

M2-BIG DATA

GPGPU - Chapter 4

Exercice 1



Réalisé par:
Gaby Maroun

Encadré par:
Dr. Etancelin JM

December 13, 2020

Objectives

Implement a basic dense matrix multiplication kernel. This is the first version of several optimizations to follow (see next exercises). The program computes :

$$C = AB$$

where A , B and C are general rectangular matrices.

Instructions

Edit the given code to perform the following :

```
1 #include <stdio.h>
2 #include <cuda.h>
3 #include <time.h>
4
5 #include "matmul_utils.hpp"
6
7
8 // Cuda kernel
9 __global__ void dgemm( float *A, float *B, float *C,
10                      int numARows, int numAColumns, int numBRows, int numBColumns)
11 {
12     // @TODO@ : Complete here the kernel code
13
14     // Calculate the row index of the C element and A
15     int Row= blockIdx.y * blockDim.y + threadIdx.y;
16
17     // Calculate the column index of C and B
18     int Col= blockIdx.x * blockDim.x + threadIdx.x;
19
20     if ( (Row< numARows ) && (Col < numBColumns ) ) {
21         float Cvalue =0;
22         // each thread computes one element of the block sub-matrix
23         for(int k=0; k< numAColumns ; k++) {
24             Cvalue += A[Row*numAColumns +k]* B[k*numBColumns +Col];
25         }
26         C[Row*numBColumns + Col ] = Cvalue;
27     }
28 }
29
30 int main(int argc, char **argv)
31 {
32     if(argc!=4) { printf("Usage : %s [nb of rows for A] [nb of cols for A] [nb of cols for B]\n", argv[0]); exit(2); }
33     //initilize a pseudo-random number generator
34     srand( time(0) );
```

```

35  int dimension = 32;
36  int numARows, numAColumns, numBRows, numBColumns, numCRows, numCColumns;
37  // Read given dimensions
38  numARows = atoi(argv[1]);
39  numAColumns = atoi(argv[2]);
40  numBColumns = atoi(argv[3]);
41  // Compute the remaining dimensions for given ones
42  // @TODO@
43  numBRows = numAColumns;
44  // @TODO@
45  numCRows = numARows;
46  // @TODO@
47  numCColumns = numBColumns ;
48
49  printf("Matrix multiplication dimensions: [%d;%d] = [%d;%d] x [%d;%d]\n",
50         numCRows, numCColumns, numARows, numAColumns, numBRows, numBColumns);
51  // host pointers
52  float *host_a, *host_b, *host_c;
53  // Device pointers
54  float *dev_a, *dev_b, *dev_c;
55
56  // Allocations on host
57  host_a = (float *) calloc (numARows*numAColumns, sizeof(float));
58  host_b = (float *) calloc (numBRows*numBColumns, sizeof(float));
59  host_c = (float *) calloc (numCRows*numCColumns, sizeof(float));
60
61  // Initialize vectors
62  init (host_a, host_b, numARows, numAColumns, numBRows, numBColumns);
63
64  // Allocations on device
65  // @TODO@ : complete device allocations
66  cudaMalloc(&dev_a, numARows*numAColumns * sizeof(float));
67  cudaMalloc(&dev_b, numBRows*numBColumns * sizeof(float));
68  cudaMalloc(&dev_c, numCRows*numCColumns * sizeof(float));
69
70  // Copy from host to device
71  // @TODO@ : complete copy from host to device
72  cudaMemcpy(dev_a, host_a, numARows*numAColumns * sizeof(float),
73             cudaMemcpyHostToDevice);
74
75  cudaMemcpy(dev_b, host_b, numBRows*numBColumns * sizeof(float),
76             cudaMemcpyHostToDevice);
77
78  // Invoke kernel
79  // @TODO@ : complete compute grid and block dim
80  dim3 DimGrid ((numARows-1)/dimension+1, (numBColumns-1)/dimension+1, 1);
81  dim3 DimBlock (dimension, dimension, 1);
82
83  // Initialize C device data
84  cudaMemcpy(dev_c, 0, numARows * numBColumns * sizeof(float));

```

```

84 // Call the kernel
85 // @TODO@ : complete to call the kernel
86 dgemm<<< DimGrid, DimBlock >>>(dev_a, dev_b, dev_c, numARows, numAColumns,
    numBRows, numBColumns);
87
88 // Copy result from device to host
89 // @TODO@ : complete copy from device to host
90 cudaMemcpy(host_c, dev_c, numCRows*numCColumns*sizeof(float),
    cudaMemcpyDeviceToHost);
91
92 // Check result
93 check(host_a, host_b, host_c, numARows, numAColumns, numBRows, numBColumns);
94
95 // Free device memory
96 // @TODO@ : complete to deallocate memory
97 cudaFree(dev_a); cudaFree(dev_b); cudaFree(dev_c);
98 cudaFree(host_a); cudaFree(host_b); cudaFree(host_c);
99
100
101 return 0;
102 }

```

I went on and tried some different examples and the results came as follow:

The screenshot shows a terminal window with the following content:

```

gmaroun@scinfe058: /import/etud/3/gmaroun/Bureau/stockage/Semestre 3/GPGPU/Chap4/
Fichier  Edition  Affichage  Rechercher  Terminal  Aide
Matrix multiplication dimensions: [3;3] = [3;3] x [3;3]
Ok
gmaroun@scinfe058: /import/etud/3/gmaroun/Bureau/stockage/Semestre 3/GPGPU/Chap4/
Ex1$ ./1-basicMatMul 3 4 3
Matrix multiplication dimensions: [3;3] = [3;4] x [4;3]
Ok
gmaroun@scinfe058: /import/etud/3/gmaroun/Bureau/stockage/Semestre 3/GPGPU/Chap4/
Ex1$ ./1-basicMatMul 5 4 3
Matrix multiplication dimensions: [5;3] = [5;4] x [4;3]
Ok
gmaroun@scinfe058: /import/etud/3/gmaroun/Bureau/stockage/Semestre 3/GPGPU/Chap4/
Ex1$ ./1-basicMatMul 3 4 4
Matrix multiplication dimensions: [3;4] = [3;4] x [4;4]
Ok
gmaroun@scinfe058: /import/etud/3/gmaroun/Bureau/stockage/Semestre 3/GPGPU/Chap4/
Ex1$ ./1-basicMatMul 3 3 4
Matrix multiplication dimensions: [3;4] = [3;3] x [3;4]
Ok
gmaroun@scinfe058: /import/etud/3/gmaroun/Bureau/stockage/Semestre 3/GPGPU/Chap4/
Ex1$ ./1-basicMatMul 3 5 4
Matrix multiplication dimensions: [3;4] = [3;5] x [5;4]
Ok
gmaroun@scinfe058: /import/etud/3/gmaroun/Bureau/stockage/Semestre 3/GPGPU/Chap4/

```

Questions

1. Before coding : What are the relations between the three matrices dimensions to have a well defined multiplication ?

To multiply a matrix by another matrix we need to do the "dot product" of rows and columns.
For example :

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \times \begin{bmatrix} 7 & 8 \\ 9 & 10 \\ 11 & 12 \end{bmatrix} = \begin{bmatrix} 58 & 64 \end{bmatrix}$$

We apply the dot product on our matrix A of 2 rows and 3 columns and B of 3 rows and 2 columns and insert the resulted number in the matrix C

2. . How many floating operations are being performed in your matrix multiply kernel ? EXPLAIN.

There are $nARows * nBCols(2 * nACols)$ floating operations being performed. By having to read each row in the matrix A and each column in the matrix B which makes the $nARows * nBCols$ then multiplying the rows of A with the cols of B before adding it to the Cvalue , which the kernel do it $2 * nACols$ times

3. How many global memory reads are being performed by your kernel ? EXPLAIN.

There are $nARows * nBCols(2 * nACols)$ reads from A and B that are being performed by the kernel as explained in quest 2

4. How many global memory writes are being performed by your kernel ? EXPLAIN.

There are $nARows * nBCols$ writes in C that are being performed by the kernel represented by the number of times the CValue was written into the C matrix.

5. Compute the arithmetic intensity of your kernel. The arithmetic intensity is a FLOP/Byte number standing for the number of floating point operations performed per byte of global memory accessed.

$$\frac{nARows * nBCols(2 * nACols)}{(nARows * nBCols(2 * nACols) + nARows * nBCols)}$$

$$\Rightarrow \frac{nARows * nBCols(2 * nACols)}{nARows * nBCols((2 * nACols) + 1)} \Rightarrow \frac{2 * nACols}{2 * nACols + 1} (FLOP/Byte)$$

La fin.