

《汇编语言》大作业 2-矩阵乘法的优化

在本次作业中, 我一共采用了共 5 中方法来优化计算, 下面逐个讲解, 有关运行时间的计算, 我在本地运行时发现每次运算的时间方差较大, 于是我擦爱用了计算 7 次, 去掉其中的最大值和最小值, 求平均的办法来相对严谨地计算运行时间

1 Python 原始代码

```
1 import numpy as np
2 import time
3
4 N = 4096
5 A = np.random.randint(0, 10, (N, N), dtype=np.int32)
6 B = np.random.randint(0, 10, (N, N), dtype=np.int32)
7 C = np.zeros((N, N), dtype=np.int32)
8
9 start = time.time()
10 for i in range(N):
11     for j in range(N):
12         for k in range(N):
13             C[i, j] += A[i, k] * B[k, j]
14 end = time.time()
15 print("Python naive time:", end - start)
```

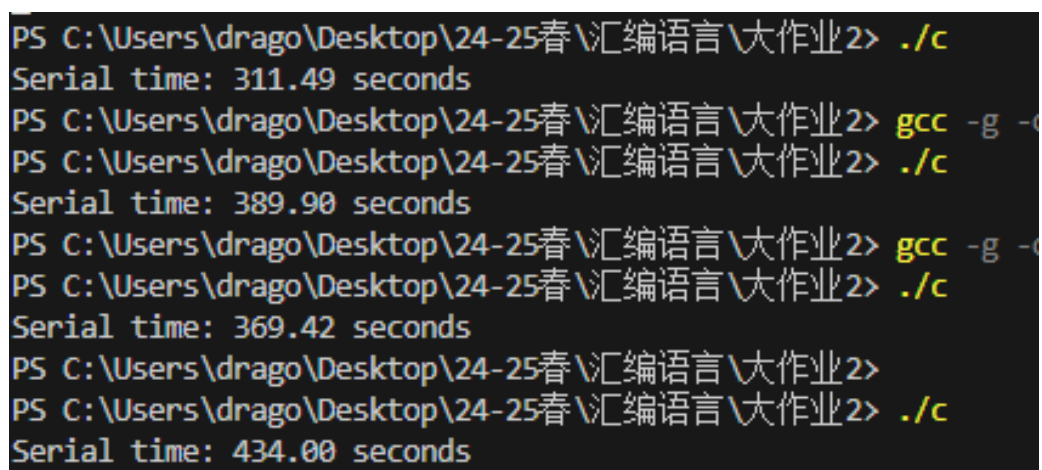
这是最原始的 python 代码, 由于这次矩阵的大小为 4096, 所以计算量非常大, 直接运行会花费很长时间, 大约花费了 35 小时, 大约 126000 秒 (这个数字是根据小规模矩阵推算得到, 可能并不准确)

2 线性计算的 C 代码

```
1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <time.h>
4
5 #define N 4096
6
7 //最基本的矩阵乘法函数, 逐个行列相乘
8 void matrixmultiply(int n, int **A, int **B, int **C) {
9     for (int i = 0; i < n; ++i)
10         for (int j = 0; j < n; ++j) {
11             int sum = 0;
12             for (int k = 0; k < n; ++k)
13                 sum += A[i][k] * B[k][j];
```

```
14     C[i][j] = sum;
15 }
16 }
17
18 int main() {
19     //为ABC分配内存空间
20     int **A = malloc(N * sizeof(int*));
21     int **B = malloc(N * sizeof(int*));
22     int **C = malloc(N * sizeof(int*));
23     for (int i = 0; i < N; ++i) {
24         A[i] = malloc(N * sizeof(int));
25         B[i] = malloc(N * sizeof(int));
26         C[i] = malloc(N * sizeof(int));
27         for (int j = 0; j < N; ++j) {
28             A[i][j] = rand() % 10;
29             B[i][j] = rand() % 10;
30         }
31     }
32     clock_t start = clock();
33     matrixmultiply(N, A, B, C);
34     clock_t end = clock();
35     printf("Serial time: %.2f seconds\n", (double)(end - start) / CLOCKS_PER_SEC);
36     //释放内存
37     for (int i = 0; i < N; ++i) {
38         free(A[i]); free(B[i]); free(C[i]);
39     }
40     free(A); free(B); free(C);
41 }
```

部分运行截图如下



```
PS C:\Users\drago\Desktop\24-25春\汇编语言\大作业2> ./c
Serial time: 311.49 seconds
PS C:\Users\drago\Desktop\24-25春\汇编语言\大作业2> gcc -g -c
PS C:\Users\drago\Desktop\24-25春\汇编语言\大作业2> ./c
Serial time: 389.90 seconds
PS C:\Users\drago\Desktop\24-25春\汇编语言\大作业2> gcc -g -c
PS C:\Users\drago\Desktop\24-25春\汇编语言\大作业2> ./c
Serial time: 369.42 seconds
PS C:\Users\drago\Desktop\24-25春\汇编语言\大作业2>
PS C:\Users\drago\Desktop\24-25春\汇编语言\大作业2> ./c
Serial time: 434.00 seconds
```

图 1: 线性计算的 C 代码部分运行截图

以下是 7 次运行中选取的 5 次有效时间的表格

运行编号	1	2	3	4	5
运行时间 (s)	311.50	389.90	369.42	401.27	352.89
平均时间 (s)	364.99				

这个程序是很平凡的想法, 在此不多赘述

3 矩阵分块计算

由于之后的几个哈按树中 main 函数的差意不大, 所以在此将相同的 main 函数和宏定义省去

```

1  #define BLOCK 128
2
3  void matrixmultiply(int n, int **A, int **B, int **C) {
4      for (int ii = 0; ii < n; ii += BLOCK)
5          for (int jj = 0; jj < n; jj += BLOCK)
6              for (int kk = 0; kk < n; kk += BLOCK)
7                  for (int i = ii; i < ii + BLOCK && i < N; ++i)
8                      for (int j = jj; j < jj + BLOCK && j < N; ++j) {
9                          int sum = 0;
10                         for (int k = kk; k < kk + BLOCK && k < N; ++k)
11                             sum += A[i][k] * B[k][j];
12                         C[i][j] += sum;
13                     }
14 }
```

在这个代码中我们可以看到, 我们将整个大的矩阵拆分为几个小的块, 每次计算一个小块的矩阵, 这样可以减少缓存未命中的次数, 提高计算速度, 这个 BLOCK 的大小是根据我们的 cache 自定的, 所以这个 BLOCK 的大小要根据不同的输入规模和不同的计算机进行调整, 在我的计算机上, 我测试了几个不同的 BLOCK 值 (16,32,64,256), 我发现 128 和 256 的加速效果最好, 而且二者相差不多, 我选取了 128

以下是部分代码运行截图

```
PS C:\Users\drago\Desktop\24-25春\汇编语言\大作业2> ./blocking
blocking time: 125.38 seconds
PS C:\Users\drago\Desktop\24-25春\汇编语言\大作业2>
PS C:\Users\drago\Desktop\24-25春\汇编语言\大作业2> ./blocking
blocking time: 129.72 seconds
PS C:\Users\drago\Desktop\24-25春\汇编语言\大作业2> ./blocking
blocking time: 127.62 seconds
PS C:\Users\drago\Desktop\24-25春\汇编语言\大作业2> ./blocking
blocking time: 123.56 seconds
```

图 2: 矩阵分块计算的 C 代码部分运行截图

以下是 7 次运行中选取的 5 次有效时间的表格

运行编号	1	2	3	4	5
运行时间 (s)	125.38	129.72	127.42	123.56	120.29
平均时间 (s)	125.27				

4 多核多线程加速

```
1 #include <omp.h>
2
3 void matrixmultiply(int n, int **A, int **B, int **C) {
4     #pragma omp parallel for
5     for (int i = 0; i < n; ++i)
6         for (int j = 0; j < n; ++j) {
7             int sum = 0;
8             for (int k = 0; k < n; ++k)
9                 sum += A[i][k] * B[k][j];
10            C[i][j] = sum;
11        }
12 }
```

我们的并行多线程中, 对于 C 程序没有改变, 就是多了如下指令

```
1 #pragma omp parallel for
```

这条指令的作用是告诉编译器创建一组等于我计算机核心数的线程, 并且将接下来的 for 循环中的迭代空间自动划分给这些线程, 实际上就是将串行的运作转为多个线程并行的工作, 我们在命令

```
1 gcc -S -fopenmp -mavx2 -mfma -o openmp.s openmp.c
```

下所得到的.s 文件中可以看到如下命令

```
1 call omp_get_num_threads
2 call omp_get_thread_num
```

代表调用多核多线程来解决串行的循环

以下是部分代码运行截图



```
PS C:\Users\drago\Desktop\24-25春\汇编语言\大作业2> ./openmp
OpenMP time: 24.89 seconds
PS C:\Users\drago\Desktop\24-25春\汇编语言\大作业2> ./openmp
OpenMP time: 41.01 seconds
PS C:\Users\drago\Desktop\24-25春\汇编语言\大作业2> ./openmp
OpenMP time: 37.82 seconds
PS C:\Users\drago\Desktop\24-25春\汇编语言\大作业2> ./openmp
OpenMP time: 30.51 seconds
PS C:\Users\drago\Desktop\24-25春\汇编语言\大作业2> ./openmp
OpenMP time: 35.78 seconds
PS C:\Users\drago\Desktop\24-25春\汇编语言\大作业2> ./openmp
OpenMP time: 32.45 seconds
PS C:\Users\drago\Desktop\24-25春\汇编语言\大作业2> █
```

图 3: 多核多线程加速的 C 代码部分运行截图

以下是 7 次运行中选取的 5 次有效时间的表格

运行编号	1	2	3	4	5
运行时间 (s)	37.82	30.51	35.78	32.45	31.98
平均时间 (s)	33.71				

5 SIMD 指令集和多核多线程加速

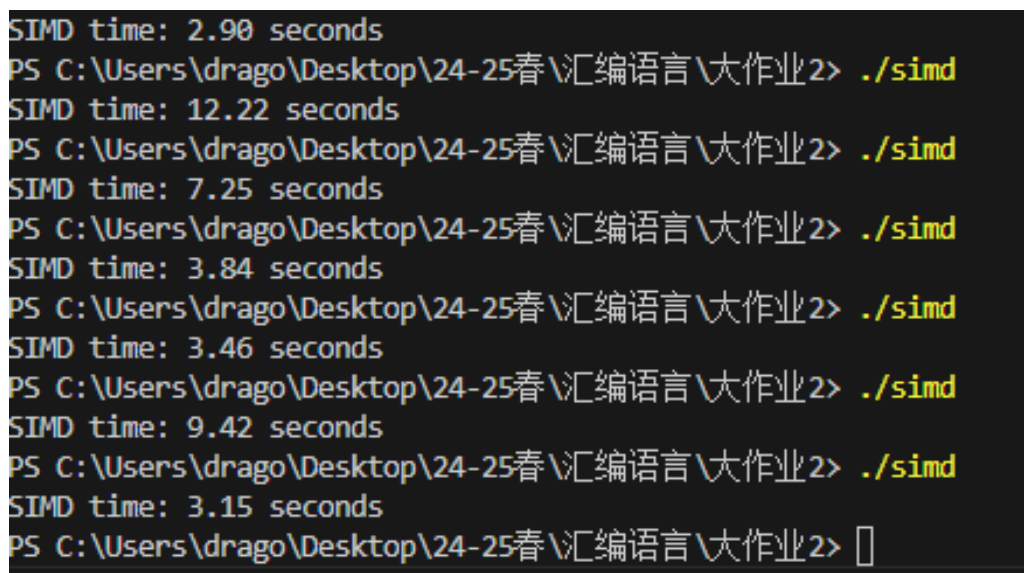
```
1 #include <immintrin.h>
2 #include <stdio.h>
3 #include <stdlib.h>
4 #include <time.h>
5 #include <omp.h>
6
7 #define N 4096
8
9 void matrixmultiply(int n, float **A, float **B, float **C) {
10     #pragma omp parallel for
11     for (int i = 0; i < n; ++i)
12         for (int j = 0; j < n; ++j) {
13             __m256 sum = _mm256_setzero_ps();
14             for (int k = 0; k < N; k += 8) {
15                 __m256 a = _mm256_loadu_ps(&A[i][k]);
```

```
16     __m256 b = _mm256_loadu_ps(&B[k][j]);
17     sum = _mm256_fmadd_ps(a, b, sum);
18 }
19 float tmp[8];
20 _mm256_storeu_ps(tmp, sum);
21 float total = 0;
22 for (int x = 0; x < 8; ++x) total += tmp[x];
23 C[i][j] = total;
24 }
25 }
26
27 int main() {
28     float **A = (float**)malloc(N * sizeof(float*));
29     float **B = (float**)malloc(N * sizeof(float*));
30     float **C = (float**)malloc(N * sizeof(float*));
31     for (int i = 0; i < N; ++i) {
32         A[i] = (float*)malloc(N * sizeof(float));
33         B[i] = (float*)malloc(N * sizeof(float));
34         C[i] = (float*)calloc(N, sizeof(float));
35         for (int j = 0; j < N; ++j) {
36             A[i][j] = (float)(rand() % 10);
37             B[i][j] = (float)(rand() % 10);
38         }
39     }
40     clock_t start = clock();
41     matrixmultiply(N, A, B, C);
42     clock_t end = clock();
43     printf("SIMD time: %.2f seconds\n", (double)(end - start) / CLOCKS_PER_SEC);
44
45     for (int i = 0; i < N; ++i) {
46         free(A[i]); free(B[i]); free(C[i]);
47     }
48     free(A); free(B); free(C);
49
50     return 0;
51 }
```

在这里由于使用了 SIMD 指令集, 因此为我们改用了 float 型, 其中比较关键的是标识符 `__m256`, 这个代表我们使用了 256 位的向量寄存器, 可以同时处理 8 个 float 类型的数据, 从而实现数据的并行处理. 我们应该用碧嘉澳查昂的寄存器来尽可能的用我们并行所得到的提升来覆盖生成寄存器的开销, 这样才能得到更好的性能提升, 在这里我们使用了 `_mm256_loadu_ps` 和 `_mm256_storeu_ps` 来加载和存储数据, 同时使用 `_mm256_fmadd_ps` 来进行乘加操作, 这样可以减少指令的数量, 提高计算效率. 在汇编中我们发现了形如

```
1  vmovups (%rax), %ymm0
2  vmovups %ymm0, 128(%rbx)
3  vmovups 160(%rbx), %ymm0
4  vmovups %ymm0, 64(%rbx)
5  vmovups 128(%rbx), %ymm0
6  vmovups %ymm0, 32(%rbx)
7  vmovups 192(%rbx), %ymm0
8  vmovups %ymm0, (%rbx)
9  vmovups 32(%rbx), %ymm1
10 vmovups (%rbx), %ymm0
11 vfmadd231ps 64(%rbx), %ymm1, %ymm0
```

的关键汇编指令, 我们并行地同时处理了很多组数据, 大大提升了效率
以下是部分代码运行截图



```
SIMD time: 2.90 seconds
PS C:\Users\drago\Desktop\24-25春\汇编语言\大作业2> ./simd
SIMD time: 12.22 seconds
PS C:\Users\drago\Desktop\24-25春\汇编语言\大作业2> ./simd
SIMD time: 7.25 seconds
PS C:\Users\drago\Desktop\24-25春\汇编语言\大作业2> ./simd
SIMD time: 3.84 seconds
PS C:\Users\drago\Desktop\24-25春\汇编语言\大作业2> ./simd
SIMD time: 3.46 seconds
PS C:\Users\drago\Desktop\24-25春\汇编语言\大作业2> ./simd
SIMD time: 9.42 seconds
PS C:\Users\drago\Desktop\24-25春\汇编语言\大作业2> ./simd
SIMD time: 3.15 seconds
PS C:\Users\drago\Desktop\24-25春\汇编语言\大作业2> □
```

图 4: SIMD 指令集加速的 C 代码部分运行截图

以下是 7 次运行中选取的 5 次有效时间的表格

运行编号	1	2	3	4	5
运行时间 (s)	7.25	3.84	3.46	9.42	3.15
平均时间 (s)	5.82				

我们发现 SIMD 程序所得到的运行时间的方差比较大, 个人猜测是 simd 指令这种运算方法和多核多线程的方式对于输入的数据比较敏感

6 其他加速方法-Strassen 算法

Strassen 算法的核心思想是将大矩阵分解为小矩阵, 来通过减少乘法次数来降低计算复杂度, 最终可以达到 $O(n^{\log_2 7}) \approx O(n^{2.81})$ 的复杂度, 但是代码实现很复杂, 代码如下

```
1  #define THRESHOLD 64 // 小于此尺寸时用普通乘法
2
3  void classic_multiply(int n, int **A, int **B, int **C) {
4      for (int i=0; i<n; ++i)
5          for (int j=0; j<n; ++j)
6              for (int k=0; k<n; ++k)
7                  C[i][j] += A[i][k] * B[k][j];
8  }
9
10 void add(int n, int **A, int **B, int **C) {
11     for (int i=0; i<n; ++i)
12         for (int j=0; j<n; ++j)
13             C[i][j] = A[i][j] + B[i][j];
14 }
15
16 void sub(int n, int **A, int **B, int **C) {
17     for (int i=0; i<n; ++i)
18         for (int j=0; j<n; ++j)
19             C[i][j] = A[i][j] - B[i][j];
20 }
21
22 void strassen(int n, int **A, int **B, int **C) {
23     if (n <= THRESHOLD) {
24         classic_multiply(n, A, B, C);
25         return;
26     }
27     int k = n/2;
28     // Allocate submatrices
29     int **A11 = malloc(k*sizeof(int*));
30     int **A12 = malloc(k*sizeof(int*));
31     int **A21 = malloc(k*sizeof(int*));
32     int **A22 = malloc(k*sizeof(int*));
33     int **B11 = malloc(k*sizeof(int*));
34     int **B12 = malloc(k*sizeof(int*));
35     int **B21 = malloc(k*sizeof(int*));
36     int **B22 = malloc(k*sizeof(int*));
37     int **C11 = malloc(k*sizeof(int*));
38     int **C12 = malloc(k*sizeof(int*));
39     int **C21 = malloc(k*sizeof(int*));
40     int **C22 = malloc(k*sizeof(int*));
41     for (int i=0; i<k; ++i) {
42         A11[i] = A[i];
```



```
43     A12[i] = A[i] + k;
44     A21[i] = A[i+k];
45     A22[i] = A[i+k] + k;
46     B11[i] = B[i];
47     B12[i] = B[i] + k;
48     B21[i] = B[i+k];
49     B22[i] = B[i+k] + k;
50     C11[i] = C[i];
51     C12[i] = C[i] + k;
52     C21[i] = C[i+k];
53     C22[i] = C[i+k] + k;
54 }
55 // Allocate temp matrices
56 int **M1 = malloc(k*sizeof(int*));
57 int **M2 = malloc(k*sizeof(int*));
58 int **M3 = malloc(k*sizeof(int*));
59 int **M4 = malloc(k*sizeof(int*));
60 int **M5 = malloc(k*sizeof(int*));
61 int **M6 = malloc(k*sizeof(int*));
62 int **M7 = malloc(k*sizeof(int*));
63 int **T1 = malloc(k*sizeof(int*));
64 int **T2 = malloc(k*sizeof(int*));
65 for (int i=0; i<k; ++i) {
66     M1[i] = calloc(k, sizeof(int));
67     M2[i] = calloc(k, sizeof(int));
68     M3[i] = calloc(k, sizeof(int));
69     M4[i] = calloc(k, sizeof(int));
70     M5[i] = calloc(k, sizeof(int));
71     M6[i] = calloc(k, sizeof(int));
72     M7[i] = calloc(k, sizeof(int));
73     T1[i] = calloc(k, sizeof(int));
74     T2[i] = calloc(k, sizeof(int));
75 }
76 // M1 = (A11 + A22)*(B11 + B22)
77 add(k, A11, A22, T1);
78 add(k, B11, B22, T2);
79 strassen(k, T1, T2, M1);
80 // M2 = (A21 + A22)*B11
81 add(k, A21, A22, T1);
82 strassen(k, T1, B11, M2);
83 // M3 = A11*(B12 - B22)
84 sub(k, B12, B22, T2);
85 strassen(k, A11, T2, M3);
86 // M4 = A22*(B21 - B11)
87 sub(k, B21, B11, T2);
88 strassen(k, A22, T2, M4);
89 // M5 = (A11 + A12)*B22
90 add(k, A11, A12, T1);
91 strassen(k, T1, B22, M5);
92 // M6 = (A21 - A11)*(B11 + B12)
```

```
93     sub(k, A21, A11, T1);
94     add(k, B11, B12, T2);
95     strassen(k, T1, T2, M6);
96     // M7 = (A12 - A22)*(B21 + B22)
97     sub(k, A12, A22, T1);
98     add(k, B21, B22, T2);
99     strassen(k, T1, T2, M7);
100    // C11 = M1 + M4 - M5 + M7
101    for (int i=0; i<k; ++i)
102        for (int j=0; j<k; ++j)
103            C11[i][j] = M1[i][j] + M4[i][j] - M5[i][j] + M7[i][j];
104    // C12 = M3 + M5
105    for (int i=0; i<k; ++i)
106        for (int j=0; j<k; ++j)
107            C12[i][j] = M3[i][j] + M5[i][j];
108    // C21 = M2 + M4
109    for (int i=0; i<k; ++i)
110        for (int j=0; j<k; ++j)
111            C21[i][j] = M2[i][j] + M4[i][j];
112    // C22 = M1 - M2 + M3 + M6
113    for (int i=0; i<k; ++i)
114        for (int j=0; j<k; ++j)
115            C22[i][j] = M1[i][j] - M2[i][j] + M3[i][j] + M6[i][j];
116    // Free allocated memory
117    for (int i=0; i<k; ++i) {
118        free(M1[i]); free(M2[i]); free(M3[i]); free(M4[i]);
119        free(M5[i]); free(M6[i]); free(M7[i]);
120        free(T1[i]); free(T2[i]);
121    }
122    free(A11); free(A12); free(A21); free(A22);
123    free(B11); free(B12); free(B21); free(B22);
124    free(C11); free(C12); free(C21); free(C22);
125    free(M1); free(M2); free(M3); free(M4); free(M5); free(M6); free(M7);
126    free(T1); free(T2);
127 }
```

这个算法的实现是网上查阅资料所得到的, 核心思想就是利用矩阵分块, 来一定程度上减少乘法的次数

以下是部分代码运行截图

7.2 加速瓶颈

综合我们上面的几种加速方法的原理和效果, 我认为有如下加速瓶颈

1. 矩阵分块的大小选择对性能影响较大, 需要根据具体硬件进行调优
2. 在多核多线程中, 线程的个数以及线程之间的通信和同步开销会影响性能
3. SIMD 指令集能否用更长的寄存器?