

General Matrix Multiplication Optimization

Apple Zhang

College of Computer Science and Software Engineering, Shenzhen University

December 5, 2022



Table of Contents

1 Optimization Methods

2 Experiments



深圳大学
SHENZHEN UNIVERSITY

Table of Contents

1 Optimization Methods

2 Experiments

Trivial Matrix Multiplication

$$C_{ij} = \sum_{r=1}^n A_{ir} B_{rj}. \quad (1)$$

```

1 void gemm_trivial(int m, int n, int k, double *a, int lda,
2                   double *b, int ldb,
3                   double *c, int ldc)
4 {
5     // TRIVIAL MATRIX MULTIPLICATION
6     for (int j = 0; j < n; j++) {
7         for (int i = 0; i < m; i++) {
8             for (int r = 0; r < k; r++) {
9                 C(i, j) += A(i, r) * B(r, j);
10            }
11        }
12    }
13 }

```



Simple cache optimization

Change the iteration order: $jir \rightarrow jri$.

```

1 void gemm_cache(int m, int n, int k, double *a, int lda,
2                 double *b, int ldb,
3                 double *c, int ldc)
4 {
5     // OPTIMIZE THE ORDER OF MATRIX MULTIPLICATION
6     for (int j = 0; j < n; j++) {
7         for (int r = 0; r < k; r++) {
8             register double brj = B(r, j);
9             for (int i = 0; i < m; i++) {
10                 C(i, j) += A(i, r) * brj;
11             }
12         }
13     }
14 }
```

Simple cache optimization

Scan the matrix in column order, and reduce the cache miss.

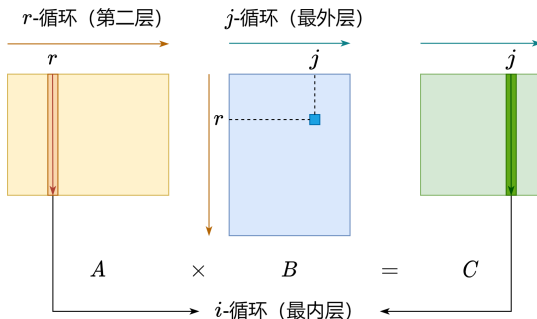


Figure: Simple cache optimization

Compute 1x4 vector

Avoid reading the same elements from matrix A repeatedly.

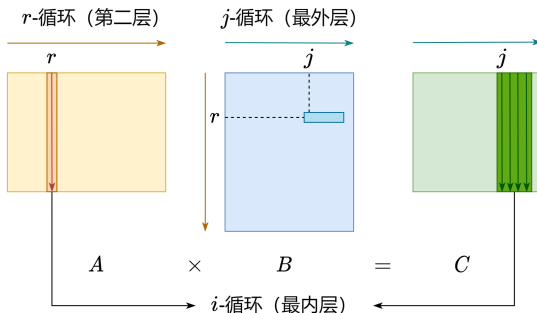


Figure: Compute four elements in each iteration.

Compute 4x4 block

Scan four rows of A and four columns of B .

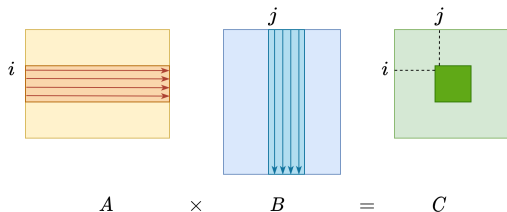


Figure: Compute a 4x4 matrix in each iteration.

Compute 4x4 block with AVX2

Consider SIMD, use AVX2 instructions for acceleration.

```

1      for (int r = 0; r < k; r++) {
2          ar_reg.v = _mm256_load_pd(&A(0, r));
3
4          b0_reg.v = _mm256_broadcast_sd(br0_ptr++);
5          b1_reg.v = _mm256_broadcast_sd(br1_ptr++);
6          b2_reg.v = _mm256_broadcast_sd(br2_ptr++);
7          b3_reg.v = _mm256_broadcast_sd(br3_ptr++);
8
9          c0_reg.v = _mm256_add_pd(c0_reg.v,
10                                 _mm256_mul_pd(ar_reg.v, b0_reg.v));
11         c1_reg.v = _mm256_add_pd(c1_reg.v,
12                                 _mm256_mul_pd(ar_reg.v, b1_reg.v));
13         c2_reg.v = _mm256_add_pd(c2_reg.v,
14                                 _mm256_mul_pd(ar_reg.v, b2_reg.v));
15         c3_reg.v = _mm256_add_pd(c3_reg.v,
16                                 _mm256_mul_pd(ar_reg.v, b3_reg.v));
17     }

```



Block matrices optimization

Since for matrix multiplication, we have

$$C = AB = \begin{bmatrix} A_1 & A_2 & \cdots & A_K \end{bmatrix} \begin{bmatrix} B_1 \\ B_2 \\ \vdots \\ B_K \end{bmatrix} \quad (2)$$

we can partition the matrices for better **spatial locality**.

Block matrices optimization

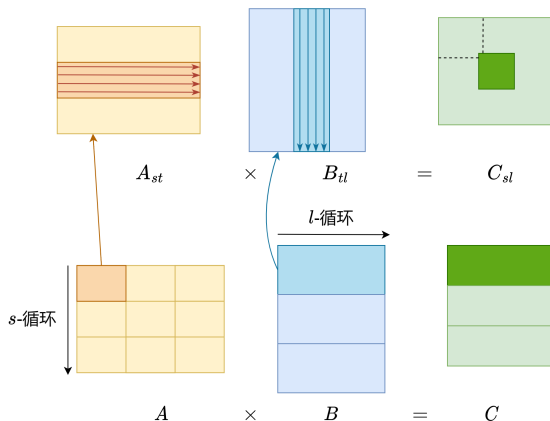


Figure: Block matrix multiplication

Reorganize the matrices

Rearrange the elements of the matrices at beginning for better spatial locality.

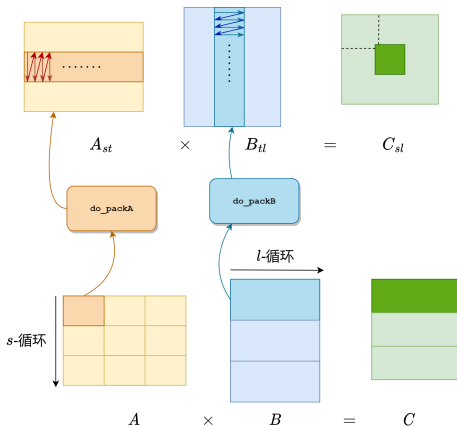


Figure: Rearrange the elements of matrices

Multi-thread with OpenMP

Use OpenMP to implement multi-thread.

```

1 void gemm_pack_memory(int m, int n, int k, double *a, int lda,
2                       double *b, int ldb,
3                       double *c, int ldc)
4 {
5     for (int l = 0; l < k; l += ikk) {
6         int lb = min(k - l, ikk);
7         #ifdef _OPENMP // MULTI-THREAD VIA OPENMP
8             #pragma omp parallel for schedule(dynamic)
9         #endif
10        for (int s = 0; s < m; s += ikm) {
11            int sb = min(m - s, ikm);
12            inner_kernel_packAB(sb, n, lb, &A(s, l), lda, &B(l, 0),
13                               ldb, &C(s, 0), ldc);
14        }
15    }

```

Table of Contents

1 Optimization Methods

2 Experiments

Evaluation

- OS: Linux (WSL)
- CPU: Intel i5-11400H

```
• (base) apple@LAPTOP-INDIFA6M:~/codes/csexp1-gemm$ ./gemm_test 1024 6
#thread -> 6/12
Trivial:      3.006579 sec
Cache:        0.328501 sec
1x4:          0.229344 sec
4x4:          0.510782 sec
4x4-avx:      0.327548 sec
Block:        0.120619 sec
packOpenMP:   0.032284 sec

max elem-wise error = 1.63425e-13
CORRECT
final speedup: 93.129073x
• (base) apple@LAPTOP-INDIFA6M:~/codes/csexp1-gemm$
```

Figure: The results of the program.



Evaluation

$$\text{GFLOPS} = \frac{2 \times 10^{-9} mnk}{\text{time}} \quad (3)$$

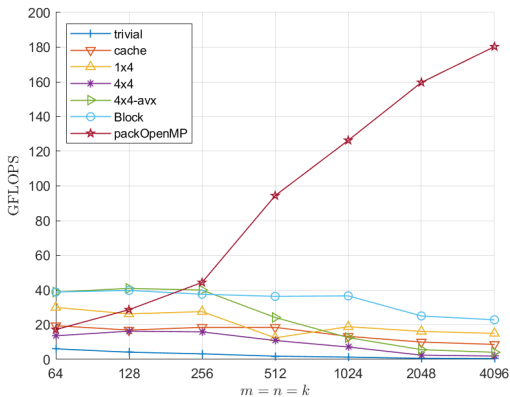


Figure: GFLOPS of each optimization.



深圳大学
SHENZHEN UNIVERSITY

Evaluation

But compared with Matlab (MKL)...

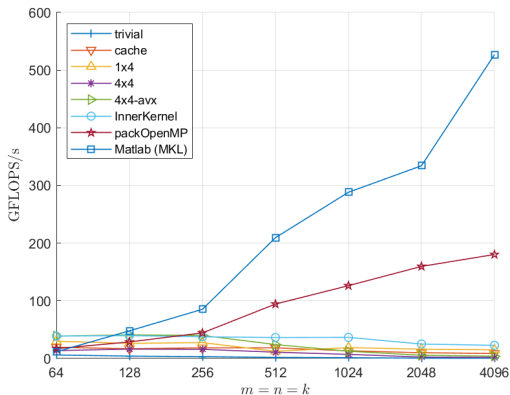


Figure: GFLOPS of each optimization.

Thanks!

My full codes can be viewed in
<https://github.com/Apple-Zhang/optimize-gemm>

