

Computing Performance Benchmarks among CPU, GPU, and FPGA

MathWorks

Authors:

**Christopher Cullinan
Christopher Wyant
Timothy Frattesi**

Advisor:

Xinming Huang

Abstract

In recent years, the world of high performance computing has been developing rapidly. The goal of this project was to conduct computing performance benchmarks on three major computing platforms, CPUs, GPUs, and FPGAs. A total of 66 benchmarks were evaluated. GPUs outperformed the other platforms in terms of execution time. CPUs outperformed in overall execution combined with transfer time. FPGAs outperformed for fixed algorithms using streaming. The team made several recommendations for further research in this area.

Acknowledgements

The successful completion of this project would not have been feasible without the help of several key individuals. First, we would like to express our gratitude to our advisor, Xinming Huang, who met with us weekly throughout the course of the project offering advice and wisdom. In addition, we would like to thank our sponsor at MathWorks for their financial support of this project.

Authorship

Christopher Cullinan

Christopher C. was responsible for the CPU Multicore portion of this project. This included gathering and testing 34 benchmarks on the AMAX machine. In addition, he wrote the CPU sections for the background, benchmark, and results along with the future work and executive summary sections of this report.

Timothy Frattesi

Timothy was responsible for the GPU portion of this project. This included gathering and testing 24 benchmarks on the GeForce GTX 460 and GeForce 9800 GTX+ NVIDIA graphics cards. In addition, he wrote the GPU sections for the background, benchmark, and results of this report along with the future work and executive summary sections.

Christopher Wyant

Christopher W. was responsible for the FPGA portion of this project. This included gathering and testing of 8 benchmarks using the ISim program of Xilinx on a Virtex-5 board. Christopher W. wrote the FPGA sections of the background, benchmark, and results section of the report. In addition, he was responsible for writing the abstract, acknowledgements, introduction and conclusion of this report.

Executive Summary

Ever since the beginning of modern day computing, engineers and developers have been trying to squeeze every ounce of performance into their devices. How do we test for performance though? What quantifiable measurements can be made to justify the superiority of one device over another? These are the types of questions this project aims to answer. In this project, we investigated the performance abilities of current generation multiprocessing hardware. Through looking at multicore CPUs, general purpose GPU computing and an FPGA, we compared device capabilities to determine which platform future investments should be focused towards and why.

To begin our endeavor, we used benchmarking to accentuate the strengths and weaknesses of these three devices. Benchmarking is the technique of using crafted programs in order to attach quantifiable performance metrics to targeted computer subsystems. By using cross platform, as well as individual, benchmarks developed across a plethora of computational necessities, we determined which device would be best suited towards specific tasks. For this project, we tested our benchmarks across two Intel Xeon 5650 CPUs, the Virtex-5 FPGA and NVIDIA's GeForce GTX460 and 9800 GTX+ GPUs.

Realizing the sophistication early on in this project, we decided to use already written benchmarking suites to conduct our tests. A benchmarking suite is nothing more than a compilation of individual benchmarks with specific intent. In total, we used seven benchmarking suites. For the FPGA, we used cores designed by Xilinx Core Generator and MATLAB Simulink HDL Coder, which contained

benchmarks encompassing mathematical algorithms and encryption. For the CPUs, we used three benchmarking suites; SPEC CPU2006, Rodinia, and John Burkardt. Lastly, for the GPUs, we used the Parboil, Rodinia, and SHOC benchmarking suites. The CPU and GPU suites tested mathematical algorithms, high performance simulation, and common computational necessities such as compression and sorting.

To further enhance the findings of this project, we discussed several future recommendations at the end of this report. These recommendations include testing a broader spectrum of benchmarks capable of running across all three platforms. Another possibility is to look into newer technologies, such as an Accelerated Processing Unit (APU). We believe that by including these recommendations, a conclusion of greater impact can be reached.

Table of Contents

Abstract.....	i
Acknowledgements.....	ii
Authorship.....	iii
Executive Summary	iv
Table of Contents	vi
Table of Figures	ix
Table of Tables.....	x
1. Introduction	1
2. Background.....	2
2.1. CPU.....	2
2.1.1. A Brief History	2
2.1.2. CPU Design	4
2.1.3. CPU Multicore	6
2.2. Graphics Processing Unit.....	9
2.2.1. GPU History.....	9
2.2.2. GPU Architecture and Parallelism.....	12
2.3. FPGA Background.....	17
2.3.1. History of FPGAs	17
2.3.2. Early Programmable Devices.....	18
2.3.3. FPGA Architecture	19
3. Benchmarks.....	24
3.1. FPGA Benchmarks.....	26
3.2. SPEC CPU2006.....	28
3.2.1. CPUINT2006	28
3.2.2. CFP2006	34
3.3. Rodinia Suite	41
3.4. John Burkardt Benchmarks.....	44
3.5. SHOC Suite.....	45
3.6. Parboil Suite	48

Breadth.....	48
4. Results	52
4.1. Devices.....	52
4.1.1. GPU.....	52
4.1.2. CPU	53
4.1.3. FPGA	53
4.2. ALL Devices.....	54
4.3. CPU & GPU.....	55
4.4. Individual Results.....	57
4.4.1. GPU Specific	57
4.4.2. CPU Results	65
4.4.3. FPGA Results.....	72
5. Future Work.....	74
6. Conclusion.....	75
7. Appendices.....	76
7.1. Parboil Results 9800 GTX+	76
7.2. Parboil Results GTX 460	78
7.3. Rodinia Results 9800 GTX+	80
7.4. Rodinia Results GTX 460	81
7.5. SHOC Max Flops GTX 460.....	84
7.6. SHOC Bus Download Speed GTX 460.....	85
7.7. SHOC Device Memory GTX 460.....	86
7.8. SHOC SPMV GTX 460.....	87
7.9. SHOC MD GTX 460	88
7.10. SHOC Reduction GTX 460	89
7.11. SHOC S3D GTX 460.....	90
7.12. SHOC Scan GTX 460.....	91
7.13. SHOC SGEMM GTX 460	92
7.14. SHOC Sort GTX 460.....	93
7.15. SHOC Stencil 2D GTX 460.....	94
7.16. SHOC Triad GTX 460	95

7.17.	SHOC Max Flops 9800 GTX+	96
7.18.	SHOC Bus Download Speed 9800 GTX+	97
7.19.	SHOC SPMV 9800 GTX+	98
7.20.	SHOC MD 9800 GTX+	99
7.21.	SHOC Reduction 9800 GTX+	100
7.22.	SHOC S3D 9800 GTX+	101
7.23.	SHOC SGEMM 9800 GTX+	102
7.24.	SHOC Sort 9800 GTX+	103
7.25.	SHOC Stencil 2D 9800 GTX+	104
7.26.	SHOC Triad 9800 GTX+	105
7.27.	SPEC CPU2006 Integer Results (No Auto-Parallel)	106
7.28.	SPEC CPU2006 Integer Results (Auto-Parallel Enabled)	106
7.29.	Speedup of SPEC CPU2006 Integer Results	107
7.30.	SPEC CPU2006 Floating Point Results (No Auto-Parallel)	107
7.31.	SPEC CPU2006 Floating Point Results (Auto-Parallel Enabled)	108
7.32.	Speedup of SPEC CPU2006 Floating Point Results	108
7.33.	Rodinia/Burkardt Benchmarks Average Execution Times	109
7.34.	Rodinia/Burkardt Benchmarks Speedup between Thread Count	109
7.35.	FPGA Results	110
	Bibliography	111

Table of Figures

Figure 1 - AMD K10 Architecture	6
Figure 2 - Intel I7 950 Quad Core Processor Design	7
Figure 3 - S386C911 (1991)	10
Figure 4 - NVIDIA GeForce 256 (1998)	10
Figure 5 - GeForce 6600 (2004)	11
Figure 6 - GeForce GTX 560 (2011)	11
Figure 7 - CUDA Software Architecture	12
Figure 8 - Improvements in CUDA GPU Architecture	14
Figure 9 - Thread Hierarchy	16
Figure 10 - Stream Multiprocessor	16
Figure 11 - PLA Architecture	18
Figure 12 - CPLD Architecture	19
Figure 13 - Three Input LUT	20
Figure 14 - Switch Block	21
Figure 15 - FPGA Routing	22
Figure 16 - Specialized Slices	23
Figure 17 - Virtex-7 Architecture	23

Table of Tables

Table 1 - FFT Results	54
Table 2 - CPU & GPU Results.....	56
Table 3 - GeForce Specification	58
Table 4 - Data Size vs. Data Rate	59
Table 5 – Global vs. Local Memory Speeds	60
Table 6 – Parboil Suite Timings.....	63
Table 7 - Rodinia Particle Filter Results	64
Table 8 - SPEC CPUINT2006 w/ GCC, G++, GFortran compiler w/o auto parallel	65
Table 9 - SPEC CPUINT2006 run with Intel compiler in auto parallel.....	66
Table 10 - Speedup percentages from the GNU run to the Intel run	66
Table 11 - SPEC CFP2006 w/ GCC, G++, GFortran compiler w/o auto parallel.....	68
Table 12 - SPEC CFP2006 run with Intel compiler in auto parallel	68
Table 13 - Speedup percentages from the GNU run to the Intel run	69
Table 14 - Rodinia/Burkardt average execution time on 2, 4, 8, 12, 24 threads	70
Table 15 - Speedup between thread counts	71
Table 16 - FPGA Results	72

1. Introduction

The world of high performance computing is a rapidly evolving field of study. Many options are open to businesses when designing a product. GPUs can provide astonishing performance using the hundreds of cores available. On the other hand, FPGAs can provide computational acceleration to many signal and data processing applications. The question arises as to what level each platform performs at for different benchmark algorithms.

To determine computing levels of each device we implemented existing benchmark suites for each device. We tested several applications to see which computing method was fastest for the various applications. While the benchmarks did not completely lineup between the different processors the information gathered laid a good foundation between different devices.

In order to explain the information in a clear manner, we broke up the information into several sections. Presented first is the background of the devices, both general and specific. The second section outlines all of the benchmarks we tested so that the reader can understand the limitations of the project. The third section discusses the results we gathered as well as a discussion of what they mean. Lastly, we discuss some recommendations we would make for similar projects in the future and some closing remarks.

2. Background

The background gives an overview of the different device we used in our project. Within each section are the history of the device and a general overview of how they work. This information gives a brief but necessary background into the platforms used for this project.

2.1. CPU

2.1.1. A Brief History

The history of the Central Processing Unit (CPU) is in all respects a relatively short one, yet it has revolutionized almost every aspect of our lives. In the early 1970, if I were to ask someone what a CPU was, they would have most likely responded “A what!” Yet just over 40 years later, CPUs have become an integral part of our lives. From desktop computers to cell phones, most of us do not go more than a few hours without somehow interacting with a CPU. Despite its indisputable popularity, most do not know how this hype all started.

In 1971, the 4-bit Intel 4004 was the first in the legacy of the CPU. It was the first commercially available CPU on chip, made possible by the all-new silicon gate technology. The max CPU clock rate of this revolutionary hardware was 740 kHz, an astonishing speed at the time. This little guy could execute 92,000 instructions per second with a single instruction cycle of 10.8 microseconds and a transistor count of 2,300. At the time, this device was truly a feat in computing technology, which paved the road for much more innovation to come. [1]

Intel dominated CPU infancy, coming out with several subsequent CPU designs including the 8086 (1978), 8088 (1979), 80186 (1980). Then, in 1993, one of the most popular names in the history of the CPU surfaced, the Intel Pentium processor. This legendary device operated at a whopping 60 MHz and 100 Millions of Instructions per Second (MIPS). The trend of innovation from Intel continued for several years until another major competitor in today's market made their first competitive appearance with the AMD AM5x86 in 1995. A fierce competition between Intel and AMD has continued since. [2]

The next milestone in the CPU history was the commercial release of the first 1GHz processor. This achievement was reached by the AMD Athlon in 1999 and then by the Intel Pentium III just two days later after. For this reason, "Athlon" was fitting name for AMD's milestone processor because it is the Greek word for "Champion/trophy of the games". The AMD Athlon is an x86-compatible processor containing 22 million transistors in a slim size of 184 mm². [3]

Nowadays, it is a common occurrence to see CPUs clocked well above 1GHz in devices as small as our cell phones. In just over 40 years, we have gone from 740 kHz to the GHz level (over a 1300 % increase) and increased the count of on chip transistors from 2,300 to more than a billion (over a 434,000 % increase). We are now producing CPUs with multiple cores on the same chip, which are capable of support an increasingly important feature known as parallel computing, which we will talk about in more detail later in this report. [2]

2.1.2. CPU Design

In a nutshell, CPU design is the design engineering process of creating a central processing unit to be used in a computing system. Many factors go into designing a CPU, especially with the level of sophistication in modern day CPUs. There are six primary focuses that designers must account for when creating a CPU, and they are; data paths, control unit, memory components, clock circuitry, pad transceiver circuitry, and logic gate cell library. [4]

A data path by definition is, “A collection of functional units, such as arithmetic logic units or multipliers, that perform data processing operations”. [5] Intuitively from its name, data paths provide routes for data to traverse between the components of a CPU. These routes are typically known as “Buses”. The majority of CPUs include both a data path and a control unit, where the control unit specializes in regulating data path and main memory interaction. [5]

Most modern day CPUs have several types of memory modules on chip. Two of the most popular are register memory and cache, both of which are normally high speed SRAM. Registers are the memory cells built directly into the CPU since they contain specific data vital to CPU operations. Cache is the next portion of memory in a CPU and is usually, in more complex processors, divided into L1 (level one) and L2 (level two) cache. Both L1 and L2 cache are there to store data that is most often used by the CPU and is typically SRAM as well. [6]

The CPU clock is the sinusoidal frequency reference signal typically created by a crystal oscillator. This sinusoidal waveform is first translated into a square waveform of the same frequency by internal circuitry and then used to synchronize

the internal components of the CPU. The clock signal traverses to the various CPU components via a clock distribution network. [7]

Lastly, the logic gate cell library is the collection of logic gates used to implement computational logic in the CPU. The logic library collection consists of low-level logic functions including AND, OR, INVERT gates as well as flip-flops, latches, and buffers. A vital feature of these libraries is that they are fixed height and variable width, meaning they can be placed in organized rows. This makes the process of automated digital layout of these components possible and efficient. [8]

To give you an idea of a simple CPU design, the following figure shows the AMD K10 Architecture of 2007. This processor is slightly outdated but is good for our purposes to show the anatomy of a modern day CPU.

AMD K10 Architecture
 Red: Difference between K8 and K10 Architecture
 (Die Änderungen zwischen der K8- und K10-Architektur sind rot markiert)

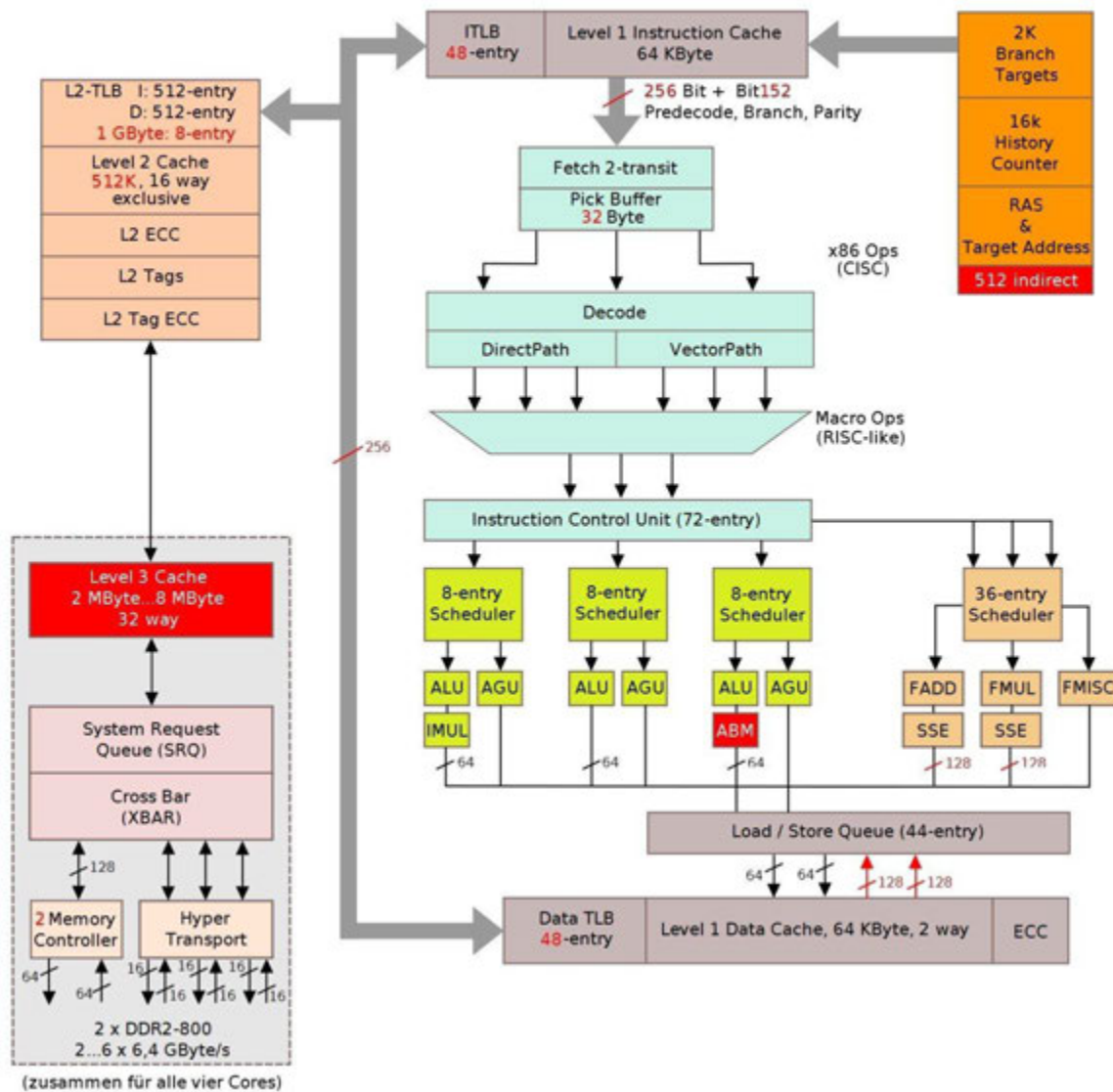


Figure 1 - AMD K10 Architecture

2.1.3. CPU Multicore

The term “multi-core” refers to a multiple core processor that is simply an integrated circuit where two or more processors have been attached for increased performance via parallel processing. [9] Parallel processing is a type of computation

where many calculations are performed simultaneously. This method of computation is based on the principle that large problems can be solved faster by breaking them down into smaller pieces and then solving those pieces concurrently. Because of this basic principle, parallel computing has become the dominant standard in computer architecture in the most popular form of multicore processing. As an example of a multicore processor design, the following figure shows the internals to an Intel I7 950 Quad Core Processor.

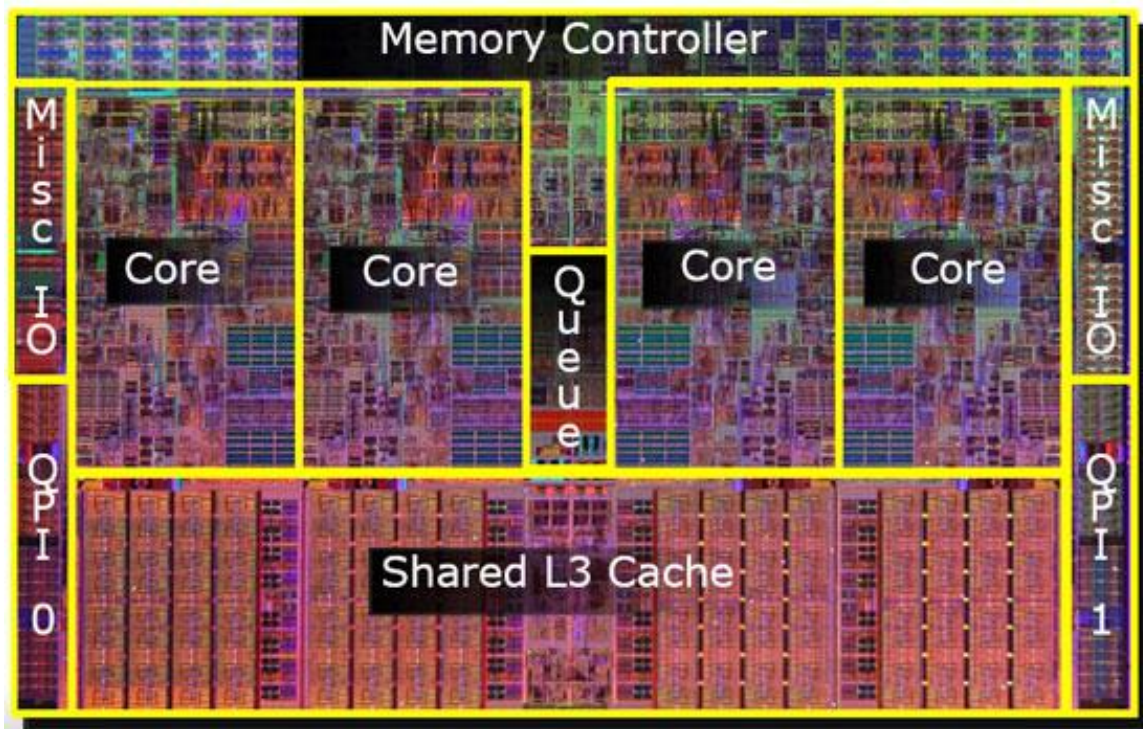


Figure 2 - Intel I7 950 Quad Core Processor Design

In the real world, parallel processing is not limited to integrated circuits. Virtually everything in our natural universe uses the principles of parallel processing. A few examples are galaxy formation, planetary movements, weather and ocean patterns, automobile assembly lines, rush hour traffic, and even ordering a hamburger at a fast food restaurant. Within all of these phenomena, numerous

complex and interrelated events occur simultaneously to achieve a common goal.
[10]

In the past, parallel computing was attributed mainly to high end computing and was used to solve complex mathematical problem in various areas of study. A few of these areas included Atmospheric studies, Physics, Mechanical Engineering, Electrical Engineering, and Seismology. Parallel computing is still used in those areas today, but the rise of commercial applications has been a major contributor to both the need for faster computing and the dispersion of parallel computing into common electronic devices such as phones, desktops, laptops, etc. Database mining, web search engines, medical imaging, and advanced graphics are just a few of the applications that utilize parallel computing. [10]

At the beginning of this section, we briefly discussed the incentives to use parallel computing; now we will delve deeper to justify the use of parallel computing. Parallel computing saves time and money by both shortening the time to the outcome and because parallel components are cheap. Second, larger problems can be solved through the use of parallel computing that are not possible by using a single computing resources. Finally, there are many limitations to serial computing. These limitations include transmission speeds, limits to miniaturization, and economic limitations. To get away from these confines, modern computer architectures are heavily relying on multiple execution units, pipelined instructions, and multi-core at the hardware level to increase performance. [10]

2.2. Graphics Processing Unit

As stated by Prof. Jack Dongarra, "GPUs have evolved to the point where many real-world applications are easily implemented on them and run significantly faster than on multi-core systems. Future computing architectures will be hybrid systems with parallel-core GPUs working in tandem with multi-core CPUs." [11] From this comment one can see one of the many thoughts on where the future or high performance computing is headed.

2.2.1. GPU History

Much like computers in general, GPU's have progressed rapidly over the last 30 years since their introduction to the market. As GPU's have progressed over the years their core functions have remained the acceleration and processing of images.

The introduction of graphics units came early in the 1980's where both Intel and IBM brought specialized products to the market. Other companies such as Commodore and Texas Instruments also added simple graphics capabilities either on chip or using an external card. These cards had simplistic functionality and were relatively expensive. Functions such as filling an area, shape drawing, and modification of simple images were all that these early processors could support.



Figure 3 - S386C911 (1991)



Figure 4 - NVIDIA GeForce 256 (1998)

The 1990's were the real beginning as far as the takeoff of GPUs. From the beginning in 1991, S3 rolled out their 86C911 card which was one of the first standards for the GPU industry. Two dimensional graphic processing had made its way into almost every system by the mid 90's and the race was on to move towards 3D processing. Two notable chip sets in the race for dedicated 3D graphics include the 3dfx Voodoo and Rendition's Verite. Until the late 90's all 3D rendering was done with the assistance of CPUs, also known as hardware assisted 3D graphics which we still see in lower end laptops today. [12]

To assist in the commonality of graphics processing several “languages” were brought about in the late 90's including both OpenGL and Direct. Throughout the 90's OpenGL prospered as the software's capability was usually ahead of Direct and it was capable of being used across cards and platforms. Towards the end of the 90's these two API's introduced support for transform and lighting (T&L) which provided a huge jump in GPU processing. T&L allowed for easier mapping of 3D images to a 2D plane while incorporating the lighting all into one. By this time there were only a few competing companies; NVIDIA, ATI, 3dfx, and S3. The end of the

90's saw the NVIDIA GeForce 256, the first readily available commercial card, bringing 3D graphics, NVIDIA, and Direct to their own level. [12]



Figure 5 - GeForce 6600 (2004)



Figure 6 - GeForce GTX 560 (2011)

Through 2010 and to today we continue to see significant advancements in the 3D rendering abilities of the GPU. On the front of programming the most notable improvements include programmable shading and floating point abilities. ATI and NVIDIA hold the majority of today's market share in graphics processing and thus have been major forces in shaping how these units improve.

One of the most significant advancements of the past decade is general purpose computing for GPU's. Due to the highly parallel structure of modern graphics cards it is possible to use them to perform research and analysis, often times competing or surpassing modern CPUs. While this can be done with almost any modern card, NVIDIA's introduction of the Compute Unified Device Architecture (CUDA) from NVIDIA this idea has become standardized. OpenCL is also a common language for performing GPU computation, but it does not support as many programming languages or have the same amount of industry support. CUDA

architecture is the main advancement that is allowing our society to take high performance computing from CPUs and FPGAs and move it to a quicker paralleled set of computations on thousands of threads instead of tens of threads. [13]

2.2.2. GPU Architecture and Parallelism

Here we will look at the both the CUDA architecture and the hardware architecture that corresponds to it. NVIDIA has created this specialized architecture to achieve the massively parallel systems that we have today.

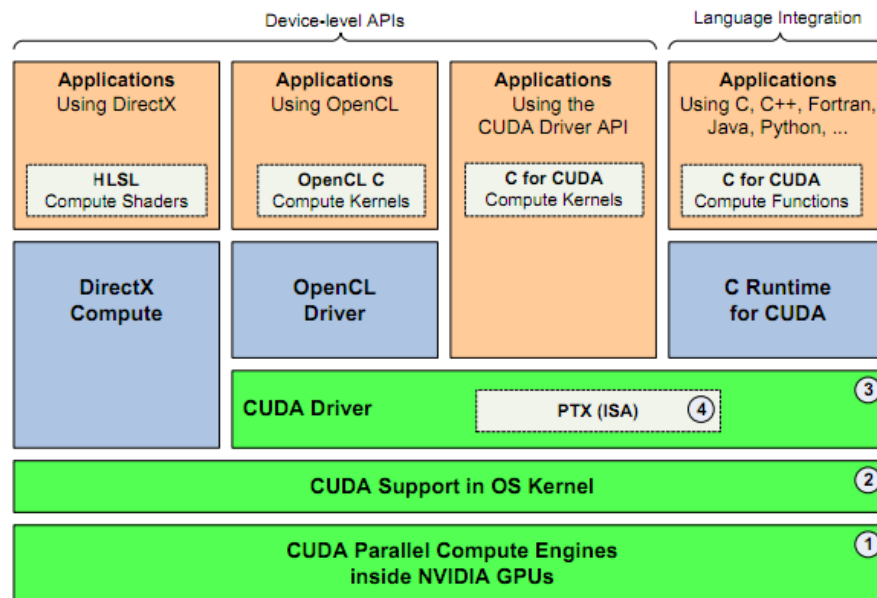


Figure 7 - CUDA Software Architecture

Looking backwards from computer to device we start with the Parallel Thread Execution (PTX) instruction set. This specially optimized instruction set allows for special optimization specifically designed for parallel processing on NVIDIA cards. The PTX instruction set allows NVIDIA to set a standard across multiple generations of GPUs as well as provide a common set of instructions for

both optimization and developer programming. This instruction set rests within the CUDA drivers which are provided dating back to NVIDIA's GeForce 8800 series. [14]

The next level up, CUDA Support in OS is provided through the CUDA Toolkit which is currently on version 4.1. This toolkit provides all of the necessary components to write and run CUDA code in an IDE such as Visual Studio or simply in a text editor. Currently supported are both a LLVM based and the standard nvcc compiler. The move to the LLVM based compiler has shown 10% increases in speed, but is geared more towards the new Fermi architecture and beyond. With the release of the new version of the CUDA toolkit, NVIDIA also provided a much larger base of functions for image processing.

The base level is the actual CUDA Parallel Computing engines that are the basis of all modern NVIDIA cards. The specialized hardware architecture in these cards is what allows for such successful parallelism in general purpose computing. On the current generation of CUDA capable hardware, Fermi, one can setup and process 65,535 simultaneous threads in grids of up to 1024x1024x64. These grid sizes have doubled and allowed for a third dimension since CUDAs initial release in 2006. The number of cores capable of providing floating point and integer functionality has also increased six fold. In the following figure the hardware advancements can be seen in overview as the architecture has changed. Further improvements can all be seen in the tables attached in Appendix X. [15] [16]

GPU	G80	GT200	Fermi
Transistors	681 million	1.4 billion	3.0 billion
CUDA Cores	128	240	512
Double Precision Floating Point Capability	None	30 FMA ops / clock	256 FMA ops /clock
Single Precision Floating Point Capability	128 MAD ops/clock	240 MAD ops / clock	512 FMA ops /clock
Special Function Units (SFUs) / SM	2	2	4
Warp schedulers (per SM)	1	1	2
Shared Memory (per SM)	16 KB	16 KB	Configurable 48 KB or 16 KB
L1 Cache (per SM)	None	None	Configurable 16 KB or 48 KB
L2 Cache	None	None	768 KB
ECC Memory Support	No	No	Yes
Concurrent Kernels	No	No	Up to 16
Load/Store Address Width	32-bit	32-bit	64-bit

Figure 8 - Improvements in CUDA GPU Architecture

Parallelism is generated through the use of threads on the GPU. CUDA has the ability to run a particular kernel across multiple threads and multiple cores. Picture the newest card in NVIDIA's lineup has 1024 cores per card. If each core can provide up to 1024 threads per 65535 blocks, there is a possibility to run upwards of 60 billion instances of a single kernel per card in a system. Spanning this number across multi-card systems, or even multiple systems connected together, it is easy to see how parallelism prevails and allows for faster processing. This type of parallelism is known as single instruction multiple data (SIMD). CUDA devices are also capable of running thousands of small programs simultaneously as well.

In order to achieve this number of parallel threads there is a complex setup of memory. Each and every thread an individual piece of memory that is unshared. This contains data such as program counters and individual registers. From there we move up to thread blocks which share memory amongst themselves and then up to sixty five thousand blocks sharing memory per core. These sets of blocks are

called grids and can share a set of application memory for use by smaller threads with global memory. All of the threads running a particular instance of a kernel are kept in synchronization through the use of special code functions that wait for all threads to be completed before reading or writing large changes from global memory. The memory hierarchy described here is important because each smaller level allows for quicker access to data; thus like the L cache of a CPU, there is less need to constantly read and write to slower global memory. [16] [17]

To handle all of these threads processing at the same time there is a unique feature called a warp handler. A warp is a group of 32 threads; the smallest data size for a SIMD setup. When programming in CUDA, users work with blocks so it is up to the warp handler to determine how to divide the instructions. The Fermi architecture has a dual warp scheduler, allowing it to process and divide up two sets of instructions at a time. With 32 bit mathematics it is also possible to dispatch two of a single type of instruction or a mix at one time. Since this setup works in sets of 32 in order to achieve the peak performance on CUDA capable GPUs is to run kernels in sets divisible by 32. [17]

Along with the efficient thread hierarchy, the Fermi architecture relies on NVIDIA's third generation of stream multiprocessing (SM) for its hardware architecture. In this revision of SM, there are 32 CUDA cores per multiprocessor, giving each card 16 to 32 SMs; with each core having it's out floating point and arithmetic logic units. Figure 9 below is a simplified example of third generation SM. Pictured is a single multiprocessor, with the hardware available in each SM. In the figure the 32 individual processing cores can be seen along with the 16

load/store units and the 4 special functions units. These units are accessed once the warp schedulers divide tasks amongst the cores. Following these processing units are the standard graphical processing hardware units such as those that perform tessellations and texturing. It is also noted that in this version each individual core has its own integer and floating point units. [16] [17]

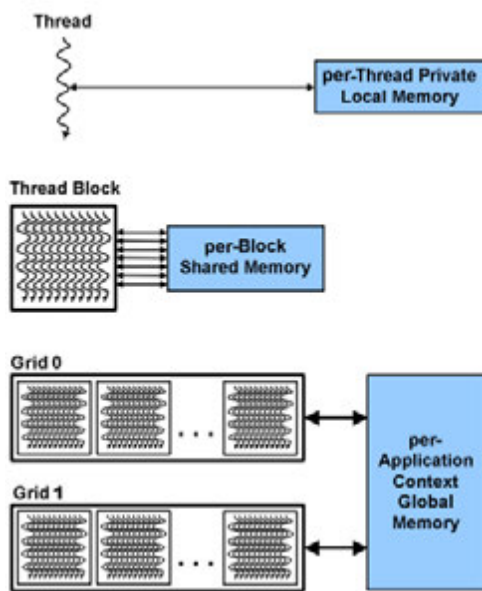


Figure 9 - Thread Hierarchy



Figure 10 - Stream Multiprocessor

Added support for the new IEEE floating point standard has given these new cards the ability to use a fused add and multiply in one step. Data precision has also been improved so that the integer units provide 64 bit support while the floating point units finally provide full 32 bit support. The sixteen load and store and four special function units allow for up to sixteen thread address to be calculated at a

time and four instances of functions such as sine or square roots. With most of these instructions executing with one instruction per clock cycle having so many SMs on a single card allows for distribution of complex equations and faster execution. [16]

2.3. FPGA Background

2.3.1. History of FPGAs

Ross Freeman, co-founder of the company Xilinx, invented Field Programmable Gate Arrays (FPGA) in 1984 while working for the company Zilog. After inventing the FPGA Freeman left Zilog with his Patent (patent 4,870,302) to found Xilinx. While Xilinx is a multi-billion dollar company today, Freeman did not live to see this become a reality passing away in 1989. He was honored in 2009 by being inducted into the National Inventor's Hall of Fame for his work on FPGAs. [18]

The first FPGAs released to the market had only several thousand gates and had several disadvantages to their counterparts, ASICs. They were slower, consumed more power, and had limited functionality. The industry of FPGAs grew slowly through the 1990s. In 1992 the U.S. Naval Surface Warfare department completed a project on FPGAs that implemented 600,000 logic gates. During this time, the main applications for FPGAs were networking and telecommunications.

By the late 90s, the number of gates on a single FPGA reached the millions and many of the disadvantages compared to ASICs were diminishing. FPGAs began entering many other industries because of the low time from development to market introduction. Much money could be generated by being the first to the market. [19]

Today, FPGAs can cater to many different applications. Different series and families are application specific and have additional logic to support faster processes. FPGAs have a high capacity for parallelization and pipelining processes. Often, they are used as peripherals to CPUs to carry out specific processes that a CPU has trouble handling.

2.3.2. Early Programmable Devices

The idea behind FPGAs originated from two devices, Programmable Logic Arrays (PLA) and Complex Programmable Logic Devices (CPLDs). PLAs were introduced during the early 1970s as one-time programmable chips to implement logic functions. The AND gates and OR gates were connected with a communication matrix that could be programmed by burning fuses to implement a truth table. The limiting factors were the number of inputs, AND gates, and OR gates. [20]

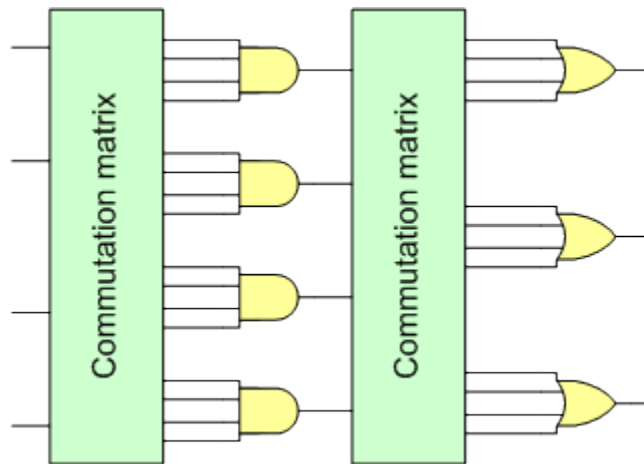


Figure 11 - PLA Architecture

CPLDs built upon the idea of PDAs with an interconnection matrix connecting all of the inputs and outputs. The connection matrix was formed of on-chip Flash

memory to configure macrocells. These macrocells are very similar in structure to a PLA. CPLDs are very similar to FPGAs as the main difference is in the underlying architecture. [20]

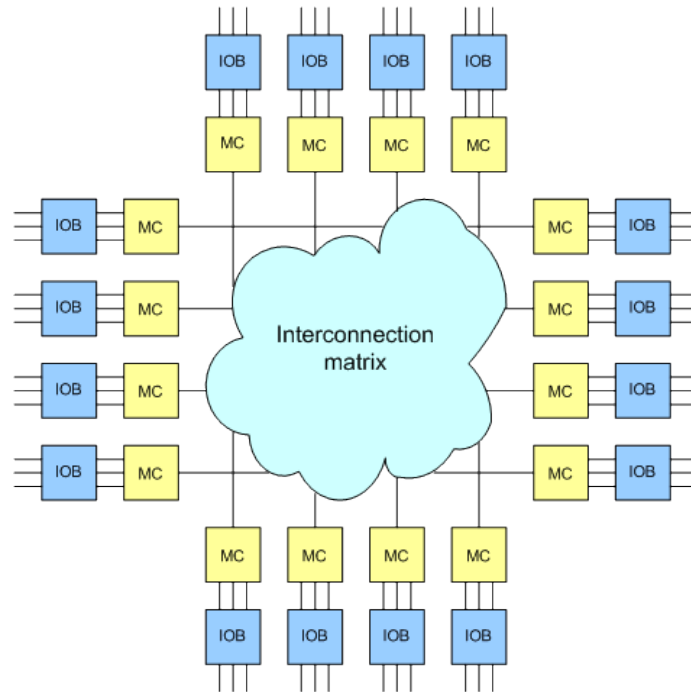


Figure 12 - CPLD Architecture

2.3.3. FPGA Architecture

Field Programmable Gate Array is a semiconductor device comprised of many logic blocks with configurable interconnections between these. The logic blocks are capable of acting as simple logic gates, such as AND and XOR. In addition to the logic gates, there are routing channels that run between each logic block. These channels are programmable and enable different logic blocks to talk to each other. In recent years, more specific circuits are implemented on FPGAs for

application specific purposes. These can include multipliers and DSP circuits, which speed up processing for those applications.

The main component of FPGAs is the logic block. Millions of these are replicated in a network throughout the chip. They are implemented in a Lookup Table (LUT) usually consisting of four input pins. The LUTs have small piece of memory attached that is programmed for output logic depending on the input. Essentially, a truth table defined for that piece of logic. Some designers are increasing the number of input pins to six to increase speed.

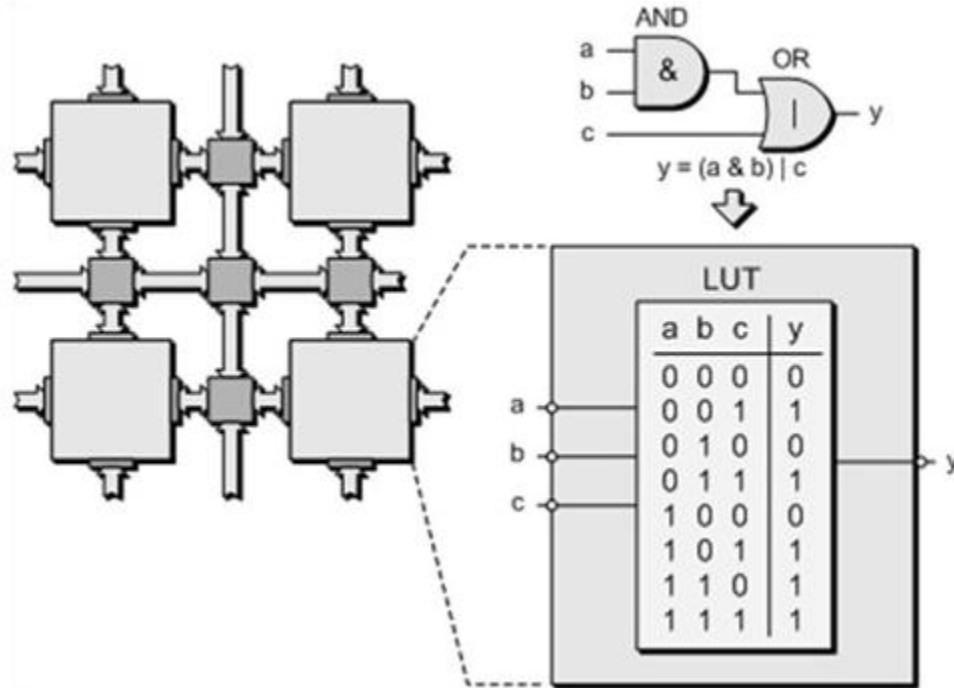


Figure 13 - Three Input LUT

Each LUT has only one output. This output can then be stored in a flip flop to preserve values over a clock cycles, or it can run to other LUTs to further implement logic. The Virtex-5, which the benchmarks in the report are based on, uses six input LUTs.

The routing channels, which run between logic blocks, are used to connect various LUTs together. The routing channels are controlled by switch blocks that control connections between crossing wires. An example of this is seen below in Figure 14.

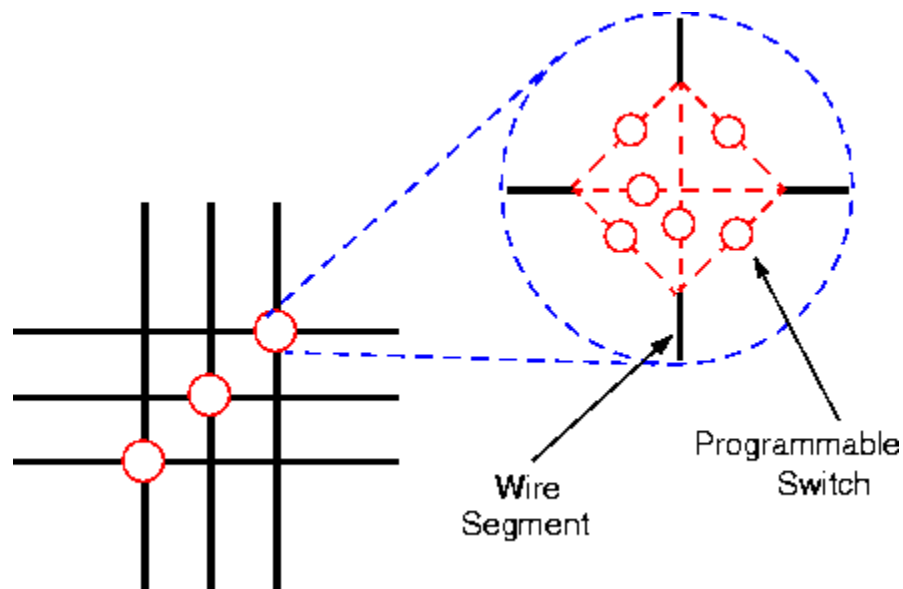


Figure 14 - Switch Block

These connections allow for an immense amount of configurable logic allowing an FPGA to carry out its functionality. Figure 15 below shows the layout of blocks throughout a network. [21]

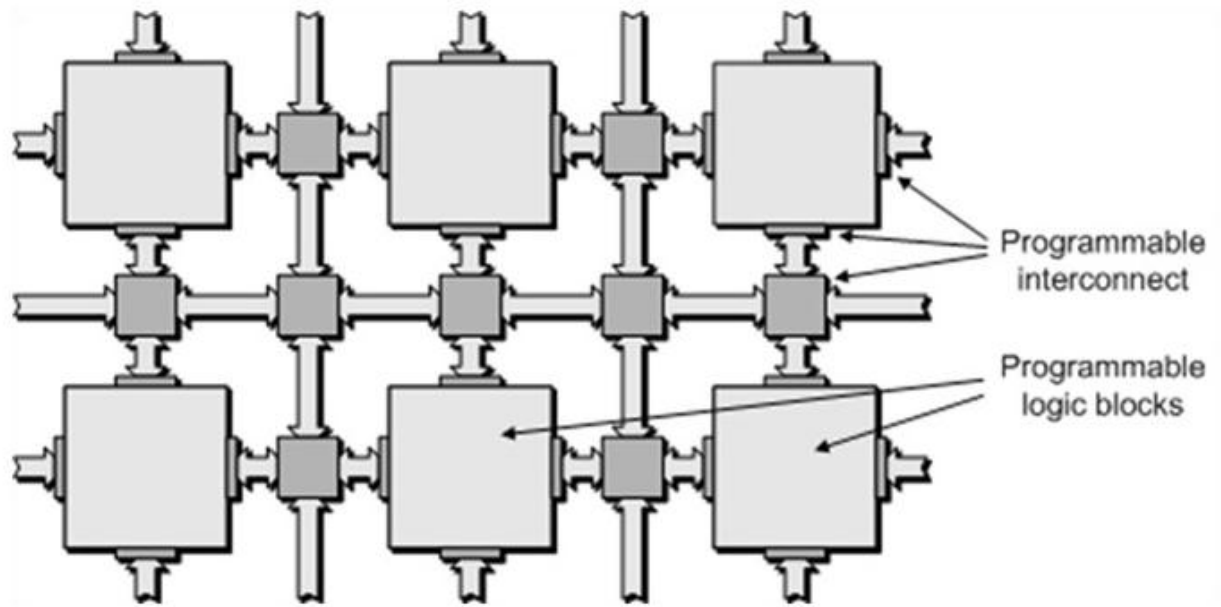


Figure 15 - FPGA Routing

In newer FPGA series other blocks are implemented for application specific functionality. These specialized blocks, or slices, run their specific functionality faster than can be implemented using LUTs and routing wires. These slices also drastically cut down the number of LUTs used. The two main blocks are the multipliers and DSP slices. To implement the multiplication of two 32 bit numbers would require more than 2000 operations for a single multiply. Different FGPA series have different types and quantities of specialized blocks based on designed applications. [22]

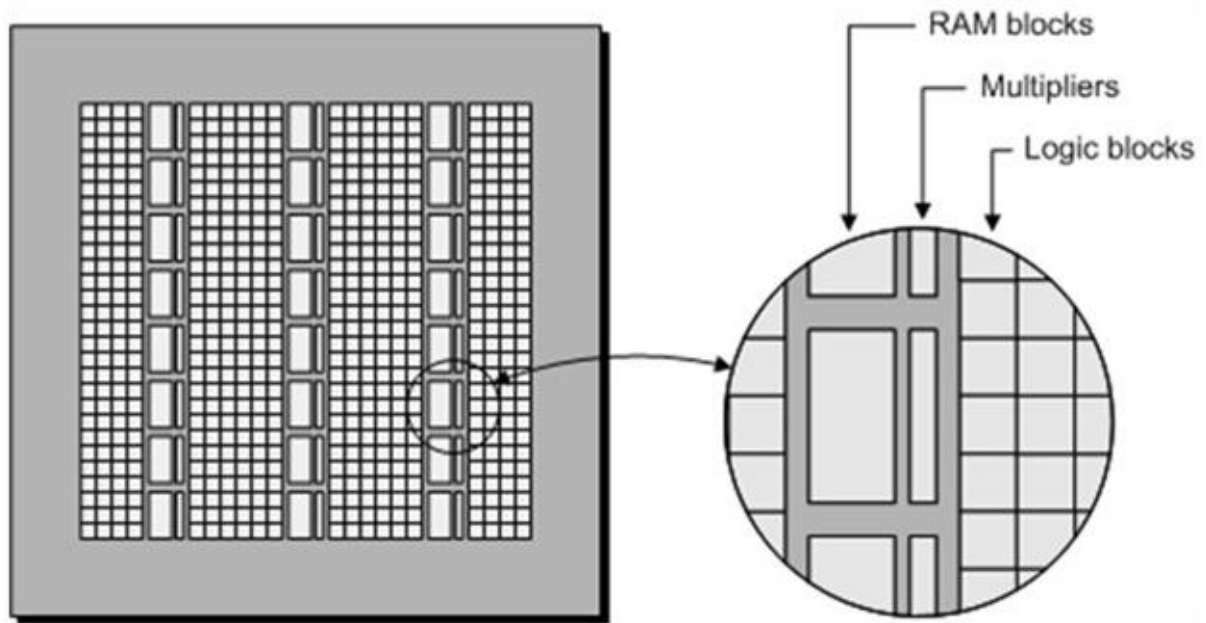


Figure 16 - Specialized Slices

The latest FPGA is the Virtex-7 from Xilinx. This model has increased computing power and efficiency. The architecture for this model can be seen below in Figure 17.

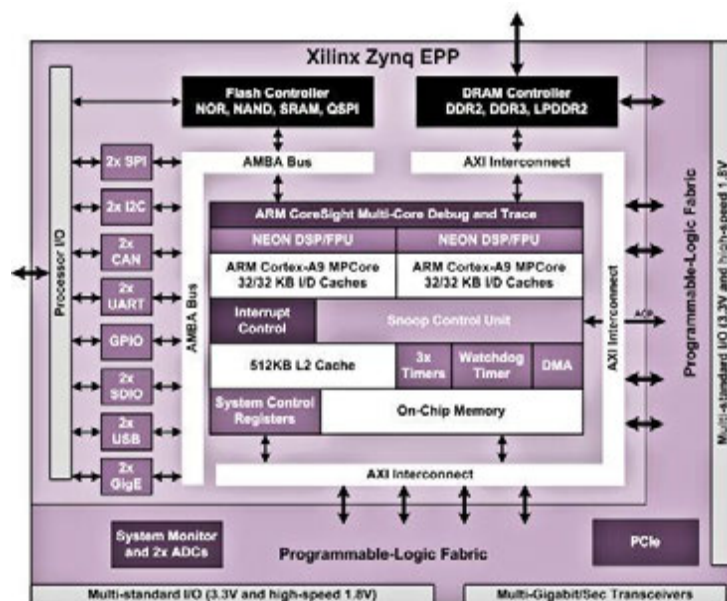


Figure 17 - Virtex-7 Architecture

3. Benchmarks

Benchmarking has been around for years and is widely used as the standard to which we base our computing processing power today. Programs available from many different packages are used to measure the metrics of new processors daily. In the beginning, performance was measured by various system specifications, such as clock rate. Many designers and consumer realized that while this may give some indication of the processing power of a machine, it did not incorporate the entire scope of the situation. Many benchmarks are used today to apply stress to processors through different processes for different applications. Benchmarks come in four different types, each having its own strength and weaknesses: real applications, small benchmarks, benchmark suites, and synthetic benchmarks.

Real applications are benchmarks based pre-existing programs. They are comprised of a typical user's workload during the day. The advantages to running real applications are that they directly translate over to improved performance times on programs and very accurately mirror everyday workloads. However, these benchmarks are usually very big and require more time to run and transfer to other machines. In addition, it can be very hard to pinpoint a processing bottleneck to discover what instruction types need to be improved for the greatest performance increase.

The next type of benchmark is the small benchmark. It consists of a very small code segment that exists in many other applications. For example, the following C code could be a small benchmark.

```
for (j = 0; j<8; j++)
```

```
    S = S + Aj × Bi-j;
```

This code runs very fast and does not take very long to compile or transfer between systems. In addition, it can give developers a very good idea about what portion of their process is bottlenecking the systems in order to make improvements and can be easy to simulate during design to test functionality. However, many designers can take advantage of the limited instructions used to design a system specialized for that particular loop. This abuses the benchmark and does not report accurate data to both designers and customers.

Benchmark suites are compilations of different benchmarks from different industries that together represent a variety of computing loads on a machine. Standard Performance Evaluation Corporation (SPEC) provides one of the main suites available to developers. SPEC began in 1989 with SPEC89 CPU intensive benchmark suite. Many companies came together and agreed on a set of programs that represented a user's typical workload. Suites are very useful in covering a diverse set of parameters and characteristics; however, they are still susceptible to abuse. In addition, they require periodic updates to change the applications as typical workloads change.

Lastly, some programmers advocate the use of synthetic benchmarks. These programs attempt to mirror the characteristics of other applications while using much less space and processing time. In reality, they do not perform any functional task on the processor, but do give an accurate representation of processing power.

It is always important to run many iterations of a benchmark. Timings will change and an average of all iterations should be used to provide a more accurate picture for comparisons. Often manufacturers can find one or more benchmarks that their platform particularly excels at and abuse this standard for advertisement. The best benchmarks to look at for different computers are those posted by non-profit organizations which are unbiased. [23]

3.1. FPGA Benchmarks

Benchmarking is not traditionally done for FPGAs. Rather, the datasheets for different FPGAs have details about maximum clock cycle that they can operate at. The user then looks at the program that they want to run and they can see number of clock cycles it takes to produce a result. The throughput does not change between iterations on the same machine. The following benchmarks that were run on the Virtex-5 series FPGA were compiled using either the Xilinx Core Generator or MATLAB Simulink HDL Coder.

The Xilinx CORE Generator is where most of the applications are from. It is built into the Xilinx program. When opening a core you have several options that can be selected based on your needs such as extra pin I/O, processing type, and size. The Simulink HDL coder is a new addition to the program. Rather than designing our own code, the current demos available in the Simulink program were used to generate HDL code, which is then tested in the Xilinx program. The HDL Coder is also capable of generating testbenches for the application.

Fast Fourier Transform (FFT)

A discrete Fourier transform that operates with reduced computational power. The number of computations needed is reduced from $2N^2$ to $2N \log_2 N$ where N is the number of points necessary for the computation. The FFT used in this study is a 1024 point FFT.

Finite Impulse Response Filter (FIR)

FIR Filters are designed to be simple to implement for Digital Signal Processing (DSP). FIR filters are more commonly used than IIR filters because of the advantages offered such as fractional arithmetic and fewer practical problems. The FIR filter used for performance testing in this study is a low pass filter.

Advanced Encryption Standard (AES)

The Advanced Encryption Standard (AES) specifies a FIPS-approved cryptographic algorithm that can be used to protect electronic data. The AES algorithm is a symmetric block cipher that is capable of encrypting (encipher) and decrypting (decipher) information. Encryption converts data to an unintelligible form called cipher text; decrypting the cipher text converts the data back into its original form, called plaintext. The AES algorithm is capable of using cryptographic keys of 128, 192, and 256 bits to encrypt and decrypt data in blocks of 128 bits. (<http://csrc.nist.gov/publications/fips/fips197/fips-197.pdf>)

Double Precision Floating Point Multiplication

This block takes in two double floating point numbers for multiplication. This is similar to a small benchmark for CPUs. A single FPGA is able to implement multiple instances of this application allowing for large amounts of parallel processing.

3.2. SPEC CPU2006

Standard Performance Evaluation Corporation (SPEC) CPU 2006 is a benchmark suite that contains over twenty-five different benchmarks. These industry-standardized benchmarks were created to stress a system's processor, memory architecture, and compilers. This particular suite is divided into two sub categories, CINT2006 and CFP2006. The CINT2006 portion consists of twelve total benchmarks that measure integer operation performance, while the CFP2006 portion consists of seventeen benchmarks that measure floating point operation performance. A brief description of each SPEC CPU2006 benchmark can be seen below. [24]

3.2.1. CPUINT2006

400.perlbench

"400.perlbench is a cut-down version of Perl v5.8.7, the popular scripting language. SPEC's version of Perl has had most of OS-specific features removed. In addition to the core Perl interpreter, several third-party modules are used:

- SpamAssassin v2.61
- Digest-MD5 v2.33
- HTML-Parser v3.35

- MHonArc v2.6.8
- IO-stringy v1.205
- MailTools v1.60
- TimeDate v1.16" [25]

401.bzip2

"401.bzip2 is based on Julian Seward's bzip2 version 1.0.3. The only difference between bzip2 1.0.3 and 401.bzip2 is that SPEC's version of bzip2 performs no file I/O other than reading the input. All compression and decompression happens entirely in memory. This is to help isolate the work done to only the CPU and memory subsystem." [25]

403.gcc

"403.gcc is based on gcc Version 3.2. It generates code for an AMD Opteron processor. The benchmark runs as a compiler with many of its optimization flags enabled.

403.gcc has had its inlining heuristics altered slightly, so as to inline more code than would be typical on a UNIX system in 2002. It is expected that this effect will be more typical of compiler usage in 2006. This was done so that 403.gcc would spend more time analyzing its source code inputs, and use more memory. Without this effect, 403.gcc would have done less analysis, and needed more input workloads to achieve the run times required for CPU2006." [25]

429.mcf

“429.mcf is a benchmark that is derived from MCF, a program used for single-depot vehicle scheduling in public mass transportation. The program is written in C. The benchmark version uses almost exclusively integer arithmetic.

The program is designed for the solution of single-depot vehicle scheduling sub-problems occurring in the planning process of public transportation companies. It considers one single depot and a homogeneous vehicle fleet. Based on a line plan and service frequencies, so-called timetabled trips with fixed departure/arrival locations and times are derived. Each of these timetabled trips has to be serviced by exactly one vehicle. The links between these trips are so-called dead-head trips. In addition, there are pull-out and pull-in trips for leaving and entering the depot.” [25]

445.gobmk

“The program plays Go and executes a set of commands to analyze Go positions.” [25]

456.hmmr

“Profile Hidden Markov Models (profile HMMs) are statistical models of multiple sequence alignments, which are used in computational biology to search for patterns in DNA sequences.

The technique is used to do sensitive database searching, using statistical descriptions of a sequence family's consensus. It is used for protein sequence analysis.” [25]

458.sjeng

“458.sjeng is based on Sjeng 11.2, which is a program that plays chess and several chess variants, such as drop-chess (similar to Shogi), and 'losing' chess.

It attempts to find the best move via a combination of alpha-beta or priority proof number tree searches, advanced move ordering, positional evaluation and heuristic forward pruning. Practically, it will explore the tree of variations resulting from a given position to a given base depth, extending interesting variations but discarding doubtful or irrelevant ones. From this tree the optimal line of play for both players ("principle variation") is determined, as well as a score reflecting the balance of power between the two.

The SPEC version is an enhanced version of the free Sjeng 11.2 program, modified to be more portable and more accurately reflect the workload of current professional programs.” [25]

462.libquantum

“Libquantum is a library for the simulation of a quantum computer. Quantum computers are based on the principles of quantum mechanics and can solve certain computationally hard tasks in polynomial time. In 1994, Peter Shor discovered a polynomial-time algorithm for the factorization of numbers, a problem of particular interest for cryptanalysis, as the widely used RSA cryptosystem depends on prime factorization being a problem only to be solvable in exponential time. An implementation of Shor's factorization algorithm is included in libquantum.

Libquantum provides a structure for representing a quantum register and some elementary gates. Measurements can be used to extract information from the system. Additionally, libquantum offers the simulation of decoherence, the most important obstacle in building practical quantum computers. It is thus not only possible to simulate any quantum algorithm, but also to develop quantum error correction algorithms. As libquantum allows adding new gates, it can easily be extended to fit the ongoing research, e.g. it has been deployed to analyze quantum cryptography.” [25]

464.h264ref

“464.h264ref is a reference implementation of H.264/AVC (Advanced Video Coding), the latest state-of-the-art video compression standard. The standard is developed by the VCEG (Video Coding Experts Group) of the ITU (International Telecommunications Union, <http://www.itu.int>) and the MPEG (Moving Pictures Experts Group, <http://mpeg.chiariglione.org>) of the ISO/IEC (International Standardization Organization, <http://www.iso.ch>). This standard replaces the currently widely used MPEG-2 standard, and is being applied for applications such as the next-generation DVDs (Blu-ray and HD DVD) and video broadcasting.” [25]

471.omnetpp

“The benchmark performs discrete event simulation of a large Ethernet network. The simulation is based on the OMNeT++ discrete event simulation system (www.omnetpp.org), a generic and open simulation framework. OMNeT++'s

primary application area is the simulation of communication networks, but its generic and flexible architecture allows for its use in other areas such as the simulation of IT systems, queuing networks, hardware architectures or business processes as well. The Ethernet model used in this benchmark is publicly available from the address given in the References.” [25]

473.astra

“473.astar (pronounced: A-star) is derived from a portable 2D path-finding library that is used in game's AI. This library implements three different path-finding algorithms: First is the well-known A* algorithm for maps with passable and non-passable terrain types. Second is a modification of the A* path finding algorithm for maps with different terrain types and different move speed. Third is an implementation of A* algorithm for graphs. This is formed by map regions with neighborhood relationship. The library also includes pseudo-intellectual functions for map region determination.” [25]

483.xalancbmkg

“This program is a modified version of Xalan-C++, an XSLT processor written in a portable subset of C++. You use the XSLT language to compose XSL style sheets. An XSL style sheet contains instructions for transforming XML documents from one document type to another document type (XML, HTML, or other). In structural terms, an XSL style sheet specifies the transformation of one tree of nodes (the XML input) into another tree of nodes (the output or transformation result).” [25]

3.2.2. CFP2006

410.bwaves

“410.bwaves numerically simulates blast waves in three dimensional transonic transient laminar viscous flow.

The initial configuration of the blast waves problem consists of a high pressure and density region at the center of a cubic cell of a periodic lattice, with low pressure and density elsewhere. Periodic boundary conditions are applied to the array of cubic cells forming an infinite network. Initially, the high pressure volume begins to expand in the radial direction as classical shock waves. At the same time, the expansion waves move to fill the void at the center of the cubic cell. When the expanding flow reaches the boundaries, it collides with its periodic images from other cells, thus creating a complex structure of interfering nonlinear waves. These processes create a nonlinear damped periodic system with energy being dissipated in time. Finally, the system will come to an equilibrium and steady state.

The algorithm implemented is an unfactored solver for the implicit solution of the compressible Navier-Stokes equations using the Bi-CGstab algorithm, which solves systems of non-symmetric linear equations iteratively.” [26]

416.gamess

“A wide range of quantum chemical computations are possible using GAMESS. The benchmark 416.gamess does the following computations for the reference workload:

- Self-consistent field (SCF) computation (type: Restricted Hartree-Fock) of cytosine molecule using the direct SCF method
- SCF computation (type: Restricted open-shell Hartree-Fock) of water and cu2+ using the direct SCF method
- SCF computation (type: Multi-configuration Self-consistent field) of triazolium ion using the direct SCF method” [26]

433.milc

“The program generates a gauge field, and is used in lattice gauge theory applications involving dynamical quarks. Lattice gauge theory involves the study of some of the fundamental constituents of matter, namely quarks and gluons. In this area of quantum field theory, traditional perturbative expansions are not useful. Introducing a discrete lattice of space-time points is the method of choice.” [26]

434.zeusmp

“434.zeusmp is based on ZEUS-MP, a computational fluid dynamics code developed at the Laboratory for Computational Astrophysics (NCSA, University of Illinois at Urbana-Champaign) for the simulation of astrophysical phenomena. ZEUS-MP solves problems in three spatial dimensions with a wide variety of boundary conditions.

The program solves the equations of ideal (non-resistive), non-relativistic, hydrodynamics and magnetohydrodynamics, including externally applied gravitational fields and self-gravity. The gas can be adiabatic or isothermal, and the

thermal pressure is isotropic. Boundary conditions may be specified as reflecting, periodic, inflow, or outflow.” [26]

435.gromacs

“435.gromacs is derived from GROMACS, a versatile package that performs molecular dynamics, i.e. simulation of the Newtonian equations of motion for systems with hundreds to millions of particles.

The benchmark version performs a simulation of the protein Lysozyme in a solution of water and ions. The structure of a protein is normally determined by experimental techniques such as X-ray crystallography or NMR spectroscopy. By simulating the atomic motions of these structures, one can gain significant understanding of protein dynamics and function, and, in some cases, it might even be possible to predict the structure of new proteins.” [26]

436.cactusADM

“CactusADM is a combination of Cactus, an open source problem solving environment, and BenchADM, a computational kernel representative of many applications in numerical relativity (ADM stands for ADM formalism developed by Arnowitt, Deser and Misner). CactusADM solves the Einstein evolution equations, which describe how space-time curves as response to its matter content, and are a set of ten coupled nonlinear partial differential equations, in their standard ADM 3+1 formulation. A staggered-leapfrog numerical method is used to carry out the update.” [26]

437.leslie3d

“437.leslie3d is derived from LESlie3d (*Large-Eddy Simulations with Linear-Eddy Model in 3D*), a research-level Computational Fluid Dynamics (CFD) code. It is the primary solver used to investigate a wide array of turbulence phenomena such as mixing, combustion, acoustics and general fluid mechanics.

For CPU2006, the program has been set up to solve a test problem which represents a subset of such flows, namely the temporal mixing layer. This type of flow occurs in the mixing regions of all combustors that employ fuel injection (which is nearly all combustors). Also, this sort of mixing layer is a benchmark problem used to understand physics of turbulent mixing.” [26]

444.namd

“The 444.namd benchmark is derived from the data layout and inner loop of NAMD, a parallel program for the simulation of large biomolecular systems. Although NAMD was a winner of a 2002 Gordon Bell award for parallel scalability, serial performance is equally important to the over 10,000 users who have downloaded the program over the past several years. Almost all of the runtime is spent calculating inter-atomic interactions in a small set of functions. This set was separated from the bulk of the code to form a compact benchmark for CPU2006. This computational core achieves good performance on a wide range of machines, but contains no platform-specific optimizations.” [26]

447.dealll

“The SPEC CPU2006 benchmark 447.dealII is a program that uses deal.II, a C++ program library targeted at adaptive finite elements and error estimation. The library uses state-of-the-art programming techniques of the C++ programming language, including the Boost library. It offers a modern interface to the complex data structures and algorithms required for adaptivity and enables use of a variety of finite elements in one, two, and three space dimensions, as well as time-dependent problems.

The main aim of deal.II is to enable development of modern finite element algorithms, using among other aspects sophisticated error estimators and adaptive meshes. Writing such programs is a non-trivial task, and successful programs tend to become very large and complex.” [26]

450.soplex

“450.soplex is based on SoPlex Version 1.2.1. SoPlex solves a linear program using the Simplex algorithm.” [26]

453.povray

“POV-Ray is a ray-tracer. Ray-tracing is a rendering technique that calculates an image of a scene by simulating the way rays of light travel in the real world but it does so backwards. In the real world, rays of light are emitted from a light source and illuminate objects. The light reflects off of the objects or passes through transparent objects. This reflected light hits the human eye or a camera lens. As the

vast majority of rays never hit an observer, it would take forever to trace a scene. Thus, ray-tracers like POV-Ray start with their simulated camera and trace rays backwards out into the scene. The user specifies the location of the camera, light sources, and objects as well as the surface textures and their interiors.” [26]

454.calculix

“454.calculix is based on CalculiX, which is a free software finite element code for linear and nonlinear three-dimensional structural applications. It uses the classical theory of finite elements described in books such as the work by O.C. Zienkiewicz and R.L. Taylor, "The Finite Element Method", Fourth Edition, McGraw Hill, 1989. CalculiX can be used to solve a variety of problems such as static problems (bridge and building design), buckling, dynamic applications (crash, earthquake resistance) and eigenmode analysis (resonance phenomena).” [26]

459.GemsFDTD

“GemsFDTD solves the Maxwell equations in 3D in the time domain using the finite-difference time-domain (FDTD) method. The radar cross section (RCS) of a perfectly conducting (PEC) object is computed. GemsFDTD is a subset of the code GemsTD developed in the General ElectroMagnetic Solvers (GEMS) project.” [26]

465.tonto

“Tonto is an open source quantum chemistry package, designed by Dylan Jayatilaka and Daniel J. Grimwood. Objectives include simplicity and portability;

aspects not seen in many quantum chemistry codes. The code is easily extendable by chemists with limited programming skills and time, and is easy to understand and use.

Tonto is written within an object oriented design, in Fortran 95. It uses derived types and modules to represent classes. Classes range from integers and text files, through to atoms, space groups and molecules. There is a "self" variable in most routines, which should be familiar from many OO languages. Tonto uses dynamic memory instead of common blocks, and uses array operations where possible." [26]

470.lbm

"This program implements the so-called "Lattice Boltzmann Method" (LBM) to simulate incompressible fluids in 3D as described in. It is the computationally most important part of a larger code which is used in the field of material science to simulate the behavior of fluids with free surfaces, in particular the formation and movement of gas bubbles in metal foams. For benchmarking purposes and easy optimization for different architectures, the code makes extensive use of macros which hide the details of the data access. A visualization of the results of the submitted code can be seen below (flow through a porous medium, grid size 150x150x150, 1000 time steps)." [26]

481.wrf

“481.wrf is based on the Weather Research and Forecasting (WRF) Model, which is a next-generation mesoscale numerical weather prediction system designed to serve both operational forecasting and atmospheric research needs.

WRF features multiple dynamical cores, a 3-dimensional variational (3DVAR) data assimilation system, and a software architecture allowing for computational parallelism and system extensibility. The parallel portions of the code have been turned off for SPEC CPU2006 as the interest here is in single processor performance.” [26]

482.sphinx3

“Sphinx-3 is a widely known speech recognition system from Carnegie Mellon University. The 482.sphinx3 benchmark focuses on the CPU-intensive portion of this speech recognition system.” [26]

3.3. Rodinia Suite

The following benchmarks were taken from the Rodinia Suite.

Leukocyte

The leukocyte application detects and tracks rolling leukocytes (white blood cells) in in vivo video microscopy of blood vessels. The velocity of rolling leukocytes provides important information about the inflammation process, which aids biomedical researchers in the development of anti-inflammatory medications.

In the application, cells are detected in the first video frame and then tracked through subsequent frames. Detection is accomplished by computing for every pixel in the frame the maximal Gradient Inverse Coefficient of Variation (GICOV) score across a range of possible ellipses. The GICOV score for an ellipse is the mean gradient magnitude along the ellipse divided by the standard deviation of the gradient magnitude. The matrix of GICOV scores is then dilated to simplify the process of finding local maxima. For each local maximum, an active contour algorithm is used to more accurately determine the shape of the cell. [27]

LU Decomposition

LU Decomposition is an algorithm to calculate the solutions of a set of linear equations. The LUD kernel decomposes a matrix as the product of a lower triangular matrix and an upper triangular matrix. [27]

SRAD

SRAD (Speckle Reducing Anisotropic Diffusion) is a diffusion method for ultrasonic and radar imaging applications based on partial differential equations (PDEs). It is used to remove locally correlated noise, known as speckles, without destroying important image features. SRAD consists of several pieces of work: image extraction, continuous iterations over the image (preparation, reduction, statistics, computation 1 and computation 2) and image compression. The sequential dependency between all of these stages requires synchronization after each stage (because each stage operates on the entire image). [27]

K-means

K-means are a statistical analysis of clusters where each value is placed into a group where it has the most similar mean. Initial data points are chosen and the subsequent iterations shuffle the data around until there is convergence. K-means is often considered a complex calculation, but over time modified algorithms have improved speed. [27]

Heart Wall

The heart wall benchmark uses several types of processing to create a benchmark testing “braided parallelism”. By using image despeckling and edge detection a new image is produced to detect shapes. Once this process is complete ellipses are added to the process and finally the entire image is tracked from frame to frame. This allows for testing parallelism of both multiple tasks and massive amounts of data. [27]

Hot Spot

This benchmark runs a simulation of processor power and temperature and how cells affect their neighbors. These calculations are done by a series of differential equations. The differential equations are run on a temperature map by taking power usage into account until the entire map has been normalized. [27]

Needleman

This benchmark is a nonlinear optimization of DNA. It finds the optimal path to particular cells. Based on cells surrounding elements, initially filled by the program, backwards calculation is done to find proper alignment. The larger the calculated score the closer to a match there is. Scores are calculated by looking at the north, west, and north western cells and points deducted to missing elements. [27]

Particle Filter

This benchmark estimated locations in noisy environments by looking at the location and path of an object. This is done by making guesses, checking their probability, normalizing the guesses, and updating the location of the object. This implementation of a Particle Filter looking at the speed up provided by GPU parallelism in order to make this application possible in real time applications. For our purposes, trials were run on 1, 2 and 5 million data points in 16, 32 and 64 sized processing grids. [27]

3.4. John Burkardt Benchmarks

John Burkardt is a computer programmer that has been working in the computer science field for many years. His first job was at the Pittsburgh Supercomputing Company from 1988 – 1992. Up until current day, he has had a robust career and worked at several places including Bell Helicopter, the Mathematics Department at Iowa State University, Virginia Tech, and is currently working at the Department of Scientific Computing at Florida State University. Mr.

Burkardt has written many applications, mostly for educational purposes, and two of them were used in this project for their parallel computing attributes using the OpenMP library. These two applications are an FFT benchmark and a Prime Number Counting benchmark. [28]

FFT

This benchmark is a C program that computes a Fast Fourier Transform using parallelism via OpenMP. Included in this program is the ability to change the number of threads used. This allows the user to compare the execution time of different numbers of threads. [28]

Prime Number Counting

This benchmark is a C program that counts the number of primes between 1 and N . In the default case, N is set to 131,072 but can be changed in the source code. For the sake of this project, we only used the default case. This program also allowed the user to change the number of active threads running for comparison of execution times. [28]

3.5. SHOC Suite

Scalable Heterogeneous Computing Benchmark (SHOC) provides benchmarks that are setup to run across GPUs, CPUs, and cluster computing through Message Passing Interface (MPI). While originally designed for use with OpenCL, the current version 1.1.2 supports NVIDIA's CUDA language and MPI. SHOC

provides both stress and performance tests, covering areas from mathematical problem and linear algebra to image processing.

Lennord Jones

This benchmark, based on molecular dynamics, performs nbbody computations using the Lennord Jones potential. N atoms are spread randomly over a cubic domain for this particular execution. [29]

Reduction

As the name suggests, this benchmark is massive sum reduction performed on floating point data. Sets of data points are reduced in individual threads before being looked over again until a single number is reached. [29]

Chemical Modeling (S3D)

Using a three dimensional grid, with one thread per point, this benchmark is a massively parallel calculation of chemical rates using S3D. S3D is a direct numerical solver based on the Navier-Stokes equation. This benchmark relies heavily on floating point calculations, an estimated 10 kFLOPS per thread. [29]

Parallel Prefix Sum

Also known as a scan, this benchmark uses the addition of previous smaller summations to yield a new overall result. With the addition of multiple parallel threads, increasing the number of starting values make this harder to run on modern CPU based systems. [29]

Radix Sort

This benchmark is a radix sorting algorithm where sorting is performed based on individual locations within a number. For example sorting can be based first on the highest digit, such as the hundreds place for example, and then proceed to look at the next level down. Sorting as the algorithm processes each level allows for the list to be sorted quickly in parallel. [29]

2D Stencil Operation

A 9 point stencil computation where each value is updated in turn based on the X-point pattern defined by the algorithm. The standard pattern that updates this algorithm remains the same during this benchmark; however, multiple iterations show speedup in larger systems where multiple points can be processed at once. [29]

Vector Dot Product

This benchmark calculates and checks the bandwidth of a given device while performing a vector dot product operation. Vector dot product results in calculations to determine if the vectors are orthogonal by performing normalization and then looking at the angle between the two vectors. [29]

Device Information

These benchmarks include calculating data about the particular card that is being used to execute other benchmarks. These statistics include the download speed of the device, the devices memory capabilities, and the max flops possible on the device. Since these numbers could all be limiting factors in the calculation of the performance of any one device it is important to know where these numbers lie so that they can be taken into account. [29]

3.6. Parboil Suite

Parboil is a suite designed by the IMPACT group, Illinois Microarchitecture Project utilizing Advanced Compiler Technology, mainly for GPUs. Designed to focus on the massively parallelism systems of today's GPUs, modeling and complex mathematics form the basis of the benchmarks. While the IMPACT group also provides a CPU simulation compiler for CUDA, it is used solely for GPU testing in the context of our project.

Breadth First Search

Breadth First Searching a style of graph searching algorithm that goes through every node in a tree until the correct result is found. This algorithm basically looks through the tree layer by layer, finding children for its next round as it goes. It will not skip to any of the found children until every node on the current level is finished being examined. With this type of exhaustive searching large trees can require large amounts of space and time to compute. Benchmarking with a

breadth first searching algorithm will allow us to see the advantages, if any, to more modern multiprocessor multithreaded systems. [30]

Dense Matrix Multiplication

In matrix multiplication, the dot product of a column by a row is performed. The limitations on this idea are that the two matrices must have similar dimensions, or at least have the x dimension of A equal to the y dimension of B. Dense matrix multiplication follows this same principle, except that the properties of a dense matrix are different than those of a normal matrix. Consider for example that every block in a matrix is connected to every other block in the matrix, so that if you change one you affect the entire matrix. This make multiplication of dense matrices a time consuming task for simple non parallel processors and an excellent way to benchmark the idea of parallel processing in modern hardware. [30]

FFT

The Fast Fourier Transform (FFT) is a very important algorithm for Digital Signal Processing (DSP). Simply put, FFT is a computationally efficient method for calculating the Discrete Fourier Transform (DFT). The DFT is a mathematical operation of transforming a time domain function of finite (discrete) size into a frequency domain function. There are several ways of computing the DFT, FFT being one of them. Without FFT, calculating the DFT can be a long and tedious process. The FFT simplifies this process by breaking the set of data into smaller and smaller chunks and then calculating the DFT. For example, say there were 32 data

points. Trying to calculate the DFT on all 32 points at once would be very difficult. FFT breaks down that data set into smaller and smaller sets until, in this case, there are 16 sets of 2. It is much faster and easier to calculate the DFT at a data size of two than at a data size of 30, especially when parallel computing is involved and several of those 16 calculations can be done at the same time.

Sum of Absolute Differences

Sum of absolute differences (SAD) is an algorithm often used in video encoding as it is a comparison between the current image and the next occurring frame. The algorithm simply finds the absolute value of the given image minus the next before proceeding to add them all together. The smallest set is the closest to the same image. SAD is often used across multiple platforms for compressing video, simple animation, and object recognition. [30]

Distance Cutoff Coulombic Potential

This benchmark computes the Coulombic Potential of each grid point in a 3D matrix. This benchmark relies on the calculation of the unequal dipole forces of a water molecule. Speedup is shown by splitting up the calculations among parallel threads. [30]

Saturating Histogram

Using a 2 dimensional matrix, this benchmark calculates a saturated histogram based on the input data. For our particular set of tests, the input data is that of a silicon wafer with a Gaussian representation of data. [30]

Lattice Boltzman

The Lattice Boltzman benchmark is a fluid dynamics simulation within an enclosed container. Within this benchmark, particle collision and general interaction is calculated using a discrete equation. This method for fluid dynamic simulation is easily ported to GPU benchmarking for its ease in parallelism. [30]

Sparse Matrix Dense Vector Multiplication

The SPMV benchmark is similar to the dense matrix multiplication in that all points are related to one another, thus when one point is acted upon it effects the results of other points. In this case the matrix itself is sparsely populated so there are not as many points to work with. This particular execution uses JDS format so it allows for padding with zeros and multiple alignments. [30]

Two Angular Correction Function

The TPACF benchmark performs a statistical analysis of a spatial distribution, often used when measuring astronomical bodies. It calculates a histogram of distances between every set of objects within the data set. Completion of this benchmark places distances, normally on an exponential curve, within reach of a straight sloped line. Parallelism allows for processing of a multitude of points at one time. [30]

4. Results

4.1. Devices

The following is the exact specifications for the devices that these benchmarks were tested on. This information includes: clock rates, system memory, and host operating system among other specifications.

4.1.1. GPU

All of the GPU benchmarks were run on the same system. Since transfer time and driver commands are issued by the CPU we must look at the whole system. The test system is running Ubuntu 10.04 with the Linux kernel version 2.6.32-38-generic. Other system specifications are as follows:

- CPU: AMD Athlon 64 x2 5200+ running at 2.611GHz dual core
- 3GB of system memory, DDR2 800MHz
- GPU 1: NVIDIA GeForce GTX 460
 - CUDA Capability 2.1
 - 336 CUDA capable cores
 - 1024MB of global onboard memory
 - GPU clocks are as follows:
 - Graphics 675MHz
 - Memory 1800MHz
 - Processor 1350MHz
- GPU 2: NVIDIA GeForce 9800 GTX+
 - CUDA Capability 1.1
 - 128 CUDA capable cores
 - 512MB of global onboard memory
 - GPU clocks are as follows:

- Graphics 738MHz
- Memory 1100MHz
- Processor 1836MHz
- System PCIE version 1.1

4.1.2. CPU

All of the CPU benchmarks were also run on the same system known at Worcester Polytechnic Institute as the AMAX machine. This system has the following specifications:

- 2x – Intel Xeon X5650 processors running at 1.6 GHz. The processor is capable by factory default to run at 2.67GHz. These processors have 6 cores and 12 threads a piece, totaling 12 cores and 24 threads using Hyper Threading.
- The total memory of this machine is 24.6 GB.
- The operating system running on this machine is Linux version 2.6.18-274.7.1.el5.
- The compiler is GCC, G++, GFortran versions 4.3.4.
- Intel C++ Compiler XE version 12.0.1.116 build 2010116

4.1.3. FPGA

The FPGA benchmarks were performed using simulations on a Virtex-5 family board. All simulations were performed using the ISIM program through the Xilinx program. Virtex-5 boards come in many different series and are specialized for the following applications taken from the Xilinx datasheet:

- Virtex-5 LX: High-performance general logic applications
- Virtex-5 LXT: High-performance logic with advanced serial connectivity

- Virtex-5 SXT: High-performance signal processing applications with advanced serial connectivity
- Virtex-5 TXT: High-performance systems with double density advanced serial connectivity
- Virtex-5 FXT: High-performance embedded systems with advanced serial connectivity [31]

4.2. All Devices

As a comparison between all three platforms we looked to the Fast Fourier Transform. The FFT benchmark was chosen since it was easily portable between all three devices. For the CPU the benchmark was acquired from John Burkardt, while the GPU one came from the Parboil Suite and the FPGA core came from MATLAB HDL Coder. While the benchmark came from different sources, the FFT implementation was similar. A factor for concern could come from the use of the different style compilers, making different optimizations to the code. The timings from this benchmark can be seen here in Table 1.

FFT Benchmark 262,155 Points							
Bench mark	CPU					GPU	FPGA
	2 Threads	4 Threads	8 Threads	12 Threads	24 Threads	Max Threads	Virtex-5
FFT	76.17 ms	45.41 ms	31.63 ms	27.85 ms	31.36 ms	8.13 us (Execution)	2.59 ms

Table 1 - FFT Results

Using these timings it is obvious that the GPU processed the FFT the fastest. The GPU throughput at an FFT this size is 311 times more than that of the FPGA and 3351 times faster than the CPU Multicore with 12 threads. What is missing from this data is the transfer time for the GPU to send and receive the data. This time is 69 ms for this size data. With this time added in the GPU actually becomes the slowest of all three platforms. Now, the FPGA is the fastest at 26.67 times faster than the GPU and 10.76 times faster than the CPU.

What these results can infer is that depending on the hookup of your GPU to minimize transfer time you could be better off using another platform. As GPU transfer times become smaller and smaller, it is evident that in terms of speed, your best choice will be a GPU.

4.3. CPU & GPU

In this section, we will discuss how the CPU and GPU compared against each other. The three benchmarks that we were able to test across both of these platforms were all from the Rodinia benchmark suite. The Leukocyte, LU Decomposition, and Speckle Reducing Anisotropic Diffusion (SRAD) benchmarks were the three that were executed and the average timings from three runs can be seen as a comparison in Table 2 below.

Benchmark	CPU					GPU
	2 Threads	4 Threads	8 Threads	12 Threads	24 Threads	Max Threads
Leukocyte	11.70 s	6.11 s	3.65 s	2.16 s	1.67 s	0.157 s
LU Decomposition	242.82 ms	130.00 ms	76.40 ms	69.23 ms	145.20 ms	7.38 ms
Speckle Reduction	638.14 ms	415.07 ms	306.41 ms	283.93 ms	492.26 ms	282.93 ms

Table 2 - CPU & GPU Results

Looking at the above data in most cases the GPUs beats out CPUs in raw processing. Raw processing is something to note because, for the GPU iterations, data is all moved to the GPU's global memory before testing. This keeps the GPU from having to ask for data over the slower system bus.

Looking at the times required for the Leukocyte benchmark there is a 10.6x speedup between the 24 Hyper threads of the AXAM machine and the CUDA based GTX 460. Again this is looking at raw processing as even in the benchmark papers it is stated that "Although the overall kernel executed in slightly less than a second, the memory allocation and copying overheads add more than eleven seconds to the overall runtime." [32]

The ultimate success of CUDA results from the many code optimizations and compacting the code into a single kernel; thus, removing much of the extra overhead that decreases speed.

Moving onto the LU Decomposition, it is easy to see that the GPU was much faster than the CPU running this benchmark. This is most likely due to the functionality of LU Decomposition, which is purely linear algebra mathematics.

Since this benchmark is purely mathematical based, the 336 cores and consequently the numerous ALUs of the GPU would be must faster in this type of computation than the 12 cores of the CPU. The only place that the CPU would be able to make up for its slower computation would be in the data transfer portion of this benchmark. In this study though, we are only looking at raw execution and not data transfer, so the GPU is much faster than the CPU in the computation of the LU Decomposition benchmark.

The Speckle Reduction benchmark uses mostly partial differential equations in its computations. As described in the background section, this benchmark consists of several, sequential parts, making synchronization of these parts a necessary feature. This necessity for synchronization is what most likely allowed the CPU execution time to equal that of the GPU. In the flat out computational portions of this benchmark, the GPU would most likely beat the CPU due to its sheer number of available cores. Advancements in CPU pipelining, however, allow the CPU to begin the next instruction while the current one is executing. This basically means that the CPU has the ability to synchronize and start its next instruction at the same time. We believe that is the uniqueness of CPU pipelining that allows it to compete with the GPU in execution time for this benchmark.

4.4. Individual Results

4.4.1. GPU Specific

For any GPU today the biggest area of a setback is in the area of data transfer. Until recently, and still the majority, GPUs were only on dedicated cards requiring a

link to the CPU and memory over a data bus. The limitation of these busses to transfer the computed data severely limits the abilities of GPUs in high performance computing. While bus speeds have come a long way since the original PCI and AGP graphics interfaces, even the newest version of PCIe limits speeds to 32 GB/sec assuming full utilization of every channel both upload and download.

PCIe Version	Raw Bitrate (GT/sec)	Total Bandwidth (GB/sec)
1.0a	2.5	8
1.1	2.5	8
2.0	5	16
2.1	5	16
3.0	8	32
4.0 (theoretical)	16	X

Table 3 - GeForce Specification

Just looking at the specifications for the GeForce GTX460 we can see that its maximum memory bandwidth is clocked at 115.2 GB/sec which is far in excess of the 32 GB/sec we have with the current transfer standard. Looking at the SHOC Benchmark Suite we can view download speeds and the speed of the onboard memory of the tested cards as well. In the following table it is clear that larger file sizes allow for use of more bandwidth, but the graphics cards simply are not capable of utilizing the full bandwidth. The maximum data rate we see below is 3.278 GB/sec whereas the host systems capability was 4 GB/sec in both the upload and download channels. With small file sizes, such as the 1kB file, the graphics card barely reaches a tenth of its potential. As the file size increases past the 512kB point, the time required to move the data increases in a linear fashion with a factor of two.

Focusing on the transfer capabilities of new generations of GPUs would allow for much faster processing of large data sets by significantly reducing transfer times.

Size of Data Chunk (kB)	GB / second	Time (milliseconds)
1	0.0968	0.01057
2	0.1919	0.01067
4	0.3399	0.01204
8	0.6074	0.01348
16	1.0987	0.01491
32	1.5839	0.02069
64	2.1322	0.03074
128	2.5777	0.05084
256	2.8825	0.09094
512	3.0872	0.1698
1024	3.1711	0.3306
2048	3.1773	0.6602
4096	3.1769	1.3211
8192	3.2272	2.5995
16384	3.2552	5.1541
32768	3.2357	10.369
65536	3.2696	20.525
131072	3.2733	41.004
262144	3.2781	81.887
524288	3.2783	163.768

Table 4 - Data Size vs. Data Rate

In contrast to moving data to and from the device, moving data around between the internal memory levels is a faster process. Looking below at Table 5 gives us another look at how slow global memory transfers are in comparison to the internal capabilities of the GPU. The internal memory movement is on the order of a hundred to three hundred GB/sec.

Read / Write	Memory Block Size	Speed (GB/sec)
Local	32	190 / 181
Local	64	288 / 328
Local	128	299 / 388
Local	256	289 / 385
Local	512	277 / 368
Global	32	7.4 / 3.7
Global	64	5.8 / 3.5
Global	128	4.8 / 3.4
Global	256	4.3 / 3.4
Global	512	4.1 / 3.3

Table 5 – Global vs. Local Memory Speeds

As we can see in the table of timings below, from the Parboil Suite of Benchmarks, large portions of time during any given benchmark is transferring data to the GPU or the result back to main memory. While transferring data is not necessarily the bulk of the time, a more efficient method for transferring data is needed to help speed up the overall computation times. As seen in

arboil	Average - 460	Average - 9800	Percent Diff
CUTCP			
GPU	0.037495333	0.036806	-0.92775365
Copy	0.006528333	0.006477667	-0.38956379
FFT			
GPU	0.000812667		
Copy	0.069062		
LBM - long			
GPU	30.15880433	93.86147667	51.36472182
Copy	0.354906667	0.334240333	-2.99882802
LBM - short			
GPU	1.006141333	3.124029333	51.2784621
Copy	0.340840333	0.332132333	-1.29396043
MM - long			
GPU	0.008865333	0.010750667	9.61120174
Copy	0.010738333	0.010793667	0.256981856
SAD			
GPU	0.001495667	0.002197667	19.00722022
Copy	0.115957	0.064151	-28.7638528
SPVM - large			
GPU	0.000219333	0.000319333	18.56435644
Copy	0.076647333	0.054017333	-17.3191426
SPVM - medium			
GPU	0.000117	0.000114	-1.2987013
Copy	0.051963667	0.047360667	-4.6343125
SPVM - small			
GPU	4.73333E-05	4.93333E-05	2.068965517
Copy	0.049963667	0.045377	-4.81081875
TPACF			
GPU	1.361817	1.343129333	-0.69087015
Copy	0.081934	0.051611333	-22.7058976

Table 6 the smallest transfer time is still on the order of milliseconds (6.528 ms); magnitudes larger than a single cycle of execution for today's modern CPUs and GPUs (0.5 nanosec). An interesting note comes from the comparison of the older 9800 card to the current 460; the copy times actually increased by an average of

9.18% for the new architecture. This may be attributed to the added memory banks or new hierarchy.

Parboil	Average - 460	Average - 9800	Percent Diff
CUTCP			
GPU	0.037495333	0.036806	-0.92775365
Copy	0.006528333	0.006477667	-0.38956379
FFT			
GPU	0.000812667		
Copy	0.069062		
LBM - long			
GPU	30.15880433	93.86147667	51.36472182
Copy	0.354906667	0.334240333	-2.99882802
LBM - short			
GPU	1.006141333	3.124029333	51.2784621
Copy	0.340840333	0.332132333	-1.29396043
MM - long			
GPU	0.008865333	0.010750667	9.61120174
Copy	0.010738333	0.010793667	0.256981856
SAD			
GPU	0.001495667	0.002197667	19.00722022
Copy	0.115957	0.064151	-28.7638528
SPVM - large			
GPU	0.000219333	0.000319333	18.56435644
Copy	0.076647333	0.054017333	-17.3191426
SPVM - medium			
GPU	0.000117	0.000114	-1.2987013
Copy	0.051963667	0.047360667	-4.6343125
SPVM - small			
GPU	4.73333E-05	4.93333E-05	2.068965517
Copy	0.049963667	0.045377	-4.81081875
TPACF			
GPU	1.361817	1.343129333	-0.69087015
Copy	0.081934	0.051611333	-22.7058976

Table 6 – Parboil Suite Timings

Continuing to look at the Parboil Suite, there are other noticeable improvements with the new Fermi architecture. The full table of results can be seen in the appendices. Using the new architecture, the time spent on interactions between the CPU and GPU decreased by a noticeable amount. CPU computation time decreased by 8.77% and the time the CPU spent handling GPU commands

decreased by 2.21%. These increases show how newer GPUs are capable of handling more commands by themselves and, while still reliant on CPUs, are moving towards being able to compute by themselves.

Moving to the data from the Rodinia suite other observations on GPU processing can be made. Once again, as with the Parboil suite, we can see that the time taken to move memory around with the newer GTX 460 is still slower. Overall this effect could be due to the test system's PCIE bus. We did not have an intermediate GPU to test floating point speedup with the new architecture, but using the 9800 GTX+ we could observe native increase. Looking at the table below we see that the older GPU beat the new architecture by a very slight margin; an average increase of -0.97% for the GTX 460. Since the older architecture was optimized for standard arithmetic and did not support floating point, this result is reasonable.

Rodinia	Average - 460	Average - 9800	Percent Diff
Particle Float (non-float)			
A			
GPU execution	0.000117667	0.000116889	-0.3315964
B			
GPU execution	0.000129333	0.000129444	0.042936883
C			
GPU execution	0.000126333	0.000123444	-1.15658363
D			
GPU execution	0.000120333	0.000116111	-1.78571429
E			
GPU execution	0.000125333	0.000122111	-1.30220027
F			
GPU execution	0.000123	0.000121	-0.81967213
G			
GPU execution	0.000123667	0.000119889	-1.55109489
H			
GPU execution	0.000129667	0.000127556	-0.82073434

Table 7 - Rodinia Particle Filter Results

4.4.2. CPU Results

SPEC CPUINT2006 Results

There were two instances of the CPUINT2006 benchmark suite run on the CPU for this project. The first run through was done on only a single core and single thread using the GCC, G++, and GFortran 4.3.4 (GNU) compiler. The second was done using the Intel C++ Compiler XE. The GNU compiler run through had the auto parallel feature off, meaning it used only a single thread. The Intel compiler run through had the auto parallel feature enabled, meaning it was using multicore parallelism to complete the benchmark. The results of these two runs can as well as the overall speedup percentages can be seen in the following tables.

	Iteration #1 [s]	Iteration #2 [s]	Iteration #3 [s]
400.perlbench	398.00	396.00	396.00
401.bzip2	582.00	582.00	581.00
403.gcc	375.00	375.00	375.00
429.mcf	373.00	373.00	371.00
445.gobmk	510.00	511.00	511.00
456.hmmer	869.00	869.00	869.00
458.jeng	595.00	612.00	595.00
462.libquantum	508.00	506.00	506.00
464.h264ref	710.00	706.00	709.00
471.omnetpp	374.00	373.00	374.00
473.atar	494.00	495.00	494.00
483.xalancbmk	268.00	260.00	259.00

Table 8 - SPEC CPUINT2006 w/ GCC, G++, GFortran compiler w/o auto parallel

	Iteration #1 [s]	Iteration #2 [s]	Iteration #3 [s]
400.perlbench	410.00	411.00	411.00
401.bzip2	539.00	539.00	538.00
403.gcc	330.00	331.00	334.00
429.mcf	184.00	184.00	184.00
445.gobmk	600.00	601.00	607.00
456.hmmer	213.00	214.00	214.00
458.jeng	508.00	508.00	507.00
462.libquantum	12.20	14.30	13.00
464.h264ref	959.00	959.00	1031.00
471.omnetpp	339.00	340.00	339.00
473.atar	350.00	349.00	349.00
483.xalancbmk	214.00	214.00	214.00

Table 9 - SPEC CPUINT2006 run with Intel compiler in auto parallel

	Iteration Speedup	#1	Iteration Speedup	#2	Iteration Speedup	#2
400.perlbench	-3.02%		-3.79%		-3.79%	
401.bzip2	7.39%		7.39%		7.40%	
403.gcc	12.00%		11.73%		10.93%	
429.mcf	50.67%		50.67%		50.40%	
445.gobmk	-17.65%		-17.61%		-18.79%	
456.hmmer	75.49%		75.37%		75.37%	
458.jeng	14.62%		16.99%		14.79%	
462.libquantum	97.60%		97.17%		97.43%	
464.h264ref	-35.07%		-35.84%		-45.42%	
471.omnetpp	9.36%		8.85%		9.36%	
473.atar	29.15%		29.49%		29.35%	
483.xalancbmk	20.15%		17.69%		17.37%	
Average Increase Per Benchmark	21.72%		21.51%		20.37%	
Total Average Increase			21.20%			

Table 10 - Speedup percentages from the GNU run to the Intel run

From the tables above, you can see that the multicore, auto parallel run through generally yielded a quicker execution time than the run through with only a single thread. There were, however, a few of the benchmarks that actually ran better using only a single thread. This is probably due to the fact that the work load of these specific benchmarks was better optimized for only a single thread. Overall, the multicore run through produced a 21.2% increase in speed over than of single thread run.

SPEC CFP2006 Results

Once again, there were two instances of the SPEC CPU2006 floating point suite benchmarks run on the CPU. There were two benchmarks in this suite, *bwaves* and *wrf*, which were left out of these runs because of an invalid run error. Both of these benchmarks would build successfully, but every time we tried to run them, we would get this invalid run error that we could not figure out how to fix. Otherwise, all the other benchmarks in the floating point suite ran fine and there timings can be seen in the following two tables. The first table shows the timings for the run on the GCC, G++, and GFortran compiler with the auto parallel feature off (single threaded). The subsequent table shows the timings of the run on the Intel compiler with the auto parallel feature on (multi-threaded).

	Iteration #1 [s]	Iteration #2 [s]	Iteration #3 [s]
416.gamess	937	940	938
433.milc	479	463	489
435.gromacs	579	579	578
436.cactusADM	1372	1441	1338
437.leslie3d	604	604	603
444.namd	496	497	496
447.dealII	430	429	430
450.soplex	270	270	283
453.povray	236	235	237
454.calculix	1484	1484	1484
459.GemsFDTD	517	517	517
465.tonto	652	649	652
470.lbm	378	379	378
482.sphinx3	632	630	632
434.zeusmp	623	625	625

Table 11 - SPEC CFP2006 w/ GCC, G++, GFortran compiler w/o auto parallel

	Iteration #1 [s]	Iteration #2 [s]	Iteration #3 [s]
416.gamess	1238	1185	1197
433.milc	190	189	190
435.gromacs	485	489	482
436.cactusADM	60.9	53.1	50.3
437.leslie3d	87.7	90.5	95.8
444.namd	457	456	457
447.dealII	293	293	293
450.soplex	296	263	286
453.povray	191	191	190
454.calculix	382	292	375
459.GemsFDTD	119	122	121
465.tonto	469	462	469
470.lbm	49.9	49.9	50.1
482.sphinx3	528	544	514
434.zeusmp	93.9	93	90.8

Table 12 - SPEC CFP2006 run with Intel compiler in auto parallel

Though there were a few benchmarks that produced a negative speed increase from the single threaded to the multi-threaded run, most benchmarks

displayed a significant increase in speed. The speedup of each benchmarking iteration, as well as the average total speedup of all the run-throughs, can be seen in the following table.

	Iteration #1 Speedup	Iteration #2 Speedup	Iteration #3 Speedup
416.gamess	-32.12%	-26.06%	-27.61%
433.milc	60.33%	59.18%	61.15%
435.gromacs	16.23%	15.54%	16.61%
436.cactusADM	95.56%	96.32%	96.24%
437.leslie3d	85.48%	85.02%	84.11%
444.namd	7.86%	8.25%	7.86%
447.dealII	31.86%	31.70%	31.86%
450.soplex	-9.63%	2.59%	-1.06%
453.povray	19.07%	18.72%	19.83%
454.calculix	74.26%	80.32%	74.73%
459.GemsFDTD	76.98%	76.40%	76.60%
465.tonto	28.07%	28.81%	28.07%
470.lbm	86.80%	86.83%	86.75%
482.sphinx3	16.46%	13.65%	18.67%
434.zeusmp	84.93%	85.12%	85.47%
Average Increase Per Benchmark	42.81%	44.16%	43.95%
Total Average Increase	43.64%		

Table 13 - Speedup percentages from the GNU run to the Intel run

As you can see, there was an average total speedup of 43.64% for the SPEC CPU2006 floating point suite. This is over two times the speedup we saw from the CPU2006 integer suite. We believe that such an increase in speedup is due to the complexity of floating point operations. A single thread would be able to run simpler integer operations faster than complex floating point operations. Thus, you would see a more significant speed increase when using the multi-threaded

capabilities of a CPU to calculate floating point operations as opposed to integer operations.

Rodinia / John Burkardt Results

Each of the Rodinia and Burkardt benchmarks was run five different times at three iterations a piece using 2, 4, 8, 12, and 24 threads. Since the CPU we were testing on has two processors, totaling 12 cores, we believed that this range of thread counts would cover the most practical multithreading situations. The following table shows the average execution time for three iteration of each benchmark on the various thread counts. The first five benchmarks in this table are from the Rodinia suite and the other two are from Burkardt suite.

	2 threads	4 threads	8 threads	12 threads	24 threads
Leukocyte (s)	11.70	6.11	3.65	2.16	1.67
LU Decomposition (ms)	242.82	130.00	76.40	69.23	145.20
Speckle Reduction (ms)	638.14	415.07	306.41	283.93	492.26
Means (s)	3.35	3.67	2.91	2.16	1.67
Stream Clusters (s)	47.08	25.18	14.61	11.91	10.40
FFT (ms)	76.17	45.41	31.63	27.85	31.36
Primes (s)	2.04	1.17	0.64	0.44	0.36

Table 14 - Rodinia/Burkardt average execution time on 2, 4, 8, 12, 24 threads

As you can see, the execution times change when transitioning to a different number of threads. The total speedup percentage change between thread counts can be seen in the following table.

	2-4 threads	4-8 threads	8-12 threads	12-24 threads
Leukocyte (s)	47.78%	40.26%	40.82%	22.69%
LU Decomposition (ms)	46.46%	41.23%	9.38%	-109.74%
Speckle Reduction (ms)	34.96%	26.18%	7.34%	-73.37%
kmeans (s)	-9.55%	20.71%	25.77%	22.69%
Stream Clusters (s)	46.52%	41.98%	18.48%	12.68%
FFT (ms)	40.38%	30.35%	11.95%	-12.60%
Primes (s)	42.65%	45.30%	31.25%	18.18%
Average Increase Between Threads	35.60%	35.14%	20.71%	-17.07%

Table 15 - Speedup between thread counts

Several interesting observations can be made from the information in the table above. First, we saw that three out of the seven benchmarks had a decrease in speed when transitioning between 12 and 24 threads. This is most likely due to the hardware limitations of our processor. At twelve threads, our processor can dedicate one core to each thread because it has a total of 12 cores. Once we transition to 24 threads though, 12 of the 24 threads become “virtual” threads that are implemented at the software level. This basically means that each core is handling two threads apiece. While this technology may be good for some applications, running two threads on a single core can sometimes produce slower execution times than using a single thread per core.