

Set 9 - GPUS II

Issued: May 8, 2017

Question 1: Diffusion on GPUs - Part II

In this exercise we will improve the 2D diffusion code on GPUs.

- a) Improve the 2D diffusion kernel of the previous exercise by employing the shared memory of the GPU.
- b) Write a kernel for `GetMoment()` to further reduce the memory transfer needed between GPU and CPU by calculating (at least) partial sums on the GPU.

Hint: An important part of the operation is a reduction.

While you can spend a full talk on how to optimize reductions¹, a simple GPU reduction will do here. It is called only every 100th update and will not dominate the execution time of our code. A simple reduction scheme within a block using the shared memory can be achieved by selecting only specific thread indices:

```
1  __shared__ float data[8];
2
3  ...
4  if(threadIdx.x < 4)
5      data[threadIdx.x] += data[threadIdx.x + 4];
6  __syncthreads();
7  if(threadIdx.x < 2)
8      data[threadIdx.x] += data[threadIdx.x + 2];
9  __syncthreads();
10 if(threadIdx.x < 1)
11     data[threadIdx.x] += data[threadIdx.x + 1];
```

- c) Think about good choices for `blocksPerGrid` and `threadsPerBlock` of your kernels and explain your choice.

¹M. Harris, Optimizing Parallel Reduction in CUDA, 2007,
<http://developer.download.nvidia.com/assets/cuda/files/reduction.pdf>