

Report II

Large data structures

Memory Managed utility and how to protect program from too large data

Paweł Lipiór, Michał Żoczek

October 30, 2019

1 How large data structures are handled

During computer calculations often deal with large data structures. Computer is a great tool to do these calculations but we have to remember that it also has limited hardware resources. During our test with `vectorAdd` sample we have encountered some of these problems. Firstly, vector length has been limited by `Int` maximum size so we have been forced to use `long unsigned Int`. In some cases we could be limited by system memory capacity but we have been using simple dynamically allocated table of floats so our resources were enough.

```
1 unsigned long long numElements = 1<<30;
2 unsigned long long size = numElements * sizeof(float);
```

However the most limited resource often is device memory capacity. In case we suspect data structure is too large to handle it whole at the same time in GPU memory we can split data into smaller ones and do calculations in steps. If we have access to Pascal or later GPU we can also use management memory to automate this process.

2 Using Managed Memory

While programming in CUDA we can see that memory host and device management is tedious and in some cases hard to do. It takes time to write a right code, test and debug it. Fortunately it can be avoided to some extent by using unified memory – a memory address space that is accessible by both host and device. We can just replace all functions responsible for allocate host and device memory and copying data between these by on simple function `cudaMallocManaged`. Below there is comparison of these two methods.

```

1  // Allocate the host input vector A
2  float *h_A = (float *)malloc(size);
3
4  // Allocate the host input vector B
5  float *h_B = (float *)malloc(size);
6
7  // Allocate the host output vector C
8  float *h_C = (float *)malloc(size);
9
10 // Allocate the device input vector A
11 float *d_A = NULL;
12
13 // Allocate the device input vector B
14 float *d_B = NULL;
15
16 // Allocate the device output vector C
17 float *d_C = NULL;
18
19 // Copy the host input vectors A and B in host memory to the device
    input vectors in
20 // device memory
21 err = cudaMemcpy(d_A, h_A, size, cudaMemcpyHostToDevice);
22 }

```

```

1  // Allocate the vector A
2  float *A;
3  checkCudaErrors(cudaMallocManaged(&A, size));
4
5  // Allocate the vector B
6  float *B;
7  checkCudaErrors(cudaMallocManaged(&B, size));
8
9  // Allocate the vector C
10 float *C;
11 checkCudaErrors(cudaMallocManaged(&C, size));

```

This way it is possible to simplify the code and avoid mistakes and reduce time needed to release final and fully-featured calculation tool. On the other hand `cudaMallocManaged` can not be smarter than the programmer and we have to take care that all process are fine. E.g we can not allocate more memory than the device hardware has.

3 Protecting program from too large data

Another troublesome aspect of working with GPU is managing resources, especially memory. While using managed memory makes it easier, it does not mean that it takes care of everything. One such aspect that is left for programmer is to make sure that data used during kernel execution will not exceed available

GPU memory. In program used for this exercise it was achieved by first checking GPU global memory using `cuDeviceGetProperties` function and then comparing obtained number of bytes with size of data.

```
1 if(deviceProp.totalGlobalMem < size)
2     {
3         fprintf(stderr, "Data size too large\n");
4         exit(EXIT_FAILURE);
5     }
```

Owing to these simply if statement there is no possibility to exceed the size of device memory. There is also limitation of number of blocks that can be used at once, but for modern GPUs its so big ($2^{31} - 1$) that it is rarely a problem.