer(generateXml()));
      }


    for (int i = 0; i < 1000; i++) {
        executor.submit(() -> sendXmlToServer(generateXml()));
    }
```

Figure 6: Main function of SocketClient.java