

Wojciech Pawlak (s091820)

Final Report

A GPU-accelerated Navier-Stokes Solver using OpenCL

Special Course at GPUlab, Scientific Computing Section, DTU
Supervisor: Allan P. Engsig-Karup Ph.D.

July 26, 2012

Contents

Contents	i
1 Introduction	1
2 Survey of GPU programming	2
2.0.1 Automation of porting to GPU process	2
2.1 OpenCL	3
2.1.1 Differences to CUDA	3
2.2 Previous research	3
2.3 Recent developments in CUDA and OpenCL	3
2.3.1 OpenCL	3
2.3.2 CUDA	3
2.3.3 Other Technologies	3
C++ AMP (C++ Accelerated Massive Parallelism)	3
OpenACC	3
AMD Accelerated Parallel Processing (APP) SDK	4
2.4 Current advances in NVIDIA and AMD Architectures	4
2.4.1 NVIDIA	4
Tesla	4
Fermi	4
Kepler	4
2.4.2 AMD	4
APU	4
2.4.3 Intel	4
Intel SDK for OpenCL Applications 2012	4
3 Navier-Stokes Solver	5
3.1 Computational fluid dynamics	5
3.2 Description of Navier-Stokes Solver	5
3.3 Other attempts of implementing the solver	5
4 Design and Implementation	6
4.1 Design	6
4.2 Implementation	6
4.2.1 Benchmarking code	7
4.2.2 Visualisation code	7
4.2.3 Naive kernels	7

<i>CONTENTS</i>	ii
4.2.4 Shared memory kernels	7
5 Performance Analysis	8
5.1 Computed results	8
6 Conclusions	9
References	10

Introduction

The goal of this report is to document the work done for individual special course “A GPU-accelerated Navier-Stokes Solver using OpenCL”. The project took place in GPUlab of Scientific Computing Section of Department of Informatics and Mathematical Modeling (DTU Informatics) at Technical University of Denmark. The project work for this special course took place between late March until early August 2012 and was supervised by Allan P. Engsig-Karup Ph.D. The course was worth 5 ECTS points.

The practical goal of this project was to design and implement a scientific computing application for execution on GPUs. The implemented application was then verified and validated using standard benchmarks.

The performance of implemented PDE solver was analyzed.

Apply optimization techniques for improving performance of PDE solver on GPUs.

Understand how to write a parallel program using OpenCL for heterogenous computing on many-core architectures.

Apply basic principles for numerical approximation/discretization.

Section 2 is

Section 3 is

Section 4 is

Section 5 is

Survey of GPU programming

For the exploitation of the high performance of GPUs, codes must express sufficient fine-grained parallelism, minimize GPU–CPU data transfer, and achieve coalesced memory access. One option for writing the GPU code is to do so directly in CUDA [1] or OpenCL [27]. For a large Fortran code written well before the existence of GPUs, this entails a massive effort to rewrite, debug, and maintain a separate branch of the code. An alternative to performing such a monolithic code rewrite is to follow the ‘accelerator design’, as so called by Cohen and Molemaker, of porting bottleneck loops or subroutines in isolation, with memory transfer calls made just before and after these bottleneck loops or subroutines. Although this allows for the possibility of an incremental porting effort, it suffers the drawback of introducing the crippling bottleneck of a large amount of data transfer across the system bus, which has at least an order of magnitude lower bandwidth than that of the internal GPU memory, and thus severely restricts the possible gain in performance when porting

Heterogeneous: Developers leverage AMD GPUs and CPUs for optimal application performance and user experience

Industry Standards: OpenCL™ and DirectCompute enable cross-platform development

High performance: Massively parallel, programmable GPU architecture enables superior performance and power efficiency

With OpenCL

Leverage CPUs and GPUs to accelerate parallel computation Get dramatic speedups for computationally intensive applications Write accelerated portable code across different devices and architectures

AMD APP Acceleration

2.0.1 Automation of porting to GPU process

<http://stackoverflow.com/questions/1126989/what-future-does-the-gpu-have-in-computing>

2.1 OpenCL

OpenCL (Open Computing Language) is an open standard defined by Khronos Group [16]. In contrast to the CUDA library supported mainly by NVIDIA it is a vendor-independent framework that enables the programmer to develop software that can be executed on heterogeneous platforms consisting of CPUs, GPUs and potentially other processing units like FPGAs.

Kernel Basic unit of executable code -similar to a C function
 Data-parallel or task-parallel
 Program Collection of kernels and other functions Analogous to a dynamic library
 Applications queue kernel execution instances Queued in-order Executed in-order or out-of-order

Pipeline - table: find platform find device create context and queue build program write device buffers execute kernel in range read device buffers

2.1.1 Differences to CUDA

2.2 Previous research

2.3 Recent developments in CUDA and OpenCL

2.3.1 OpenCL

What applications use OpenCL GPU-acceleration? OpenCL in Photoshop CS6 WinZip 16.5

OpenCL Studio

<http://www.youtube.com/user/OpenCLStudio>

2.3.2 CUDA

2.3.3 Other Technologies

C++ AMP (C++ Accelerated Massive Parallelism)

<http://msdn.microsoft.com/en-us/library/hh265137>

OpenACC

Directives like OpenMP for multicore CPU programming

<http://openacc.org/>

AMD Accelerated Parallel Processing (APP) SDK

formerly ATI Stream

<http://developer.amd.com/sdks/AMDAPPSDK/Pages/default.aspx>

<http://www.amd.com/us/products/technologies/amd-app/Pages/eyespeed.aspx>

2.4 Current advances in NVIDIA and AMD Architectures

today we have two viable and competitive product lines, Nvidia and Advanced Micro Devices (AMD) GPUs, with support for a wide range of programming languages.

GPUs on super-computers

2.4.1 NVIDIA

Tesla

Fermi

Kepler

2.4.2 AMD

APU

2.4.3 Intel

Intel SDK for OpenCL Applications 2012

Navier-Stokes Solver

3.1 Computational fluid dynamics

3.2 Description of Navier-Stokes Solver

Reynolds number < 1000

3.3 Other attempts of implementing the solver

Design and Implementation

4.1 Design

A set of kernels

Modularity

All code should be executed on GPU for valid comparison.

4.2 Implementation

Code in project is based on structure from Griebels code[1]. Least as possible was changed in structure as to keep it comparable with the original code.

Code is written in C and OpenCL C.

Two-dimensional arrays in original One-dimensional arrays for OpenCL code. 1D arrays cannot be passed to OpenCL kernels.

Worksize

Implementation consists of sets of kernels as

Griebels cod intrinsics: Memory allocation

Code is portable and works under both on Windows and Unix platforms. Tested on Windows 7 and Linux clusters. Compilers used are nvcc and Visual Studio on Windows and gcc on Linux.

Table Computer specs.

Table

4.2.1 Benchmarking code

Code is timed with standard library time.h clock function. No OpenCL timer code because on CPU. Code timed with buffer allocation/memory access and without.

4.2.2 Visualisation code

Matlab scripts to visualise code. Tests.

4.2.3 Naive kernels

Straightforward port of functions to kernels. Getting rid of for loops. Ensuring that boundaries are not crossed.

4.2.4 Shared memory kernels

Performance Analysis

5.1 Computed results

Conclusions

References

- [1] Michael Griebel, Thomas Dornseifer, and Tilman Neunhoeffler. *Numerical Simulation in Fluid Dynamics: A Practical Introduction*. Siam Monographs on Mathematical Modeling and Computation. Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 1998. 4.2