

体系结构Lab5

CPU

baseline核心代码：

▼

C++ | 复制代码

```
1 void gemm_baseline(float *A, float *B, float *C)
2 {
3     for (int i = 0; i < N; ++i)
4     {
5         for (int j = 0; j < N; ++j)
6         {
7             for (int k = 0; k < N; ++k)
8             {
9                 C[i * N + j] += A[i * N + k] * B[k * N + j];
10            }
11        }
12    }
13 }
```

AVX核心代码：

▼

C++ | 复制代码

```
1 void gemm_avx(float *A, float *B, float *C)
2 {
3     __m256 tmp;
4     for (int i = 0; i < N; ++i)
5     {
6         for (int j = 0; j < N; j += 8)
7         {
8             tmp = _mm256_loadu_ps(C + i * N + j);
9             for (int k = 0; k < N; ++k)
10            {
11                tmp = _mm256_add_ps(tmp, _mm256_mul_ps(_mm256_broadcast_ss
12                (A + i * N + k), _mm256_loadu_ps(B + k * N + j)));
13                _mm256_storeu_ps(C + i * N + j, tmp);
14            }
15        }
16    }
```

AVX_BLOCK核心代码：

```
1 void gemm_avx_one_block(int i0, int j0, int k0, float *A, float *B, float *C)
2 {
3     __m256 tmp;
4     for (int i = i0; i < i0 + BLOCKSIZE; ++i)
5     {
6         for (int j = j0; j < j0 + BLOCKSIZE; j += 8)
7         {
8             tmp = _mm256_loadu_ps(C + i * N + j); //取出16个float
9             for (int k = k0; k < k0 + BLOCKSIZE; ++k)
10            {
11                tmp = _mm256_add_ps(tmp, _mm256_mul_ps(_mm256_broadcast_ss
12                (A + i * N + k), _mm256_loadu_ps(B + k * N + j)));
13            }
14            _mm256_storeu_ps(C + i * N + j, tmp);
15        }
16    }
17
18 void gemm_avx_block(float *A, float *B, float *C)
19 {
20     for (int i = 0; i < N; i += BLOCKSIZE)
21     {
22         for (int j = 0; j < N; j += BLOCKSIZE)
23         {
24             for (int k = 0; k < N; k += BLOCKSIZE)
25             {
26                 gemm_avx_one_block(i, j, k, A, B, C);
27             }
28         }
29     }
30 }
```

不同规模矩阵结果

三种矩阵乘法实现见 `cpu.c` 源文件。

取 `BLOCKSIZE=1<<3` ,分别取 `N=1<n,n={6,7,8,9,10}` 三种实现的耗时及优化效率如下：

定义优化效率为 $\eta = \frac{t - t_0}{t_0} * 100\%$

```

1  N=64,BLOCKSIZE=8
2  0.00017
3  0.00006 65.06%
4  0.00003 79.52%
5  N=128,BLOCKSIZE=8
6  0.00145
7  0.00036 75.12%
8  0.00033 77.47%
9  N=256,BLOCKSIZE=8
10 0.01446
11 0.00372 74.30%
12 0.00341 76.39%
13 N=512,BLOCKSIZE=8
14 0.12085
15 0.03496 71.07%
16 0.01957 83.81%
17 N=1024,BLOCKSIZE=8
18 2.14926
19 0.53419 75.15%
20 0.27251 87.32%

```

不难看出AVX和AVX_BLOCK都能显著提高矩阵乘法速率，且AVX_BLOCK比AVX更高效，但提升率并不完全与规模成正比，即N增大，提升率不一定增大。当然，上述数据可能明确说明AVX和AVX_BLOCK能提高效率，至于提升了多少，上面的数值由于存在较大误差，只能作为简单参考。

原因分析

在单次乘法中，AVX能够同时计算8个单精度或4个双精度浮点数，大大减少了计算的时间。因此，使用AVX指令集能够提高单次矩阵乘法的计算性能。

采用分块AVX算法可以进一步提高性能，原因是分块AVX算法将矩阵乘法划分为多个小矩阵，每个小矩阵使用AVX指令集进行计算，这种算法在计算时可以使用CPU缓存，因此可以大大减少访问内存的次数，从而进一步减少运算时间。同时，由于分块操作实现了更好的数据局部性，能够更好的使矩阵缓存，从而能够更快的进行计算，提高矩阵计算的速度。

不同分块大小结果

接下来，取 $N=1 \ll 8$ ，分别取 $BLOCKSIZE=1 \ll b, b=\{3, 4, 5, 6, 7\}$ ，三种方式的耗时及优化率如下所示：（注意，BLOCKSIZE不能小于8，否则将导致结果计算错误）

理论上，BLOCKSIZE参数只会影响第三种方式（AVX_BLOCK），但从下面的数据可以看出，即使N不变，前两种乘法方式的耗时依旧有变化，这是因为计算机内部的复杂机制造成的无法消除的误差，即使采用多次计算取平均值的做法依然只能小幅减少该误差。

即便有上述所述的误差影响，我们依旧能看出，BLOCKSIZE对乘法计算性能的影响不完全是正相关的。在一定范围内，BLOCKSIZE的增大能提高性能，但当BLOCKSIZE进一步增大时，性能反而降低。就本次实验而言，在 `N=256` 的情况下，当 `BLOCKSIZE=32` 时，程序性能最佳。

▼

C++ | 复制代码

```
1  N=256,BLOCKSIZE=8
2  0.01430
3  0.00374 73.85%
4  0.00202 85.87%
5  N=256,BLOCKSIZE=16
6  0.01447
7  0.00233 83.93%
8  0.00169 88.31%
9  N=256,BLOCKSIZE=32
10 0.01485
11 0.00380 74.43%
12 0.00249 83.23%
13 N=256,BLOCKSIZE=64
14 0.01537
15 0.00226 85.27%
16 0.00191 87.58%
17 N=256,BLOCKSIZE=128
18 0.01517
19 0.00231 84.78%
20 0.00236 84.47%
```

更改N的大小，发现在 `N={1<n,n={8,9,10}}` 时，都是取 `BLOCKSIZE=32` 可以取得最佳结果，当N进一步增大时，耗时明显增大，实验成本增加，不在继续。不过目前测试结果已经表明，`BLOCKSIZE=32` 是一个不错的选择。

原因分析

当BLOCKSIZE 较小时，能够利用CPU缓存，减少了访问内存的频率，因此性能有所提高。但是，每个块中的计算量并不是特别稠密，所以它在利用AVX指令时也会浪费一些计算能力。

当BLOCKSIZE 适中时，采用的块的大小适中，可以利用一个256位向量并行计算32个乘积，显著提高了程序并行计算的效率。同时，块的大小也适合CPU缓存，可以更好的利用CPU缓存，从而减少了对内存的访问，同时具有较高的并行度和稠密的计算量，因此效果最好。

当BLOCKSIZE继续增加时，块中的计算量将变得更加稀疏，这样在利用AVX指令时也会浪费一些计算能力，同时由于块的数量减少，程序的并行度也会降低，因此性能并没有得到显著提升。

总之，BLOCKSIZE的取值太大也不能太小，但没有严格的确切数值可以保证BLOCKSIZE在该值时取得最佳结果，实际应用中，应采用实验测试的方式确定BLOCKSIZE的取值。

其它矩阵乘法优化方法

1. 预取数据。通过预取矩阵乘法中需要用到的数据到CPU缓存中，可以减少对内存的频繁访问，从而提高了程序的性能。
2. 循环重排序。通过调整循环顺序来优化空间访问局部性，从而提高效率
3. 循环展开。通过循环展开来加速循环过程。这一内容已在LAB4中详细探究
4. 使用多线程。对于特别大的数组，可以使用多线程计算分别计算各个部分，利用同步技术，互斥锁技术等保证结果正确性即可。

GPU

baseline

核心代码：

```
1 ▾ __global__ void gemm_baseline(float *A, float *B, float *C, int N)
2 ▾ {
3     int idx_x = blockIdx.x * blockDim.x + threadIdx.x;
4     int idx_y = blockIdx.y * blockDim.y + threadIdx.y;
5
6     float sum = 0.f;
7     if (idx_x < N && idx_y < N)
8     {
9         for (int kk = 0; kk < N; ++kk)
10        {
11            sum += A[idx_y * N + kk] * B[kk * N + idx_x];
12        }
13        C[idx_y * N + idx_x] = sum;
14    }
15 }
```

取 `N={128,256,512,1024}`, `SIZE=32`, `blockSize=(SIZE,SIZE,1)`, `gridSize=(5,5,1)` 时的结果如下：


```

1 SIZE=32,N=128
2 blockSize=(8,8,1)
3 Time(%)      Total Time  Instances      Average      Minimum
   Maximum Name
4 -----
5 100.0         4128         1      4128.0      4128
   4128 gemm_baseline
6 SIZE=32,N=128
7 blockSize=(16,16,1)
8 Time(%)      Total Time  Instances      Average      Minimum
   Maximum Name
9 -----
10 100.0         5600         1      5600.0      5600
   5600 gemm_baseline
11 SIZE=32,N=128
12 blockSize=(32,32,1)
13 Time(%)      Total Time  Instances      Average      Minimum
   Maximum Name
14 -----
15 100.0         14880         1      14880.0      14880
   14880 gemm_baseline
16

```

取 `blockSize=(8,8,1)` 时效果最佳。值得注意的是，该程序效率影响因子较多，当取不同的N和不同的gridSize时得到的结果不同。上述结果并不能说明在任何情况下都是取(8,8,1)最好。

仅对上述情况进行分析，当blockSize设置较大时，会导致单个线程块中的线程数变多，从而使线程数超过一个流处理器可容纳的最大线程数。这就会导致无法充分利用GPU资源，从而出现效率下降的情况。

取 `N=128,blockSize=(32,32,1),gridSize=(G,G,1),G={6,7,8,9,10,11,12}`

结果如下：

1 ▾

G=5

2

Time(%)

Total Time

Instances

Average

Minimum

Maximum

Name

3

4 ▾

100.0

14880

1

14880.0

14880

14880

gemm_baseline

5

G=6

6

Time(%)

Total Time

Instances

Average

Minimum

Maximum

Name

7

8 ▾

100.0

14880

1

14880.0

14880

14880

gemm_baseline

9

G=7

10

Time(%)

Total Time

Instances

Average

Minimum

Maximum

Name

11

12 ▾

100.0

15296

1

15296.0

15296

15296

gemm_baseline

13

G=8

14

Time(%)

Total Time

Instances

Average

Minimum

Maximum

Name

15

16 ▾

100.0

15040

1

15040.0

15040

15040

gemm_baseline

17

g=9

18

Time(%)

Total Time

Instances

Average

Minimum

Maximum

Name

19

20 ▾

21	100.0	14944	1	14944.0	14944
22	14944	gemm_baseline			
	G=10				
	Time(%)	Total Time	Instances	Average	Minimum
23	Maximum	Name			
	-----	-----	-----	-----	-----
24	-----	-----	-----	-----	-----
	-----	-----	-----	-----	-----
25	100.0	14944	1	14944.0	14944
26	14944	gemm_baseline			
	G=11				
	Time(%)	Total Time	Instances	Average	Minimum
27	Maximum	Name			
	-----	-----	-----	-----	-----
28	-----	-----	-----	-----	-----
	-----	-----	-----	-----	-----
29	100.0	15360	1	15360.0	15360
30	15360	gemm_baseline			
	G=12				
	Time(%)	Total Time	Instances	Average	Minimum
31	Maximum	Name			
	-----	-----	-----	-----	-----
32	-----	-----	-----	-----	-----
	-----	-----	-----	-----	-----
33	100.0	15455	1	15455.0	15455
	15455	gemm_baseline			

分析可知，G=5和6时效率最高，而G=9和10时效率相同且高于上下的结果。可以发现， gridSize与效率并非完全正相关。实际使用过程过，应该结合N,SIZE,blockSize等的取值，灵活决定。

block

核心代码：

```

1  ▾ __global__ void gemm_block(float *A, float *B, float *C, int N) {
2
3      __shared__ float s_a[SIZE][SIZE];
4      __shared__ float s_b[SIZE][SIZE];
5
6      int idx_x = blockIdx.x * blockDim.x + threadIdx.x;
7      int idx_y = blockIdx.y * blockDim.y + threadIdx.y;
8
9      float sum = 0.0;
10 ▾  for (int bk = 0; bk < N; bk += SIZE) {
11      s_a[threadIdx.y][threadIdx.x] = A[idx_y * N + bk + threadIdx.x];
12      s_b[threadIdx.y][threadIdx.x] = B[(bk + threadIdx.y) * N + idx_x];
13      __syncthreads();
14
15 ▾      for (int i = 0; i < SIZE; ++i) {
16          sum += s_a[threadIdx.y][i] * s_b[i][threadIdx.x];
17      }
18      __syncthreads();
19  }
20
21 ▾  if (idx_x < N && idx_y < N) {
22      C[idx_y * N + idx_x] = sum;
23  }
24  }

```

取 `N=512, SIZE={8, 16, 32}, blockSize=(32, 32, 1), gridSize=(5, 5, 1)` 时结果如下:
(SIZE>32将无法获得正确结果)

Markdown | 复制代码

1	SIZE=8					
2	Time(%)	Total Time	Instances	Average	Minimum	
	Maximum	Name				
3	-----	-----	-----	-----	-----	-----
	-----	-----	-----	-----	-----	-----
4	100.0	18592	1	18592.0	18592	
	18592	gemm_block				
5	SIZE=16					
6	Time(%)	Total Time	Instances	Average	Minimum	
	Maximum	Name				
7	-----	-----	-----	-----	-----	-----
	-----	-----	-----	-----	-----	-----
	-----	-----	-----	-----	-----	-----
8	100.0	17760	1	17760.0	17760	
	17760	gemm_block				
9	SIZE=32					
10	Time(%)	Total Time	Instances	Average	Minimum	
	Maximum	Name				
11	-----	-----	-----	-----	-----	-----
	-----	-----	-----	-----	-----	-----
	-----	-----	-----	-----	-----	-----
12	100.0	39424	1	39424.0	39424	
	39424	gemm_block				

从上述结果可以发现，分块不能过大，一方面，分块过大可能导致计算速率降低，另一方面，当分块大于32时，计算结果将不在正确。

取 `N=512,SIZE=32,blockSize=(SIZE,SIZE,1),gridSize=(G,G,1),G={4,8,10,12}` 时，结果如下：

▼

Markdown | 复制代码

1	G=4					
2	Time(%)	Total Time	Instances	Average	Minimum	
	Maximum	Name				
3	<div>-----</div> <div>-----</div> <div>-----</div>					
4	100.0	47211	1	47211.0	47211	
	47211	gemm_block				
5	G=8					
6	Time(%)	Total Time	Instances	Average	Minimum	
	Maximum	Name				
7	<div>-----</div> <div>-----</div> <div>-----</div>					
8	100.0	38880	1	38880.0	38880	
	38880	gemm_block				
9	G=10					
10	Time(%)	Total Time	Instances	Average	Minimum	
	Maximum	Name				
11	<div>-----</div> <div>-----</div> <div>-----</div>					
12	100.0	39168	1	39168.0	39168	
	39168	gemm_block				
13	G=12					
14	Time(%)	Total Time	Instances	Average	Minimum	
	Maximum	Name				
15	<div>-----</div> <div>-----</div> <div>-----</div>					
16	100.0	46431	1	46431.0	46431	
	46431	gemm_block				
17						

较大的gridSize会导致较大的内存带宽负载，从而降低性能。每个线程块都需要访问内存来读取和写入数据，当gridSize变大时，需要访问内存的频率也会随之增加，从而使内存成为性能瓶颈。故实际应用中 gridSize不能取得过大。

较小的gridSize会导致无法有效利用GPU资源，同样会导致效率降低。

在上述条件下， gridSize=(8,8,1)时结果最佳。

取 N=512,SIZE=32,blockSize=(B,B,1),B={4,8,16,32},gridSize=(5,5,1) 时结果如下:

Markdown | 复制代码

1	B=4					
2	Time(%)	Total Time	Instances	Average	Minimum	
	Maximum	Name				
3	<div></div>					
4	100.0	7936	1	7936.0	7936	
	7936	gemm_block				
5	B=8					
6	Time(%)	Total Time	Instances	Average	Minimum	
	Maximum	Name				
7	<div></div>					
8	100.0	7808	1	7808.0	7808	
	7808	gemm_block				
9	B=16					
10	Time(%)	Total Time	Instances	Average	Minimum	
	Maximum	Name				
11	<div></div>					
12	100.0	13792	1	13792.0	13792	
	13792	gemm_block				
13	B=32					
14	Time(%)	Total Time	Instances	Average	Minimum	
	Maximum	Name				
15	<div></div>					
16	100.0	39008	1	39008.0	39008	
	39008	gemm_block				

结果显示，blockSize=(8,8,1)时结果最佳。

blockSize过小，无法有效利用GPU资源，blockSize过大，线程块需要频繁地访问内存，从而导致内存带宽负载较高，使GPU的效率降低。