
中山大学计算机院本科生实验报告

(2023 学年春季学期)

课程名称：并程序序设计

批改人：

实验	3-Pthreads 并行矩阵乘法与数组求和	专业（方向）	计算机科学与技术计科一班
学号	21307099	姓名	李英骏
Email	liyj323@mail2.sysu.edu.cn	完成日期	2024 年 4 月 28 日

目录

1	实验目的	2
2	实验过程和核心代码	2
2.1	矩阵乘法	3
2.1.1	main.cpp	3
2.1.2	Pthread_Matmul.cpp / .hpp	4
2.1.3	性能测试	7
2.2	数组求和	8
2.2.1	main.cpp	8
2.2.2	IntArray	11
2.2.3	Linear Aggregation	12
2.2.4	Tree Aggregation	12
3	实验结果	13
3.1	矩阵乘法	13
3.2	数组求和	15
4	实验感想	17

1 实验目的

1. 使用 Pthreads 实现并行矩阵乘法, 并分析其性能。(尝试不同数据/任务划分方式)
2. 使用 Pthreads 实现并行数组求和, 并分析其性能。(可分析不同聚合方式的影响)

2 实验过程 and 核心代码

矩阵乘法中尝试了两种数据划分方式:

1. 划分 1: 将 A 按行均分, 在每个线程中与完整的 B 相乘, 得出结果矩阵的部分行
2. 划分 2: 将 A 按行均分, **B 按列均分**, 在每个线程中相乘, 得出结果矩阵的一块

数组求和中尝试了两种聚合方式:

1. 所有线程算出结果后, 在主线程中统一求和
2. 树形规约

代码文件结构如图:

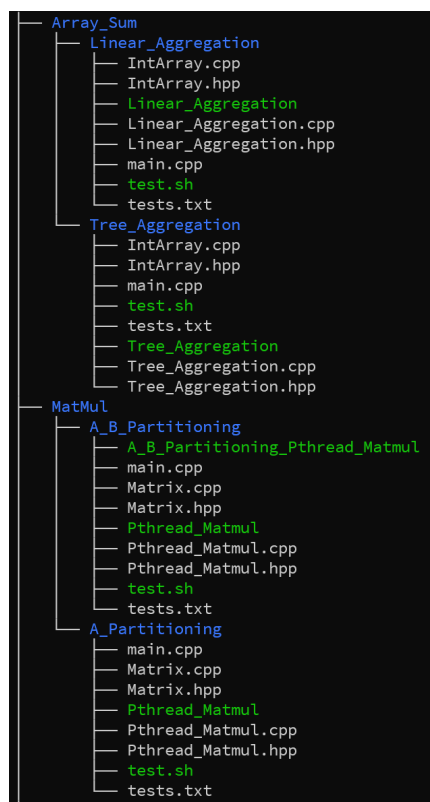


图 1: main.cpp

2.1 矩阵乘法

2.1.1 main.cpp

```
MatMul > A_Partitioning > main.cpp > main(int, char * [])
1  #include "Pthread_Matmul.hpp"
2  #include <chrono>
3  #include <stdlib.h>
4
5  long int thread_count;
6  int m, n, k;
7  Matrix A, B, C;
8  int main(int argc, char *argv[])
9  {
10
11     // 参数数量不正确则在线程中打印报错并return
12     if (argc != 5) {
13         std::cerr << "Usage: " << argv[0] << " m n k" << std::endl;
14         return 1;
15     }
16
17     // 读取矩阵参数
18     thread_count = std::strtol(argv[1], NULL, 10);
19     m = std::stoi(argv[2]);
20     n = std::stoi(argv[3]);
21     k = std::stoi(argv[4]);
22
23     // 初始化矩阵
24     A = Matrix(m, n, true);
25     B = Matrix(n, k, true);
26     C = Matrix(m, k, false);
27     // 开始计时
28     std::chrono::duration<double, std::milli> elapsed;
29     auto start_time = std::chrono::high_resolution_clock::now();
30     pthread_t *thread_handles = new pthread_t[thread_count];
31     ThreadTasks thread_tasks(thread_count);
32     // 计算矩阵乘法
33     for (long thread = 0; thread < thread_count; ++thread) {
34         pthread_create(&thread_handles[thread], NULL, pthread_matmul, (void *)&thread_tasks[thread]);
35     }
36     for (long thread = 0; thread < thread_count; ++thread) {
37         pthread_join(thread_handles[thread], NULL);
38     }
39     auto end_time = std::chrono::high_resolution_clock::now();
40     // 打印矩阵A B C
41     std::cout << "Matrix A:\n";
42     A.print_matrix();
43     std::cout << "\nMatrix B:\n";
44     B.print_matrix();
45     std::cout << "\nMatrix C:\n";
46     C.print_matrix();
47
48     delete [] thread_handles;
49     // 打印耗时
50     elapsed = end_time - start_time;
51     std::cout << "Time taken for matrix calculation: " << elapsed.count() / 1000 << " seconds\n";
52
53     return 0;
54 }
```

图 2: main.cpp

详见注释.

2.1.2 Pthread_Matmul.cpp / .hpp

划分 1

```
MatMul > A_Partitioning > Pthread_Matmul.hpp > ...
1  #ifndef P_MATMUL
2  #define P_MATMUL
3
4  #include <pthread.h>
5  #include "Matrix.hpp"
6  extern Matrix A, B, C;
7
8  struct ThreadTask {
9      int row_start;
10     int row_end;
11 };
12
13 class ThreadTasks
14 {
15 private:
16     ThreadTask *thread_tasks;
17
18 public:
19     ThreadTasks() : thread_tasks(nullptr) {}
20     ThreadTasks(long int thread_count)
21     {
22         thread_tasks = new ThreadTask[thread_count];
23         int m = A.get_rows();
24         int rows_per_thread = m / thread_count;
25         for (long int thread = 0; thread < thread_count - 1; ++thread) {
26             thread_tasks[thread].row_start = thread * rows_per_thread;
27             thread_tasks[thread].row_end = (thread + 1) * rows_per_thread - 1;
28         }
29         thread_tasks[thread_count - 1].row_start = (thread_count - 1) * rows_per_thread;
30         thread_tasks[thread_count - 1].row_end = m - 1;
31     }
32     ~ThreadTasks()
33     {
34         delete[] thread_tasks;
35         thread_tasks = nullptr;
36     }
37     ThreadTask& operator[](int index)
38     {
39         return this->thread_tasks[index];
40     }
41 };
42
43 void *pthread_matmul(void *arg);
44
45 #endif
```

图 3: Pthread_Matmul.hpp

注意 20-31 行, 记总线程数为 N , 三个矩阵分别为 $A_{m \times n} B_{n \times k} = C_{m \times k}$. 我们将矩阵 A 几乎平均地按行分成 N 份, 分发给 N 个 threads, 由于 $A B C$ 是全局的, 我们只需要计算每个线程需要计算的行号, 传入即可.

矩阵乘法:

```
MatMul > A_Partitioning > Pthread_Matmul.cpp > pthread_matmul(void *)
1  # include "Pthread_Matmul.hpp"
2
3  void *pthread_matmul(void *arg)
4  {
5      ThreadTask *task = (ThreadTask *)arg;
6
7      int row_start = task->row_start;
8      int row_end = task->row_end;
9      const double *const local_A = &(A.MAT[row_start * A.cols]);
10     const double *const local_B = B.MAT;
11     double *local_C = &(C.MAT[row_start * C.cols]);
12     const int m = row_end - row_start + 1;
13     const int n = A.cols;
14     const int k = B.cols;
15
16     sequential_matmul(local_A, local_B, local_C, m, n, k);
17     return nullptr;
18 }
```

图 4: Pthread_Matmul.cpp

在每个线程中计算地址偏移, 再进行矩阵乘法即可.

划分 2

```

20 ThreadTasks(long int thread_count)
21 {
22     thread_tasks = new ThreadTask[thread_count];
23     int m = A.get_rows();
24     int k = B.get_cols();
25     int A_partitions = 0;
26     int B_partitions = 0;
27     int sqrt_n = static_cast<int>(std::sqrt(thread_count + 0.1)); // 计算平方根并向下取整
28     for (int factor = sqrt_n; factor >= 1; --factor) {
29         if (thread_count % factor == 0) {
30             A_partitions = thread_count / factor;
31             B_partitions = factor;
32             break;
33         }
34     }
35     int A_rows_per_thread = m / A_partitions;
36     int B_cols_per_thread = k / B_partitions;
37     //std::cout << "MMM" << m << " " << A_partitions << " " << A_rows_per_thread << " " << (int)m / A_partitions << std::endl;
38     //std::cout << "A B PER THREAD " << A_rows_per_thread << " " << B_cols_per_thread << std::endl;
39     for (long int thread = 0; thread < thread_count; ++thread) {
40         int i = thread / B_partitions;
41         int j = thread % B_partitions;
42         if (i == A_partitions - 1) {
43             thread_tasks[thread].A_row_start = i * A_rows_per_thread;
44             thread_tasks[thread].A_row_end = m - 1;
45         } else {
46             thread_tasks[thread].A_row_start = i * A_rows_per_thread;
47             thread_tasks[thread].A_row_end = (i + 1) * A_rows_per_thread - 1;
48         }
49
50         if (j == B_partitions - 1) {
51             thread_tasks[thread].B_col_start = j * B_cols_per_thread;
52             thread_tasks[thread].B_col_end = k - 1;
53         } else {
54             thread_tasks[thread].B_col_start = j * B_cols_per_thread;
55             thread_tasks[thread].B_col_end = (j + 1) * B_cols_per_thread - 1;
56         }
57     }
58 }

```

图 5: Pthread_Matmul.hpp

在第二种划分中, 我们将线程总数 N 拆分为 $N = a * b$, 将矩阵 A 按行分成 a 份, B 按列分成 b 份, 从而在每个线程中计算出 C 的一块。

```

1 #include "Pthread_Matmul.hpp"
2 void *pthread_matmul(void *arg)
3 {
4     ThreadTask *task = (ThreadTask *)arg;
5     const int m = A.get_rows();
6     const int n = A.get_cols();
7     const int k = B.get_cols();
8     std::cout << "A: " << task->A_row_start << " " << task->A_row_end << std::endl;
9     std::cout << "B: " << task->B_col_start << " " << task->B_col_end << std::endl;
10    const int A_row_start = task->A_row_start;
11    const int A_row_end = task->A_row_end;
12    const int B_col_start = task->B_col_start;
13    const int B_col_end = task->B_col_end;
14    partitioned_sequential_matmul(A.MAT, B.MAT, C.MAT, A_row_start, A_row_end, B_col_start, B_col_end, m, n, k);
15    return nullptr;
16 }

```

图 6: Pthread_Matmul.cpp

其中 partitioned_sequential_matmul 的作用如图:

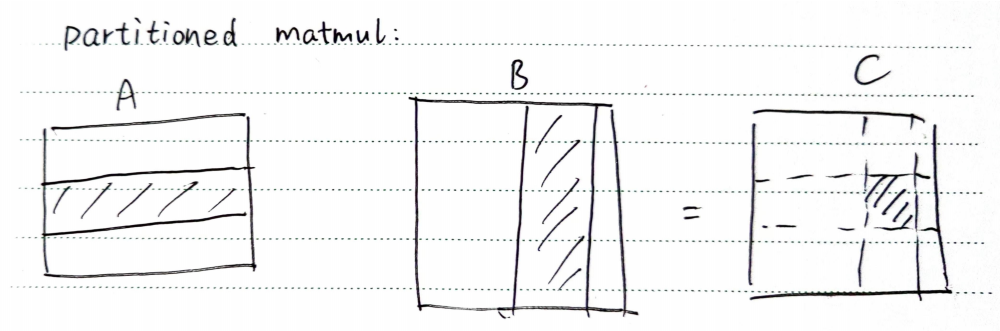


图 7: Pthread_Matmul.cpp

2.1.3 性能测试

我们使用如下脚本进行性能测试

```
MatMul > A_B_Partitioning > $ testsh
1  #!/bin/bash
2
3  g++ -g -Wall -o A_B_Partitioning_Pthread_Matmul main.cpp Matrix.cpp Pthread_Matmul.cpp -lpthread
4
5  # 矩阵维度
6  declare -a inputs=("128 128 128" "256 256 256" "512 512 512" "1024 1024 1024" "2048 2048 2048")
7  # 进程数量
8  declare -a procs=(1 2 4 8 16)
9
10 for proc in "${procs[@]}"
11 do
12     echo "worker线程数量: ${proc}"
13     for input in "${inputs[@]}"
14     do
15         echo "    对于参数m,n,k: $input"
16         total_time=0
17
18         # 每组运行3次
19         for i in {1..3}
20         do
21             echo "        运行 #${i}"
22             # 使用mpirun运行
23             output=$(./A_B_Partitioning_Pthread_Matmul $proc $input)
24
25             time_taken=$(echo "$output" | grep -o -E 'Time taken for matrix calculation: [0-9.]+ seconds' | grep -o -E '[0-9.]+')
26             echo "            time_taken: $time_taken"
27             total_time=$(echo "$total_time + $time_taken" | bc)
28         done
29
30         average_time=$(echo "scale=6; $total_time / 3" | bc)
31         echo "        平均运行时间: $average_time 秒"
32     done
33 done
34
```

图 8: test.sh

2.2 数组求和

2.2.1 main.cpp

聚合 1

```
1  #include "Linear_Aggregation.hpp"
2  #include <chrono>
3
4  long long thread_count;
5  long long N;
6  IntArray A;
7  int main(int argc, char *argv[])
8  {
9      if (argc != 3) {
10         std::cerr << "Usage: " << argv[0] << "Threads N" << std::endl;
11         return 1;
12     }
13     thread_count = std::strol(argv[1], NULL, 10);
14     N = std::strol(argv[2], NULL, 10);
15     A = IntArray(N);
16     // A.test();
17     // 开始计时
18     std::chrono::duration<double, std::milli> elapsed;
19     auto start = std::chrono::high_resolution_clock::now();
20
21     // 线程handles
22     pthread_t *thread_handles = new pthread_t[thread_count];
23     ThreadTasks thread_tasks(thread_count);
24     for (long thread = 0; thread < thread_count; ++thread) {
25         pthread_create(&thread_handles[thread], NULL, pthread_array_sum, (void *)&thread_tasks[thread]);
26     }
27     long long totalSum = 0;
28     for (int i = 0; i < thread_count; ++i) {
29         pthread_join(thread_handles[i], nullptr);
30         totalSum += thread_tasks[i].result;
31     }
32
33     // 停止计时
34     auto stop = std::chrono::high_resolution_clock::now();
35
36     // 打印A 和 sum(A)
37     std::cout << "Array A:\n";
38     A.print();
39     std::cout << "Sum:\n"
40         << totalSum << '\n';
41
42     // 检查求和结果是否正确
43     /*long long realsum = A.partitioned_sum(0, A.get_size() - 1);
44     std::cout << "Delta:\n"
45         << realsum - totalSum << '\n';*/
46     elapsed = stop - start;
47     delete[] thread_handles;
48
49     // 打印总耗时
50     std::cout << "Time taken for matrix calculation: " << elapsed.count() / 1000 << " seconds\n";
51
52     return 0;
53 }
54
```

图 9: main.cpp

聚合 2(树形聚合)

```
Array_Sum > Tree_Aggregation > G main.cpp > ...
3
4 long long thread_count;
5 long long N;
6 IntArray A;
7
8 // 全局数组和路障
9 std::vector<long long int> results;
10 pthread_barrier_t barrier;
11 int main(int argc, char *argv[])
12 {
13     if (argc != 3) {
14         std::cerr << "Usage: " << argv[0] << "Threads N" << std::endl;
15         return 1;
16     }
17     thread_count = std::strtol(argv[1], NULL, 10);
18     N = std::strtol(argv[2], NULL, 10);
19     A = IntArray(N);
20     // 全局数组,用于线程间通信
21     results.resize(thread_count, 0);
22     //A.test();
23     // 开始计时
24     std::chrono::duration<double, std::milli> elapsed;
25     auto start = std::chrono::high_resolution_clock::now();
26
27     // 线程 handles
28     pthread_t *thread_handles = new pthread_t[thread_count];
29     ThreadTasks thread_tasks(thread_count);
30     pthread_barrier_init(&barrier, NULL, thread_count);
31     for (long long thread = 0; thread < thread_count; ++thread) {
32         pthread_create(&thread_handles[thread], NULL, pthread_array_sum, (void *)&thread_tasks[thread]);
33     }
34
35     for (int i = 0; i < thread_count; ++i) {
36         pthread_join(thread_handles[i], nullptr);
37     }
38     long long totalSum = results[0];
39     // 停止计时
40     auto stop = std::chrono::high_resolution_clock::now();
41
42     // 打印A 和 sum(A)
43     std::cout << "Array A:\n";
44     //A.print();
45     std::cout << "Sum:\n"
46             << totalSum << '\n';
47
48     // 检查求和结果是否正确
49     /* long long realsum = A.partitioned_sum(0, A.get_size() - 1);
50     std::cout << "Delta:\n"
51             << realsum - totalSum << '\n'; */
52     elapsed = stop - start;
53     delete[] thread_handles;
54
55     // 打印总耗时
56     std::cout << "Time taken for matrix calculation: " << elapsed.count() / 1000 << " seconds\n";
57
58     return 0;
59 }
60
```

图 10: main.cpp

区别主要在于:

1. 树形规约为了实现进程间的同步和交流,使用了一个全局数组和一个 barrier.
2. 线性规约需要对所有线程的 results 求和,树形规约只要取数组的 0 号元素作为结果即可.

注释掉的部分用于检查算法正确性,如图:

```
10 94 16 50 44 15 55 2 17 57 43 8 21 81 34 66 77 22 61 61 19 49 100 11 4
9 15 77 3 46 18 75 20 99 28 45 20 49 45 77 82 86 13 6 41 31 44 12 30 92
97 27 91 3 54 54 77 51 15 62 92 88 84 3 86 48 80 2 85 48 2 72 13 82 69 1
65 32 92 14 75 57 54 29 16 23 73 35 18 20 41 44 48 60 28 28 1 9 75 46 9
6 67 32 52 15 89 57 64 71 76 42 98 53 87 27 5 29 93 2 89 91 84 20 2 19 8
5 79 87 66 29 96 46 58 13 47 81 65 34 81 76 66 11 36 69 80 81 1 88 68 7
58 68 23 23 41 36 51 6 80 41 53 95 74 43 66 43 82 2 27 27 61 75 79 15 91
10 85 5 79 57 85 71 53 5 40 68 9 21 58 90 72 52 5 76 76 4 90 23 10 25 2
82 75 76 28 70 58 8 17 25 55 91 27 10 8 63 66 14 86 85 21 85 51 99 5 53
32 30 50 83 88 48 26 93 35 42 43 52 89 20 61 72 1 30 13 1 60 69 81 67 9
0 24 82 38 10 55 76 46 98 60 77 51 100 77 17 2 58 50 69 13 12 65 29 28 1
5 87 84 87 92 9 81 24 37 26 97 57 89 8 99 7 29 63 97 32 49 6 15 88 62 63
65 27 51 3 71 16 30 77 6 77 51 96 27 87 55 73 96 24 94 4 29 23 85 79 28
97 46 53 2 7 92 26 44 44 35 59 100 37 3 11 44 34 87 74 50 99 4 86 82 25
69 49 63 68 78 6 65 45 74 94 89 90 76 56 2 69 67 23 87 25 91 30 49 60 6
9 58 43 8 70 71 81 36 61 31 24 36 36 5 27 12 43 16 61 3 7 75 79 44 21 50
41 60 80 32 97 66 7 51 91 96 52 80 7 47 80 28 23 34 76 46 69 74 79 ]
Sum:
6471821
Delta:
0
Time taken for matrix calculation: 0.0014349 seconds
```

图 11: check

上图 Delta 为 0 表示并行求和的结果与单线程的循环求和结果一致, 说明算法正确.

2.2.2 IntArray

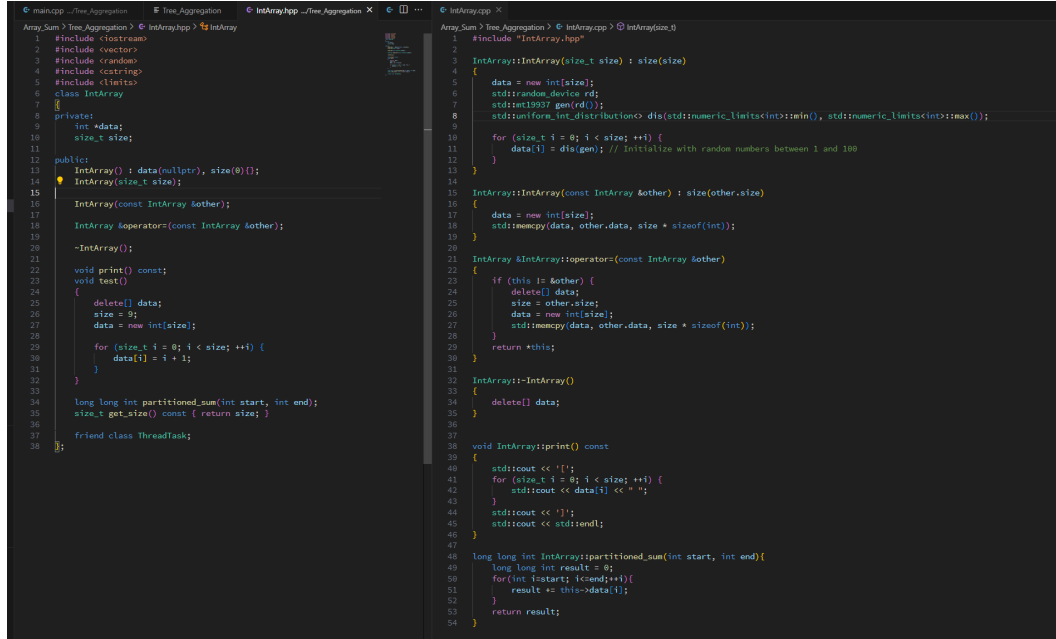


图 12: IntArray

第 3-13 行为初始化函数, 可以看到使用均匀分布函数, 将每个元素初始化为 `int`.
48-53 行为计算数组的部分和.

2.2.3 Linear Aggregation

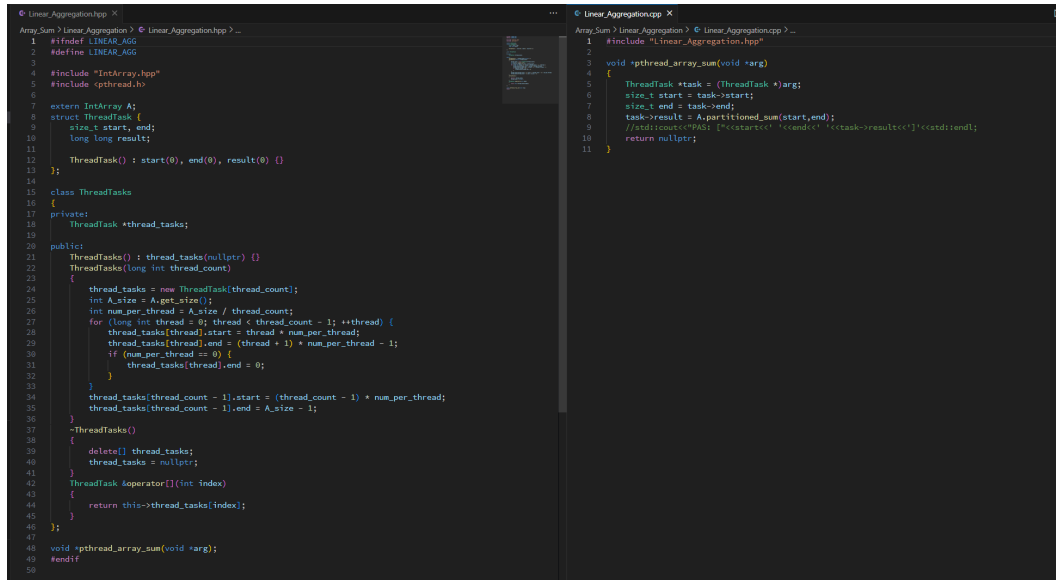


图 13: Linear Aggregation

简单计算部分和即可。

2.2.4 Tree Aggregation

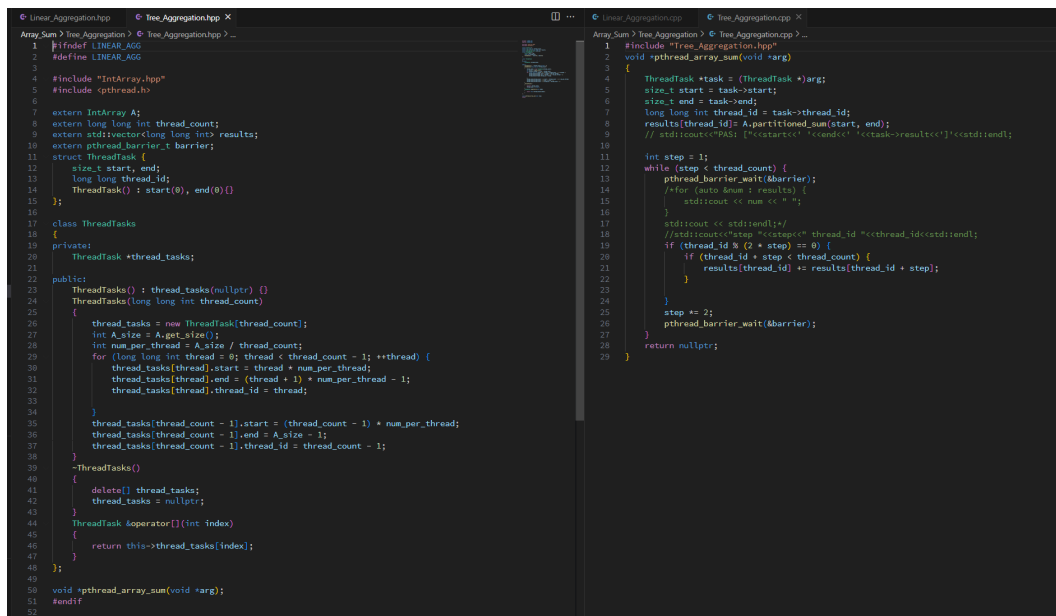


图 14: Tree Aggregation

右侧 11-27 行为树形规约, 为简便只考虑了线程数是 2 的幂次的情形. 实现了下图的二叉树:

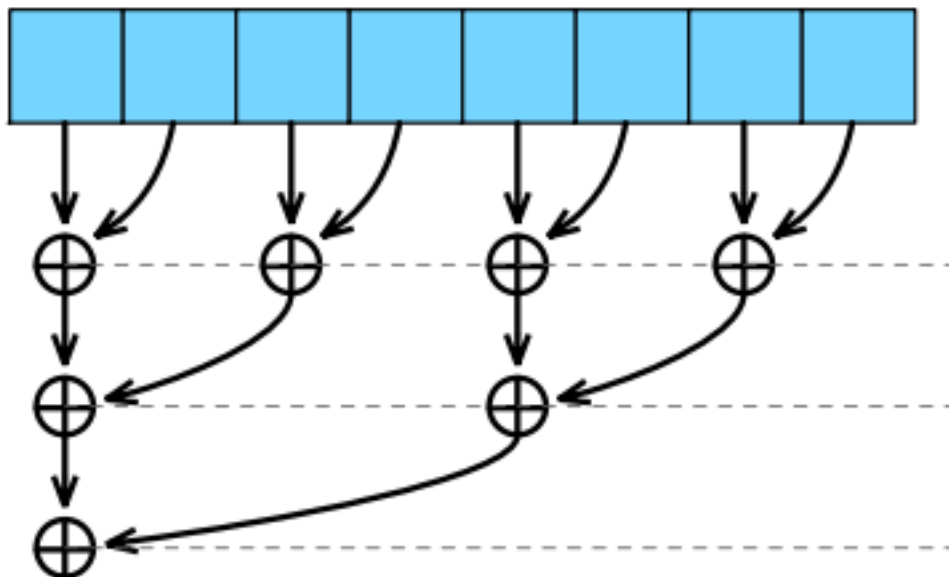


图 15: Tree Aggregation

3 实验结果

3.1 矩阵乘法

划分 1:

表 1: 划分 1 线程数与矩阵规模的性能数据 (秒)

worker 进程数 \ 矩阵规模	128	256	512	1024	2048
1	.003431	.026890	.219175	1.711090	13.725566
2	.002038	.014694	.109977	.853172	6.847813
4	.001028	.007010	.056423	.418382	3.421770
8	.000829	.004758	.041580	.243998	2.082076
16	.000796	.003455	.028399	.210076	1.474280

划分 2:

表 2: 划分 2 线程数与矩阵规模的性能数据 (秒)

worker 进程数 \ 矩阵规模	128	256	512	1024	2048
1	.003523	.026857	.210081	1.726030	13.360266
2	.001849	.013418	.106289	.843146	6.833443
4	.001564	.010202	.058136	.448828	3.569666
8	.001168	.007704	.031460	.243478	1.864370
16	.001078	.006368	.035039	.192912	1.255803

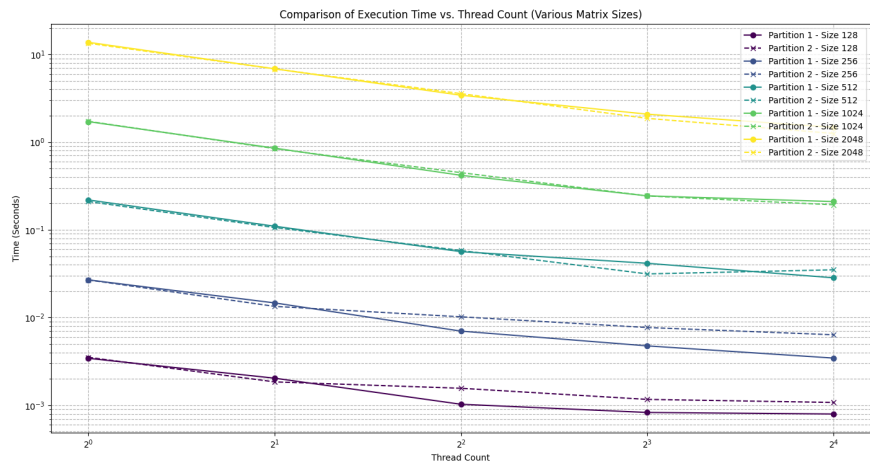


图 16: result

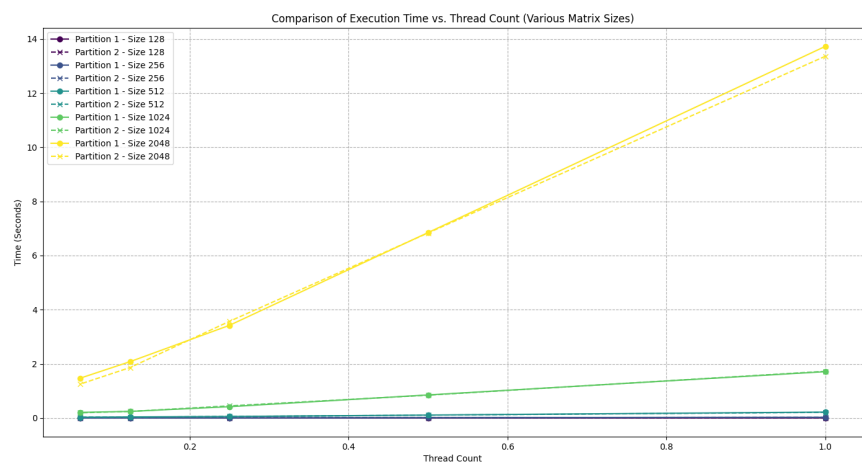


图 17: result

观察出以下结论:

1. 时间开销随着线程数量的增加几乎反比例地下降 (图中近似直线).
2. 两种划分算法性能相仿. 第二种算法在矩阵较大时略优于第一种算法.

3.2 数组求和

划分 1:

表 3: 聚合 1 线程数与数组规模的性能数据 (秒)

worker 进程数 \ 数组规模 (M)	1	2	4	8	16	32	64	128
1	.001481	.002983	.005647	.011060	.021771	.044333	.088551	.178929
2	.000806	.001530	.002885	.005695	.011572	.023171	.047254	.092074
4	.000623	.000949	.001713	.004073	.007763	.014384	.023820	.045619
8	.000762	.001147	.001613	.003309	.006813	.013369	.013515	.041619
16	.001367	.001779	.002284	.004367	.005254	.011448	.018628	.033708

划分 2:

表 4: 聚合 2 线程数与数组规模的性能数据 (秒)

worker 进程数 \ 数组规模 (M)	1	2	4	8	16	32	64	128
1	.001505	.002958	.005740	.011594	.022237	.043872	.085593	.173632
2	.000893	.001481	.002779	.005510	.011183	.022117	.044044	.092859
4	.000681	.001270	.002966	.002889	.005716	.011321	.023251	.052073
8	.000629	.001289	.001648	.002874	.004138	.009129	.019201	.035495
16	.000517	.000816	.001202	.002754	.004697	.008007	.015590	.026738

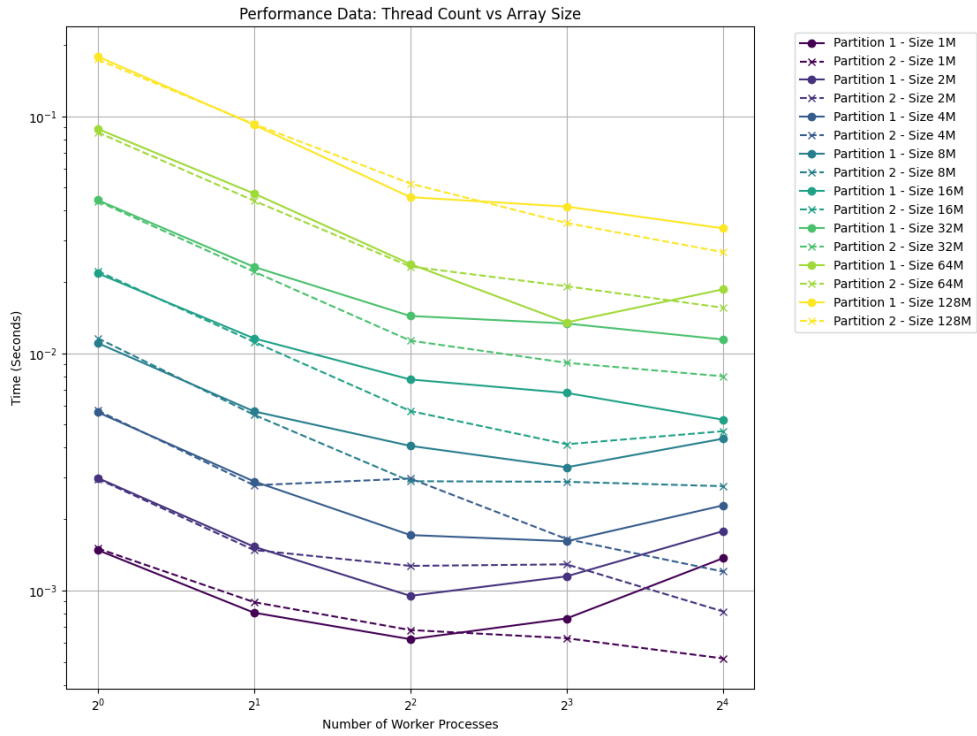


图 18: result

观察出以下结论:

1. 树形规约在线程数较多时, 显著优于线性规约, 而在线程数较少时, 由于 barrier 的影响, 可能劣于线性规约.
2. 数组维度较小时, 线程数增加反而会降低求和速度. 这是由于线程池创建和销毁线程的开销, 以及上下文切换的开销较大.

4 实验感想

通过本次实验，我深入了解了并行程序设计中矩阵乘法和数组求和的实现方式，特别是如何利用 pthread 实现高效的并行计算。实验中，我尝试了两种不同的数据划分方式，即按行均分和行列均分，这不仅加深了我对并行计算中数据划分策略的理解，也让我认识到在实际应用中，针对不同的计算任务选择合适的划分策略是多么重要。

同时，我还了解了树形规约的实现方式。

此外，性能测试也是本次实验的重要部分。通过对比不同划分方式下的执行时间，我更加直观地理解了并行计算的性能优化。