

ENERGY-EFFICIENT MANAGEMENT OF RECONFIGURABLE COMPUTERS

by

Rathiijit Sen

A dissertation submitted in partial fulfillment of
the requirements for the degree of

Doctor of Philosophy

(Computer Sciences)

at the

UNIVERSITY OF WISCONSIN-MADISON

2016

Date of final oral examination: 05/13/16

The dissertation is approved by the following members of the Final Oral Committee:

David A. Wood, Professor, Computer Sciences

Mark D. Hill, Professor, Computer Sciences

Mikko H. Lipasti, Professor, Computer Sciences

Gurindar S. Sohi, Professor, Computer Sciences

Michael M. Swift, Associate Professor, Computer Sciences

© Copyright by Rathijit Sen 2016

All Rights Reserved

To my parents Meena Sen and Ranjit Kumar Sen

ACKNOWLEDGMENTS

I am honored to have Prof. David Wood as my advisor. He has patiently guided me and supported me over so many years. David is an erudite scholar and outstanding professor. I have learned so much from him about computer architecture and analytical modeling. He encouraged me to consider the theoretical underpinnings in addition to the practical feasibility of problem solutions. He helped me explore ideas, refine drafts of writeups, improve my presentation skills, and network with other researchers through attending conferences. Thank you, David, for your time, support, and help.

I would like to thank the other members of my committee, Profs. Mark Hill, Mikko Lipasti, Guri Sohi, and Mike Swift, for giving me valuable feedback on my work. I have learned a lot from their vast experience and deep insight. Mark and Mike gave me great inputs for improving my papers and posters. Mike pointed me to many interesting papers. I have learned a lot from Mark about queueing models and how to present my work better. I had insightful discussions with Mikko about computer architecture research. Guri gave me power meters that I used to measure server power for this work.

I am also grateful to the Multifacet research group for providing me with necessary infrastructure to do the work in this dissertation. The server clusters, Lapis full-system simulator, Ruby memory hierarchy simulator, and Wisconsin commercial workloads have been invaluable to me. Technical discussions at group meetings were highly instructive.

Dan Gibson helped me a lot with setting up of the simulation infrastructure. Somayeh Sardashti and Hamid Reza Ghasemi also helped with creating simulation checkpoints.

This work was supported in part by the National Science Foundation (CNS-0916725, CCF-1017650, CNS-1117280, CCF-1218323, CNS-1302260, CCF-1438992, CCF-1533885), Microsoft Corporation (MSN140822), Sandia National Labs/DOE (MSN123960/DOE890426), and a University of Wisconsin Vilas award.

I would like to thank Prof. Somesh Jha and Prof. Jignesh Patel for their collaboration and help on multiple projects in the areas of static program analysis and databases. Jignesh gave me an opportunity to contribute to the Quickstep project and generously supported me during my final summer at UW-Madison. I look forward to continuing our collaboration. I would also like to thank Prof. Tom Reps for supporting me during my first year at UW-Madison and for giving me an opportunity to work on projects involving static analysis of executables.

My current and former office mates—Jayneel Gandhi, Derek Hower, Lena Olson, and Somayeh Sardashti—have given me company and have always been ready to help me. I have enjoyed countless hours chatting with them. I will miss Lena’s delicious home-baked muffins. Derek helped me understand many details of the Ruby simulator.

I have enjoyed working with Gagan Gupta on a long project involving static analysis and having many conversations with him about computer architecture, emerging technologies, reading groups, paper reviews and rebuttals. I also explored other projects with Arkaprava Basu, Asim Kadav, Nilay Vaish, Cong Wang, and Jianqiao Zhu. Discussions with Siddharth Barman helped me appreciate different perspectives on theory and systems research. Jason Power’s elevator pitch workshop was useful and fun.

I have also benefited from interacting with many other excellent colleagues including Muhammad Shoaib Bin Altaf, Newsha Ardalani, Piramanayagam Arumuga Nainar, Raghuraman Balasubramanian, Emily Blem, Jayaram Bobba, Evan Driscoll, Polina Dudnik, Yasuko Eckert, Christopher Feilbach, Venkatraman Govindaraju, Swapnil Haria, Bill Harris, Joel Hestness, Nick Kidd, Marc De Kruijf, Akash Lal, Junghee Lim, Jaikrishnan Menon, Sanketh Nalli, Tony Nowatzki, Marc Orr, Sankaralingam Pannerselvam, Matt Sinclair, Srinath Sridharan, Swaminathan Sundararaman, Aditya Thakur, Aditya Venkataraman, Haris Volos, and Hongil Yoon.

A number of external researchers, including UW-Madison computer architecture alumni, gave useful feedback on my work. I would like to acknowledge Luiz Barroso, Edouard Bugnion, Lisa Hsu, Mike Marty, Kathryn McKinley, and Ravi Rajwar for helpful suggestions and comments. The Architecture Affiliates meetings helped to get useful feedback also from Alaa Almeldeen, Brad Beckmann, John Davis, Dan Gibson, Peter Hsu, Konrad Lai, Kevin Moore, Steve Reinhardt, Greg Thorson, Greg Wright, and others.

I learned about new applications for analytical models while working with Trishul Chilimbi, Wei Huang, Srilatha Manne, and Indrani Paul during internships at Microsoft and AMD. I also enjoyed interacting with Manish Arora, Joseph Greathouse, and fellow interns Lavanya Subramanian and Zhe Wang. Weekend trips around and near Seattle with Arkaprava Basu, Tushar Krishna, and Shekhar Srikantaiah were enjoyable.

Angela Thorp, our graduate program coordinator, always helped me with understanding and following graduate school and departmental procedures. The Computer Systems Lab staff, and in particular Tim Czerwonka, promptly resolved all issues with the computing infrastructure that I ran into while doing my work.

Prof. Y. N. Srikant (IISc) and Prof. Reinhard Wilhelm (Saarland University) have greatly encouraged me. I wish to also acknowledge the friendship of some wonderful people, among them being Prof. Jan Reineke (Saarland University), Oindrilla Gupta, Kuntal Dey, Mohamed Abdel Maksoud, S. V. N. Narayana Rao, and Arpan Sen.

My sister Rakhee, brother-in-law Sandip, and niece Swagnita have brought me joy with their kindness, encouragement, and arrangements for many memorable excursions.

It has been a long and difficult journey, one on which I could not have persevered without the steadfast support of my parents. They have been a pillar of strength for me. I dedicate this dissertation to them. Thank you for your boundless love, encouragement, and good wishes throughout the years.

CONTENTS

Contents	v
List of Tables	ix
List of Figures	x
Abstract	xiv
1 Introduction	1
1.1 <i>Iron Law of Energy</i>	5
1.1.1 Load Management	6
1.1.2 Configuration Management	7
1.2 <i>Service-Level Agreements (SLA)-aware Governors</i>	8
1.3 <i>Contributions</i>	10
1.4 <i>Implications</i>	13
2 Energy Efficiency Ideals and the Iron Law	15
2.1 <i>Overview</i>	15
2.2 <i>Terminology and Infrastructure</i>	18
2.3 <i>Inadequacy of Conventional Energy Efficiency Ideals</i>	20
2.4 <i>Redefining EP and Dynamic EP</i>	24
2.5 <i>Power-Performance Pareto Frontier (Dynamic EO)</i>	27
2.6 <i>Computational PUE</i>	30
2.7 <i>Load and Configuration Management</i>	32
2.8 <i>The Π-dashboard</i>	35
2.9 <i>Conclusion</i>	38

3	Pareto Governors	39
3.1	<i>Overview</i>	39
3.2	<i>Infrastructure</i>	42
3.3	<i>Governors in Linux</i>	43
3.4	<i>Two-level governor design</i>	44
3.5	<i>Deployment Scenarios</i>	46
3.6	<i>SLAee: Maximize energy efficiency</i>	48
3.7	<i>SLAee: Adding L2 Prefetch Control</i>	56
3.8	<i>SLAee: Adding Control for Wall Power</i>	59
3.9	<i>SLApower: Maximize performance within a power cap/budget</i>	62
3.10	<i>SLAperf: Maximize power savings given a performance target</i>	65
3.10.1	Governing for absolute performance targets	66
3.10.2	Governing for relative performance targets	70
3.10.3	Governing to minimize idle time	74
3.11	<i>Limitations</i>	77
3.11.1	Socket-Wide Control	78
3.11.2	Intrusive Profiling	79
3.11.3	Sampling Inconsistency and Non-representativeness	79
3.11.4	Non-Zero Reaction Times	80
3.12	<i>Conclusion</i>	80
4	Cache Reuse Models	82
4.1	<i>Overview</i>	82
4.2	<i>Infrastructure</i>	85
4.3	<i>Measures of Temporal Locality</i>	89
4.3.1	Reuse Distance Distributions	90

4.3.2	$T \rightarrow r(T)$ is a Lossy Transformation	90
4.3.3	$d(T)$ Estimation	91
4.4	<i>Per-set Locality</i>	93
4.4.1	$r(S')$ Estimation	94
4.4.2	Matrix dimension and Truncation of r	97
4.4.3	Poisson approximation to Binomial	98
4.5	<i>Cache Hit Functions</i>	99
4.5.1	Estimating $\phi(\text{LRU})$	102
4.5.2	Estimating $\phi(\text{RANDOM})$	103
4.5.3	Estimating $\phi(\text{NMRU})$	106
4.5.4	Estimating $\phi(\text{PLRU})$	107
4.5.5	Estimation Accuracy and Computation Time	113
4.6	<i>Hardware Support</i>	114
4.6.1	New hardware support to estimate reuse distributions	115
4.6.2	Set-Counters, Way-Counters, and Shadow Tags	130
4.7	<i>Index Hashing</i>	133
4.8	<i>Limitations</i>	136
4.9	<i>Conclusion</i>	137
5	Cache Power Budgeting	139
5.1	<i>Overview</i>	139
5.2	<i>Infrastructure</i>	142
5.3	<i>Cache Resizing Opportunities</i>	145
5.4	<i>Operations overview</i>	150
5.5	<i>Cache miss rate prediction</i>	153
5.6	<i>Performance and Power prediction models</i>	158

5.7	<i>Model-driven Power Budgeting</i>	164
5.7.1	Basic model	165
5.7.2	On-chip power-budgeting model	166
5.7.3	System power-budgeting model	168
5.7.4	Results and Limitations	169
5.8	<i>Conclusion</i>	174
6	Related Work	176
6.1	<i>Overview</i>	176
6.2	<i>Energy Efficiency Characterization</i>	177
6.3	<i>Power-Performance States</i>	179
6.4	<i>Optimization Goals</i>	180
6.5	<i>Cache Models</i>	181
6.6	<i>Reconfiguration Knobs</i>	185
6.6.1	Classification	197
7	Conclusion	204
A	SPECpower power-performance	209
	Bibliography	210

LIST OF TABLES

2.1	R^2 values for polynomial fits to SPECpower Pareto frontier.	37
4.1	Number of sets and associativity for different cache sizes.	85
4.2	System configuration.	86
4.3	Workload characteristics.	86
4.4	Relative miss ratios for difference cache sizes and replacement policies. . . .	88
4.5	Average absolute values of prediction errors over all cache configurations. . .	114
4.6	Average of absolute errors with different filters.	122
4.7	Average of relative errors with different filters.	122
4.8	Number of samples selected with different sampling configurations.	123
4.9	Number of entries in the Histogram array for different sampling configurations.	124
4.10	Plain vs hashed indexing.	136
5.1	System configuration.	143
5.2	Workloads.	144
6.1	Classification, by semantic types, of system reconfiguration capabilities. . . .	201
6.2	Classification (cont.), by semantic types, of system reconfiguration capabilities.	202
6.3	Classification (cont.), by semantic types, of system reconfiguration capabilities.	203

LIST OF FIGURES

1.1 Trends in Processor power-performance profiles.	4
2.1 Power-Performance profile with conventional server configuration.	21
2.2 Conventional efficiency model of servers.	21
2.3 Power-Performance profile for super-proportional systems.	22
2.4 Performance (Load) vs Efficiency for super-proportional systems.	22
2.5 EOP and Dynamic EO models.	25
2.6 CPUE(c, l) and LUE(l).	33
2.7 Resource Usage Effectiveness.	34
2.8 Coordination architecture.	36
3.1 State transitions to Dynamic EO for meeting SLAs.	39
3.2 Example power-performance profile with Wall Power.	46
3.3 Example power-performance profile with Socket + Mem Power.	47
3.4 Power-Performance traces for applu.	49
3.5 Power-Performance traces for graph500.	50
3.6 BIPS-per-Watt on HS with different policies.	53
3.7 R(10) freq. distribution for applu (0.8–3.5 GHz).	55
3.8 Average energy efficiency of applu as a function of the number of instructions executed and processor frequency.	55
3.9 Prefetch Impact.	57
3.10 L2 Prefetch mode distribution by RF(10).	58
3.11 BIPS-per-watt of governors with (RF(10)) and without (P, PF, R(10)) dynamic control for L2 Prefetching.	58

3.12 RAPL and Wall Power correlation.	60
3.13 BIPS-per-watt of governors P and RF(10) with wall (full-system) power.	61
3.14 Power-performance profiles for graph500 and md for SLApower.	64
3.15 Power-performance profiles for md and graph500 for SLAperf.	67
3.16 Execution profiles for graph500 with SLA = 3.5 BIPS and 5.5 BIPS.	68
3.17 Power-performance profiles for SPECpower with Linux governors.	70
3.18 Power-performance profiles for SPECpower with RF_SLAperf(10) and R_SLAperf(10)	73
3.19 Distributions of settings for SPECpower with RF_SLAperf(10)	73
3.20 Power-performance profiles for SPECpower with RF_Active	76
3.21 Distributions of settings for SPECpower with RF_Active(10,2)	76
3.22 Distributions of settings for SPECpower with RF_Active(100,20)	77
3.23 Distributions of settings for SPECpower with RF_Active(500,20)	77
4.1 “Instantaneous” and cumulative miss ratios.	87
4.2 Absolute reuse distance visualization.	91
4.3 Model vs Estimated $\mathbf{d}(\mathbf{T})$	92
4.4 Effect of the number of sets (S) on per-set locality for <code>oltp</code>	93
4.5 Reuse distribution transformations with stochastic Binomial Matrices.	95
4.6 LRU prediction with limited reuse information.	97
4.7 Equation 4.3 pseudo-code with Poisson approximation.	99
4.8 Representative hit ratio functions (Φ_k).	101
4.9 Actual vs estimated miss ratios with RANDOM replacement policy.	106
4.10 Actual vs estimated miss ratios with NMRU replacement policy.	108
4.11 PLRU subtrees for $A' = 4$	110
4.12 PLRU subtree decomposition.	112

4.13 Actual vs estimated miss ratios with PLRU replacement policy.	113
4.14 Schematic of new hardware support.	115
4.15 Probability of false hit in a 1024-bit Bloom filter with 2 hash functions.	117
4.16 Online estimation of miss ratios using E filters.	119
4.17 Online estimation errors with E filters.	119
4.18 Online estimation of miss ratios using B filters.	120
4.19 Online estimation errors with B filters.	120
4.20 Online estimation of miss ratios using CB filters.	121
4.21 Online estimation errors with CB filters.	121
4.22 Online estimation of miss ratios using E filters and with a set sample randomly chosen at the start of every sample.	125
4.23 Online estimation errors using E filters and CLT criteria.	129
4.24 Online estimation errors using B filters and CLT criteria.	129
4.25 Online estimation errors using CB filters and CLT criteria.	130
4.26 PLRU trees demonstrating non-inclusion.	132
4.27 Miss ratio reduction with hashed indexing.	135
5.1 Power-performance for blackscholes with DVFS and cache resizing.	140
5.2 MPKI vs cache size.	146
5.3 Cache power budgeting opportunities and pitfalls.	148
5.4 Operations Overview.	151
5.5 Training and Prediction intervals.	152
5.6 Workload reuse distributions.	154
5.7 LLC access dependence on associativity for small caches.	156
5.8 LLC MPKI dependence on associativity for small caches.	156
5.9 Averages of Absolute Error for miss ratio estimation.	157

5.10 Model Error for miss ratio estimation.	159
5.11 Phase Error for miss ratio estimation.	160
5.12 CPI regression for commercial workloads.	161
5.13 Combined CPI regression for commercial workloads.	162
5.14 Model Error for CPI estimation.	163
5.15 Phase Error for CPI estimation.	163
5.16 Model Error for system power estimation.	164
5.17 Phase Error for system power estimation.	164
5.18 First-order power-budgeting model.	165
5.19 System power budgeting results.	170
5.20 Comparison of performance gains between the oracle and our model.	171
5.21 Comparison of power savings between the oracle and our model.	171
5.22 Comparison matrix between the oracle and our model.	172
A.1 SPECpower power-performance with different configurations.	209

ABSTRACT

Power and energy consumption are first-order constraints on the design and operation of computer systems today. Improving energy efficiency reduces the amount of energy needed to perform a given computation as well as enables more computation to be performed for the same amount of energy. This saves operational costs to use these systems as well as capital costs to provision for them.

Conventionally, energy proportionality (energy consumption in proportion to the work done, or equivalently, power consumption in proportion to utilization/performance/load served) as proposed by Barroso and Hölzle, has been the gold standard of an ideal system's energy efficiency. While this model is valid for fixed-resource systems, modern systems are reconfigurable in many aspects, allowing them to adapt to changing workload characteristics. Smart reconfigurability increases energy efficiency. However, we show that reconfigurability invalidates the conventional notions of ideal energy proportionality if the system starts to behave super-proportionally. Super-proportional systems provide more performance (or work) in proportion to the power (or energy) used. We propose a new ideal model, Energy Optimal Proportional (EOP), that subsumes the conventional model and improves upon it by also accounting for super-proportional systems.

EOP can guide system designers to improve the maximum efficiency attainable over the operating range and forms a basis for comparisons of energy efficiency across systems. Power-performance Pareto optimality, on the other hand, can guide system operators to manage load and configure resources appropriately to make the current system execute efficiently. We propose a new intellectual framework that interrelates these two complementary energy efficiency goals.

The rest of this dissertation focuses on energy-efficient management. We develop new reactive governors that coordinate processor frequency (and voltage) and hardware

prefetching to improve energy efficiency on a real (Haswell) server. We also propose a space-efficient hardware mechanism to estimate temporal locality (reuse) in cache accesses. The estimated distributions can be used by our new analytical models for cache performance to drive resizing decisions of the last-level cache.

Finally, we propose a new classification system for system reconfiguration capabilities. The classification is based on the semantics of what the reconfiguration affects—computation, communication, storage, scheduling, speculation. We hope that this classification will be insightful to future researchers while exploring the space of reconfigurable systems, in categorizing existing work and in identifying coordination options that have been less well explored.