



Министерство науки и высшего образования Российской Федерации
Федеральное государственное бюджетное образовательное
учреждение высшего образования
«Московский государственный технический университет имени
Н.Э. Баумана
(национальный исследовательский университет)»
(МГТУ им. Н.Э. Баумана)

ФАКУЛЬТЕТ «Информатика и системы управления»

КАФЕДРА «Программное обеспечение ЭВМ и информационные технологии»

Отчет по лабораторной работе №2 по дисциплине "Анализ алгоритмов"

Тема Параллельные вычисления

Студент Байрамгалин Я.Р.

Группа ИУ7-53Б

Преподаватели Волкова Л.Л., Строганов Ю.В.

Москва — 2022 г.

Оглавление

Введение	3
1 Аналитическая часть	4
1.1 Протокол HTTP	4
1.2 Умножение матриц	6
1.3 Серверное приложение	7
1.4 Рапараллеливание обработки запросов	7
2 Конструкторская часть	10
2.1 Серверное приложение	10
2.2 Очередь задач	10
2.3 Потребители	10
2.4 Умножение матриц	11
2.5 Парсинг http запросов	13
2.6 Модель оценки эффективности серверного приложения . . .	13
3 Технологическая часть	14
3.1 Требования к программному обеспечению	14
3.2 Выбор средств реализации	15
3.3 Реализация алгоритмов	15
3.4 Тестирование	21
4 Исследовательская часть	23
4.1 Технические характеристики	23
4.2 Анализ метрик	23
Заключение	25
Список использованных источников	26

Введение

Ранее на всех предметах в МГТУ им. Баумана рассматривались лишь алгоритмы, которые выполняются последовательно. В данной лабораторной работе предстоит погрузиться в безграничный мир параллельных вычислений.

Многопоточность в современных информационных системах используется повсеместно, так как во многих случаях позволяет обрабатывать большие массивы данных за то же процессорное время.

Целью данной лабораторной работы является создание серверного приложения, способного эффективно обрабатывать множество запросов параллельно, нахождение оптимального количества потоков для получения наилучшего результата.

Для достижения цели ставятся следующие задачи:

- разработать простое серверное приложение;
- распараллелить обработку поступающих запросов;
- сравнить время обработки запросов и максимальное количество запросов в секунду, которое способно обработать серверное приложение;
- оценить оптимальное количество обслуживающих потоков;
- провести сравнительный анализ на основе полученных экспериментально данных.

1 Аналитическая часть

Запросом будем называть сообщение (под сообщением подразумевается конечный набор байтов) в формате http запроса. Формат http выбран как самый популярный и простой формат клиент-серверного взаимодействия.

Ответом будет называть сообщение в формате http ответа.

В рамках данной лабораторной работы поставлена задача разработки серверного приложения. Серверным приложением будем называть endpoint, совершающий обработку запроса и возвращающий некоторый ответ.

Обработчик должен перемножать 2 матрицы, полученные в теле запроса, и отправлять ответ - результат перемножения матриц или же уведомлять, что перемножение выполнить невозможно.

1.1 Протокол HTTP

HTTP — протокол прикладного уровня передачи данных, изначально — в виде гипертекстовых документов в формате HTML, в настоящее время используется для передачи произвольных данных.

Формат http запроса состоит из:

- Версии протокола.
- HTTP-метода, обычно глагола подобно GET, POST или существительного, как OPTIONS или HEAD, определяющее операцию, которую клиент хочет выполнить. Обычно, клиент хочет получить ресурс (используя GET) или передать значения HTML-формы (используя POST), хотя другие операция могут быть необходимы в других случаях.
- Пути к ресурсу.
- Заголовков с метаданной (опционально).
- Тела запроса, предворяемого пустой строкой (опционально).

Пример http запроса представлен в листинге 1.1.

Формат http ответа состоит из:

- Версии протокола.
- Код состояния. Наиболее популярными являются: 200 (запрос выполнен успешно), 4xx (плохой запрос), 5xx (внутренняя ошибка сервера).
- Сообщение состояния, например OK, Bad request, Internal Server Error.
- Заголовки с метаданной (опционально).
- Тела ответа, предворяемого пустой строкой (опционально).

Пример http ответа представлен в листинге 1.2.

```
1 POST /echo/post/json HTTP/1.1
2 Host: reqbin.com
3 Authorization: Bearer mt0dgHmLJMVQhvjpNXDyA83vA_PxH23Y
4 Content-Type: application/json
5 Content-Length: 80
6
7 {
8   "Id": 12345,
9   "Customer": "John Smith",
10  "Quantity": 1,
11  "Price": 10.00
12 }
```

Листинг 1.1: Пример http запроса

```

1 HTTP/1.1 201 Created
2 Location: http://localhost/objectserver/restapi/alerts/status/kf/12481%3
  ANCOMS
3 Cache-Control: no-cache
4 Server: libnhttpd
5 Date: Wed Jul 4 15:31:53 2012
6 Connection: Keep-Alive
7 Content-Type: application/json; charset=UTF-8
8 Content-Length: 304
9
10 {
11   "entry": {
12     "affectedRows": 1,
13     "keyField": "12481%3ANCOMS",
14     "uri": "http://localhost/objectserver/restapi/alerts/status/kf/12481%3
  ANCOMS"
15   }
16 }

```

Листинг 1.2: Пример http ответа

1.2 Умножение матриц

В этом, а также во всех последующих разделах матрицей размера $m \cdot n$ будем называть прямоугольную таблицу, содержащую m строк и n столбцов.

Для начала дадим определение умножения матриц. Пусть есть матрица A размера $m \cdot n$ и матрица B размера $n \cdot l$. Тогда матрицей $C = A \cdot B$ будем называть матрицу

$$C = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1l} \\ a_{21} & a_{22} & \dots & a_{2l} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{ml} \end{pmatrix} \quad (1.1)$$

где

$$c_{ij} = \sum_{r=1}^m a_{ir} b_{rj} \quad (i = \overline{1, m}; j = \overline{1, l}) \quad (1.2)$$

Для умножения матриц будет использован стандартный алгоритм умно-

жения матриц, который заключается в подсчете значения для каждого c_{ij} согласно формуле 1.2.

1.3 Серверное приложение

Socket — название программного интерфейса для обеспечения обмена данными между процессами. Процессы при таком обмене могут исполняться как на одной ЭВМ, так и на различных ЭВМ, связанных между собой только сетью. Сокет — абстрактный объект, представляющий конечную точку соединения.

Для взаимодействия между машинами используются адреса и порты. Адрес представляет собой 32-битную структуру для протокола IPv4, 128-битную для IPv6. Номер порта — целое число в диапазоне от 0 до 65535. Эта пара и определяет сокет («гнездо», соответствующее адресу и порту).

Для выполнения задачи построения серверного приложения необходимо читать запросы, поступающие на определенный сокет, посылать их в обработчик и затем записывать их по другому сокету, созданному для принятия ответа.

1.4 Рапараллеливание обработки запросов

Для решения поставленной задачи обработки запросов существует несколько способов организации параллельности вычислений:

- Выполнять все в одном потоке. Такой подход является неэффективным, так как пока не будет выполнен первый запрос, поступивший на обработчик следующий запрос не будет взят в работу.
- Обработка каждого запроса в отдельном потоке. Такой подход также не является эффективным, так как при высокой загрузке такого серверного приложения количество потоков в системе может стать слишком большим и значительная часть процессорного времени будет тратиться на организацию работы потоков, а не на обработку самих запросов.

- Создание системы с фиксированным числом потоков, где один поток занимается лишь чтением запросов, поступающих на сокет, а остальные обрабатывают запросы и посылают ответы. Данный подход кажется оптимальным в данной ситуации, однако требует дополнительного анализа для поиска оптимального числа потоков.

Последний пункт может быть формализован как реализация модели пула потоков. Данная модель параллельных вычислений состоит в организации работы следующих сущностей:

- производитель, роль которого заключается в отправке задач в общую очередь задач;
- очередь задач, роль которой заключается в накоплении задач производителя и их передаче потребителю;
- потребители, роли которых заключаются в забираании задач из очереди и их обработке.

Схему работы пула потоков более наглядно можно проследить на рисунке 1.1.

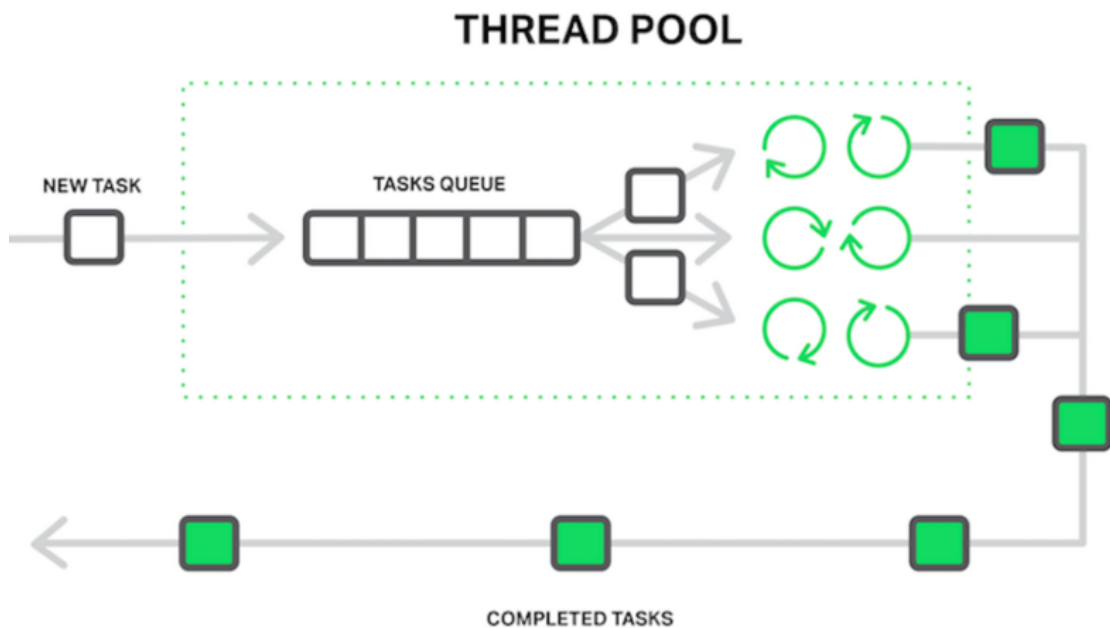


Рис. 1.1: Схема работы пула потоков

Вывод

Рассмотрены теоретические построения серверного приложения. Получены знания, достаточные для перехода к программной реализации задачи.

2 Конструкторская часть

2.1 Серверное приложение

Алгоритм работы серверного приложения можно разделить на 2 части: инициация сервера и обработка запросов.

Для инициация сервера необходимо сделать следующие:

- создать endpoint для комуникации системным вызовом `socket`;
- установить настройки сокета;
- соединить созданный дескриптор сокера с ассоциированным портом.

Для обработки запросов после инициации сервера необходимо бесконечно (до получения сигнала о завершении) выполнять следующие действия:

- слушать сокет системным вызовом `listen`;
- при получении запроса создать новое соединение системным вызовом `accept`;
- отправить полученный запрос и дескриптор нового соединения в очередь задач.

2.2 Очередь задач

Очередь задач — программная сущность, поддерживающая следующие операции: добавить задачу в очередь, получить задачу из очереди.

2.3 Потребители

Потребители — фиксированное число потоков, которые:

- забирают задачу из очереди;

- обрабатывают задачу (см. далее Умножение матриц);
- посылают ответ по переданному дескриптору соединения;

забирают задачи из очереди задач

2.4 Умножение матриц

Время выполнения стандартного алгоритма не зависит от каких-либо характеристик выбранных матриц, из-за чего рассматривать лучший и худший случай будет излишне.

На рисунке 2.1 приведена схема стандартного алгоритма умножения матриц.

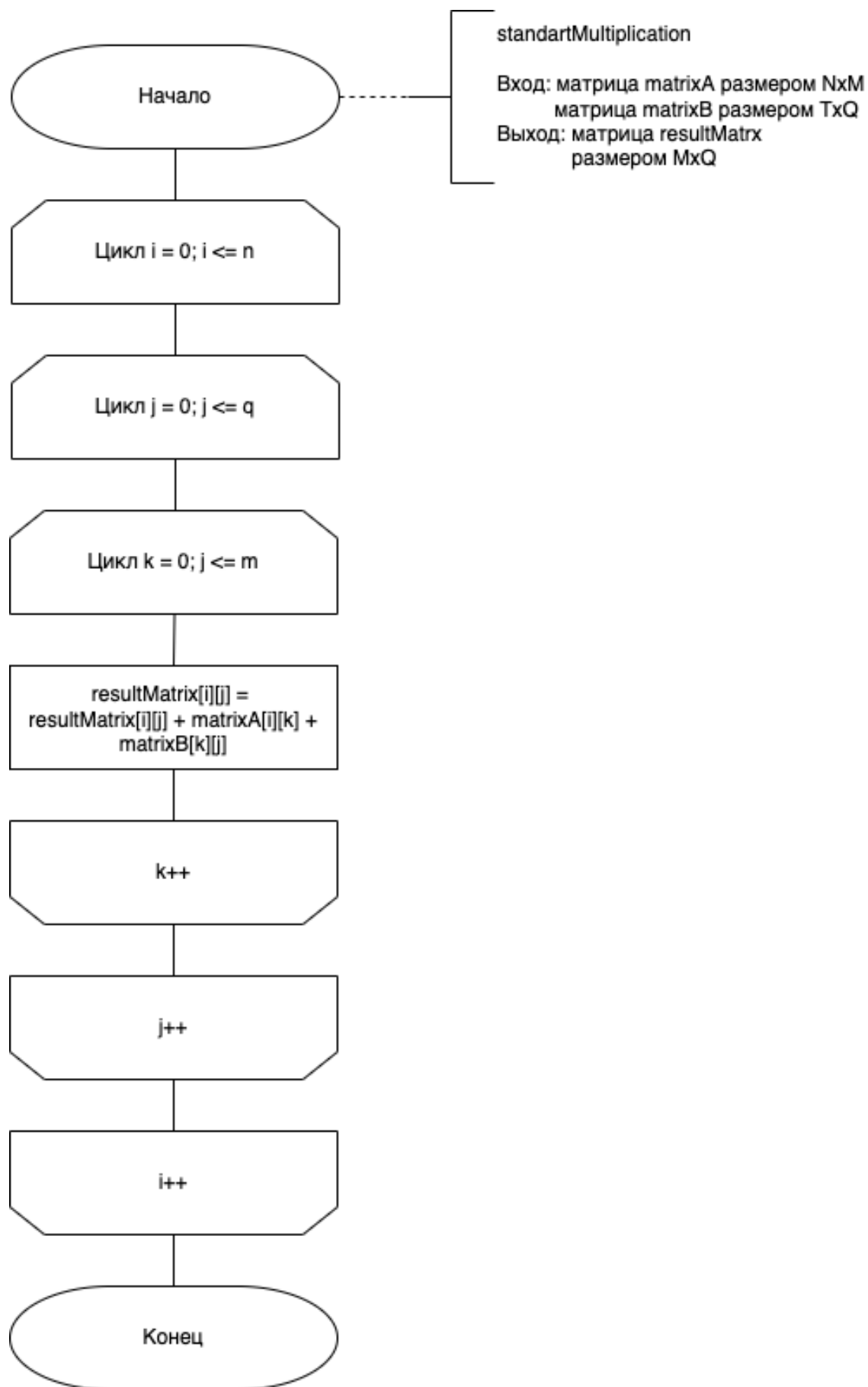


Рис. 2.1: Схема стандартного алгоритма умножения матриц

2.5 Парсинг http запросов

Задача парсинга http запросов не относится к сути выполненной лабораторной работы и поэтому подробное описание алгоритма парсинга http запросов не приводится.

2.6 Модель оценки эффективности серверного приложения

Для оценки эффективности серверного приложения, разрабатываемого в рамках данной лабораторной работы, предложены следующие метрики: время обработки запросов в персентилях 50, 66, 75, 80, 90, 95, 98, 99, 100; количество запросов в секунду, при котором сервер может стабильно работать (rps).

3 Технологическая часть

В данном разделе будут приведены требования к программному обеспечению, используемые технологии и реализации алгоритмов.

3.1 Требования к программному обеспечению

К программе предъявляется ряд требований:

- запрос поступает по заданному адресу сетевого соединения на порт 5002;
- запрос поступает в формате http запроса, версия протокола 1.1;
- в запросе HTTP метод — POST;
- в запросе путь к ресурсу — `/api/multiply_matrix`;
- в запросе указан заголовок Content-Type, его значение — text;
- в запросе указан заголовок Content-Length, его значение — длина тела запроса;
- в теле запроса располагаются через пробел (также допустимо использовать символ переноса строки вместо пробела) целые числа следующим образом: первые 2 числа m_1 , n_1 — размеры первой матрицы, далее расположены $m_1 \cdot n_1$ чисел — значения первой матрицы, далее 2 числа m_2 , n_2 — размеры второй матрицы, далее $m_2 \cdot n_2$ чисел — значения второй матрицы;
- ответ поступает в формате http ответа.
- ответ содержит заголовок Content-Type, его значение — text;
- ответ содержит заголовок Content-Length, его значение — длина тела запроса;

- при успешной обработке запроса в ответе содержится код выполнения 200, статус ОК, в теле ответа через пробел (или знак переноса строки) располагаются целые числа, два первых из которых показывают размер полученной матрицы, остальные - содержание полученной матрицы;
- при ошибке соответствия формата запроса протоколу http запроса в ответе содержится код выполнения 400, статус Bad Request;
- при невозможности умножения указанных в запросе матриц или несоответствия тела запроса формату, указанному выше в ответе содержится код выполнения 400, статус Bad Request.

3.2 Выбор средств реализации

В качестве языка программирования для реализации данной лабораторной работы был выбран язык программирования C++[1]. Язык позволяет управлять всеми ресурсами компьютера и позволяет делать необходимые системные вызовы.

Тестирование корректности работы серверного приложения производилось с помощью утилиты curl[2].

Метрики, изложенные в подразделе 2.6 замерены утилитой ab - Apache HTTP server benchmarking tool[3]

3.3 Реализация алгоритмов

Для параллельной обработки запросов были разработаны два класса: BlockingQueue, ThreadPool. Их реализация представлена в листингах 3.1, 3.2 соответственно.

Непосредственно реализация серверного приложения представлена в классе Server в листинге 3.3.

```

1 #pragma once
2
3 #include <deque>
4 #include <mutex>

```

```

5 #include <thread>
6 #include <chrono>
7 #include <condition_variable>
8
9 template<typename Task>
10 class BlockingQueue {
11     public:
12     void Add(Task task) {
13         std::lock_guard guard(mutex_);
14         queue_.push_back(std::move(task));
15         queue_not_empty_.notify_one();
16     }
17
18     Task Take() {
19         std::unique_lock lock(mutex_);
20         while (true) {
21             if (!queue_.empty()) {
22                 return TakeLocked();
23             } else {
24                 queue_not_empty_.wait(lock);
25             }
26         }
27     }
28
29     private:
30     Task TakeLocked() {
31         Task oldest = std::move(queue_.front());
32         queue_.pop_front();
33         return oldest;
34     }
35
36     std::deque<Task> queue_; // Guarded by mutex_
37     std::mutex mutex_;
38     std::condition_variable queue_not_empty_;
39 };

```

Листинг 3.1: Реализация блокирующей очереди

```

1 #pragma once
2
3 #include <thread>
4 #include <functional>
5 #include <deque>
6 #include <vector>
7 #include <stdexcept>
8
9
10 #include "pool/blocking_queue.hpp"
11

```



```

12 template<typename Task>
13 class StaticThreadPool {
14     public:
15     explicit StaticThreadPool(int thread_count)
16     : thread_count_(thread_count) {
17         StartWorkers();
18     }
19
20     void Submit(Task task) {
21         blocking_queue_.Add(std::move(task));
22     }
23
24     void Join() {
25         for (size_t i = 0; i < thread_count_; ++i) {
26             blocking_queue_.Add({});
27         }
28     }
29
30     private:
31     void StartWorkers() {
32         if (!pool_.empty()) {
33             throw std::runtime_error("Bad start workers call");
34         }
35
36         for (int i = 0; i < thread_count_; ++i) {
37             std::cout << "Starting worker " << i << "\n";
38             pool_.emplace_back([this]() {
39                 WorkerRoutine();
40             });
41         }
42     }
43
44     void WorkerRoutine() {
45         while(true) {
46             auto task = blocking_queue_.Take();
47             if (!task) {
48                 break;
49             }
50             task();
51         }
52     }
53
54     const int thread_count_;
55     BlockingQueue<Task> blocking_queue_;
56     std::vector<std::thread> pool_;
57 };

```

Листинг 3.2: Реализация пула потоков

```

1 #pragma once
2
3 #include <iostream>
4 #include <string>
5 #include <functional>
6
7 #include "pool/thread_pool.hpp"
8 #include "components.hpp"
9
10 class Server {
11     using Task = std::function<void()>;
12     public:
13     Server();
14
15     void StartListening();
16
17     private:
18     void ProcessRequest(int new_socket_descriptor) const;
19
20     int socket_descriptor_;
21     components::EndpointProvider endpoint_provider_;
22     StaticThreadPool<Task> thread_pool_;
23 };

```

Листинг 3.3: Реализация сервера из файла server/server.hpp

```

1 #include "server.hpp"
2
3 #include <iostream>
4 #include <stdexcept>
5 #include <string>
6 #include <functional>
7 #include <memory>
8
9 #include <unistd.h>
10 #include <sys/socket.h>
11 #include <netinet/in.h>
12
13 #include "models/request.hpp"
14 #include "models/response.hpp"
15 #include "utils/parse.hpp"
16 #include "views/multiply_matrix/handler.hpp"
17 #include "views/handler.hpp"
18 #include "pool/thread_pool.hpp"
19 #include "components.hpp"
20
21 namespace {
22

```

```

23     int kServerPortNumber = 5002;
24
25 }
26
27 Server::Server()
28 : socket_descriptor_(socket(AF_INET, SOCK_STREAM, 0)),
29 endpoint_provider_(components::EndpointProvider()),
30 thread_pool_(24) {
31     endpoint_provider_.Register(
32         views::Handler::HandleRequest,
33         RequestType::kGet, "/");
34     endpoint_provider_.Register(
35         views::multiply_matrix::Handler::HandleRequest,
36         RequestType::kPost, "/api/multiply_matrix");
37
38     if (socket_descriptor_ == -1) {
39         throw std::runtime_error("Socket error");
40     }
41
42     int enable = 1;
43     const auto opt_rc = setsockopt(socket_descriptor_, SOL_SOCKET,
44         SO_REUSEADDR,
45         &enable, sizeof(int));
46     if (opt_rc == -1) {
47         throw std::runtime_error("Setopt error");
48     }
49
50     sockaddr_in address{};
51     address.sin_family = AF_INET;
52     address.sin_addr.s_addr = INADDR_ANY;
53     address.sin_port = htons(kServerPortNumber);
54
55     auto* converted_address = (sockaddr*)&address;
56     if (bind(socket_descriptor_, converted_address, sizeof(address)) < 0) {
57         throw std::runtime_error("Bind error");
58     }
59
60     std::cout << "Server is ready to accept connections\n";
61 }
62
63 void Server::StartListening() {
64     listen(socket_descriptor_, 100);
65
66     sockaddr_in client_address{};
67     socklen_t client_length = sizeof(client_address);
68     std::cout << "Listening on port " << kServerPortNumber << "\n";
69
70     while (true) {

```

```

70     const auto new_sfd = accept(
71         socket_descriptor_, (sockaddr*)&client_address, &client_length);
72     if (new_sfd < 0) {
73         throw std::runtime_error("Accept error");
74     }
75     std::this_thread::sleep_for(std::chrono::milliseconds(1));
76     thread_pool_.Submit([this, new_sfd]() {
77         ProcessRequest(new_sfd);
78         close(new_sfd);
79     });
80 }
81 }
82
83 void Server::ProcessRequest(int new_socket_descriptor) const {
84     constexpr const size_t buffer_size = 524288;
85     std::unique_ptr<char[]> buffer = std::make_unique<char[]>(buffer_size);
86     bzero(buffer.get(), buffer_size);
87
88     if (recv(new_socket_descriptor, buffer.get(),
89         buffer_size, 0) == -1) {
90         throw std::runtime_error("Read error");
91     }
92
93     Request request{};
94     Response response{};
95     try {
96         request = utils::Parse(buffer.get(), utils::To<Request>());
97         const auto endpoint = endpoint_provider_.Get(request.type, request.
98             path);
99         response = endpoint(request);
100     } catch (std::exception& error) {
101         std::cerr << "Error occurred while handling request : " << error.what
102             ();
103         response = Response{500, "Server error"};
104     }
105
106     const auto response_string = BuildResponseString(response);
107     const auto written_bytes = write(
108         new_socket_descriptor, response_string.c_str(), response_string.size());
109
110     if (written_bytes < 0) {
111         std::cerr << "Write error";
112     }
113 }

```

Листинг 3.4: Реализация сервера из файла server/server.cpp

3.4 Тестирование

Рассмотрены несколько тестовых случаев, указанных ниже.

- запрос из листинга 3.5 - ответ из листинга 3.6

```
1  post /api/multiply_matrix HTTP/1.0
2  Host: localhost:5002
3  User-Agent: curl/7.84.0
4  Accept: */*
5  Content-Type: text
6  Content-Length: 25
7
8  2 2
9  5 7
10 9 1
11 2 2
12 7 4
13 1 -5
14
```

Листинг 3.5: Запрос, проверяющий правильность ответа на корректных данных

- запрос из листинга 3.5 - ответ из листинга 3.6

```
1  HTTP/1.1 200 OK
2  Content-Length: 25
3  Content-Type: text
4
5  2 2
6  42 -15
7  64 31
8
```

Листинг 3.6: Ожидаемый ответ на запрос из листинга 3.5

Вывод

Разработано серверное приложение, удовлетворяющее изложенным требованиям

4 Исследовательская часть

В данном разделе будут приведены примеры работы программ, постановка эксперимента и сравнительный анализ алгоритмов на основе полученных данных.

4.1 Технические характеристики

Технические характеристики устройства, на котором выполнялось тестирование, следующие.

- Операционная система: macOS Monterey 12.4[2].
- Память: 16 Гбайт.
- Процессор: 2,6 ГГц 6-ядерный процессор Intel Core i7[3].

Тестирование проводилось на ноутбуке, включенном в сеть электропитания. Во время тестирования ноутбук был нагружен только встроенными приложениями окружения, а также непосредственно системой тестирования.

4.2 Анализ метрик

В таблице 4.1 представлена зависимость времени выполнения запросов в персентелях от количества потоков в системе.

В таблице 4.2 представлено количество запросов в секунду, которое сервер может успешно обработать.

Таблица 4.1: Количество запросов в секунду, которое может обработать сервер

Потоков, ед.	50%, мкс	90%, мкс	95%, мкс	98%, мкс	100%, мкс
2	3956	4011	4162	4254	4287
6	904	947	964	1002	1075
12	734	792	901	932	1154
23	710	790	810	854	859
31	726	885	900	1023	1119

Таблица 4.2: Количество запросов в секунду, которое может обработать сервер

Потоков, ед.	tps
2	20
6	32
12	111
23	198
31	131

Вывод

Оптимальным является количество потоков, соответствующее количеству виртуальных ядер.

Заключение

В ходе выполнения лабораторной работы цель достигнута. Решены следующие задачи:

- разработано серверное приложение;
- проанализирована зависимость времени ответа от числа потоков в системе;
- проанализирована зависимость количества обрабатываемых запросов от числа потоков в системе;
- оценено оптимальное количество обслуживающих потоков;
- проведен сравнительный анализ на основе полученных экспериментально данных.
- подготовлен отчет о лабораторной работе;

СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ

- [1] C++ language [Электронный ресурс]. Режим доступа: <https://en.cppreference.com/w/cpp/language> (дата обращения: 04.10.2022).
- [2] macos Monterey 12.03 [Электронный ресурс]. Режим доступа: <https://www.apple.com/ru/macos/monterey/> (дата обращения: 04.10.2022).
- [3] Процессор Intel® Core™ i5-1135G7 [Электронный ресурс]. Режим доступа: <https://www.intel.ru/content/www/ru/ru/products/sku/208658/intel-core-i51135g7-processor-8m-cache-up-to-4-20-ghz/specifications.html> (дата обращения: 04.10.2022).