

组成原理矩阵乘法作业2实验报告

学生姓名：刘修铭 学号：2112492 指导老师：董前琨 班级：李涛老师

一、实验要求

在Taishan服务器上使用vim+gcc编译环境完成：

1. 实现矩阵乘法；
2. 完成矩阵乘法的优化，并在1024~4096规模的矩阵上进行测试；
3. 总结出不同层次，不同规模下的矩阵乘法优化对比，对比指标包括计算耗时、运行性能、加速比等；
4. 对比Taishan服务器和自己个人电脑上程序运行时间等相关指标，分析不同电脑上的运行差异的原因；
5. 总结优化过程中遇到的问题和解决方式。

二、矩阵乘法的实现及优化

```
1  #include <iostream>
2  #include <vector>
3  #include <chrono>
4  #include <thread>
5
6  using namespace std;
7  using namespace std::chrono;
8
9  // 矩阵乘法（按照定义计算）
10 vector<vector<int>>> matrixMultiplication(const vector<vector<int>>& matrix1, const vector<vector<int>>&
    matrix2) {
11     int rows1 = matrix1.size();
12     int cols1 = matrix1[0].size();
13     int rows2 = matrix2.size();
14     int cols2 = matrix2[0].size();
15
16     if (cols1 != rows2) {
17         cout << "Error: The number of columns in matrix1 should be equal to the number of rows in
matrix2." << endl;
18         return {};
19     }
20
21     vector<vector<int>>> result(rows1, vector<int>(cols2, 0));
22
23     for (int i = 0; i < rows1; i++) {
24         for (int j = 0; j < cols2; j++) {
25             for (int k = 0; k < cols1; k++) {
26                 result[i][j] += matrix1[i][k] * matrix2[k][j];
27             }
28         }
29     }
30
31     return result;
32 }
33
34 // 优化的矩阵乘法（分块矩阵乘法）
35 vector<vector<int>>> optimizedMatrixMultiplication(const vector<vector<int>>& matrix1, const
    vector<vector<int>>& matrix2) {
36     int rows1 = matrix1.size();
37     int cols1 = matrix1[0].size();
38     int rows2 = matrix2.size();
39     int cols2 = matrix2[0].size();
40
41     if (cols1 != rows2) {
```

```

42     cout << "Error: The number of columns in matrix1 should be equal to the number of rows in
matrix2." << endl;
43     return {};
44 }
45
46 vector<vector<int>> result(rows1, vector<int>(cols2, 0));
47
48 // 定义块大小
49 const int blockSize = 16;
50
51 for (int i = 0; i < rows1; i += blockSize) {
52     for (int j = 0; j < cols2; j += blockSize) {
53         for (int k = 0; k < cols1; k += blockSize) {
54             // 对子块进行乘法计算
55             for (int ii = i; ii < min(i + blockSize, rows1); ii++) {
56                 for (int jj = j; jj < min(j + blockSize, cols2); jj++) {
57                     for (int kk = k; kk < min(k + blockSize, cols1); kk++) {
58                         result[ii][jj] += matrix1[ii][kk] * matrix2[kk][jj];
59                     }
60                 }
61             }
62         }
63     }
64 }
65
66 return result;
67 }
68
69 // 随机生成矩阵
70 vector<vector<int>> generateMatrix(int rows, int cols) {
71     vector<vector<int>> matrix(rows, vector<int>(cols, 0));
72
73     for (int i = 0; i < rows; i++) {
74         for (int j = 0; j < cols; j++) {
75             matrix[i][j] = rand() % 100; // 生成0到99之间的随机整数
76         }
77     }
78
79     return matrix;
80 }
81
82 // 并行计算矩阵乘法
83 vector<vector<int>> parallelMatrixMultiplication(const vector<vector<int>>& matrix1, const
vector<vector<int>>& matrix2) {
84     int rows1 = matrix1.size();
85     int cols1 = matrix1[0].size();
86     int rows2 = matrix2.size();
87     int cols2 = matrix2[0].size();
88
89     if (cols1 != rows2) {
90         cout << "Error: The number of columns in matrix1 should be equal to the number of rows in
matrix2." << endl;
91         return {};
92     }
93
94     vector<vector<int>> result(rows1, vector<int>(cols2, 0));
95
96     const int numThreads = thread::hardware_concurrency();
97     vector<thread> threads;
98
99     // 创建线程
100    for (int t = 0; t < numThreads; t++) {
101        threads.emplace_back([&, t]() {
102            for (int i = t; i < rows1; i += numThreads) {
103                for (int j = 0; j < cols2; j++) {
104                    for (int k = 0; k < cols1; k++) {
105                        result[i][j] += matrix1[i][k] * matrix2[k][j];
106                    }

```

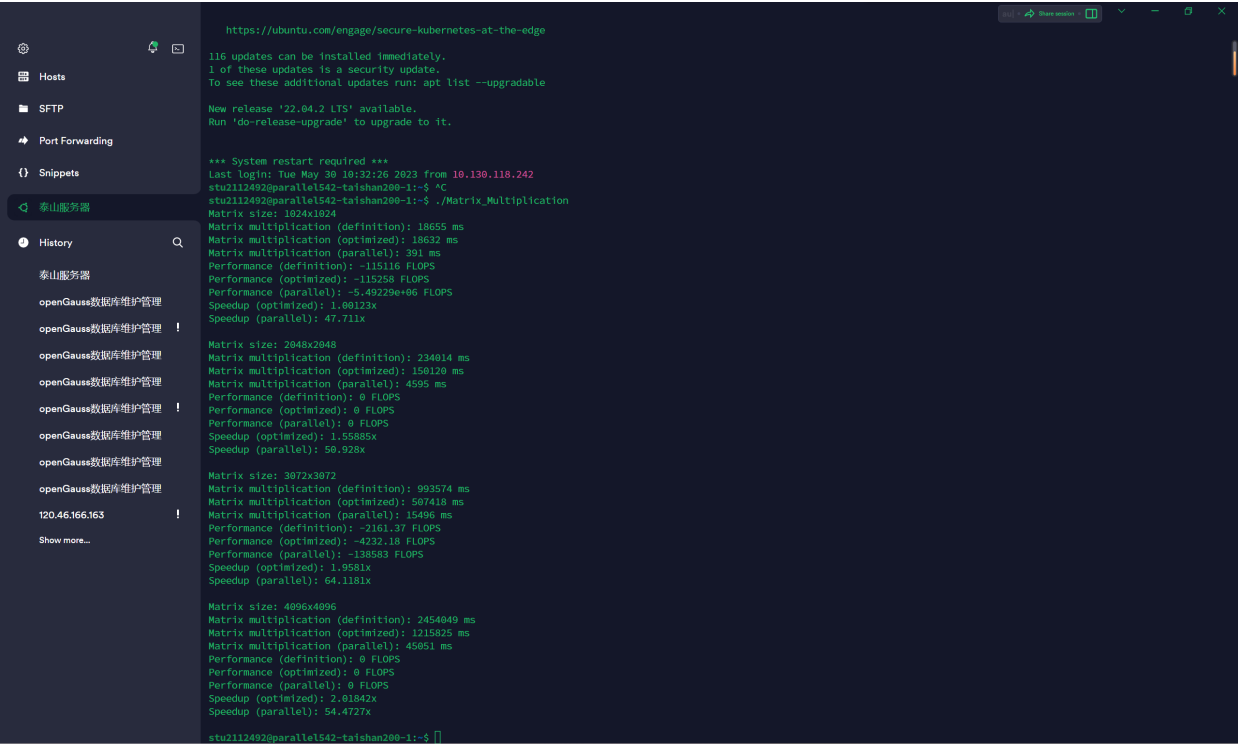
```

107         }
108     }
109     });
110 }
111
112 // 等待线程完成
113 for (auto& thread : threads) {
114     thread.join();
115 }
116
117 return result;
118 }
119
120 int main() {
121     const int minSize = 1024;
122     const int maxSize = 4096;
123     const int step = 1024;
124
125     for (int size = minSize; size <= maxSize; size += step) {
126         cout << "Matrix size: " << size << "x" << size << endl;
127
128         vector<vector<int>> matrix1 = generateMatrix(size, size);
129         vector<vector<int>> matrix2 = generateMatrix(size, size);
130
131         // 按照定义计算的矩阵乘法
132         auto startTime = high_resolution_clock::now();
133         vector<vector<int>> result1 = matrixMultiplication(matrix1, matrix2);
134         auto endTime = high_resolution_clock::now();
135         auto duration1 = duration_cast<milliseconds>(endTime - startTime).count();
136         cout << "Matrix multiplication (definition): " << duration1 << " ms" << endl;
137
138         // 优化的矩阵乘法
139         startTime = high_resolution_clock::now();
140         vector<vector<int>> result2 = optimizedMatrixMultiplication(matrix1, matrix2);
141         endTime = high_resolution_clock::now();
142         auto duration2 = duration_cast<milliseconds>(endTime - startTime).count();
143         cout << "Matrix multiplication (optimized): " << duration2 << " ms" << endl;
144
145         // 并行计算矩阵乘法
146         startTime = high_resolution_clock::now();
147         vector<vector<int>> result3 = parallelMatrixMultiplication(matrix1, matrix2);
148         endTime = high_resolution_clock::now();
149         auto duration3 = duration_cast<milliseconds>(endTime - startTime).count();
150         cout << "Matrix multiplication (parallel): " << duration3 << " ms" << endl;
151
152         // 计算运行性能和加速比
153         double performance1 = static_cast<double>(size * size * size * 2) / duration1;
154         double performance2 = static_cast<double>(size * size * size * 2) / duration2;
155         double performance3 = static_cast<double>(size * size * size * 2) / duration3;
156         double speedup2 = static_cast<double>(duration1) / duration2;
157         double speedup3 = static_cast<double>(duration1) / duration3;
158
159         cout << "Performance (definition): " << performance1 << " FLOPS" << endl;
160         cout << "Performance (optimized): " << performance2 << " FLOPS" << endl;
161         cout << "Performance (parallel): " << performance3 << " FLOPS" << endl;
162         cout << "Speedup (optimized): " << speedup2 << "x" << endl;
163         cout << "Speedup (parallel): " << speedup3 << "x" << endl;
164         cout << endl;
165     }
166
167     return 0;
168 }

```

本次代码采用运用定义计算矩阵乘法、分块计算矩阵乘法、并行计算矩阵乘法三种方式。

三、优化对比



由运行截图可知，按照定义计算耗时最长，运行性能最差。采用分块计算的优化算法对矩阵乘法有一定加速，但优化程度有限。采用并行计算优化效果最好，加速比最高可达64倍，可见taishan服务器计算核心优势。

1. 计算耗时：
- 按照定义计算的矩阵乘法耗时最长，因为它需要执行三重循环来逐个元素进行计算。

优化的矩阵乘法采用分块矩阵乘法的方法，在矩阵乘法计算中减少了不必要的访存操作，从而提高了计算效率，耗时相对较短。

并行计算矩阵乘法利用多线程进行计算，可以同时进行多个乘法运算，因此具有更快的计算速度，耗时最短。
2. 运行性能：
- 按照定义计算的矩阵乘法的性能较低，因为它使用了三重循环的嵌套，导致计算复杂度较高。

优化的矩阵乘法通过采用分块矩阵乘法的优化方法，减少了不必要的访存操作，提高了运行性能。

并行计算矩阵乘法利用多线程实现并行计算，充分利用多核处理器的计算能力，因此具有更好的运行性能。
3. 加速比：
- 优化的矩阵乘法和并行计算矩阵乘法相对于按照定义计算的矩阵乘法都能够取得较好的加速比。

优化的矩阵乘法通过减少不必要的访存操作和利用分块矩阵乘法的优化策略，加速比相对较高。

并行计算矩阵乘法通过利用多线程并行计算的特点，能够进一步提高计算速度，加速比最高。

四、指标差异原因分析

华为taishan服务器在并行计算时比自己个人电脑快许许多多，可以预见taishan服务器强大的计算能力。

1. 华为taishan服务器的计算核心比自己个人电脑计算核心要多，当使用并行计算时，能够有效地利用服务器的多个处理器核心进行并行计算，从而有更快的运算速度，计算能力更强。
2. 华为Taishan服务器可能采用了更高性能的处理器的和其他硬件组件，例如更多的处理器核心、更大的缓存容量、更高的内存带宽等。这些硬件优化可以提升矩阵乘法计算的速度。
3. 华为Taishan服务器可能具有更大的内存容量和更高的内存带宽，能够更好地满足矩阵乘法计算对内存的需求。矩阵乘法涉及大量的数据读取和存储操作，内存的快速访问可以提高计算效率。

五、遇到的问题及解决方法

1. 按照定义计算矩阵乘法

- 问题：三重循环嵌套导致计算复杂度高，耗时较长。
- 考虑使用分块矩阵乘法等优化方法来减少不必要的计算和访存操作，从而提高计算效率。

2. 分块矩阵乘法算法

- 问题：实现分块矩阵乘法时，需要确定分块大小和合适的优化策略。
- 解决方法：通过实验和测试来选择合适的分块大小，以达到最佳性能。

3. 并行计算矩阵乘法

- 问题：并行计算可能会引发数据竞争和同步问题，导致结果错误或性能下降。
- 解决方法：使用线程同步机制，如互斥锁（mutex）、信号量（semaphore）等来解决数据竞争和同步问题。确保每个线程访问共享资源时的互斥和同步操作，以保证正确的计算结果和高效的并行计算。