

Hardware-oriented optimisation in CUDA

University of Regensburg



Tobias Sizmann

day month year

Abstract

Contents

1	Introduction	4
1.1	From single-core CPUs over multi-core CPUs to GPUs	4
1.2	What is CUDA	4
2	Tree Reduction on GPUs	7
2.1	The importance of reductions	7
2.2	The tree reduction algorithm	7
2.3	Naive implementation with CUDA	8
2.3.1	Serial tree reduction on a CPU	8
2.3.2	Tree reduction in CUDA for small arrays	10
2.3.3	Tree reduction in CUDA for arbitrary array sizes	12
2.4	Conclusion	14
A	Appendix One	15

Chapter 1

Introduction

1.1 From single-core CPUs over multi-core CPUs to GPUs

With the invention of the MOS transistor in 1959 and more specifically the silicon-gate MOS transistor in 1968 the development of single-chip microprocessors started. The first processors were single-threaded and their bottleneck was (besides other things) their clock rate. However, clock rates increased exponentially over time and subsequently did the computing power. In the 2000s this development was eventually brought to a halt at around 3.4 GHz due to thermodynamical limits of silicon. While pushing past that limit is possible, the extra costs for cooling usually outweigh the increase in computing power. This is the beginning of the multiprocessing era. The idea is simple: When one thread cannot run faster, just increase the number of threads. In 2006 the first desktop PCs with two cores were sold and ever since the number of threads is steadily increasing. There was one problem in particular, however, that CPUs (Central Processing Unit) were not efficient in solving, namely, rendering graphics. This task required a lot of independent and small calculations that needed to be done in real time - a prime example for a massively parallelisable computation. Since rendering graphics required such a different type of computing, the first GPUs (Graphics Processing Unit) were developed. These GPUs featured less single thread computation power but have a higher number of threads. State of the art GPUs have thousands of threads. This being said, the number of threads cannot be easily compared between a CPU and a GPU or even between different GPUs as they follow different design paradigms. How the GPU threads work exactly will be explored in this work. Generally, one can say that GPUs perform well in massively parallelisable computations where the single computations are not complex while CPUs excel at complicated single threaded problems (for example running the event loop of a desktop application). This shift to a higher number of threads rather than thread quality has been greatly motivated by scientific calculations, machine learning and graphics.

1.2 What is CUDA

Writing code for a GPU is not as straight forward as for a CPU since it is more dependent on its hardware. Different GPU developers use different application pro-

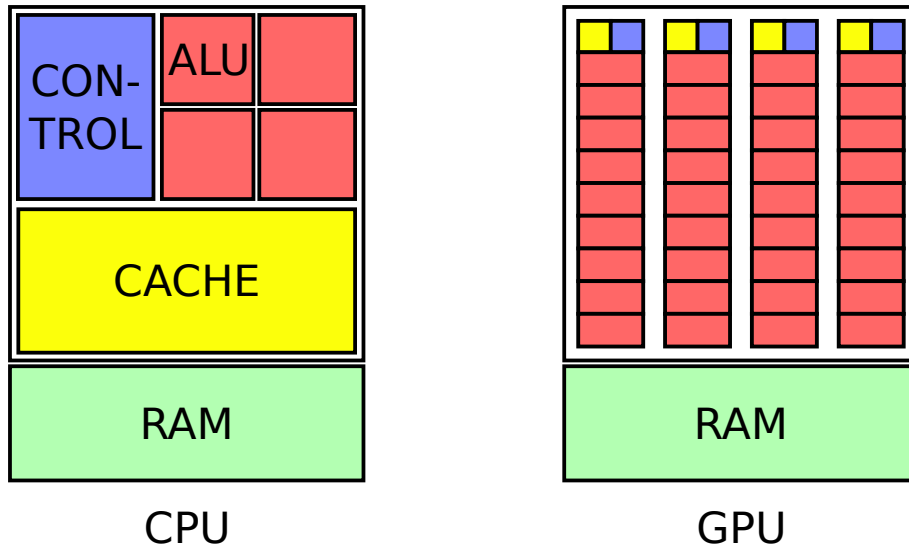


Figure 1.1: Comparison of the basic architectural differences of a GPU and a CPU. In green the random access memory (RAM), in yellow the cache, in blue the flow control unit and in red the arithmetic-logic unit (ALU). Both designs feature a RAM that all threads have access to. The GPU is split into many small "CPUs" (vertical groups) with their own flow control and cache. These are called streaming multiprocessors (SMPs or SMs). Generally the cache hierarchy is more complicated as depicted in the figure. However, the defining property here is that there exists a non-global cache level that is assigned to a group of ALUs, namely the SMs. Note that the notion of an ALU has slightly different meaning for a CPU and a GPU. For a CPU one ALU usually corresponds to one thread. For a GPU one ALU usually corresponds to a group of threads (for modern GPUs: 32), called a warp.

gramming interfaces (APIs). The biggest GPU designer Nvidia developed a framework called CUDA, which will be used in this work. It is simple to learn, offers efficient implementations and a plethora of literature and support. The downside is that it only supports Nvidia GPUs and is not open-source. An alternative would be OpenCL, which is open-source but harder to learn.

CUDA works as an extension to the C programming language and comes with its own compiler. Sections of code are distributed to either the CPU (called host) or the GPU (called device). When only writing code for the host normal C is used. Writing device code is more sophisticated. Here, the programmer first needs to write a so called kernel, which can be thought of the interior of a for-loop. Then the host must transfer required data to the memory of the GPU and call the kernel. When calling the kernel, the boundaries of the for-loop are set. The loop is then executed in parallel on the GPU. This API only allows parallelisation of for-loops. While this might seem like a very strict limitation, it closely relates to how a Nvidia GPU works.

The GPU can be thought of being organised in streaming multiprocessors (SMPs or SMs), warps and threads. Warps will be explained in more detail later. A SM groups together a set of (hardware) threads and allows synchronisation between them. Threads of different SMs cannot be synchronised. Also, threads within a SM share a cache, which cannot be accessed by threads of another SM. The domain of the for-loop is organised in blocks and (CUDA) threads, where the threads are grouped in blocks. The blocks are executed by the SMs, which means that only threads within a block can be synchronised and have access to the same cache. Note the differentiation of hardware and CUDA threads. This is an unfortunate ambiguity in the terminology of CUDA. While each CUDA thread is mapped to one hardware thread, a hardware thread can be mapped to (i.e. execute) several, none or one CUDA threads. This will become clearer in the next chapter.

Chapter 2

Tree Reduction on GPUs

2.1 The importance of reductions

A reduction in terms of parallel programming is an operation, where a large dataset of entries (e.g. numbers) is reduced to one entry. The simplest example is the sum of a set of numbers and will be used for the rest of this work. More generally, a reduction is defined by an operation $\circ : X \times X \rightarrow X$ on two entries with the following properties:

$$a \circ b = b \circ a \quad (\text{commutativity}), \quad (2.1)$$

$$a \circ (b \circ c) = (a \circ b) \circ c \quad (\text{associativity}). \quad (2.2)$$

This guarantees that the result is (mathematically) independent of the order in which the elements are reduced. Note, that these properties do not ensure numerical stability.

The reduction operation, first and foremost the sum, plays a crucial role in all of numerics. Simple linear algebra operations like the scalar product or the matrix multiplication already include a reduction: The entries of two vectors are multiplied element wise and the summed up. In machine learning, reductions are present as a key step in feed-forward networks: Again all the values of the node one layer below are multiplied by weights elementwise and the summed up to calculate the value of a singular node above. Reductions are a very basic and fundamental operation and, therefore, an efficient implementation is required.

2.2 The tree reduction algorithm

The most efficient algorithm is heavily dependent on the underlying hardware. For example, a node network with a lot of parallelisation overhead and a complicated communication topology might use a ring algorithm ("add value and pass to next"). In this work the focus is on single GPU reductions. Here, the tree reduction approach is the most successful. All available threads are called to reduce two entries to one in parallel, effectively halving the size of the dataset in one step. This is repeated until only one entry remains.

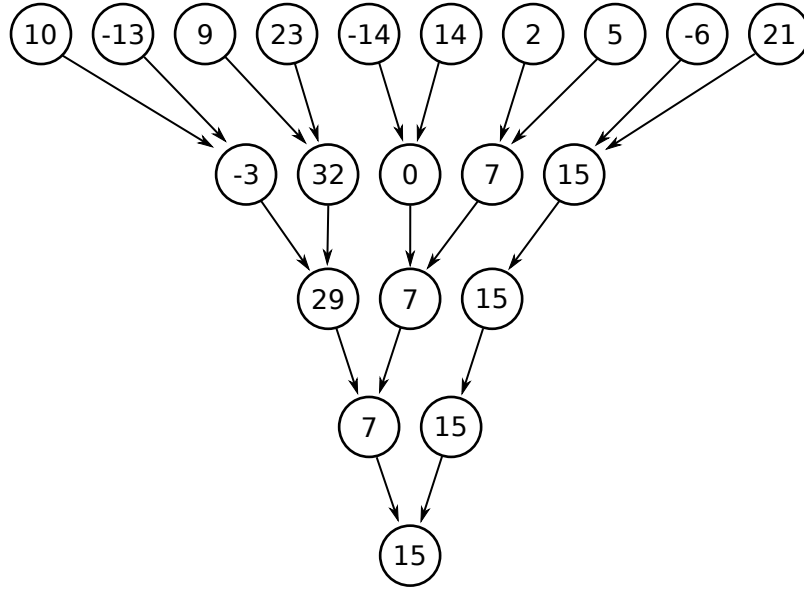


Figure 2.1: Example of a tree reduction of ten elements. Note how there is a leftover after the second reduction. These are usually handled by zeropadding (i.e. adding zeros) after each step to make the number of elements divisible by 2.

2.3 Naive implementation with CUDA

For a better understanding of the framework of CUDA and tree reduction algorithm an unoptimized code is presented which implements the algorithm for the case of addition of signed 32-bit integers. Note, that even though this is unoptimized GPU code, it runs magnitudes faster than on a CPU (exact speedup depending on the problem size and the available hardware).

When approaching a problem using CUDA the big challenge is to map the for-loop that one wants to parallelize to blocks and CUDA-threads. Often, the best starting point is to write serial CPU code to get a feeling for the problem and to have a working solution to test the optimized solutions against later. In this case, we will first implement a serial tree reduction in C and port this code to CUDA in a second step.

2.3.1 Serial tree reduction on a CPU

Our starting point is an integer array `h_in`. The `h_` denotes data stored on the host (CPU RAM), which is a useful convention, since CUDA does not distinguish between pointers to host and device memory. For simplicity, we assume that the size of the array is a power of two. If this is not the case, one could simply zeropad the data and would still be able to use the code.

Lets write a function `reduce` that executes the tree reduction. This function takes an input array `h_in`, performs the reduction and writes the result to an output `h_out`. We need two for-loops: One to iterate over the step of the reduction (vertically in fig. 2.3.1) and one to iterate over the (remaining) data set (horizontally). A temporary array is used to do the reduction on, so the input array stays

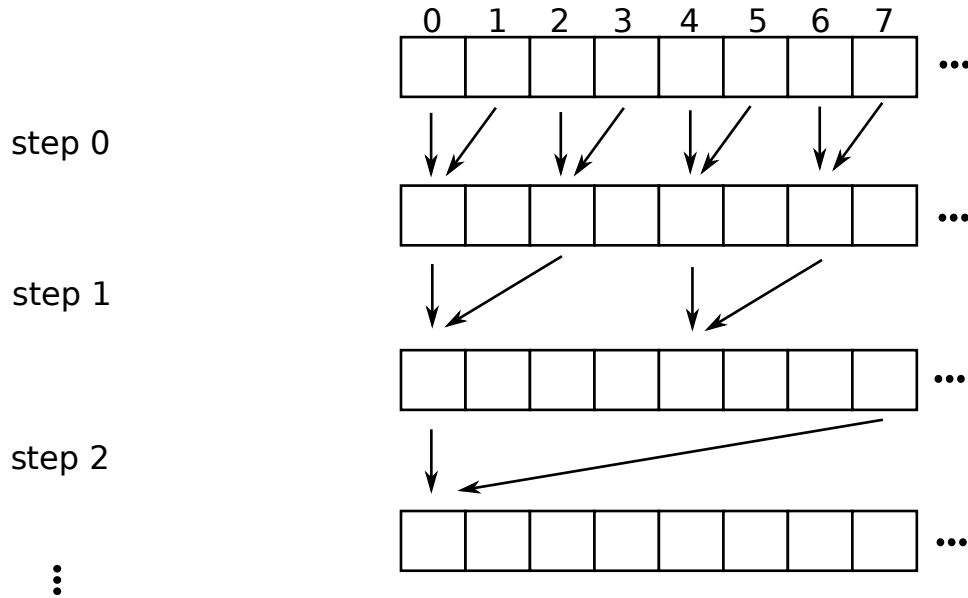


Figure 2.2: Sketch of the implementation of the tree reduction algorithm presented in this section. The cells denote entries in the data array. Two arrows pointing into a cell denote addition of the source cells values and writing of the result into the target cell. The final result is found in the first cell of the array.

untouched.

```

1 void reduce(int * h_in, int * h_out, int len) {
2     // copy data into temporary array
3     int * temp = (int *) malloc(len*sizeof(int));
4     memcpy(temp, h_in, len*sizeof(int));
5
6     int stride = 1;
7     // iterate over reduction steps (vertically)
8     for (int step = 0; step < log2(len); step++) {
9
10        // iterate over dataset (horizontally)
11        for (int i = 0; i < len; len += 2*stride) {
12            temp[i] = temp[i] + temp[i + stride];
13        }
14        stride *= 2;
15    }
16
17    // return the result
18    *h_out = temp[0];
19    free(temp);
20    return;
21 }

```

The inner loop of this implementation could be parallelised easily, since each step of the inner loop is independent. This means, however, that with each step of the outer loop all threads are created and destroyed, which leads to a large parallelisation overhead. In general and independent of CPU or GPU, one should always aim to parallelize the outermost loop. Here, the outer loop cannot be parallelized, since each step is dependent on the result of the step before. To fix this, one can swap the

two loops. While this requires slightly more code, it greatly simplifies all examples in the following. The rewritten function looks like this:

```

1 void reduce(int * h_in, int * h_out, int len) {
2     // copy data into temporary array
3     int * temp = (int *) malloc(len*sizeof(int));
4     memcpy(temp, h_in, len*sizeof(int));
5
6     // iterate over dataset (horizontally)
7     for (int i = 0; i < len; len += 1) {
8         int stride = 1;
9
10        // iterate over reduction steps (vertically)
11        for (int step = 0; step < log2(len); step++) {
12            if (i % (2 * stride) == 0) {
13                temp[i] = temp[i] + temp[i + stride];
14            }
15        }
16
17        stride *= 2;
18    }
19
20    // return the result
21    *h_out = temp[0];
22    free(temp);
23    return;
24 }

```

Note that an additional if-statement is required. This is a much better basis for parallelisation, since now the outer loop can be parallelised. However, one needs to be very careful as this is now prone to a race-condition. The parallelisation with CUDA is done in the next section.

It should be noted, that there are much faster ways to implement a reduction algorithm as a single-thread CPU application. This code has been written with the intent of parallelisation on a GPU and serves solely this purpose.

2.3.2 Tree reduction in CUDA for small arrays

Once the algorithm has been successfully implemented serially, the parallelisation of the target for-loop is always the same. One needs to map the domain of the for-loop (here integers from 0 to `len-1`) to threads. The maximum number of threads n_{Threads} a block can contain depends on the hardware, but is always a power of two. Modern Nvidia GPUs will allow for $n_{\text{Threads}} = 1024$, but it is not necessarily optimal to use all threads. This will be explored later. For now we will restrict the maximum size of our input array to 1024, such that only one block is required. With this, one can write down the kernel (terminology for a function running on the GPU).

```

1 __global__ void reduce(int * d_in, int * d_out, int len) {
2
3     // prepare shared data allocated by kernel invocation
4     // and copy input array
5     extern __shared__ int temp[];
6     temp[threadIdx.x] = d_in[threadIdx.x];
7     __syncthreads();
8
9     // do treereduction in interleaved addressing style

```

```

10     int stride = 1;
11     for (int step = 0; step < log2(len); step++) {
12
13         if (threadIdx.x % (2*stride) == 0) {
14             temp[threadIdx.x] += temp[threadIdx.x+stride];
15         }
16         __syncthreads();
17         stride *= 2;
18     }
19
20     // export result to global memory
21     if (threadIdx.x == 0) {
22         *d_out = temp[0];
23     }
24 }

```

There are several things to uncover here. First, the `__global__` void declares the function as a kernel, which runs on the device but can be invoked from the host. Secondly there is a new constant `threadIdx.x` available within the kernel. This is an identifier of the thread that is executing the kernel, is unique for all threads within a block and ranges from 0 to `numThreadsPerBlock - 1`. This replaces the index which the for-loop iterated over. The temporary array can be replaced by the ultra-fast cache shared by all threads of one block, which is allocated during kernel invocation and declared within the kernel by the line

```

1 extern __shared__ int temp[];

```

As mentioned earlier, there is a race condition between threads, which can be solved by the `__syncthreads()` method. This method acts as a wait-for-all barrier within a block (synchronisation between blocks is not possible!). Finally, the result needs to be exported. For this, only one thread is required, hence the if-statement. Note that the input and output pointer names start with a `d_`, which is not required but is a convention to mark, that these pointers live in the address-space of the device.

The kernel can be invoked from the host with the following code:

```

1 // allocate and copy to memory of device
2 int arraySize = 1024;
3 int * d_in, * d_out;
4 cudaMalloc(&d_in, sizeof(int)*arraySize);
5 cudaMalloc(&d_out, sizeof(int));
6 cudaMemcpy(d_in, h_in, sizeof(int)*arraySize,
7             cudaMemcpyHostToDevice);
8
9 // invoke kernel with the correct amount of threads and cache space
10 int numThreadsPerBlock = len;
11 int numBlocks = 1;
12 reduce <<< numBlocks, numThreadsPerBlock, sizeof(int)*
13             numThreadsPerBlock >>> (d_in, d_out, len);
14
15 // copy result to host and free memory
16 cudaMemcpy(h_out, d_out, sizeof(int), cudaMemcpyDeviceToHost);
17 cudaFree(d_in);
18 cudaFree(d_out);

```

First the data is copied to the device memory. Then the kernel is invoked by the triple angled brackets syntax. Within the brackets, the number of blocks n_{Blocks} (so far exactly one), the number of threads per block n_{Threads} and the amount of

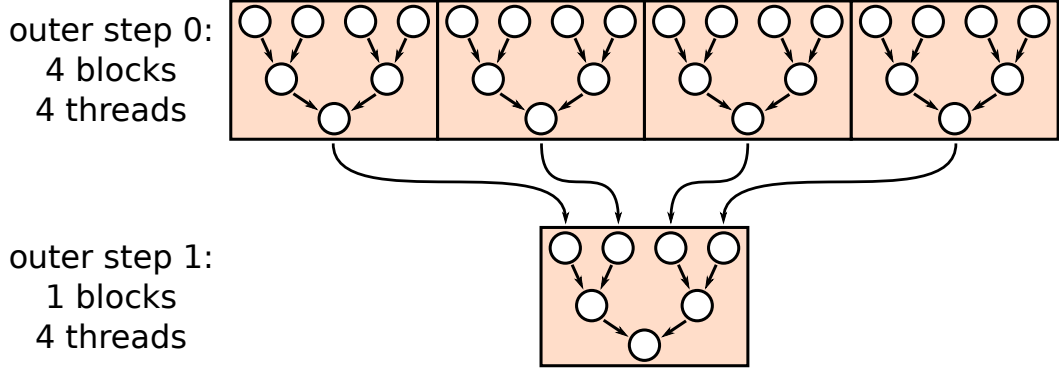


Figure 2.3: Depicted is the schematic tree reduction in block structure for an array of length 16 and 4 threads per block. The first step requires 4 blocks and the second 1 block. The blocks run on the SMs of the GPU. The first step allows for four SMs and a total of 16 physical threads to be used in parallel. Between the two steps the output array is used as new input array. The second step can only make use of a single SM.

cache space are specified. The kernel arguments are specified in parenthesis. Finally the result is copied back to the host and the device memory is free'd. Note that in practice one would write more code to error-check every step, assure that the correct upstream is used and optimize the copy procedure. Nonetheless, this is a minimal working example.

The big problem with this implementation is, that it is limited to an array size of 1024, since we only use one block. This is solved in the next section.

2.3.3 Tree reduction in CUDA for arbitrary array sizes

The limitation to an array size n_{Data} of maximum block size can be solved in the following way: Split the array into equal chunks and run the kernel on each of these chunks individually. This procedure returns an array with size $n_{\text{Data}}/n_{\text{Threads}}$ (Assuming that n_{Data} is a power of 2, otherwise round up is required). The procedure is then again used recursively, until the array is reduced to a singular value. We refer to such a reduction step as "outer step", which reduces n_{Threads} elements, compared to an "inner step" reducing 2 elements.

While this could be implemented already with the existing kernel, there is a hardware structure available to do kernel invocations with a certain amount of threads in parallel, namely the streaming multiprocessors (SMs). The SMs are a second layer of parallelisation. One block of threads is always mapped to one SM, so several blocks can be executed in parallel on several SMs. Unlike the number of threads per block, the number of blocks is unlimited and blocks are queued until a SM is free to work on it. In our case this has the following implication. Lets assume that the initial array size is $n_{\text{Data}} = n \times 1024$. Then one could use n blocks of 1024 threads to do the first step of the outer reduction and use the full computational power of the GPU this way.

This leaves us with a new problem, however. Namely, how does one calculate the position of the array, one specific thread is supposed to work on? To this end,

CUDA offers two more constants in the kernel: `blockDim.x` and `blockIdx.x`. The first one is simply the number of threads in a block and the second one is a unique identifier of the current block, ranging from 0 to $n_{\text{Blocks}} - 1$. The array position can then be calculated with `threadIdx.x + blockIdx.x * blockDim.x`. The required amount of blocks for the kernel invocation can be calculated with rounding up division. This can be neatly done with $n_{\text{Blocks}} = (n_{\text{Data}} + n_{\text{Threads}} - 1) / n_{\text{Threads}}$, where the `"/` denotes integer division.

To enable the use of blocks with our kernel three small modifications need to be done. First, the access of our input array needs to be rewritten with the new formula for the position. Secondly, the output is not a singular value but an array. Thirdly, the range of the loop and the input array length can be deduced directly from the block size and number of threads per block and, therefore, the length of the input array does not need to be passed to the kernel. The result from i -th block should be written into the i -th position. Also for simplicity the loop over the `step`-variable has been replaced by a loop over `stride` directly, saving an extra variable and the `log2` function call. The finished code looks like this:

```

1 __global__ void reduce(int * d_in, int * d_out) {
2
3     // prepare shared data allocated by kernel invocation
4     // and copy input array
5     extern __shared__ int temp[];
6     temp[threadIdx.x] = d_in[threadIdx.x + blockIdx.x*blockDim.x];
7     __syncthreads();
8
9     // do treereduction in interleaved addressing style
10    for (int stride = 1; stride < blockDim.x; stride *= 2) {
11        if (threadIdx.x % (2*stride) == 0) {
12            temp[threadIdx.x] += temp[threadIdx.x+stride];
13        }
14        __syncthreads();
15    }
16
17    // export result to global memory
18    if (threadIdx.x == 0) {
19        *d_out = temp[blockIdx.x];
20    }
21 }
22 }
```

The invocation changes slightly. First, we need to specify the number of required blocks as argument in the triple angled brackets. Secondly, the kernel has only two arguments now. Also the size of the output array changes.

```

1 // set the number of threads per block and calculate the required
2 // number of blocks
3 int arraySize = sizeof(int)*len;
4 int numThreadsPerBlock = 1024;
5 int numBlocks = (arraySize + numThreadsPerBlock - 1) /
6                 numThreadsPerBlock;
7
8 // allocate and copy into device memory
9 int * d_in, * d_out;
10 cudaMalloc(&d_in, sizeof(int)*arraySize);
11 cudaMalloc(&d_out, sizeof(int)*numBlocks);
12 cudaMemcpy(d_in, h_in, sizeof(int)*arraySize,
```

```
11     cudaMemcpyHostToDevice);  
12 // invoke kernel with the correct amount of threads and cache space  
13 reduce <<< numBlocks, numThreadsPerBlock, sizeof(int)*  
    numThreadsPerBlock >>> (d_in, d_out);  
14  
15 // copy result to host and free memory  
16 cudaMemcpy(h_out, d_out, sizeof(int)*numBlocks,  
    cudaMemcpyDeviceToHost);  
17 cudaFree(d_in);  
18 cudaFree(d_out);
```

Note that this code only executes a single outer step of the reduction. One would need to take the output array and feed it through this code until a single value is left. This was left out since it is mostly host code and not of particular interest for the rest of the work. It should be noted, however, that all code timings in the following were done for the full reduction.

2.4 Conclusion

This concludes the introduction to CUDA and tree reduction algorithms. From now on this work focuses on modifying the kernel to optimize it as much as possible. It should be noted, that a lot of optimization has been done already. For example, using the explicitly declared shared memory instead of DRAM or the implementation of the block structure already offer a great speedup compared to other ways one could come up with to do the tree reduction. The main goal now is to dive deeper into the hardware structure and achieve speedups by fixing problems like warp divergence and memory bank conflicts, which will be explained in full detail. Some algorithmic improvements will also be done and in the end all the optimisations will be benchmarked on several hardware systems.

Appendix A

Appendix One

this will be space for the appendix