

大規模計算特論B レポート課題

春山 棕

名古屋大学 情報学研究科

学生番号 251901088

2021 年 8 月 4 日

1 概要

行列サイズ 500, 800, 1000 の行列積に対して, Thread と ThreadBlock の形状を変化させて, それぞれ CUDA による並列化を行った. 行列データは一次元配列によって表現し, 列を変数 i , 行を変数 j で表現した. 行列積は三つのループを要する. 外側から, それぞれ列に関するループ (i), 行に関するループ (j), 総和を取る際のループ (k) である. 本稿では, 次の二通りの方法で並列化を行い, 両者の実行時間を比較した:

1. 列に関するループを ThreadBlock ごとに, 行に関するループを Thread ごとに並列化する
2. 列に関するループを Thread ごとに, 行に関するループを ThreadBlock ごとに並列化する

その結果, 1. の方が 2. よりも短い実行時間で計算を行うことができることがわかった. これは, 一次元配列のインデックスにおいては行列の行ごとに SharedMemory にキャッシュされるはずなので, メモリの局所性が有効活用され, キャッシュミスをおこしやすい 1. よりパフォーマンスが向上したためであると考えられる. 以下では, この考察を導く詳細なデータを示す.

2 実行時間データの作成方法

この章では, 実行時間データの作成方法を説明する. サンプルコード `mm5.cu` と `mm6.cu` をもとに,

1. 列に関するループを ThreadBlock ごとに, 行に関するループを Thread ごとに設定し, 行列サイズ 500, 800, 1000 の行列積を計算するプログラム `mm5_500.cu`, `mm5_800.cu`, `mm5_1000.cu` をそれぞれ作成した.

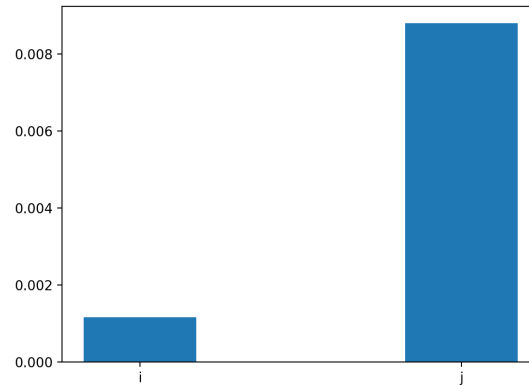


図 1: chart500.png

2. 列に関するループを Thread ごとに、行に関するループを ThreadBlock ごとに設定し、**行列サイズ 500, 800, 1000 の行列積を計算するプログラム mm6_500.cu, mm6_800.cu, mm6_1000.cu をそれぞれ作成した。**

上記のすべてのソースコードにおいて、行列サイズは変数 N で表現されている。以降では基本的に $N = 500$ と固定する。 $N = 800, 1000$ の場合も同様の操作を行った。

mm5_500.cu と mm6_500.cu をビルドしてバイナリファイル mm5_500.bin, mm6_500.bin をそれぞれ作成し、これらをそれぞれ 10 回ずつ実行し、実行時間データ data500 を作成した。このデータでは、一列目が mm5_500.cu、二列目が mm6_500.cu の実行時間データを表す（単位は秒）。これらの平均を取り、棒グラフとして比較したのが chart500.png である。なお、平均を取る際に外れ値を除外した。同様に chart800.png を作成した。また、mm5_1000.cu と mm6_1000.cu の実行結果を比較すると、行列積の結果が一致しなかったため、これらを比較するデータを作成することができなかった。

3 考察

chart500.png を分析したところ、約 7.5 倍の差があることがわかる。ここから、キャッシュミスによるオーバーヘッドは行列積の計算全体のルーチンの約 7.5 倍を占めると考えられる。

同様に chart800.png を分析したところ、約 85 倍約 8.5 倍の差があった。行列サイズを大きくすると差が縮まるというのは奇妙である。(2021 年 8 月 4 日訂正) 差が拡大したことがわかる（データの集計を取る際、平均を誤って合計で計算したため、この部分だけデータが本来の 10 倍になっていました）。

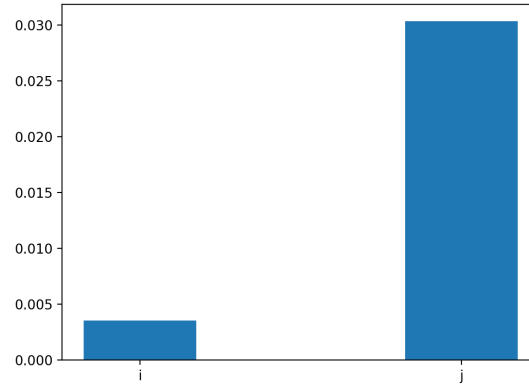


図 2: chart800.png

mm5_1000.cu と mm6_1000.cu の実行結果を比較したときに行列積の結果が一致しなかった理由は、デバイスが許容可能なスレッド数を超えて計算を行ってしまったためであると考えられる。

4 コメント

本稿では、行列サイズを大きくすると差が縮まるという奇妙な結果が出ました。もし本稿で間違っているかもしれない箇所に気づかれた場合は、こちらまで連絡をしていただけると幸いです。 haruyama.ryo@k.mbox.nagoya-u.ac.jp

5 付録

本稿で使ったプログラムとデータ、グラフを付録として添付する。

```

1 OPT = -O3 -arch=sm_70 -Xcompiler "-fopenmp"
2 all:
3     gcc -O3 -fopenmp -o mm1 mm1.c
4     nvcc ${OPT} -o mm2 mm2.cu
5     nvcc ${OPT} -o mm3 mm3.cu
6     nvcc ${OPT} -o mm4 mm4.cu
7     nvcc ${OPT} -o mm5 mm5.cu
8     nvcc ${OPT} -o mm6 mm6.cu
9
10 cu:
11     nvcc ${OPT} -o mm5_1000.bin mm5_1000.cu
12     nvcc ${OPT} -o mm5_2000.bin mm5_2000.cu
13     nvcc ${OPT} -o mm6_1000.bin mm6_1000.cu
14     nvcc ${OPT} -o mm6_2000.bin mm6_2000.cu
15
16 cum:
17     nvcc ${OPT} -o mm5_500.bin mm5_500.cu

```

```

18 nvcc ${OPT} -o mm5_800.bin mm5_800.cu
19 nvcc ${OPT} -o mm6_500.bin mm6_500.cu
20 nvcc ${OPT} -o mm6_800.bin mm6_800.cu
21
22 clean:
23 -/bin/rm ./*.o ./*~
24 -/bin/rm ./mm1 ./mm2 ./mm3 ./mm4 ./mm5 ./mm6
25 -/bin/rm *.bin *log

```

▲ Makefile

```

1 // -*- c++ -*-
2 #include <stdio.h>
3 #include <string.h>
4 #include <stdlib.h>
5 #include <cuda_runtime.h>
6 #include <omp.h>
7
8 #define frand() (double)(rand()%100)/10.0
9
10 __global__ void gpukernel
11 (int n, double *A, double *B, double *C){
12     int i, j, k;
13     i = blockIdx.x;
14     j = threadIdx.x;
15     for(k=0; k<n; k++){
16         C[i*n+j] += A[i*n+k] * B[k*n+j];
17     }
18 }
19
20 double *A, *B, *C;
21 double *d_A, *d_B, *d_C; // device memory
22 int N = 500;
23 int main(int argc, char **argv){
24     int x, y;
25     A = (double*)malloc(sizeof(double)*N*N);
26     B = (double*)malloc(sizeof(double)*N*N);
27     C = (double*)malloc(sizeof(double)*N*N);
28     for(y=0; y<N; y++){
29         for(x=0; x<N; x++){
30             A[y*N+x] = frand();
31             B[y*N+x] = frand();
32             C[y*N+x] = 0.0;
33         }
34     }
35     cudaMalloc((void**)&d_A, sizeof(double)*N*N);
36     cudaMalloc((void**)&d_B, sizeof(double)*N*N);
37     cudaMalloc((void**)&d_C, sizeof(double)*N*N);
38     cudaMemcpy(d_A, A, sizeof(double)*N*N, cudaMemcpyHostToDevice)
39     ;
40     cudaMemcpy(d_B, B, sizeof(double)*N*N, cudaMemcpyHostToDevice)
41     ;
42     cudaMemcpy(d_C, C, sizeof(double)*N*N, cudaMemcpyHostToDevice)
43     ;
44
45     double d1, d2;
46     cudaDeviceSynchronize();
47     d1 = omp_get_wtime();
48     gpukernel<<<N,N>>>(N, d_A, d_B, d_C);
49     cudaDeviceSynchronize();
50     d2 = omp_get_wtime();
51     printf("gpukernel %f sec\n", d2-d1);

```

```

50     cudaMemcpy(C, d_C, sizeof(double)*N*N, cudaMemcpyDeviceToHost)
51     ;
52     cudaFree(d_A); cudaFree(d_B); cudaFree(d_C);
53     // A
54     printf("A\n");
55     for(y=0; y<N; y++){
56         for(x=0; x<N; x++){
57             printf(" %.2f", A[y*N+x]);
58         }
59     }
60     // B
61     printf("B\n");
62     for(y=0; y<N; y++){
63         for(x=0; x<N; x++){
64             printf(" %.2f", B[y*N+x]);
65         }
66     }
67     // C
68     printf("C\n");
69     for(y=0; y<N; y++){
70         for(x=0; x<N; x++){
71             printf(" %.2f", C[y*N+x]);
72         }
73     }
74     printf("\n");
75 }
76 free(A); free(B); free(C);
77 return 0;
78 }

```

▲ mm5_500.cu

```

1 // -*- c++ -*-
2 #include <stdio.h>
3 #include <string.h>
4 #include <stdlib.h>
5 #include <cuda_runtime.h>
6 #include <omp.h>
7
8 #define frand() (double)(rand()%100)/10.0
9
10 __global__ void gpukernel
11 (int n, double *A, double *B, double *C){
12     int i, j, k;
13     i = blockIdx.x;
14     j = threadIdx.x;
15     for(k=0; k<n; k++){
16         C[i*n+j] += A[i*n+k] * B[k*n+j];
17     }
18 }
19
20 double *A, *B, *C;
21 double *d_A, *d_B, *d_C; // device memory
22 int N = 800;
23 int main(int argc, char **argv){
24     int x, y;
25     A = (double*)malloc(sizeof(double)*N*N);
26     B = (double*)malloc(sizeof(double)*N*N);
27     C = (double*)malloc(sizeof(double)*N*N);
28     for(y=0; y<N; y++){
29         for(x=0; x<N; x++){
30             A[y*N+x] = frand();

```

```

31     B[y*N+x] = frand();
32     C[y*N+x] = 0.0;
33 }
34 }
35 cudaMalloc((void**)&d_A, sizeof(double)*N*N);
36 cudaMalloc((void**)&d_B, sizeof(double)*N*N);
37 cudaMalloc((void**)&d_C, sizeof(double)*N*N);
38 cudaMemcpy(d_A, A, sizeof(double)*N*N, cudaMemcpyHostToDevice)
39 ;
40 cudaMemcpy(d_B, B, sizeof(double)*N*N, cudaMemcpyHostToDevice)
41 ;
42 cudaMemcpy(d_C, C, sizeof(double)*N*N, cudaMemcpyHostToDevice)
43 ;
44 double d1, d2;
45 cudaDeviceSynchronize();
46 d1 = omp_get_wtime();
47 gpukernel<<<N,N>>>(N, d_A, d_B, d_C);
48 cudaDeviceSynchronize();
49 d2 = omp_get_wtime();
50 printf("gpukernel %f sec\n", d2-d1);
51
52 cudaMemcpy(C, d_C, sizeof(double)*N*N, cudaMemcpyDeviceToHost)
53 ;
54 cudaFree(d_A); cudaFree(d_B); cudaFree(d_C);
55 // A
56 printf("A\n");
57 for(y=0; y<N; y++){
58     for(x=0; x<N; x++){
59         printf(" %.2f", A[y*N+x]);
60     }
61     printf("\n");
62 }
63 // B
64 printf("B\n");
65 for(y=0; y<N; y++){
66     for(x=0; x<N; x++){
67         printf(" %.2f", B[y*N+x]);
68     }
69     printf("\n");
70 }
71 // C
72 printf("C\n");
73 for(y=0; y<N; y++){
74     for(x=0; x<N; x++){
75         printf(" %.2f", C[y*N+x]);
76     }
77     printf("\n");
78 }
79 free(A); free(B); free(C);
80 return 0;
81 }

```

▲ mm5_800.cu

```

1 // -*- c++ -*-
2 #include <stdio.h>
3 #include <string.h>
4 #include <stdlib.h>
5 #include <cuda_runtime.h>
6 #include <omp.h>
7
8 #define frand() (double)(rand()%100)/10.0

```

```

9
10 __global__ void gpukernel
11 (int n, double *A, double *B, double *C){
12     int i, j, k;
13     i = blockIdx.x;
14     j = threadIdx.x;
15     for(k=0; k<n; k++){
16         C[i*n+j] += A[i*n+k] * B[k*n+j];
17     }
18 }
19
20 double *A, *B, *C;
21 double *d_A, *d_B, *d_C; // device memory
22 int N = 1000;
23 int main(int argc, char **argv){
24     int x, y;
25     A = (double*)malloc(sizeof(double)*N*N);
26     B = (double*)malloc(sizeof(double)*N*N);
27     C = (double*)malloc(sizeof(double)*N*N);
28     for(y=0; y<N; y++){
29         for(x=0; x<N; x++){
30             A[y*N+x] = frand();
31             B[y*N+x] = frand();
32             C[y*N+x] = 0.0;
33         }
34     }
35     cudaMalloc((void**)&d_A, sizeof(double)*N*N);
36     cudaMalloc((void**)&d_B, sizeof(double)*N*N);
37     cudaMalloc((void**)&d_C, sizeof(double)*N*N);
38     cudaMemcpy(d_A, A, sizeof(double)*N*N, cudaMemcpyHostToDevice);
39     ;
40     cudaMemcpy(d_B, B, sizeof(double)*N*N, cudaMemcpyHostToDevice);
41     ;
42     cudaMemcpy(d_C, C, sizeof(double)*N*N, cudaMemcpyHostToDevice);
43     ;
44     double d1, d2;
45     cudaDeviceSynchronize();
46     d1 = omp_get_wtime();
47     gpukernel<<<N,N>>>(N, d_A, d_B, d_C);
48     cudaDeviceSynchronize();
49     d2 = omp_get_wtime();
50     printf("gpukernel %f sec\n", d2-d1);
51
52     cudaMemcpy(C, d_C, sizeof(double)*N*N, cudaMemcpyDeviceToHost);
53     ;
54     cudaFree(d_A); cudaFree(d_B); cudaFree(d_C);
55     // A
56     printf("A\n");
57     for(y=0; y<N; y++){
58         for(x=0; x<N; x++){
59             printf(" %.2f", A[y*N+x]);
60         }
61     }
62     printf("\n");
63     // B
64     printf("B\n");
65     for(y=0; y<N; y++){
66         for(x=0; x<N; x++){
67             printf(" %.2f", B[y*N+x]);
68         }
69     }
70     printf("\n");

```

```

67     }
68     // C
69     printf("C\n");
70     for(y=0; y<N; y++){
71         for(x=0; x<N; x++){
72             printf(" %.2f", C[y*N+x]);
73         }
74         printf("\n");
75     }
76     free(A); free(B); free(C);
77     return 0;
78 }

```

▲ mm5_1000.cu

```

1  // -*- c++ -*-
2  #include <stdio.h>
3  #include <string.h>
4  #include <stdlib.h>
5  #include <cuda_runtime.h>
6  #include <omp.h>
7
8  #define frand() (double)(rand()%100)/10.0
9
10 __global__ void gpukernel
11 (int n, double *A, double *B, double *C){
12     int i, j, k;
13     i = threadIdx.x;
14     j = blockIdx.x;
15     for(k=0; k<n; k++){
16         C[i*n+j] += A[i*n+k] * B[k*n+j];
17     }
18 }
19
20 double *A, *B, *C;
21 double *d_A, *d_B, *d_C; // device memory
22 int N = 500;
23 int main(int argc, char **argv){
24     int x, y;
25     A = (double*)malloc(sizeof(double)*N*N);
26     B = (double*)malloc(sizeof(double)*N*N);
27     C = (double*)malloc(sizeof(double)*N*N);
28     for(y=0; y<N; y++){
29         for(x=0; x<N; x++){
30             A[y*N+x] = frand();
31             B[y*N+x] = frand();
32             C[y*N+x] = 0.0;
33         }
34     }
35     cudaMalloc((void**)&d_A, sizeof(double)*N*N);
36     cudaMalloc((void**)&d_B, sizeof(double)*N*N);
37     cudaMalloc((void**)&d_C, sizeof(double)*N*N);
38     cudaMemcpy(d_A, A, sizeof(double)*N*N, cudaMemcpyHostToDevice);
39     cudaMemcpy(d_B, B, sizeof(double)*N*N, cudaMemcpyHostToDevice);
40     cudaMemcpy(d_C, C, sizeof(double)*N*N, cudaMemcpyHostToDevice);
41
42     double d1, d2;
43     cudaDeviceSynchronize();
44     d1 = omp_get_wtime();
45     gpukernel<<<N,N>>>(N, d_A, d_B, d_C);

```



```

46  cudaDeviceSynchronize();
47  d2 = omp_get_wtime();
48  printf("gpukernel %f sec\n", d2-d1);
49
50  cudaMemcpy(C, d_C, sizeof(double)*N*N, cudaMemcpyDeviceToHost)
51  ;
52  cudaFree(d_A); cudaFree(d_B); cudaFree(d_C);
53  // A
54  printf("A\n");
55  for(y=0; y<N; y++){
56      for(x=0; x<N; x++){
57          printf(" %.2f", A[y*N+x]);
58      }
59      printf("\n");
60  }
61  // B
62  printf("B\n");
63  for(y=0; y<N; y++){
64      for(x=0; x<N; x++){
65          printf(" %.2f", B[y*N+x]);
66      }
67      printf("\n");
68  }
69  // C
70  printf("C\n");
71  for(y=0; y<N; y++){
72      for(x=0; x<N; x++){
73          printf(" %.2f", C[y*N+x]);
74      }
75      printf("\n");
76  }
77  free(A); free(B); free(C);
78  return 0;
79 }

```

▲ mm6_500.cu

```

1  // -*- c++ -*-
2  #include <stdio.h>
3  #include <string.h>
4  #include <stdlib.h>
5  #include <cuda_runtime.h>
6  #include <omp.h>
7
8  #define frand() (double)(rand()%100)/10.0
9
10 __global__ void gpukernel
11 (int n, double *A, double *B, double *C){
12     int i, j, k;
13     i = threadIdx.x;
14     j = blockIdx.x;
15     for(k=0; k<n; k++){
16         C[i*n+j] += A[i*n+k] * B[k*n+j];
17     }
18 }
19
20 double *A, *B, *C;
21 double *d_A, *d_B, *d_C; // device memory
22 int N = 800;
23 int main(int argc, char **argv){
24     int x, y;
25     A = (double*)malloc(sizeof(double)*N*N);
26     B = (double*)malloc(sizeof(double)*N*N);

```

```

27 C = (double*)malloc(sizeof(double)*N*N);
28 for(y=0; y<N; y++){
29     for(x=0; x<N; x++){
30         A[y*N+x] = frand();
31         B[y*N+x] = frand();
32         C[y*N+x] = 0.0;
33     }
34 }
35 cudaMalloc((void**)&d_A, sizeof(double)*N*N);
36 cudaMalloc((void**)&d_B, sizeof(double)*N*N);
37 cudaMalloc((void**)&d_C, sizeof(double)*N*N);
38 cudaMemcpy(d_A, A, sizeof(double)*N*N, cudaMemcpyHostToDevice)
39 ;
40 cudaMemcpy(d_B, B, sizeof(double)*N*N, cudaMemcpyHostToDevice)
41 ;
42 cudaMemcpy(d_C, C, sizeof(double)*N*N, cudaMemcpyHostToDevice)
43 ;
44 double d1, d2;
45 cudaDeviceSynchronize();
46 d1 = omp_get_wtime();
47 gpukernel<<<N,N>>>(N, d_A, d_B, d_C);
48 cudaDeviceSynchronize();
49 d2 = omp_get_wtime();
50 printf("gpukernel %f sec\n", d2-d1);
51
52 cudaMemcpy(C, d_C, sizeof(double)*N*N, cudaMemcpyDeviceToHost)
53 ;
54 cudaFree(d_A); cudaFree(d_B); cudaFree(d_C);
55 // A
56 printf("A\n");
57 for(y=0; y<N; y++){
58     for(x=0; x<N; x++){
59         printf(" %.2f", A[y*N+x]);
60     }
61     printf("\n");
62 }
63 // B
64 printf("B\n");
65 for(y=0; y<N; y++){
66     for(x=0; x<N; x++){
67         printf(" %.2f", B[y*N+x]);
68     }
69     printf("\n");
70 }
71 // C
72 printf("C\n");
73 for(y=0; y<N; y++){
74     for(x=0; x<N; x++){
75         printf(" %.2f", C[y*N+x]);
76     }
77     printf("\n");
78 }
79 free(A); free(B); free(C);
80 return 0;
81 }

```

▲ mm6_800.cu

```

1 // -*- c++ -*-
2 #include <stdio.h>
3 #include <string.h>
4 #include <stdlib.h>

```

```

5 #include <cuda_runtime.h>
6 #include <omp.h>
7
8 #define frand() (double)(rand()%100)/10.0
9
10 __global__ void gpukernel
11 (int n, double *A, double *B, double *C){
12     int i, j, k;
13     i = blockIdx.x;
14     j = threadIdx.x;
15     for(k=0; k<n; k++){
16         C[i*n+j] += A[i*n+k] * B[k*n+j];
17     }
18 }
19
20 double *A, *B, *C;
21 double *d_A, *d_B, *d_C; // device memory
22 int N = 1000;
23 int main(int argc, char **argv){
24     int x, y;
25     A = (double*)malloc(sizeof(double)*N*N);
26     B = (double*)malloc(sizeof(double)*N*N);
27     C = (double*)malloc(sizeof(double)*N*N);
28     for(y=0; y<N; y++){
29         for(x=0; x<N; x++){
30             A[y*N+x] = frand();
31             B[y*N+x] = frand();
32             C[y*N+x] = 0.0;
33         }
34     }
35     cudaMalloc((void**)&d_A, sizeof(double)*N*N);
36     cudaMalloc((void**)&d_B, sizeof(double)*N*N);
37     cudaMalloc((void**)&d_C, sizeof(double)*N*N);
38     cudaMemcpy(d_A, A, sizeof(double)*N*N, cudaMemcpyHostToDevice);
39     ;
40     cudaMemcpy(d_B, B, sizeof(double)*N*N, cudaMemcpyHostToDevice);
41     ;
42     cudaMemcpy(d_C, C, sizeof(double)*N*N, cudaMemcpyHostToDevice);
43     ;
44     double d1, d2;
45     cudaDeviceSynchronize();
46     d1 = omp_get_wtime();
47     gpukernel<<<N,N>>>(N, d_A, d_B, d_C);
48     cudaDeviceSynchronize();
49     d2 = omp_get_wtime();
50     printf("gpukernel %f sec\n", d2-d1);
51
52     cudaMemcpy(C, d_C, sizeof(double)*N*N, cudaMemcpyDeviceToHost);
53     ;
54     cudaFree(d_A); cudaFree(d_B); cudaFree(d_C);
55     // A
56     printf("A\n");
57     for(y=0; y<N; y++){
58         for(x=0; x<N; x++){
59             printf(" %.2f", A[y*N+x]);
60         }
61     }
62     printf("\n");
63     // B
64     printf("B\n");
65     for(y=0; y<N; y++){

```

```

63     for(x=0; x<N; x++){
64         printf(" %.2f", B[y*N+x]);
65     }
66     printf("\n");
67 }
68 // C
69 printf("C\n");
70 for(y=0; y<N; y++){
71     for(x=0; x<N; x++){
72         printf(" %.2f", C[y*N+x]);
73     }
74     printf("\n");
75 }
76 free(A); free(B); free(C);
77 return 0;
78 }

```

▲ mm6_1000.cu

```

1 0.001809 0.008751
2 0.001097 0.008815
3 0.001112 0.008826
4 0.001122 0.008762
5 0.001116 0.008774
6 0.001104 0.008846
7 0.001114 0.008823
8 0.001106 0.008804
9 0.001054 0.008764
10 0.001070 0.008790

```

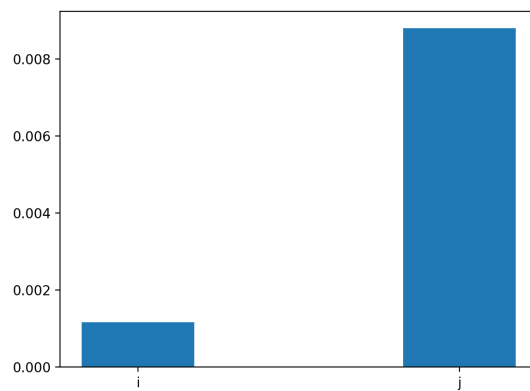
▲ data500

```

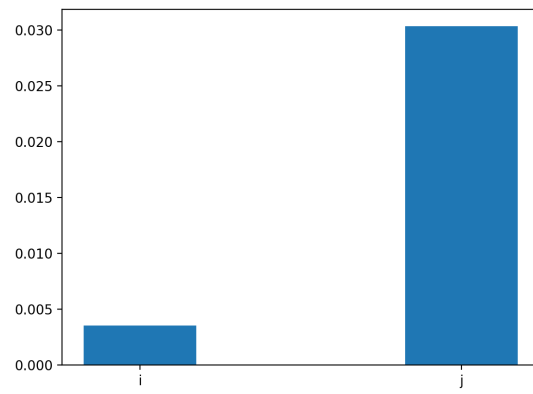
1 0.003600 0.030425
2 0.003565 0.030280
3 0.003557 0.030434
4 0.003553 0.030243
5 0.003553 0.030204
6 0.003560 0.030379
7 0.003560 0.030358
8 0.003552 0.030212
9 0.003556 0.030264
10 0.003562 0.030283

```

▲ data800



▲ chart500.png



▲ chart800.png