# Lab 1: Matrix Multiplication

姜海天 19307110022

This report is to record the performance optimization of an matrix multiplication problem.

In this problem, we are going to calculate $\alpha \cdot AB + \beta \cdot C$, where $A \in \mathbb{R}^{m \times k}, B \in \mathbb{R}^{k \times n}, C \in \mathbb{R}^{m \times n}$, and store the answer in $C$.

## 0.1 Testing Environment

All the experiments are conducted on a server equipped with dual 64-core, 128-thread `AMD EPYC 7742 64-Core` CPU. The peak double-precision performance for a single CPU is 3.48 Tflops. L1d cache: 4MiB, L1i cache: 4MiB, L2 cache: 64MiB, L3 cache: 512MiB

The compiler used for the experiments is `clang++ 14.0.6`.

All the test uses `./executable 4096 4096 4096` as the test command.

## 0.2 Computation Amount

If use the algorithm stated in the naive implementation, there will be $mn + 3mnk$ floating-point operations. If $m = n = k = 4096$ in my test case, there is 206.2G floating-point operations. So the theoretical minimum computing time is $206.2G/(2 * 3.48T/s) = 30ms$.

## 0.3 Main Results

| Optimization | Running Time(s) | Relative Speedup | Absolute Speedup | Percent of Peak |
| --- | --- | --- | --- | --- |
| Naive | 2246.3 | 1 | 1 | 0.0013% |
| + interchange loops | 224.76 | 10 | 10 | 0.013% |
| + optimization flags | 38.85 | 5.78 | 57.8 | 0.077% |
| + reduced computation | 32.36 | / | 69.4 | 0.062% |
| Parallel loops | 1.40 | 4.2 | 1609 | 2.15% |
| +tiling | 0.34 | 27.8 | 6691 | 8.94% |
| +compiler vectorization | 0.32 | 1.04 | 6941 | 9.27% |

# 1  Naive Implementation

We use the naive implementation without any compiler optimization as a baseline.

```
1   void student_gemm(int m, int n, int k,
2                     const double * A, const double * B, double * C,
3                     double alpha, double beta,
4                     int lda, int ldb, int ldc) {
5       for (int i = 0; i < m; i++) {
6           for (int j = 0; j < n; j++) {
7               C[i + j * ldc] *= beta;
8               for (int p = 0; p < k; p++) {
9                   C[i + j * ldc] += alpha * A[i + p * lda] * B[p + j * ldb];
10              }
11          }
12      }
13  }
```

The computing time for the naive method is **37min26s (2246252.1556ms)**.

Speedup: **1x**, percent of peak: **0.0013%**.

```
1   $ ~/opencilk/bin/clang++ gemm.cpp -o 1-naive
2   $ ./1-naive-1 4096 4096 4096
3   input: 4096 x 4096 x 4096
4   minimal time spent: 2246252.1556 ms
5   result: correct (err = 0.000000e+00)
```

The cache miss rate for the naive method: **4.7%** by the `valgrind` command.

## 2 Loop Order

Change the loop order from `i,j,p` in the naive implementation to `j,p,i` so that we have the locality for columns, since the matrices are stored in column major.

```cpp
void student_gemm(int m, int n, int k,
                  const double * A, const double * B, double * C,
                  double alpha, double beta,
                  int lda, int ldb, int ldc) {
    // compute β·C first
    for (int j = 0; j < n; ++j) {
        for (int i = 0; i < m; ++i) {
            C[i + j * ldc] *= beta;
        }
    }
    for (int j = 0; j < n; j++) {
        for (int p = 0; p < k; p++) {
            // put the whole column for C and A in cache
            for (int i = 0; i < m; i++) {
                C[i + j * ldc] += alpha * A[i + p * lda] * B[p + j * ldb];
            }
        }
    }
}
```

The computing time for the better loop order is **3min44s (224764.6316ms)**.

Speedup: **10x**, percent of peak: **0.013%**.

```
$ ~/opencilk/bin/clang++ gemm.cpp -o 2-loop-order
$ ./2-loop-order 4096 4096 4096
input: 4096 x 4096 x 4096
minimal time spent: 224764.6316 ms
result: correct (err = 0.000000e+00)
```

The cache miss rate for the naive method: **1.1%** by the `valgrind` command.

# 3 Compiler Optimization

The most aggressive optimization `-O3` got the best performance, and the running time is only **38.8s (38847.4110ms)** now.

Speedup: **57.8x**, percent of peak: **0.077%**.

```
1  $ ~/opencilk/bin/clang++ gemm.cpp -o 3-O1 -O1
2  $ ./3-O1 4096 4096 4096
3  input: 4096 x 4096 x 4096
4  minimal time spent: 44184.6329 ms
5  result: correct (err = 0.000000e+00)
6
7  $ ~/opencilk/bin/clang++ gemm.cpp -o 3-O2 -O2
8  $ ./3-O2 4096 4096 4096
9  input: 4096 x 4096 x 4096
10 minimal time spent: 39349.4487 ms
11 result: correct (err = 0.000000e+00)
12
13 $ ~/opencilk/bin/clang++ gemm.cpp -o 3-O3 -O3
14 $ ./3-O3 4096 4096 4096
15 input: 4096 x 4096 x 4096
16 minimal time spent: 38847.4110 ms
17 result: correct (err = 0.000000e+00)
```

# 4 Reduce Floating-point Operations

The current implementation calculates

$$c_{i,j} = \beta \cdot c_{i,j} + \sum_p \alpha \cdot a_{i,p} \cdot b_{p,j}$$

There are lots of unnecessary floating-pointoperations. The calculation with reduced floating-point operation is

$$c_{i,j} = \alpha \cdot \left( \frac{\beta}{\alpha} \cdot c_{i,j} + \sum_p a_{i,p} \cdot b_{p,j} \right)$$

Now the amount of floating-point operations have been reduced from $mn + 3mnk$ to $2mn + 2mnk$. However, it requires more read and write of memory. So there is trade-off.

This method will not be suitable for the further optimization aiming at reducing memory operations because it brings more memory operations natually.

Implementation:

```
1  void student_gemm(int m, int n, int k,
2                    const double * A, const double * B, double * C,
3                    double alpha, double beta,
4                    int lda, int ldb, int ldc) {
5      double ratio = beta / alpha;
6      // C ← (β/α)C
7      for (int j = 0; j < n; ++j) {
8          for (int i = 0; i < m; ++i) {
9              C[i + j * ldc] *= ratio;
10         }
11     }
12     // C ← C + AB
13     for (int j = 0; j < n; j++) {
14         for (int p = 0; p < k; p++) {
15             for (int i = 0; i < m; i++) {
16                 C[i + j * ldc] += A[i + p * lda] * B[p + j * ldb];
17             }
18         }
19     }
20     // C ← αC
21     for (int j = 0; j < n; ++j) {
22         for (int i = 0; i < m; ++i) {
23             C[i + j * ldc] *= alpha;
24         }
25     }
26 }
```

Experiment: The running time is **32.36s**.

Speedup: **69.4x**, percent of peak: **0.062%**. The percent of peak dropped because the amount of floating-point operations is reduced and more memory access is required.

```
1  $ ~/opencilk/bin/clang++ gemm.cpp -o 4-coeff-O3 -O3
2  $ ./4-coeff-O3 4096 4096 4096
3  input: 4096 x 4096 x 4096
4  minimal time spent: 32360.3723 ms
5  result: correct (err = 2.913225e-13)
```

# 5    Multi-Core Parallel Computation

Add parallelable for-loop for the reordered loop version. The version with excessive memory access doesn't perform better because there are too much write to memory, which is much slower than floating-point operations.

```cpp
void student_gemm(int m, int n, int k,
                  const double * A, const double * B, double * C,
                  double alpha, double beta,
                  int lda, int ldb, int ldc) {
    cilk_for (int j = 0; j < n; ++j) {
        for (int i = 0; i < m; ++i) {
            C[i + j * ldc] *= beta;
        }
    }
    cilk_for (int j = 0; j < n; j++) {
        for (int p = 0; p < k; p++) {
            for (int i = 0; i < m; i++) {
                C[i + j * ldc] += alpha * A[i + p * lda] * B[p + j * ldb];
            }
        }
    }
}
```

Experiment:   The running time is **1.4s**.

Speedup: **1609x**, percent of peak: **2.15%**.

```
$ ~/opencilk/bin/clang++ gemm.cpp -fopencilk -O3 -o 5-cilk-lo
$ ./5-cilk-lo 4096 4096 4096
input: 4096 x 4096 x 4096
minimal time spent: 1395.9282 ms
result: correct (err = 0.000000e+00)
```

## 6 Tiling

I added the tiling size `s` as an optional command line parameter. The default one is 128, which is the emperical best size.

Implementation:

```
1   // ...
2   int s = 128; // global default value for s
3   void student_gemm(int m, int n, int k,
4                       const double * A, const double * B, double * C,
5                       double alpha, double beta,
6                       int lda, int ldb, int ldc) {
7       cilk_for (int j = 0; j < n; ++j) {
8           for (int i = 0; i < m; ++i) {
9               C[i + j * ldc] *= beta;
10          }
11      }
12      cilk_for (int jh = 0; jh < n; jh += s) {
13          int J = MIN(n-jh, s);   // in case n % s ≠ 0
14          cilk_for (int ih = 0; ih < m; ih += s) {
15              int I = MIN(m-ih, s);   // in case m % s ≠ 0
16              for (int ph = 0; ph < k; ph += s) {
17                  int P = MIN(k-ph, s);   // in case p % s ≠ 0
18                  for (int jl = 0; jl < J; ++jl) {
19                      int j = jh + jl;
20                      for (int pl = 0; pl < P; ++pl) {
21                          int p = ph + pl;
22                          for (int il = 0; il < I; ++il) {
23                              int i = ih + il;
24                              C[i + j * ldc] += alpha * A[i + p * lda] * B[p + j * ldb];
25                          }
26                      }
27                  }
28              }
29          }
30      }
31  }
32  int main(int argc, const char * argv[]) {
33      if (argc ≠ 4 && argc ≠ 5) {
34          printf("Test usage: ./test m n k [s]\n");
35          exit(-1);
36      }
37      int m = atoi(argv[1]);
```

```
38      int n = atoi(argv[2]);
39      int k = atoi(argv[3]);
40      if (argc == 5) {
41          s = atoi(argv[4]);
42      }
43      // ...
44  }
```

Experiment: The best tiling step size is 128, and the running time is **0.3s(335.7088ms)**

Speedup: **6691x**, percent of peak: **8.94%**.

```
1   $ ~/opencilk/bin/clang++ gemm.cpp -fopencilk -O3 -o 6-tiling
2
3   $ ./6-tiling 4096 4096 4096 32
4   input: 4096 x 4096 x 4096
5   minimal time spent: 1455.8920 ms
6   result: correct (err = 0.000000e+00)
7
8   $ ./6-tiling 4096 4096 4096 64
9   input: 4096 x 4096 x 4096
10  minimal time spent: 708.2708 ms
11  result: correct (err = 0.000000e+00)
12
13  $ ./6-tiling 4096 4096 4096 128
14  input: 4096 x 4096 x 4096
15  minimal time spent: 335.7088 ms
16  result: correct (err = 0.000000e+00)
17
18  $ ./6-tiling 4096 4096 4096 256
19  input: 4096 x 4096 x 4096
20  minimal time spent: 554.6172 ms
21  result: correct (err = 0.000000e+00)
22
23  $ ./6-tiling 4096 4096 4096 512
24  input: 4096 x 4096 x 4096
25  minimal time spent: 1507.3377 ms
26  result: correct (err = 0.000000e+00)
```

# 7 Divide and Conquer

We only optimize for the case where the dimension is a power of 2 by recursively divide and conquer.

```
1  #define ISPOWER(x) ((x) & (-(x))) == (x)
2  void mm_base(int m, int n, int k, double alpha,
3              const double *A, const double *B, double *C,
4              int lda, int ldb, int ldc) {
5      for (int j = 0; j < n; j++) {
6          for (int p = 0; p < k; p++) {
7              for (int i = 0; i < m; i++) {
8                  C[i + j * ldc] += alpha * A[i + p * lda] * B[p + j * ldb];
9              }
10         }
11     }
12 }
13 void mm_dac(int m, int n, int k, double a,
14             const double *A, const double *B, double *C,
15             int ldA, int ldB, int ldC) {
16     assert(ISPOWER(m) && ISPOWER(n) && ISPOWER(k));
17     if (MIN(m, n) <= DAC_THOLD || k <= DAC_THOLD) {
18         mm_base(m, n, k, alpha, A, B, C, ldA, ldB, ldC);
19     } else {
20         int& R_A = m; int& R_B = k; int& R_C = m;
21         int& C_A = k; int& C_B = n; int& C_C = n;
22         #define X(M, r, c) (M + r * (R_ ## M / 2) + c * (ld ## M) * (C_ ## M / 2))
23         cilk_spawn mm_dac(m/2, n/2, k/2, a, X(A,0,0), X(B,0,0), X(C,0,0), ldA, ldB, ldC);
24         cilk_spawn mm_dac(m/2, n/2, k/2, a, X(A,0,0), X(B,0,1), X(C,0,1), ldA, ldB, ldC);
25         cilk_spawn mm_dac(m/2, n/2, k/2, a, X(A,1,0), X(B,0,0), X(C,1,0), ldA, ldB, ldC);
26                    mm_dac(m/2, n/2, k/2, a, X(A,1,0), X(B,0,1), X(C,1,1), ldA, ldB, ldC);
27         cilk_sync;
28         cilk_spawn mm_dac(m/2, n/2, k/2, a, X(A,0,1), X(B,1,0), X(C,0,0), ldA, ldB, ldC);
29         cilk_spawn mm_dac(m/2, n/2, k/2, a, X(A,0,1), X(B,1,1), X(C,0,1), ldA, ldB, ldC);
30         cilk_spawn mm_dac(m/2, n/2, k/2, a, X(A,1,1), X(B,1,0), X(C,1,0), ldA, ldB, ldC);
31                    mm_dac(m/2, n/2, k/2, a, X(A,1,1), X(B,1,1), X(C,1,1), ldA, ldB, ldC);
32         cilk_sync;
33     }
34 }
```

However, the result is that the program cannot suck all CPU power due to the synchronization, and the best result is **521.0503ms** for `DAC_THOLD=128`.

# 8    Compiler Flags

```
1  $ ~/opencilk/bin/clang++ gemm.cpp -fopencilk -O3 -o 8-avx -march=native -ffast-math
2  $ ./8-avx 4096 4096 4096
3  input: 4096 x 4096 x 4096
4  minimal time spent: 323.6070 ms
5  result: correct (err = 4.973799e-14)
```

Fast math doesn't guarantee the `err` to be 0, but it is faster.

The running time is **0.32s**.

Speedup: **6941x**, percent of peak: **9.27%**.