

PERFORMANCE ESTIMATION AND MAPPING OF APPLICATIONS ONTO GPU_s

ARUN PARAKH



DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING

INDIAN INSTITUTE OF TECHNOLOGY DELHI

AUGUST 2016

©Indian Institute of Technology Delhi (IITD), New Delhi, 2016

PERFORMANCE ESTIMATION AND MAPPING OF APPLICATIONS ONTO GPU_s

by

ARUN PARAKH

Department of Computer Science and Engineering

Submitted

in fulfillment of the requirements of the degree of Doctor of Philosophy

to the



Indian Institute of Technology Delhi
AUGUST 2016

Certificate

This is to certify that the thesis titled **Performance Estimation and Mapping of Applications onto GPUs** being submitted by **Mr. Arun Parakh** for the award of **Doctor of Philosophy** in Computer Science and Engineering is a record of bona-fide work carried out by him under our guidance and supervision at the Computer Science and Engineering Department, Indian Institute of Technology Delhi. The work presented in this thesis has not been submitted elsewhere, either in part or full, for the award of any other degree or diploma.

M. Balakrishnan

Professor

Dept. of Computer Science & Engineering

Indian Institute of Technology Delhi

New Delhi, India 110 016

Kolin Paul

Associate Professor

Dept. of Computer Science & Engineering

Indian Institute of Technology Delhi

New Delhi, India 110 016

This thesis is dedicated to..

Parent

&

Teacher

&

Student

Acknowledgments

Though only my name appears on the cover of this thesis, a great many people have contributed to its production. I owe my gratitude to all those people who have made this thesis possible and because of whom my Phd experience has been one that I will cherish forever.

My deepest gratitude is to my supervisor, Prof. M. Balakrishnan. I have been amazingly fortunate to have an guide who gave me the freedom to explore on my own, and at the same time the guidance to recover when my steps faltered. I hope that one day I would become as good a supervisor to my students as Prof. Bala has been to me.

My co-advisor, Dr. Kolin Paul, has been always there to listen and give advice. I am deeply grateful to him for the long discussions that helped me sort out the technical details of my work. Dr. Kolin always inspired me to improve my knowledge for Computer Architecture with his sound and clear concepts of this domain.

This research would not be possible without support of Mayura. She constantly monitored the work progress and feedback honestly. I am thankful to her to let me work around the clock. Writing skill of Vayam and Noopur helped me alot to prepare draft of the thesis as well as research papers so I am thankful to them.

I would like to thank my colleague Vaibhav Jain, Sarat Chandra Varma, Sohan Lal Sharma, Atul Sharma, Lava Bhargava, Rajesh Kumar Pal, Yamuna Prasad and all the other colleagues for their cooperation and support. I wish to thank Prof. Anshul Kumar, Prof. G. S. Visweswaran, and Prof. Preeti Ranjan Panada for their valuable feedback and encouragement. I would also like to thank the staff members of Computer Science department especially Mr. S. D. Sharma and Mrs. Vandana Ahluwalia for their help. They have opened their hearts and their doors, supplying me with resources I never could have found without them.

No one walks alone on the journey of life. I devote my deepest gratitude to my father and mother for their unlimited love and support. They have encouraged me throughout the years of my study. The most special thanks belong to my wife, Vineeta, for her understanding about my leaving during all these years, selfless love and support all along and encouragement and company during the preparation of this thesis. I am also grateful to my kids, Aaradhya and Aashman, for their support to keep me happy and relaxed in many tense situations. I have to give a special mention for the support given by Sanjay, Sangeeta and Manish. Their worry about me inspired me to finish this research successfully as soon as possible. I warmly appreciate the generosity and understanding of my in-laws. My father-in-law and mother-in-law has supported even during the toughest period of their life.

Arun Parakh

Abstract

Researchers strive to develop new techniques and tools to efficiently and effectively utilize hardware resources that are made available for compute intensive applications. The Graphics Processing Unit (GPU) is one such hardware, which is very popular in the domain of High Performance Computing (HPC). Our work here presents a model to estimate performance of various applications on a modern GPU with Fermi architecture from NVIDIA. First we try to estimate computation time and then follow it up with an estimation of memory access time. We have found that our average estimation errors for these applications range from -7.76% to 55%. We have analyzed the existing Map-Reduce (MR) framework and enhanced the same for new GPU architectures. Our experiments show an average of 2.5x speedup of MR framework. We have achieved performance benefits ranging from 10% to 100% for various cache sizes. Based on the above analysis, three techniques have been developed for the performance enhancement of MR framework. Using these techniques, we have saved over 32% cache miss per thread. We have reduced the number of comparisons per thread in the group phase and gained an average of 1.5x speed up in the group phase. Use of delayed writing, reduces significant cache misses and improves the execution time by 10% to 25%. The extension of this work is to reduce the performance gap between MR and native implementation of applications in CUDA (*onlyCUDA*). This work reports deployment of three complete algorithms (SW, NB and BF) and achieves penalty reduction of 5x, 6.45x and 15.87x respectively. Finally, we try to execute applications with data size larger than memory capacity, on heterogeneous platforms. This data needs to be broken into partitions of different sizes and processed in parallel on different GPUs with different characteristics. We propose a technique to improve resource utilization, faster execution and also ensure simultaneous completion of such large data applications on heterogeneous platforms.

Contents

List of Figures	v
List of Tables	ix
List of Abbreviations	xi
1 Introduction	1
1.1 Major contributions of the thesis	2
1.2 Thesis organization	7
2 Literature Survey	9
2.1 Performance modeling for GPUs	12
2.2 MR framework	14
2.3 Implementation of complete algorithms	18
2.4 Data partitioning and distribution	19
2.5 Data transfer latency	24
2.6 Summary	24
3 GPU Architecture Evolution and Programming Model	25
3.1 GPU architecture	25
3.2 Overview of MR framework	32

4	Performance Estimation of GPUs with Cache	39
4.1	Introduction	39
4.2	Methodology	40
4.2.1	Instruction count	41
4.2.2	Memory access time	47
4.2.3	Estimated execution time	54
4.3	Experiments	54
4.4	Results	56
4.5	Summary	59
5	Performance Enhancement of Map-Reduce Framework on GPU	61
5.1	Introduction	61
5.2	Applications used	62
5.3	Experiments and analysis	65
5.3.1	<i>onlyCUDA</i> execution	65
5.3.2	MR framework execution	67
5.3.3	L1 re-configuration effect	70
5.4	Enhancement in MR framework	74
5.4.1	Code restructuring	75
5.4.2	Group phase enhancement	76
5.4.3	The auxiliary functions enhancement (Delayed Writing)	78
5.4.4	Results	81
5.5	Summary	81
6	Applying Methodology to Complete Algorithms	85
6.1	Motivation	85
6.2	Optimization methodology	86

CONTENTS	iii
6.3 Implementation of proposed methodology	87
6.3.1 Application 1: SmithWaterman algorithm (SW)	89
6.3.2 Application 2: N-Body simulation (NB)	94
6.3.3 Application 3: Blowfish algorithm	97
6.4 Summary	104
7 Mapping Applications onto Heterogeneous GPU Platforms	107
7.1 Hardware and Software platform	109
7.2 Applications	110
7.3 Methodology of Mapping	110
7.3.1 Data partitioning	110
7.3.2 Performance estimation of applications for different sizes	111
7.3.3 Mapping of data partitions on heterogeneous platform	114
7.4 Experiments and Results	122
7.5 Summary	125
8 Conclusion and Future Work	127
8.1 Conclusion	127
8.2 Future scope	129
Bibliography	131
Appendix A Important Code Snippet	143
Appendix B BLOSUM50 matrix	151
Appendix C Partition Tables	153

List of Figures

3.1	Traditional (old) GPU Architecture	27
3.2	Fermi GPU Architecture	28
3.3	Case Study: Fermi GPU vs Traditional GPU	29
3.4	CUDA programming model	31
3.5	Execution flow of MARS	33
4.1	GPU execution style	40
4.2	PTX Count vs Actual Number	42
4.3	Instruction Count: 6 blocks per SMX and 8 warps per block	46
4.4	Micro-benchmarks: measured vs estimated instruction count	46
4.5	Blowfish and MM: measured vs estimated instruction count	47
4.6	Memory Latency Test	51
4.7	Computation of average memory cycles	52
4.8	Micro-benchmarks: measured vs estimated time comparison	57
4.9	Blowfish algorithm: measured vs estimated time comparison	58
4.10	MM and IS application: measured vs estimated time comparison	58
5.1	Execution time of SM, IS, MM and PVC using <i>onlyCUDA</i> on GPUs	66
5.2	Execution time of SM, IS, MM and PVC using MR framework on GPUs	68
5.3	Execution time of PVR, SS, II and WC using MR framework on GPUs	69

5.4	Effect of Memory bandwidth	70
5.5	Effect of L1 cache re-configuration on performance of SM, IS, MM and PVC	71
5.6	Effect of L1 cache re-configuration on performance of PVR, SS, II and WC	72
5.7	Effect of miss rate at L1 and L2 on average memory access cycles	73
5.8	Access patterns in MM	76
5.9	MR framework enhanced by optimizing group phase	77
5.10	<i>mapperCount</i> function enhanced by cache-aware coding	82
5.11	<i>mapperCount</i> function enhanced by cache-aware coding	83
5.12	Execution time	84
5.13	Speedup of MARS MR framework	84
6.1	Implementation of SW on GPU	90
6.2	Execution time of SW implementation	92
6.3	Performance comparison of NB Simulation	96
6.4	Execution time of NB implementation	97
6.5	Flow chart of BF algorithm	98
6.6	Time for Kernel and I/O	100
6.7	Kernel Speedup comparison	101
6.8	Performance comparison of BF implementations with Disk I/O	101
6.9	Execution time of BF implementation	103
7.1	Heterogeneous Platform	109
7.2	Measured time and Estimated time on Heterogeneous GPUs	112
7.3	Processors Characteristic	116
7.4	Mapping of IS on four GPUs	118

7.5	n-processors mapping	122
7.6	Mapping of SM on five GPUs	122
7.7	Gain using optimized Mapping	123
7.8	Speedup obtained with optimized mapping	124

List of Tables

2.1	Literature on Comparison between GPU and FPGA	12
2.2	Summary of Performance Estimation	14
2.3	Summary of MR frameworks on GPUs	17
2.4	Summary of Partitioning Literature	23
3.1	GPU Architecture Features	29
3.2	MARS APIs of user-defined functions	34
3.3	MARS APIs of functions provided by runtime	35
4.1	Address trace comparison for IS	50
5.1	Phases of MARS needed by different applications	64
5.2	Impact of number of SMXs (<i>#blocks = 16384</i>)	67
5.3	Hit/Miss counts per SMXs for 16KB and 48KB	72
5.4	Hit/Miss counts with different access pattern	76
6.1	Set of Experiment Configuration	88
6.2	Performance Penalty Reduction	103
7.1	Important terms and notations	114
7.2	Speed Ratio: GTX690/GTX590	117

B.1	Blosum50	152
C.1	Partition Table for SM on Three Generation of GPUs	153
C.2	Partition Table for BF on Two Generation of GPUs	154
C.3	Partition Table for IS on Two Generation of GPUs	155
C.4	Partition Table for MM on Two Generation of GPUs	156
C.5	Partition Table for MM on Two Generation of GPUs	156

List of Abbreviations

AC	Accelerator Cluster
AMD	Advanced Micro Devices
BF	BlowFish Algorithm
BLAST	Basic Local Alignment Search Tool
BPM	Bank-level Partition Mechanism
CPU	Central Processing Unit
CSW	CALC1 subroutine from SWIM
CUDA	Compute Unified Device Architecture
DEGIMA	DEstination for Gpu Intensive MACHine
DES	Data Encryption Standard
DG	Distribute Grep
FDTD	Finite-Difference Time-Domain
FFT	Fast Fourier Transform
FPGAs	Field Programmable Gate Arrays
FPM	Functional Performance Model
GFLOPS	Giga Floating-point Operations per Second
GPGPU	General Purpose Computing on GPU
GPMR	GPU Map Reduce
GPU	Graphics Processing Unit
GSI	Gauss-Seidel Iterative Method
HG	Histogram
HPC	High Performance Computing
II	Inverted Index
IP	Internet Protocol
JIM	Jacobi Iterative Method
KFLOPS	Kilo Floating-point Operations per Second
KM	K-Means
KNN	K-nearest Neighbor Search
LRU	Last Recent Used

LT	Laplace Transformation
MFLOPS	Mega Floating-point Operations per Second
MFPM	multi-core Functional Performance Model
MM	Matrix Multiplication
MPI	Message Passing Interface
MR	Map-Reduce
MROH	Map-Reduce overheads
MRQ	memory request queue
MT	Matrix Transpose
MV	Matrix Vector
NB	N-Body Simulation
NBC	Naive Bayes Classifier
PCA	Principle Component Analysis
PCI	Peripheral Component Interconnect
PSS	Prefix Sum Scan
PTX	Parallel Thread Execution
PVC	Page View Count
PVR	Page View Rank
rCUDA	Remote CUDA
SHM	shared memory
SIMD	Single Instruction Multiple Data
SIMT	Single Instruction Multiple Thread
SIO	Sparse Integer Occurrence
SM	String Matching
SMV	Sparse Matrix-Vector
SMX	Streaming Multiprocessors
SORI	Successive Over Relaxation Iterative
SP	Saxpy
SRC	Speed Ratio Constant
SS	Similarity Scores
SW	Smith-Waterman Algorithm
TFLOPS	Tera Floating-point Operations per Second
URL	Universal Resource Locator
WC	Word Count