# Computational Science on Many-Core Architectures: Exercise 3:

Viktor Beck, 11713110

8. November 2022

## Task 1: Strided and Offset Memory Access

### a)

Here we only sum every k-th element in the vector addition. One can see in Fig. 1 that we have an exponential decrease in the effective bandwidth. Probably this is due to spacial locality: the localities of the elements with index `i*k` are not neighbouring for subsequent ks and therefore the kernel can not be as fast as with elements located closer to each other. As expected the newer RTX3060 GPU is faster than the K40. The noise for the RTX GPU is possibly caused by itself attempting to optimize the computations.
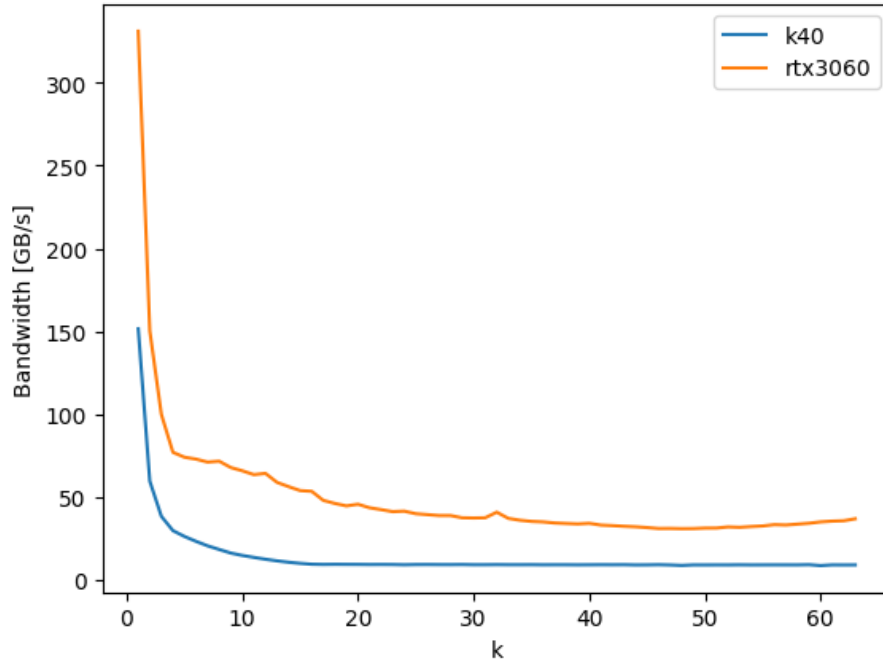


Figure 1: Strided memory access: Bandwidth for vector addition for different k. The time was measured by taking the median over 6 repetitions. Note: for computing the bandwidth the formula `3*(N/k)*sizeof(double)*t` was used, 3 because 2 reads and 1 write.

**b)**

Here we skipped the first k elements of the vector. In Fig. 2 one can see that the effective bandwidth is jumping significantly for every k which is dividable by 4 and even more by every k which is dividable by 16. These minima are obviously not a coincidence. This is probably because of the blocks and threads per block being dividable by these numbers. However, it seems that taking k = 16, 32, 48, 64, ... (I even tested it for k > 64) the effective bandwidth is the highest and therefore, one could recommend these values of k for future optimizations. Again here we can see that the RTX3060 GPU is much faster - even twice as fast as the K40 GPU.

So in general, one could recommend this approach over the one from a) where we take N/k since here we have better spacial localities when it comes to accessing the elements of the vectors.
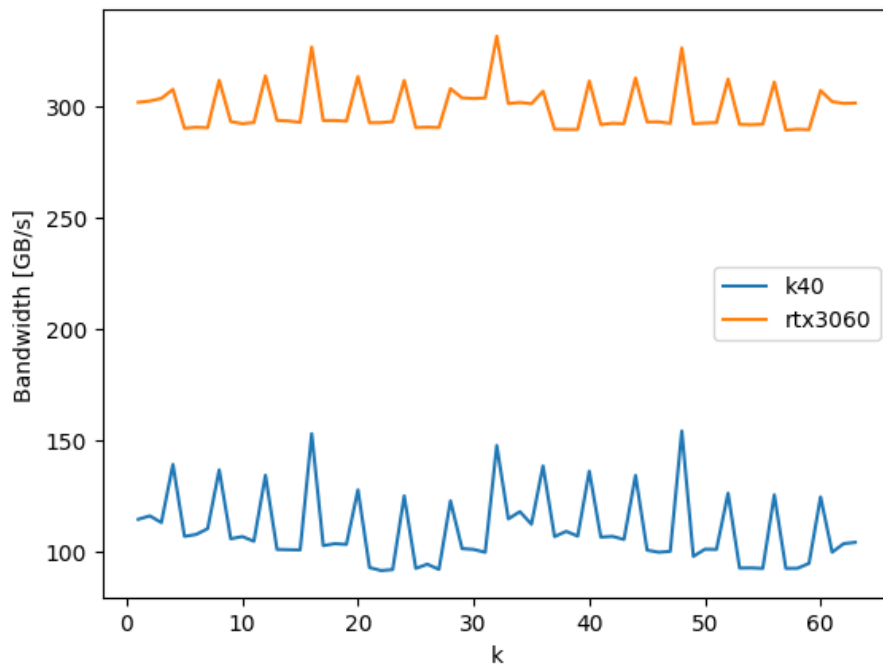


Figure 2: Offset memory access: Bandwidth for vector addition for different k

## Task 2: Dense Matrix Transpose

**a)**

The problem here was a leakage of 800 Bytes (for N=10) which was solved by simply adding `cudafree(cuda_A)` and `free(A)` at the bottom of the main function. Consequently, the leakage was due to the 'unfreed' space of the GPU kernel. Note: I guess in the description of the task you meant that the problem here is the leakage and not the wrong results of the transposition itself. I corrected this in b).

## b)

The transposed matrix of my friend's code does not produce correct results for the case that the total number of threads is smaller than N*N. If we use a for loop in the transpose function such that each thread can work again on the task we get the problem that a thread does not know which elements are already in the correct place in the matrix. E.g.: thread X writes to A[42]; after that, thread Y reads from A[42] to write to A[...] but this means that it is already reading the updated value which leads to an incorrect transposition. In an naive approach we simply make sure that the number of threads is always bigger than N*N:

```
1           transpose<<<(N*N+255)/256, 256>>>(cuda_A, N);
```

## c)

Here a new matrix B was introduced. The problem here is that we have double the amount of data that has to be passed to the GPU kernel than before with just one matrix. The naive approach was the following:

```
1  __global__ void transpose(double *A, double *B, int N)
2  {
3      int t_idx = blockIdx.x * blockDim.x + threadIdx.x;
4      // loop for the case N*N > number of threads
5      for (size_t i = t_idx; i < (N*N); i += blockDim.x*gridDim.x){
6          int row_idx = i / N;
7          int col_idx = i % N;
8          B[row_idx * N + col_idx] = A[col_idx * N + row_idx];
9      }
10 }
```

## d)

A more sophisticated version of the kernel was then implemented (with inspiration from the link you provided) as follows:

```
1  const int TILE_DIM = 16;
2  const int BLOCK_ROWS = 16;
3
4  __global__ void transposeCoalesced(double *A, double *B)
5  {
6      //create blocks (tiles) in shared memory
7      __shared__ double tile[TILE_DIM][TILE_DIM];
8
9      int x = blockIdx.x * TILE_DIM + threadIdx.x;
10     int y = blockIdx.y * TILE_DIM + threadIdx.y;
11     int width = gridDim.x * TILE_DIM;
12
13     //write from A into shared tile
14     for (int i = 0; i < TILE_DIM; i += BLOCK_ROWS){
15         tile[threadIdx.y+i][threadIdx.x] = A[(y+i)*width + x];
16     }
17
18     __syncthreads(); // wait until all blocks are finished
19     //transpose block
20     x = blockIdx.y * TILE_DIM + threadIdx.x;
21     y = blockIdx.x * TILE_DIM + threadIdx.y;
22
23     // assign blocks to right position in B
```

```
24      for (int i = 0; i < TILE_DIM; i += BLOCK_ROWS){
25          B[(y+i)*width + x] = tile[threadIdx.x][threadIdx.y + i];
26      }
27  }
```

Here we cut the matrix into tiles/blocks of 16x16 elements. Each tile is transposed separately and then assigned to the right position in matrix B.

## e)

The effective bandwidth of the naive implementations of b) and c) and the coalesced one (with loading blocks of 16x16) can be seen below. We can observe that the difference between the different approaches start out similarly but for larger N (up to $2^{14}$) there is a significant difference. It seems that the coalescent implementation of d) stabilizes for N between $10^2$ and $10^3$ while the implementation of b) drops down. Surprisingly with the GTX GPU the implementation of c) works out quite well compared to the others. As expected, there is a difference between the GTX and the K40 GPU.
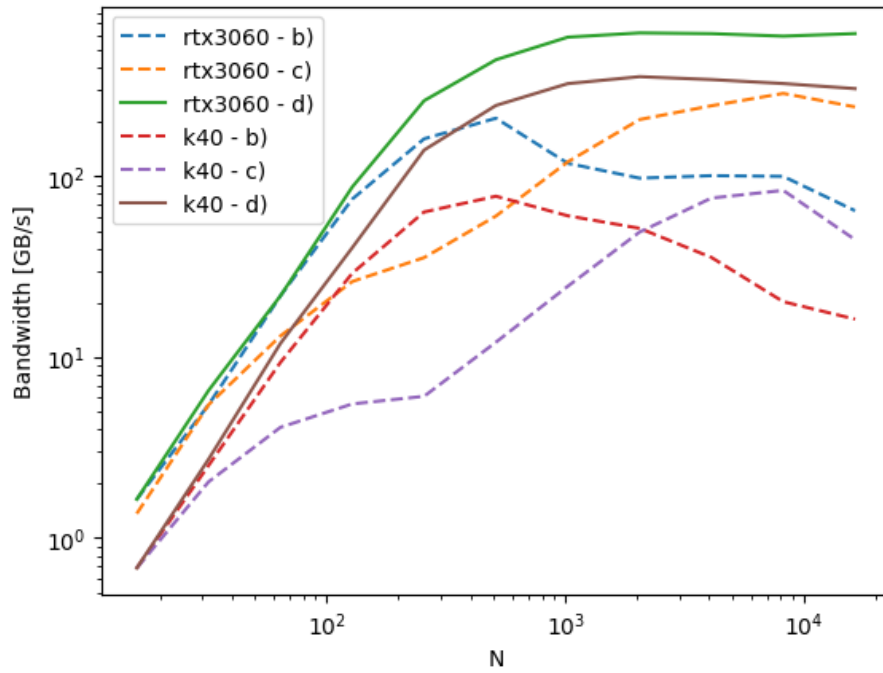


Figure 3: Effective bandwidth for different N (16 to 16384) for the K40 and the GTX 3060 GPU

# Appendix

## Task 1.ab)

```cpp
#include <iostream>
#include <stdio.h>
#include <vector>
#include <algorithm>
#include "timer.hpp"

using Vector = std::vector<double>;


__global__ void addition(int n, int k, double *x, double *y, double *z)
{
  int id = blockIdx.x*blockDim.x + threadIdx.x;
  for (size_t i = id; i < n-k; i += blockDim.x*gridDim.x){
    z[i+k] = x[i+k] + y[i+k];
    //printf("z: %g\n", z[i*k]);
  }
}

int main(void)
{

    double *x, *y, *z, *d_x, *d_y, *d_z;

    Vector bandwidth;
    Timer timer;

    //int N_min = 100000000;
    //int N_max = 100000000;
    int reps = 6; // has to be an even number

    int N = 100000000;

    // Allocate host memory and initialize
    x = (double*)malloc(N*sizeof(double));
    y = (double*)malloc(N*sizeof(double));
    z = (double*)malloc(N*sizeof(double));

    for (int i = 0; i < N; i++) {
        x[i] = i;
        y[i] = N-1-i;
        z[i] = 0;
    }

    // Allocate device memory and copy host data over
    cudaMalloc(&d_x, N*sizeof(double));
    cudaMalloc(&d_y, N*sizeof(double));
    cudaMalloc(&d_z, N*sizeof(double));

    cudaMemcpy(d_x, x, N*sizeof(double), cudaMemcpyHostToDevice);
    cudaMemcpy(d_y, y, N*sizeof(double), cudaMemcpyHostToDevice);
    cudaMemcpy(d_z, z, N*sizeof(double), cudaMemcpyHostToDevice);

    for (int k = 1; k < 64; k++){
        //std::cout << "k = " << k << std::endl;

            Vector time_vec;
```

```
58
59          // repetitions
60          for (int j = 0; j < reps; j++){
61
62              // MEASURING TIME FROM HERE
63              cudaDeviceSynchronize();
64              timer.reset();
65
66              addition<<<256, 256>>>(N, k, d_x, d_y, d_z);
67
68              cudaDeviceSynchronize();
69              time_vec.push_back(timer.get());
70              //std::cout << timer.get() << std::endl;
71
72              // TO HERE
73          }
74          std::sort(time_vec.begin(), time_vec.end());
75          bandwidth.push_back(3*(N-k)*sizeof(double)/time_vec[reps/2]*1e-9);
76          //printf("Elapsed: %g\n", time);
77      }
78      cudaMemcpy(z, d_z, N*sizeof(double), cudaMemcpyDeviceToHost);
79
80      cudaFree(d_x);
81      cudaFree(d_y);
82      cudaFree(d_z);
83      free(x);
84      free(y);
85      free(z);
86
87      printf("Gb/s:\n");
88      for (const auto& value : bandwidth){
89          std::cout << value << "," << std::endl;
90      }
91
92
93      return EXIT_SUCCESS;
94 }
```

## Task 2.ab).

```
1  #include <stdio.h>
2  #include <iostream>
3  #include "timer.hpp"
4  #include "cuda_errchk.hpp" // for error checking of CUDA calls
5  #include <vector>
6  #include <algorithm>
7
8  __global__
9  void transpose(double *A, int N)
10 {
11     int t_idx = blockIdx.x*blockDim.x + threadIdx.x;
12     int row_idx = t_idx / N;
13     int col_idx = t_idx % N;
14     if (row_idx < N && col_idx < N){
15         int tmp;
16         tmp = A[col_idx * N + row_idx];
17         A[col_idx * N + row_idx] = A[row_idx * N + col_idx];
18         A[row_idx * N + col_idx] = tmp;
19     }
20
```

```
21 | }
22 |
23 | void print_A(double *A, int N)
24 | {
25 |     for (int i = 0; i < N; i++) {
26 |         for (int j = 0; j < N; ++j) {
27 |             std::cout << A[i * N + j] << ",␣";
28 |         }
29 |     std::cout << std::endl;
30 |     }
31 | }
32 |
33 | int main(void)
34 | {
35 |     std::vector<double> N_vec = {16,32,64,128,256,512,1024,2048,4096,8192,16384};
36 |     std::vector<double> bandwidth;
37 |
38 |     double *A, *cuda_A;
39 |     Timer timer;
40 |     int reps = 6;
41 |
42 |     for (const auto &N : N_vec){
43 |         std::vector<double> time_vec;
44 |
45 |         // Allocate host memory and initialize
46 |         A = (double*)malloc(N*N*sizeof(double));
47 |
48 |         for (int i = 0; i < N*N; i++) {
49 |         A[i] = i;
50 |         }
51 |
52 |         //print_A(A, N);
53 |
54 |         // Allocate device memory and copy host data over
55 |         CUDA_ERRCHK(cudaMalloc(&cuda_A, N*N*sizeof(double)));
56 |
57 |         // copy data over
58 |         CUDA_ERRCHK(cudaMemcpy(cuda_A, A, N*N*sizeof(double), cudaMemcpyHostToDevice));
59 |
60 |         // repetitions
61 |         for (int j = 0; j < reps; j++){
62 |             // wait for previous operations to finish, then start timings
63 |             CUDA_ERRCHK(cudaDeviceSynchronize());
64 |             timer.reset();
65 |
66 |             // Perform the transpose operation
67 |             transpose<<<(N*N+255)/256, 256>>>(cuda_A, N);
68 |
69 |             // wait for kernel to finish, then print elapsed time
70 |             CUDA_ERRCHK(cudaDeviceSynchronize());
71 |             time_vec.push_back(timer.get());
72 |         }
73 |         std::sort(time_vec.begin(), time_vec.end());
74 |         bandwidth.push_back(4*N*N*sizeof(double)/time_vec[reps/2]*1e-9);
75 |
76 |         // copy data back (implicit synchronization point)
77 |         CUDA_ERRCHK(cudaMemcpy(A, cuda_A, N*N*sizeof(double), cudaMemcpyDeviceToHost));
78 |
79 |         //print_A(A, N);
80 |
81 |         cudaFree(cuda_A);
82 |         free(A);
```

```
83          }
84
85      printf("Gb/s:\n");
86      for (const auto& value : bandwidth){
87          std::cout << value << "," << std::endl;
88      }
89
90      CUDA_ERRCHK(cudaDeviceReset());  // for CUDA leak checker to work
91
92      return EXIT_SUCCESS;
93 }
```

## Task 2.c).

```
1  #include <stdio.h>
2  #include <iostream>
3  #include "timer.hpp"
4  #include "cuda_errchk.hpp" // for error checking of CUDA calls
5  #include <vector>
6  #include <algorithm>
7
8  __global__ void transpose(double *A, double *B, int N)
9  {
10     int t_idx = blockIdx.x * blockDim.x + threadIdx.x;
11     // loop for the case N*N > number of threads
12     for (size_t i = t_idx; i < (N*N); i += blockDim.x*gridDim.x){
13         int row_idx = i / N;
14         int col_idx = i % N;
15         B[row_idx * N + col_idx] = A[col_idx * N + row_idx];
16     }
17 }
18
19 void print_A(double *A, int N)
20 {
21     for (int i = 0; i < N; i++)
22     {
23         for (int j = 0; j < N; ++j)
24         {
25             std::cout << A[i * N + j] << ",␣";
26         }
27         std::cout << std::endl;
28     }
29 }
30
31 int main(void)
32 {
33     //int N = 33;
34     std::vector<double> N_vec = {16,32,64,128,256,512,1024,2048,4096,8192,16384};
35     std::vector<double> bandwidth;
36
37     // declare A and B
38     double *A, *B, *cuda_A, *cuda_B;
39     Timer timer;
40     int reps = 6;
41
42     for (const auto &N : N_vec){
43         std::vector<double> time_vec;
44
45         // Allocate host memory and initialize
46         A = (double *)malloc(N * N * sizeof(double));
```

```
47        B = (double *)malloc(N * N * sizeof(double));
48
49        for (int i = 0; i < N * N; i++)
50        {
51            A[i] = i;
52            B[i] = 0;
53        }
54
55        //print_A(A, N);
56        //print_A(B, N);
57
58        // Allocate device memory and copy host data over
59        CUDA_ERRCHK(cudaMalloc(&cuda_A, N * N * sizeof(double)));
60        CUDA_ERRCHK(cudaMalloc(&cuda_B, N * N * sizeof(double)));
61
62        // copy data over
63        CUDA_ERRCHK(cudaMemcpy(cuda_A, A, N * N * sizeof(double), cudaMemcpyHostToDevice));
64        CUDA_ERRCHK(cudaMemcpy(cuda_B, B, N * N * sizeof(double), cudaMemcpyHostToDevice));
65
66        // repetitions
67        for (int j = 0; j < reps; j++){
68            // wait for previous operations to finish, then start timings
69            CUDA_ERRCHK(cudaDeviceSynchronize());
70            timer.reset();
71
72            // Perform the transpose operation
73            transpose<<<(N + 255) / 256, 256>>>(cuda_A, cuda_B, N);
74
75            // wait for kernel to finish, then print elapsed time
76            CUDA_ERRCHK(cudaDeviceSynchronize());
77            time_vec.push_back(timer.get());
78        }
79        std::sort(time_vec.begin(), time_vec.end());
80        bandwidth.push_back(4*N*N*sizeof(double)/time_vec[reps/2]*1e-9);
81
82        // copy data back (implicit synchronization point)
83        CUDA_ERRCHK(cudaMemcpy(A, cuda_A, N * N * sizeof(double), cudaMemcpyDeviceToHost));
84        CUDA_ERRCHK(cudaMemcpy(B, cuda_B, N * N * sizeof(double), cudaMemcpyDeviceToHost));
85
86        //print_A(B, N);
87
88        cudaFree(cuda_A);
89        free(A);
90        cudaFree(cuda_B);
91        free(B);
92    }
93
94    printf("Gb/s:\n");
95    for (const auto& value : bandwidth){
96        std::cout << value << "," << std::endl;
97    }
98
99    CUDA_ERRCHK(cudaDeviceReset()); // for CUDA leak checker to work
100
101    return EXIT_SUCCESS;
102 }
```

## Task 2.d).

```
1  #include <stdio.h>
```

```cpp
#include <iostream>
#include "timer.hpp"
#include "cuda_errchk.hpp" // for error checking of CUDA calls
#include <vector>
#include <algorithm>

const int TILE_DIM = 16;
const int BLOCK_ROWS = 16;

__global__ void transposeCoalesced(double *A, double *B)
{
    //create blocks (tiles) in shared memory
    __shared__ double tile[TILE_DIM][TILE_DIM];

    int x = blockIdx.x * TILE_DIM + threadIdx.x;
    int y = blockIdx.y * TILE_DIM + threadIdx.y;
    int width = gridDim.x * TILE_DIM;

    //write from A into shared tile
    for (int i = 0; i < TILE_DIM; i += BLOCK_ROWS){
        tile[threadIdx.y+i][threadIdx.x] = A[(y+i)*width + x];
    }

    __syncthreads(); // wait until all blocks are finished
    //transpose block
    x = blockIdx.y * TILE_DIM + threadIdx.x;
    y = blockIdx.x * TILE_DIM + threadIdx.y;

    // assign blocks to right position in B
    for (int i = 0; i < TILE_DIM; i += BLOCK_ROWS){
        B[(y+i)*width + x] = tile[threadIdx.x][threadIdx.y + i];
    }
}

void print_A(double *A, int N)
{
    for (int i = 0; i < N; i++)
    {
        for (int j = 0; j < N; ++j)
        {
            std::cout << A[i * N + j] << ",␣";
        }
        std::cout << std::endl;
    }
}

int main(void)
{
    std::vector<double> N_vec = {16,32,64,128,256,512,1024,2048,4096,8192,16384};
    std::vector<double> bandwidth;

    // declare A and B
    double *A, *B, *cuda_A, *cuda_B;
    Timer timer;
    int reps = 6;

    for (const auto &N : N_vec){
        std::vector<double> time_vec;

        dim3 dimGrid(N/TILE_DIM, N/TILE_DIM, 1);
        dim3 dimBlock(TILE_DIM, BLOCK_ROWS, 1);
```

```cuda
 64          // Allocate host memory and initialize
 65          A = (double *)malloc(N * N * sizeof(double));
 66          B = (double *)malloc(N * N * sizeof(double));
 67
 68          for (int i = 0; i < N * N; i++)
 69          {
 70              A[i] = i;
 71              B[i] = 0;
 72          }
 73
 74          //print_A(A, N);
 75          //print_A(B, N);
 76
 77          // Allocate device memory and copy host data over
 78          CUDA_ERRCHK(cudaMalloc(&cuda_A, N * N * sizeof(double)));
 79          CUDA_ERRCHK(cudaMalloc(&cuda_B, N * N * sizeof(double)));
 80
 81          // copy data over
 82          CUDA_ERRCHK(cudaMemcpy(cuda_A, A, N * N * sizeof(double), cudaMemcpyHostToDevice));
 83          CUDA_ERRCHK(cudaMemcpy(cuda_B, B, N * N * sizeof(double), cudaMemcpyHostToDevice));
 84
 85          // repetitions
 86          for (int j = 0; j < reps; j++){
 87              // wait for previous operations to finish, then start timings
 88              CUDA_ERRCHK(cudaDeviceSynchronize());
 89              timer.reset();
 90
 91              // Perform the transpose operation
 92              transposeCoalesced<<<dimGrid, dimBlock>>>(cuda_A, cuda_B);
 93
 94              // wait for kernel to finish, then print elapsed time
 95              CUDA_ERRCHK(cudaDeviceSynchronize());
 96              time_vec.push_back(timer.get());
 97          }
 98          std::sort(time_vec.begin(), time_vec.end());
 99          bandwidth.push_back(4*N*N*sizeof(double)/time_vec[reps/2]*1e-9);
100
101          // copy data back (implicit synchronization point)
102          //CUDA_ERRCHK(cudaMemcpy(A, cuda_A, N * N * sizeof(double), cudaMemcpyDeviceToHost));
103          CUDA_ERRCHK(cudaMemcpy(B, cuda_B, N * N * sizeof(double), cudaMemcpyDeviceToHost));
104
105          //print_A(B, N);
106
107          cudaFree(cuda_A);
108          free(A);
109          cudaFree(cuda_B);
110          free(B);
111      }
112
113      printf("Gb/s:\n");
114      for (const auto& value : bandwidth){
115          std::cout << value << "," << std::endl;
116      }
117
118      CUDA_ERRCHK(cudaDeviceReset()); // for CUDA leak checker to work
119
120      return EXIT_SUCCESS;
121  }
```