# File server with backend authentication
## Phase 3 : Optimisation

Major Amit Pathania 163054001
Nivia Jatain 15305R007

**Department of Computer Science and Engineering**
**Indian Institute of Technology, Bombay**

November 2017

# Contents

# List of Figures

# List of Tables

# 1 Introduction

File server with backend authentication server has been implemented using TCP sockets in C++. There are two servers: Server1 and Server2. Server1 (File server) receive TCP requests from client and requests a username and password from client for either creating new account or for requesting a file. Server1 will send this username and password to Server2 (Backend authentication server) to authenticate the user or add new user account. Once authenticated, the client can then request files from the server1 and if file is present, server1 will send the file to client. **Server1 is multithreaded** to handle multiple client requests. **Server2 is multiprocess** to handle multiple authentication requests from Server1. The broad architecture of our system is described in figure1.
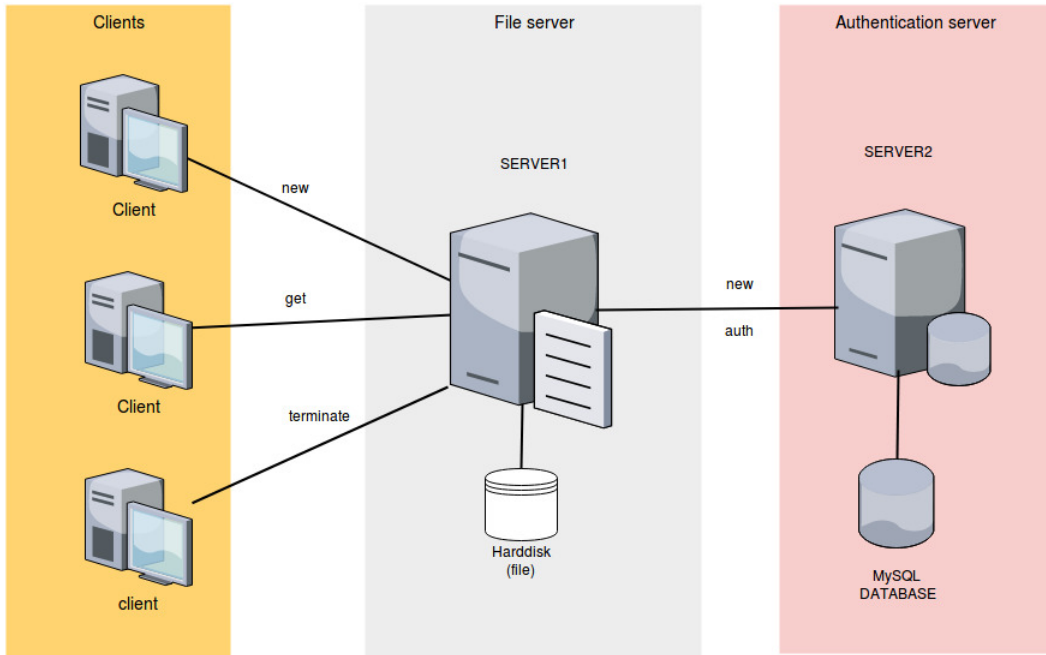


Figure 1: System Architecture

Our load generator performs close loop test using N threads.The response time, T , in a closed system is defined to be the time from when a request is submitted until it is received. In the case where the system is a single server (e.g. a web server), the server load, denoted by $\rho$ , is defined as the fraction of time that the server is busy, and is the product of the mean throughput X and the mean service demand (processing requirement) E[S].

# 2 Background

## 2.1 Design of load generator

We are performing a closed loop test where new job arrivals are only triggered by job completions.In order to simulate N concurrent users, we create N threads in our load generator. Each of the N threads emulate one user by issuing one request, waiting for it to complete, and issues the next request immediately afterwards.There is no think time between requests.
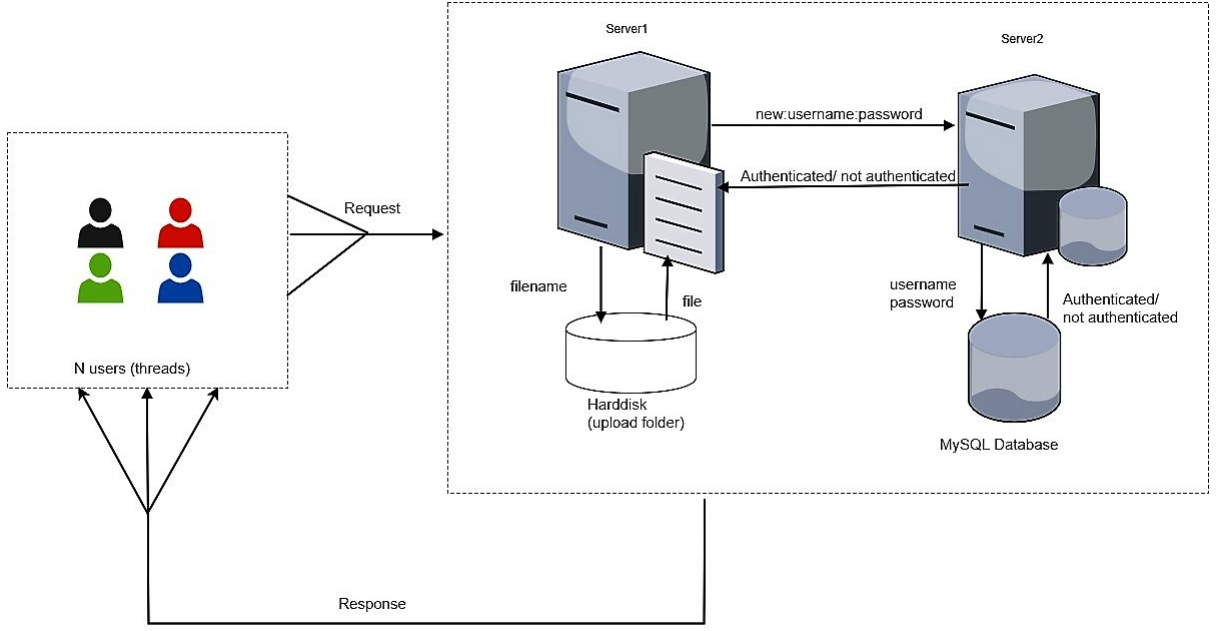
Figure 2: Design of load generator

## 2.2 Type of load/requests

Load generator can generate two types of load (or requests):

- **To create new user account**. In this case, "new" keyword followed by user-name and password will be sent to server1 (**new: uname: upassword**). Server1 will send "new" keyword followed with username and password to Server2 . Server2 will try to create user account. In this case,**there is no file transfer (disk access at Server1) and only new socket connection** is created between Client, Server 1 and server2 for each request. **Server2 will connect to mySQL database for each request** and try to create account. For creating new account, Server2 will write into the database. Hence, **disk access at Server2**.

  For each N, we defined total runtime ie 100-120 seconds.Each thread sent request to create account till runtime is over. We counted the number of requests completed and based on this time, we calculated throughput as Total Number of requests/ Actual elapsed time. For each thread,we generated the new user's username and password which was sent to server.

- **To request file from server1**. In this case, "get" keyword followed by user-name and password will be sent to server1 (**get: uname: upassword**). Server1 will send "auth" keyword followed with username and password to Server2 . Server2 will check the database for user account.If user details are correct, server1 will be informed and server1 will request filename from user. Client will then send filename with extension to server1(File server) and then file will be sent to client if that file exists at server1. In this case,**there is file transfer (ie disk access at Server1) and new socket connection** is created between Client, Server 1 and server2 for each request. **Server2 will connect to mySQL database** for each request and authenticate the user.

2

In this load test, we each user thread fetched different file of given file size everytime.We generated large number of test files of given size(1KB, 10KB,100KB, 1MB, 10MB, 100MB). In order to ensure that different files from server are fetched by different user threads with less probability of file being in cache, for each user we defined number of total fetch file requests.For instance, 2000 file request for each user. So, N=1,there were total 2000 requests,for N=2,total 4000 requests, for N=10, total 20000 requests and so on. The total time taken to complete these requests was calculated.And based on this time, we calculated throughput as Total Number of requests/ Total time taken.This approach helped us to ensure that total bytes of data downloaded by each user remained same for testing. In order to ensure that no user thread fetches the already fetched file and read the file from disk, we generated file name for each request as follows:

```
int r = (long(threadid)* 2000) +num_req+1;
strcpy(file_name,"500KB");
sprintf(count,  r);
strcat(file_name,count);
strcat(file_name,".txt");
num_req++;
```

This ensured that there is offset of 2000 between files demanded by each thread as per thread id and the same file is not requested again by some other thread. But, we had to generate 2000*N files for such test case. Also, after every one test for some value of N for given file, we cleared the buffer cache using following commands:-

```
sync; sudo sh -c "echo 3 > /proc/sys/vm/drop_caches"
```

**Generating test files**. We have to create files of varying sizes to check system load for different file sizes.Since, we don't care about contents of file and only file size matters, we used /dev/urandom option to create file. The shell code for same is as follows:

```
#! /bin/bash
for n in {1..5000}; do
  base64 /dev/urandom$|$head -c 200000000 $>$  200MB\$n.txt
done
```

## 2.3   Testing environment

Earlier in phase2, we had setup Server1 (File server) and Server2 (Authentication server) on Virtual Machine. Client is running on host machine. Client connects to server machine over bridged network between host and VM. Server contains upload folder which contains files which client can download.

However, for Phase3, we setup Servers on different system from client PC. The experiments done for load testing in Phase2 were again conducted on new system to calculate Throughput and Response Time values.

**Client and Server Hardware Configuration**. The hardware specifications of server and client machine for Phase 3 is as follows:

|  | **Load Generator** | **Server1 & 2** |
|---|---|---|
| Processor | 4x Intel(R) Core(TM) i5-4440 CPU @ 3.10GHz | Intel(R) Core(TM) i7-4790K CPU @ 4.00GHz |
| CPU(s) | 4 | 8 |
| Thread(s) per core | 1 | 2 |
| Core(s) per socket | 4 | 4 |
| Socket(s) | 1 | 1 |
| L1d cache | 32K | 32K |
| L1i cache | 32K | 32K |
| L2 cache | 256K | 256K |
| L3 cache | 6144K | 8192K |

Table 1: Client and Server Hardware Configuration

**Tuning MySQL**. By default, maximum number of simulatneous conenctions in MySQL is 150. Can be viewed using command:

show variables like 'max_connections';

The limit can be increased using command:

SET GLOBAL max_connections = 15000;

**Tuning sockets**: In order to make our server handle many concurrent connections, there is requirement to tune socket parameters. In listen system call in server, we define number of maximum incoming connections that can be handled. By default, in ubuntu this limit is 128 which is very less if we want to saturate the server.

To view maximum connections handled by socket:

cat /proc/sys/net/core/somaxconn

128(default)

**Increasing ulimit parameters**. By default the number of open files for user in Ubuntu is 1024. We have to increase it to measure maximum throughput.This is done with the ulimit command:

ulimit -a // see all the kernel parameters

ulimit -n // see the number of open files

ulimit -n 9000 // set the number open files to 9000

To make these changes permanent, add the following line to /etc/security/limits.conf:

useracct hard nofile 64000

# 3   Identification of bottleneck

The bottleneck of our system is **disk**. We monitored disk IO activity using nmon and iostat tools. For **New** request, there was **"Disk Write"** for each new user account creation in MySQL. At saturation, there was MySQL database connection error. For **Get** request, there was **"Disk Read"** at server for uploading file to the client.

The summary of results :

| Load type | File size | N* | Mean Response Time at N=N* | Max Through-put at N=N* | Bottleneck |
|---|---|---|---|---|---|
| new | | 800 | 0.219622 | 3363.49 | Disk I/O (Write) and MySQL |
| get | 1KB | 17 | 0.0054039 | 3093.7 | Disk I/O (Read) |
| get | 10KB | 10 | 0.00390524 | 2547.3 | Disk I/O (Read) |
| get | 100KB | 5 | 0.0113418 | 430.355 | Disk I/O (Read) |
| get | 1MB | 2 | 0.0265555 | 75.1187 | Disk I/O (Read) |
| get | 10MB | 1 | 0.0984429 | 10.1575 | Disk I/O (Read) |
| get | 100MB | 1 | 0.880579 | 1.1356 | Disk I/O (Read) |

Table 2: Test results before optimisation

## 3.1 Profiling the code

In order to find which portion of code/function, our server1 is spending most of the time, we used valgrind tool.

valgrind –tool=callgrind ./server1 10000

It generated a file called callgrind.out.xxxx. We used kcachegrind tool to read this file and get graphical analysis.



Figure 3: Code profiling

We can see in the profile output that our file server spends maximum time in "**upload_file**" function, ie reading the file and sending the file to client.

# 4 Optimization

## 4.1 Step1: Design changes

In order to compare performance benefits between Multi-process and multi-threaded architecture, we changed Server2 from Multi-process to multi-threaded.

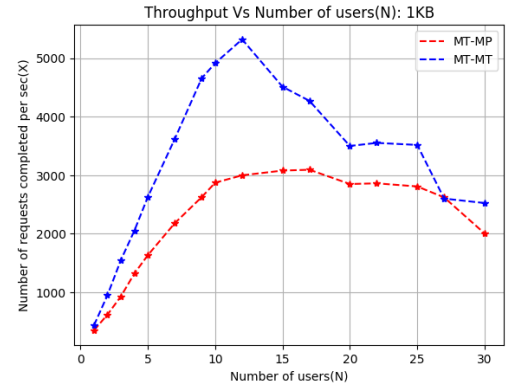### 4.1.1 Performance after design change

| Load type | File size | Before change (MT-MP) | | After Change (MT-MT) | |
|---|---|---|---|---|---|
| | | N* | Max Throughput | N* | Max Throughput |
| new | | 800 | 3363.49 | 900 | 6806.03 |
| get | 1KB | 17 | 3093.7 | 12 | 5323.32 |
| get | 10KB | 10 | 2547.3 | 9 | 3325.93 |
| get | 100KB | 5 | 430.355 | 5 | 465.524 |
| get | 1MB | 2 | 75.1187 | 1 | 81.0029 |
| get | 10MB | 1 | 10.1575 | 2 | 10.3492 |
| get | 100MB | 1 | 1.1356 | 1 | 1.14016 |

Table 3: Test results after design change

Hence, we conclude that forking a new process per connection has more overhead than creating thread per connection. As each client thread executes independently and don't have shared data/ variables, there is no overhead of synchronization using mutex/locks.



(a)



(b)



(c)



(d)

(e)



(f)



(g)

Figure 4: Comparison in performance between Multi threaded server and Multi process server for different requests

## 4.2 Step2: Disk IO optimization

When file reading system calls are used, data is usually transferred first from disk into the disk cache and then in the process memory. By prefetching the files into user space and saving it into user space using malloc, we can bypass disk access and disk buffer cache. Direct I/O is a feature of the file system whereby file reads and writes go directly from the applications to the storage device, bypassing the operating system read and write caches. We opened files in O_DIRECT mode and read every file in one I/O, i.e. create buffer in memory of the file size and read into it. We used malloc and created variables to store file contents when server started and all further client request were fetched from these without disk access. This reduced disk IO very much. This increased our throughput for 100KB file to 1154-1160 from 460-465.

However, as we increased N throughput got flattened but CPU was only 25-30% utilised. We used iperf tool to measure Bandwidth between our client and server machines and compared it with results of our load generator :

| Test | Interval | Data Transferred | Bandwidth (Mbps) | Bandwidth (MBps) |
|---|---|---|---|---|
| iperf | 10.0 sec | 1.09 GBytes | 935 | 116.875 |
| load generator | 30.0997 sec | 3.4742 GBytes | 923.384 | 115.423 |

7

Hence, we concluded **network as our bottleneck**. In order to **shift bottleneck from network to CPU**, we moved servers and load generator to same machine and pinned each server process, MySQL process and load generator to different core of the machine.

## 4.3 Step3: Code optimization

Code optimization to improve performance and reduce cache misses:

- Keeping variables that will be accessed together (or right after another) next to each other in memory by declaring them together. Also, created structure to include items accessed together.

- Removing loops to access contents of malloc(ed) files for uploading file. Used memcopy with offset to read and copy contents of file to send buffers.

### 4.3.1 Performance after disk I/O and code Optimization

**New bottleneck**: For files with size less than 1MB, MySQL was bottleneck and for files with larger file sizes, server1 process was bottleneck.



Figure 5: New Bottleneck for file size 1KB

Figure 6: New Bottleneck for file size 10KB



Figure 7: New Bottleneck for file size 100KB
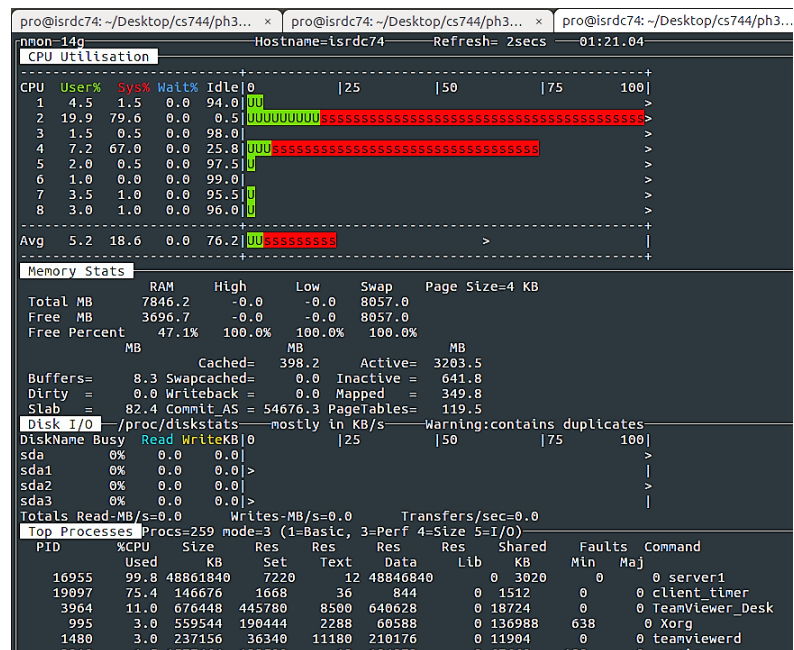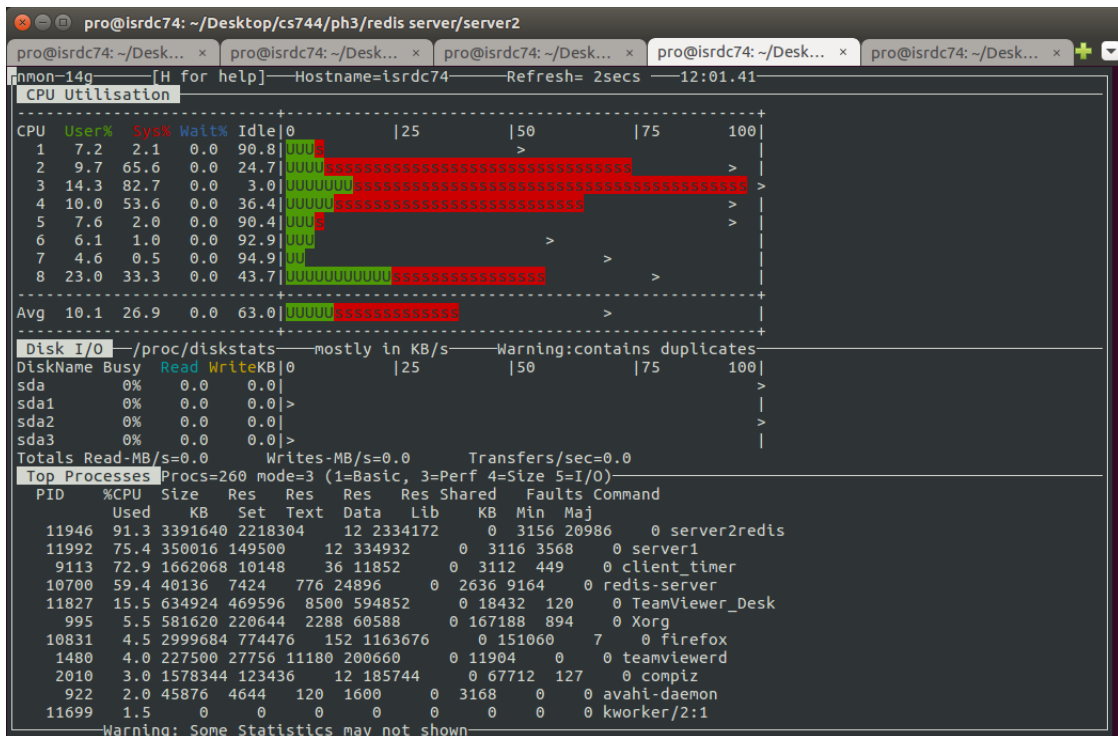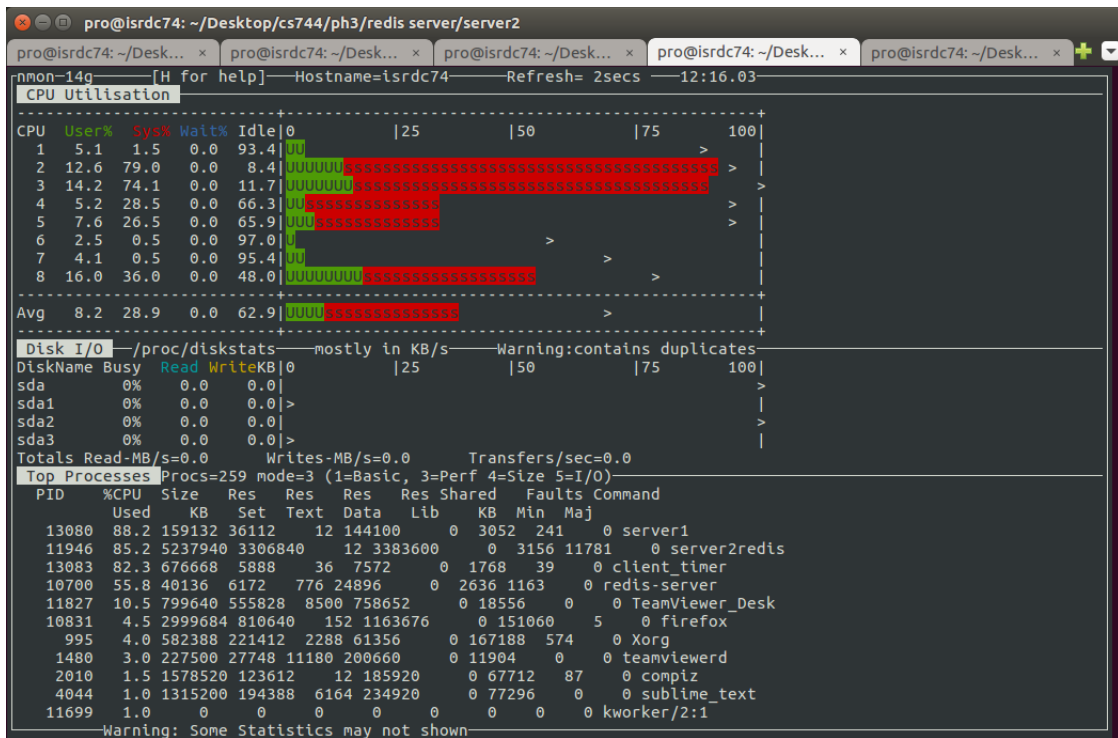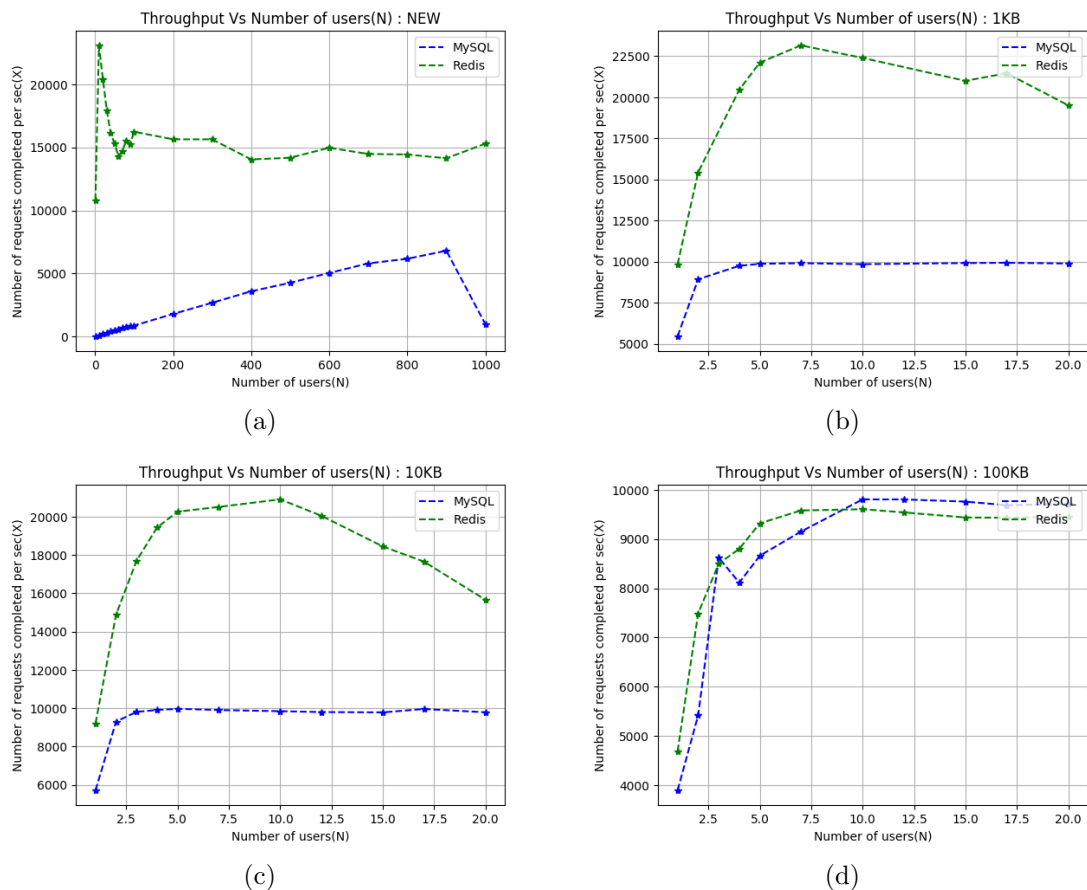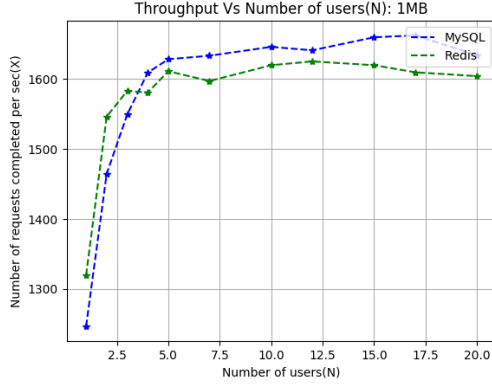
Figure 8: New Bottleneck for file size 1MB



Figure 9: New Bottleneck for file size 100MB

Hence, we replaced MySQL with Redis database.

## 4.4 Step4: Database optimisation using Redis Database

We created username-password as key-value in Redis database and changed server2 code to connect to Redis-server. For "new" request and "get" requests where file size is less than 1MB, there is increase in throughput of system and database was no longer a bottleneck.

10

Figure 10: New bottleneck for file size 1KB



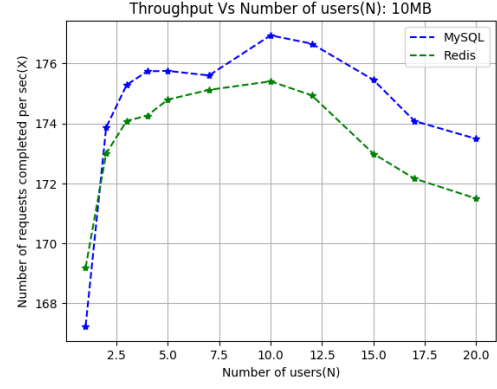Figure 11: New bottleneck for file size 10KB

Figure 12: New bottleneck for file size 100KB

## 4.4.1 Performance after database optimisation
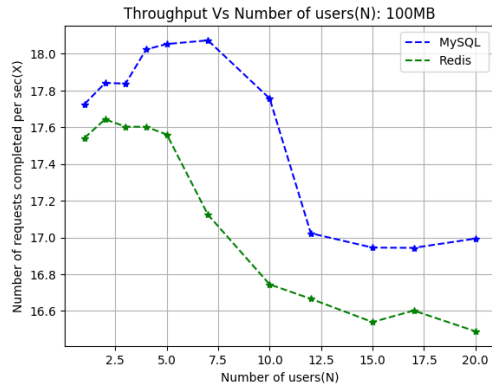


(a)



(b)
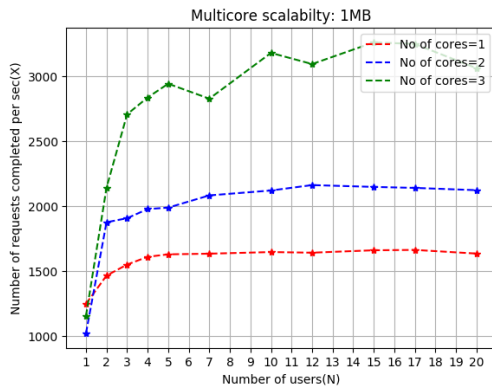
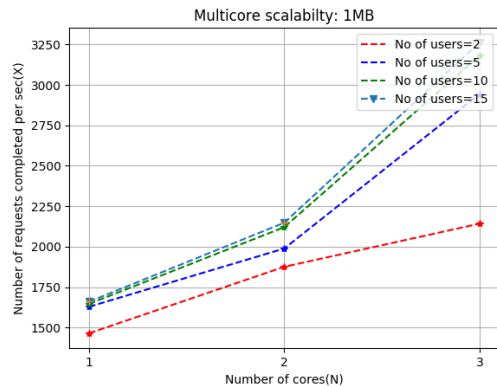

(c)



(d)

12

(e)



(f)



(g)

Figure 13: Performance after change in database

# 5 Multicore scalability

For file size=1MB, the server1 was limiting component and utilizing 96-98% CPU , so we increased number of cores using taskset. As we increased number of cores from 1 to 3, the throughput increased.



(a)



(b)

Figure 14: Multicore scalabilty
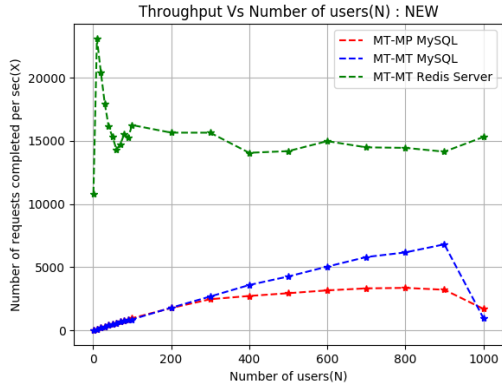
# 6 Summary of results

We first changed our server2 from multi process to multi thread and compared their performance. It was found that multi-threaded systems give better performance when there is no shared data and overhead of synchronization.

Then, in order to remove disk as bottleneck, we pre-fetched the files into user space by allocating memory to file contents using malloc. This improved the throughput a lot and we removed disk as bottleneck. We found that new bottleneck is MySQL for file sizes less than 100KB and for other cases, bottleneck was CPU core running server. So we changed our backend database from MySQL to Redis.
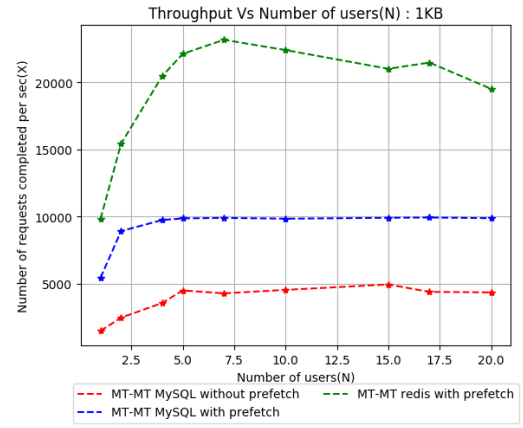
In cases where MySQL was bottleneck ie for file size less than 100KB, using Redis database increased the throughput. Also, for "new" request where disk was bottleneck for MySQL server, Redis server gave better performance because Redis is in-memory database platform and it don't write to disk continuously but saves the data to disk periodically from memory as per configuration.However for files with larger size(>100KB) the performance of systems with MySQL database and Redix Database was comparable. In such cases, bottleneck was CPU core running server.

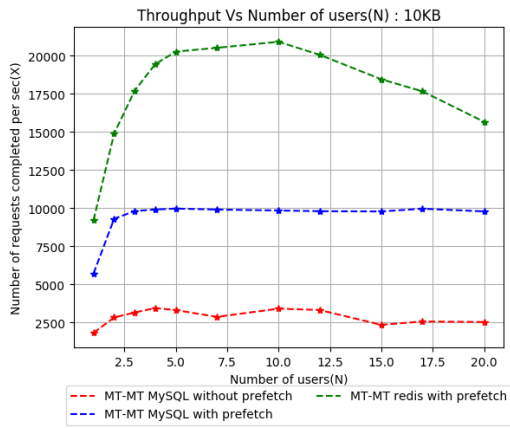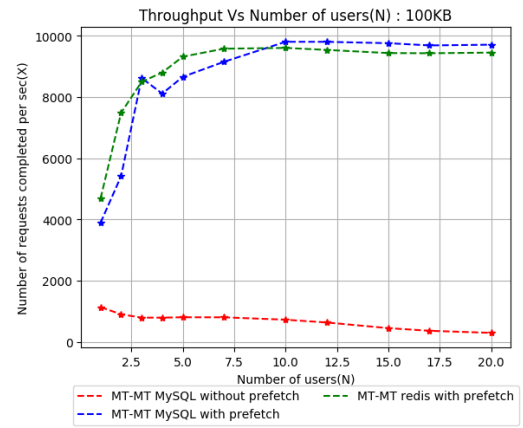| Load type | File size | Max Throughput | | |
|---|---|---|---|---|
| | | MT-MT without prefecting and MySQL DB | MT-MT with prefetching and MySQL DB | MT-MT with prefetching and RedisDB |
| new | | 6806.03 | 6806.03 | 23089.8 |
| get | 1KB | 4954.06 | 9932.58 | 23154.2 |
| get | 10KB | 3453.44 | 9973.76 | 20907.1 |
| get | 100KB | 1139.57 | 9804.96 | 9604.82 |
| get | 1MB | 135.965 | 1633.88 | 1624.99 |
| get | 10MB | 12.508 | 176.932 | 175.399 |
| get | 100MB | 1.14016 | 18.0733 | 17.6434 |

Table 4: Summary of results
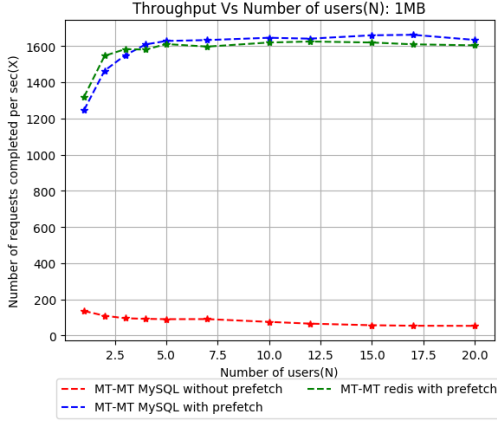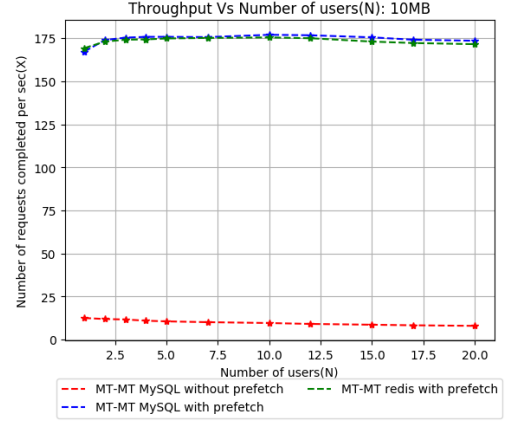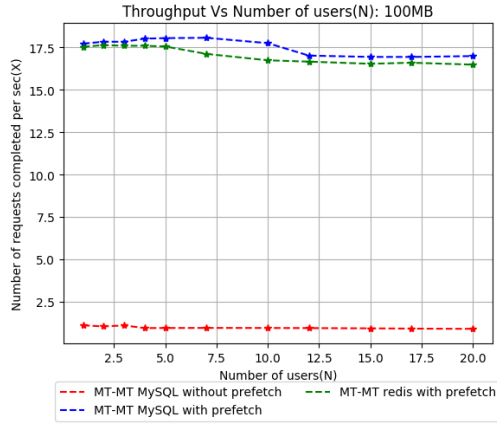
(a)



(b)



(c)



(d)

15

(e)



(f)



(g)

Figure 15: Summary of results

# 7 Submission folder

There are five folders inside : loadgen, ph2, multithreaded ,prefetching and redis.The ph2, multithreaded, prefetching and redis folder contains server1 and server2 sub-folders for respective stages of optimsations.

The folder/directory structure and files are listed below. The loadgen folder contains:

- Makefile

- client.cpp : Load generator (with Number of request as input)

- client_timer.cpp : Load generator (with run time as input)

- "downloads" folder : to store received files.

To run client/ load generator:
./client Server1_IP Server1_Listening-port request_type
Each server1 sub-folder contains:

- Makefile

16

- server1.cpp : File server code.

- "uploads" folder : to store files available for sharing/download.

- filecreate.sh : To create files with random content of specified size.

- clearbuffer.sh : To clear contents of cache buffer.

- tunesocket.sh : To tune the socket parameters.

Each server2 sub-folder of ph2, multithreaded and prefetching folder contains:

- Makefile

- server2.cpp : Authentication server code.

- installmysql.sh : script to install MySQL server and required libraries for header files.

- importdatabase.sh : script to import database and user table from cs744.sql file.

- cs744.sql : contains MySQL database and table for import.

Each server2 sub-folder of redis folder contains:

- Makefile

- server2.cpp : Authentication server code.

- installredis.sh : script to install Redis server and required libraries for header files.