

Automatic Code Optimizations on GPU Architectures

Johann Wagner
Otto-von-Guericke University
Magdeburg, Germany
johann.wagner@st.ovgu.de

Abstract—GPU programming is a hot topic. More and more tasks requiring a lot of computing power are performed on graphics cards, such as bitcoin mining. Optimizations that make these applications even more efficient to take full advantage of the power of the graphics cards are all the more important.

After a short introduction to parallel programming and NVIDIA CUDA, this paper compares different approaches to shortening the time a thread waits for memory. There are several optimizations, such as thread block merging, thread merging, data prefetching and the elimination of partition camping. These four optimizations are compared in four different categories.

At the moment there are many implementations that offer improved compilers for CUDA, but few comparisons between different optimizations and their modes of action and lines of action. This paper aims to give an overview of the most important optimizations and to compare them in interesting and simple categories.

Index Terms—GPU Compiler Optimizations

I. INTRODUCTION

General Purpose Computation on Graphics Processing Unit (GPGPU) is a frequently used term in data processing. Due to this, complex calculations can be accelerated by using GPU instead of CPU.

Today's compilers for CPU architectures are optimized a lot. There are many generally known optimizations like loop unrolling and constant precalculation. The aim is to avoid problems of CPUs, such as costly instructions like jump- and division-operations [12].

Loop unrolling tries to get rid of jump-operations, that in the best case accelerates the assembler code by unrolling a loop to procedural code without jump-operations [12]. Constant precalculation tries to evaluate a formula as best as possible without known variables [12].

In general, these optimizations try to avoid cost-intensive operations. These are either avoided or calculated at compile time, so that the actual program runs faster.

Compilers, which are used to create GPGPU applications, have to help the developer more than conventional CPU applications. The main problem with GPGPU applications is memory access. NVIDIA introduces CUDA, a programming model that allows developers to run code on graphics cards, but CUDA introduces six different types of memory, all of which have to be managed by the programmer, which is much more than traditional programming. In addition, this memory access is slow in case of incorrect handling [1] [7].

Obviously, other problems have to be solved than with CPU applications. In GPGPU applications, memory access to global memory is the main problem and a problem that can be solved in different ways. Therefore, we would like to discuss here what optimizations already exist for the GPU [1] [7].

II. BACKGROUND

This section is intended to provide an introduction to the various terms used in this paper. The terms Thread and Thread Blocks and General Purpose Computation on GPUs are explained. In addition, NVIDIA CUDA and the different types of memory introduced by CUDA will be presented to help improve understanding of the topic.

A. Threads and Thread Blocks

1) *Threads*: Threads are subroutines, so a small set of instructions, which can be managed and executed independently from other threads. Threads typically have an independent stack for local variables and share heap with all other threads. Due to access from multiple independent threads, heap access must be managed [13].

2) *Thread Blocks*: Thread Blocks is a construct, which was introduced by CUDA to organize threads in groups. Threads, that are grouped in a thread block, can access a part of the storage simultaneously [7].

3) *Hardware Implementation*: NVIDIA GTX 1080 has 20 multiprocessors with 128 threads each. This means that the threads of 20 thread blocks can be executed simultaneously. If there are more than 20 thread blocks, they are executed one after the other. That means, that 2560 threads can be executed in parallel [5].

B. General Purpose Computation on GPUs

Big Computations need a lot of calculation time. Those computations can be mostly splitted into multiple smaller workloads, which can be calculated independently [2]. While CPUs are designed to calculate bigger workloads with less threads, GPUs are designed to calculate smaller workloads on hundreds of threads. In its original form, each pixel could be calculated individually on a GPU, which implies heavy multi threading [3].

There are some advantages and disadvantages of GPGPU usage:

1) *Threading*: A problem usually contains multiple smaller workloads, that can be calculated independently on GPUs. This calculations can be processed parallel to save time. Modern GPUs have more than 2500 cores with up to 1.7 GHz, which can execute code parallel. Other papers determine speedups up to 20 times [4] [?].

2) *Data Transfer Time*: Transferring data from memory into GPU memory is usually a bottleneck. Calculations, which use data from secondary storage or produce a lot of data, have to transfer data to the host system, that moves the data to CPU memory or on secondary storage.

The problem here is the PCIe connection, which is currently specified with a maximum data rate of 32 GB per second. Thus it takes at least 250ms to fill the entire memory of a current graphics card, which is relatively slow in the context of fast calculations [4] [?].

3) *Complexity*: Complexity is a big difficulty. The greater the degree of complexity, the more difficult and slower the development process. Especially when programming GPUs, the complexity is high and the quality of the programming is decisive for the speed of the application. However, code is often not optimized when the consumer can cope with the waiting time, because the development of accelerated solutions would cost a lot of money [8].

C. NVIDIA CUDA

NVIDIA CUDA is a parallel computing platform and model, which was developed by NVIDIA to speed up applications by usage of NVIDIA GPUs. For example CUDA can be plugged as library into C programs to make use of a CUDA-capable GPU [6].

NVIDIA CUDA, further referred as CUDA, makes it easy to access GPU for general purpose programming by providing functions for different programming languages like C, C++ and Python. It provides simple functions and/or annotations to enable GPU usage.

```
// Kernel definition
__global__ void VecAdd(float* A, float* B, float* C)
{
    int i = threadIdx.x;
    C[i] = A[i] + B[i];
}

int main()
{
    ...
    // Kernel invocation with N threads
    VecAdd<<<1, N>>>>(A, B, C);
    ...
}
```

Fig. 1. Code with CUDA Support to add two vectors, accelerated with GPU - Taken from NVIDIA Programming Guide

In Figure 1, we show a simple vector addition, that is accelerated by using GPU with CUDA. You have to define a kernel function, which is a simple c function, which is called N times by N different CUDA threads. This kernel function is executed on GPU. `threadIdx.x` contains, in this simple example, just the number ($0 \leq \text{threadIdx.x} < N$) of the thread, which called the function. In Figure 2, we show the mathematical equation, which is done by the code in Figure

1. We can calculate every c_i independently in sepeate threads, which makes it possible to speed up the process. This example is trivial, however matrix operations, that are more complex, can be speed up with the same approach [7].

$$\begin{pmatrix} c_1 \\ c_2 \\ \vdots \\ c_m \end{pmatrix} = \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_m \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{pmatrix}$$

Fig. 2. Mathematical Equation for Vector Addition

D. Memory Seperation in CUDA

NVIDIA CUDA introduces six different types of memory, which are up to four more types than in other programming languages. Different types of memory has different attributes.

1) *Register*: Register Memory is quite similar to CPU Register Memory and is used as storage for local variables, etc. It is limited to 16kb and everything, which exceeds this limit, will be pushed on Local Memory [9].

2) *Local Memory*: Local Memory is an abstraction of CUDA, which makes it possible to store information in Global Memory, which is only available from current thread [9].

3) *Shared Memory*: Shared Memory is much faster than the Global Memory. It is visible to all threads, which are in the same thread block. It is intended to share data across multiple threads, that are in the same thread block [8] [7].

4) *Global Memory*: Global Memory is the slowest memory. It is visible to all threads, can be read and written from all threads and can be as General Purpose memory. It is very similar to CPU memory, which makes it easy to use to share data across multiple thread blocks [8].

5) *Constant Memory*: Constant Memory is quite similar to the Global Memory, but can only be read. It can help to reduce overhead by caching information [8] [9].

6) *Texture Memory*: Texture Memory is quite similar to the Global Memory, but can only be read. It provides some caching methods, which can be used to accelerate certain applications, but this is not used in the further explanations [7].

III. BODY

Waiting for Memory is one of the most time-consuming states in executing a program. Compared to CPU programming time can be saved in GPU programming by utilizing different memory types, which were introduced in the prior section. Fundamentally data, that is created and used by GPU threads, could be read and written in the global memory. This native approach would be slow, because only one thread could use the data at a time.

It would be better if data that is only used by the thread would also be stored in the memory that can only be used by the thread to reduce memory accesses to global memory and to shorten the wait time for other memory requests for

global memory. Unfortunately it is not possible to determine the usage of different variables at compile time, because we are not able to test a complex program deterministically, and as a consequence it is not possible to optimize the memory access time this way.

The primary objective is to keep the time for memory access as short as possible. There are various optimizations for this, such as thread block and thread merging, data prefetching and the elimination of partition camping.

A. Thread-Block Merging

Thread merging is a technique where different thread blocks are merged into one thread block to reduce memory access times.

Thread Blocks have a shared memory segment, where they can save data for execution. If multiple thread blocks use the same data from the global memory to execute its threads, we can merge them into one thread block to reduce the memory access for global memory [2].

In Figure 3, we show an example for merging two thread blocks. Both thread blocks contain a data segment, which is equivalent from global memory. After the merge the data segment is only fetched one time and the number of threads in the thread block adds up.

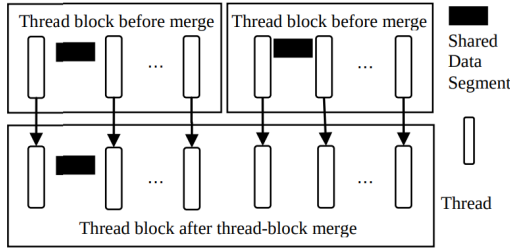


Fig. 3. Improved Memory Reuse by Thread Block Merging - Taken from A GPGPU Compiler for Memory Optimization and Parallelism Management

B. Thread Merging

Thread merging is a technique where different threads are merged into one thread to reduce memory access times. In exchange, there are less concurrent threads, which reduces the absolute computation time.

The goal of thread merging is to reuse data. In the best case, each block of data has to be fetched only once. In addition to thread block merging, which only allows for the reuse of Shared Data Segments, thread merging allows the reuse of shared registers too [2].

In Figure 4, we show an example for merging two thread blocks into one thread block by merging threads. Again, both blocks contain a shared data segment, which is equivalent. Additionally a shared register per thread is equivalent. By merging two threads into one, the shared data segment and the shared register just have to be fetched once.

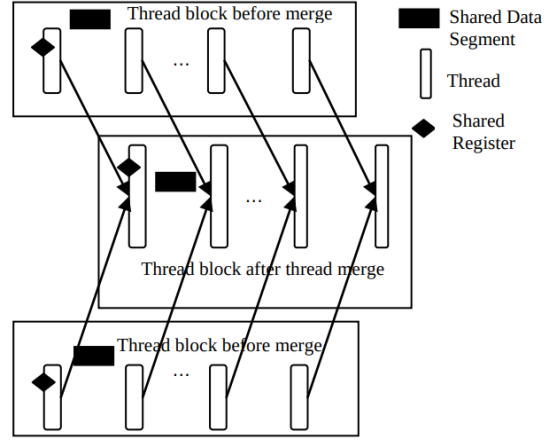


Fig. 4. Improved Memory Reuse by Thread Merging - Taken from A GPGPU Compiler for Memory Optimization and Parallelism Management

C. Data Prefetching

Data prefetching is a technique to optimize loops by prefetching the data, before they are actually needed.

Data prefetching optimizes the iteration over data structures by prefetching the data, that is needed in the following iteration step. In the current iteration step i the data from iteration step $i + 1$ gets fetched. This step allows you to overlap memory access time and computing time [2]. In exchange the actual iteration step become slower by adding some extra instructions for prefetching.

In Figure 5, we show example code with data prefetching. The code is part of the code for matrix multiplication. To understand the concept the *tmp* variable is important. In the beginning *tmp* holds the first part of data. In the following iteration steps it requests the part of data, that is used in the following step. That is an asynchronous event, that run next to the calculations of the current iteration step.

In Figure 6, we show a timeline, that shows the overlapping execution times from fetching data and executing code after the optimization. There is no time saving in the first iteration, because it is not possible to execute code without loaded data. After the first iteration, data can be loaded and code can be executed at the same time.

D. Partition Camping

Partition Camping is a problem, which is related to the memory architecture of NVIDIA CUDA. It describes the problem, that data, that is used in different threads, uses the same partition in global memory, which implies long memory access times.

The global memory of graphics cards is divided into partitions. Partitions can only be read or written by one thread at a time. If several threads want to access a partition at the same time, a queue is created and the threads cannot calculate any further. If all parts of the data set are located on a few partitions, then many threads need to access one partition, so

```

for (i=0; i<w; i=(i+16)){
    __shared__ float shared0[16];
    shared0[(0+tidx)]=a[idy][((i+tidx)+0)];
    __syncthreads();
    int k;
    for (k=0; k<16; k=(k+1)) {
        sum+=(shared0[(0+k)]*b[(i+k)][idx]);
    }
    __syncthreads();
}

```

(a) Before inserting prefetching

```

/* temp variable */
float tmp = a[idy][((0+tidx)+0)];
for (i=0; i<w; i=(i+16)) {
    __shared__ float shared0[16];
    shared0[(0+tidx)]=tmp;
    __syncthreads();
    if (i+16<w) //bound check
        tmp = a[idy][(((i+16)+tidx)+0)];
    int k;
    for (k=0; k<16; k=(k+1)) {
        sum+=(shared0[(0+k)]*b[(i+k)][idx]);
    }
    __syncthreads();
}

```

(b) After inserting prefetching

Fig. 5. Example Code for Data Prefetching - Taken from A GPGPU Compiler for Memory Optimization and Parallelism Management



Fig. 6. Example Timeline for Example Code in Figure 5 - Taken and Modified from NVIDIA CUDA Programming Guide

memory access is slower. This is called partition camping [10] [2].

If the data is evenly distributed across all available partitions, a single partition has fewer threads to access, which speeds up memory access. However, please note that data can only be efficiently redistributed if there is enough space and a good structure of the data. In addition, it can happen that redistribution makes the code faster in one place, but decelerates other places where the data is used in a different way [10] [2].

In Figure 7, we show SM-1, SM-2, ..., SM-30, which are different accessors.

In the table titled "Without Partition Camping" we show that the data for the different accessors are distributed evenly over all partitions.

In the table titled with "With Partition Camping" we show that the data is all in one partition and that all accessors would have to access one partition, resulting in long waiting times.

In conclusion, it can be said that in this example it is very easy to see that the individual partitions are used by considerably fewer accessors, which reduces the waiting time for memory.

IV. EVALUATION

In this section we compare the different optimizations for GPGPU, which we have explained before. In doing so, we will

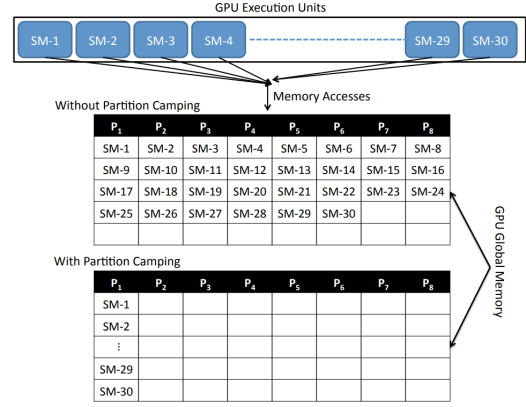


Fig. 7. Partition Camping Example - Taken from Bounding the Effect of Partition Camping in GPU Kernels

look at various points such as use cases, efficiency, difficulty and error cases. Note that all optimizations occur at compile time and do not mean any additional work for the developer of the application.

Optimization	UC	E	D	EC
Thread Block Merge	-	+	-	
Thread Merge	-	+	-	
Data Prefetching	+	-	+	
Elimination of Partition Camping	+	o	+	

Fig. 8. Comparison of different optimizations - Use Cases (UC), Efficiency (E), Difficulty (D), Error Cases (EC)

A. Use Cases

Cases for data prefetching and partition camping occur quite frequently. Nearly every loop that accesses data in global memory can be optimized with the help of data prefetching. The costs of the additional registers required for the temporary variable are relatively low.

The effect of partition camping always occurs when processing relatively large amounts of data and is therefore also optimized.

However, the use cases for thread merging and thread block merging are less frequent, since the circumstances and prerequisites are considerably more complex. Although cases in which these optimizations work can be constructed in this way, the application of this optimization is probably less in comparison to data prefetching.

B. Efficiency

The efficiency of the various optimizations can be measured by the increase in speed and number of memory accesses.

Thread merging and thread block merging reduce memory access by the number of n threads merged. This arithmetically increases the speed by a factor of n . According to REFERENCE, an acceleration of 10 times is possible in some cases.

The speed increase is not directly calculable with the elimination of partition camping and data prefetching, since it depends on many different factors.

The speed increase for the two other optimizations is much lower in detail, since only increases of the speed are possible in the same example by 2 to 4 times. These two optimizations do not reduce memory accesses.

C. Difficulty

Difficulty describes the effort of implementation for the application developer and the developer of the compiler. In this case, where the compiler automatically optimizes the application, there is no difficulty for the developer of the application, who does not have to work on the optimization.

The difficulty of data prefetching is, according to REFERENCE, easy to classify. The sample code in Figure 5 shows that the optimization is trivial to implement, apart from the bounding checks.

Implementing the elimination of partition camping is a bit more difficult to evaluate. There are several approaches to better disseminate data that produce different results. However, an initial implementation is relatively simple. Suppose there are n different partitions. Then the data of each partition can be split m times, so that it can be distributed to n partitions. Each partition contains an average of m/n data blocks.

Unfortunately, there is little information about the difficulty of implementing Thread Block Merging and Thread Merging. The detection is also not explained further, so that no conclusions can be drawn.

D. Error Cases

V. RELATED WORK

Yi Yang et. al. wrote a paper on "A GPGPU Compiler for Memory Optimization Parallelism Management" which explains how to implement the optimizations presented in this paper. Further optimizations are also presented and these optimizations are tested using examples such as matrix multiplication with large matrices [2].

Ashwin M. Aji et al. wrote an article on "Bounding the Effect of Partition Camping in GPU Kernels", which explains the effect of Partition Camping in more detail. It also explains the different patterns that can be used to avoid partition camping, which have been omitted in this paper for reasons of simplicity [10].

This paper optimizes access to shared storage areas. Greg Ruetsch et. al. from NVIDIA have carried out several optimizations using the example of matrix multiplication. An optimization was that the data is stored in such a way that no conflicts occur, which is a further development and/or modification of the partition camping solution [11].

VI. CONCLUSION

REFERENCES

- [1] CUDA Zone - <https://developer.nvidia.com/cuda-zone>
- [2] Yi Yang, Ping Xiang, Jingfei Kong, Huiyang Zhou, "A GPGPU Compiler for Memory Optimization and Parallelism Management", 2010
- [3] Victor W Lee, Changkyu Kim, Jatin Chhugani et. al "Debunking the 100X GPU vs. CPU Myth: An Evaluation of Throughput Computing on CPU and GPU", 2010
- [4] NVIDIA GeForce GTX 1080, <https://www.nvidia.com/en-us/geforce/products/10series/geforce-gtx-1080/>
- [5] NVIDIA GeForce GTX 1080 White Paper, https://international.download.nvidia.com/geforce-com/international/pdfs/GeForce_GTX_1080_Whitepaper_FINAL.pdf
- [6] GPU Accelerated Computing with C and C++ - <https://developer.nvidia.com/how-to-cuda-c-cpp>
- [7] Programming Model - Kernel <http://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#kernels>
- [8] Sain-Zee Ueng, Melvin Lathara, Sara S. Bagsorkhi, Wen-mei W. Hwu, "CUDA-Lite: Reducing GPU Programming", 2008 Complexity
- [9] Peter Bakkum, Kevin Skadron, "Accelerating SQL Database Operations on a GPU with CUDA", 2010
- [10] Ashwin M. Aji, Mayank Daga, Wu-chun Feng, "Bounding the Effect of Partition Camping in GPU Kernels", 2011
- [11] Greg Ruetsch, Paulius Micikevicius, "Optimizing Matrix Transpose in CUDA", 2009
- [12] GCC Optimize Options - <https://gcc.gnu.org/onlinedocs/gcc/Optimize-Options.html>
- [13] Lamport, Leslie (September 1979). "How to Make a Multiprocessor Computer That Correctly Executes Multiprocess Programs" (PDF). *IEEE Transactions on Computers*. C28 (9): 690691. doi:10.1109/tc.1979.1675439.