

CS 531 - HW3

Ege Alpay

May 2020

1 The Code

In this homework, we were asked to implement Futhoshiki solver which will run on GPU. Running code on GPU is not straightforward as running it on CPU. Data should be transfered between CPU and GPU, before and after the execution.

```
1  int *d_grids;
2  int *d_constraints;
3  int *d_constraint_sizes;
4
5  cudaMalloc((void **) &d_grids, sizeof(int) * no_grids *
    ↪  SIZE * SIZE);
6  cudaMalloc((void **) &d_constraints, sizeof(int) *
    ↪  NUM_CONSTRAINTS_ELEMENTS);
7  cudaMalloc((void **) &d_constraint_sizes, sizeof(int) *
    ↪  no_grids);
8
9  cudaMemcpy(d_grids, grids_1d, sizeof(int) * no_grids * SIZE
    ↪  * SIZE, cudaMemcpyHostToDevice);
10 cudaMemcpy(d_constraints, constraints_1d, sizeof(int) *
    ↪  NUM_CONSTRAINTS_ELEMENTS, cudaMemcpyHostToDevice);
11 cudaMemcpy(d_constraint_sizes, constraint_sizes,
    ↪  sizeof(int) * no_grids, cudaMemcpyHostToDevice);
12
13 int threadPerBlock = 64;
14 int blockCount = (no_grids / threadPerBlock) + 1;
```

```

15     solve << < blockCount, threadPerBlock >> > (d_grids,
    ↪     d_constraints, d_constraint_sizes, no_grids);
16     cudaDeviceSynchronize();
17
18     cudaMemcpy(grid_1d, d_grids, sizeof(int) * no_grids * SIZE
    ↪     * SIZE, cudaMemcpyDeviceToHost);
19     cudaDeviceSynchronize();
20
21     cudaFree(d_grids);
22     cudaFree(d_constraints);
23     cudaFree(d_constraint_sizes);

```

In the code above, you can find the way of running code on GPU. In the beginner code, we had *grids* variable as 3D array and *constraints* variable as 2D array. I have read that transferring these types of array are not recommended, so I have flatten these arrays into 1D. First, memory allocation is done, then data is transferred from CPU to GPU. Later on, kernel is called, you can find the implementation below, then data is transferred back to CPU from GPU. Lastly, memory is deallocated on GPU.

```

1  __global__ void solve(int *grids, int *constraints, int
    ↪  *constraint_sizes, int gridCount) {
2      int futoshiki[25];
3      int constraintSizeForPuzzle;
4
5      int globalId = blockIdx.x * blockDim.x + threadIdx.x;
6
7      if (globalId < gridCount) {
8          constraintSizeForPuzzle = constraint_sizes[globalId];
9
10         int localConstraints[60];
11
12         for (int i = 0; i < 25; i++) {
13             futoshiki[i] = grids[globalId * 25 + i];
14         }
15
16         int constraintStartIndex = 0;
17         for (int i = 0; i < globalId; i++) {

```

```

18         constraintStartIndex += constraint_sizes[i] * 4;
19     }
20
21     for (int i = 0; i < constraintSizeForPuzzle; i++) {
22         localConstraints[4 * i] =
23             ↪ constraints[constraintStartIndex + 4 * i];
24         localConstraints[4 * i + 1] =
25             ↪ constraints[constraintStartIndex + 4 * i + 1];
26         localConstraints[4 * i + 2] =
27             ↪ constraints[constraintStartIndex + 4 * i + 2];
28         localConstraints[4 * i + 3] =
29             ↪ constraints[constraintStartIndex + 4 * i + 3];
30     }
31
32     solveSingleThread(futoshiki, localConstraints,
33         ↪ constraint_sizes[globalId]);
34
35     for (int i = 0; i < 25; i++) {
36         grids[globalId * 25 + i] = futoshiki[i];
37     }
38 }
39 }

```

Function *solve* is the first function triggered on GPU. Global Id of a thread is calculated first. Then, each thread will get it's puzzle and corresponding constraints with respect to their global id. Later on, each thread will call *solveSingleThread* function by passing it's puzzle and constraints. This function is basically solves the given puzzle with backtracking. However, it does not use recursion since I have read that it is not so efficient, and memory problems occurs because of the depth. Lastly, the solved puzzle will be written back to *grids* variable. All these indexing operations are calculated based on global id of the thread.

To sum up, each thread will be working on a unique puzzle where number of the threads is equal to number of puzzles.

2 Execution times

Duration for Memory Allocation & Deallocation in GPU and data movement between GPU & CPU is not dependent to block size or thread count. You can find the duration in seconds, in the following table:

Operation	Execution Time in Seconds
Memory Allocations for GPU	0.0007
CPU to GPU Data Transfer	0.014
GPU to CPU Data Transfer	0.008
Memory Deallocation for GPU	0.004

In Cuda Programming, a block can have maximum 1024 threads. For that reason, I have tried different number of threads to observe execution time.

Number of Threads Per Block	Kernel Execution Time in Seconds
1	4.972
2	3.569
4	2.820
8	2.473
16	2.290
32	2.200
64	1.847
128	1.853
256	1.867
512	1.889
1024	1.906

Fastest execution time was obtained where 64 Threads were used in a block where number of blocks is equal to $no_grids/64$.

3 Tricks Done for Efficiency

Transferring non 1D array needs multiple *cudaMemcpy* calls. In order to reduce number of calls of *cudaMemcpy*, these types of arrays are flattened into 1D arrays. So, we were able to transfer data with only one *cudaMemcpy* call.

4 Compile and Run

For development, I have used CLion IDE to write the code. For compilation and testing the implementation, I have used gcc 9.2.0 on my Linux Mint machine. On Gandalf, gcc 4.4.7 with cuda 10.0 was used.

Following command was used to run the code on Gandalf:

```
1      module load cuda/10.0.  
2      nvcc solution.cu  
3      ./a.out input.txt
```