

**RESOURCE CONTENTION AWARE  
PERFORMANCE AND POWER  
OPTIMIZATION IN CHIP MULTIPROCESSORS**

**SOLOMON ABERA**



DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING  
INDIAN INSTITUTE OF TECHNOLOGY DELHI

MAY 2021

©Indian Institute of Technology Delhi - 2021  
All rights reserved.

# **RESOURCE CONTENTION AWARE PERFORMANCE AND POWER OPTIMIZATION IN CHIP MULTIPROCESSORS**

by

**SOLOMON ABERA**

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**

**MAY 2021**

# Dedication

I dedicate this thesis to my late Mother, Etalemahu Kassa, who could not see this journey completed.

# Certificate

This is to certify that the thesis titled **RESOURCE CONTENTION AWARE PERFORMANCE AND POWER OPTIMIZATION IN CHIP MULTIPROCESSORS** being submitted by **Mr. SOLOMON ABERA** 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 **Department of Computer Science and Engineering, 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

Department of Computer Science and Engineering

Indian Institute of Technology Delhi

New Delhi- 110016

Anshul Kumar

Emeritus Professor

Department of Computer Science and Engineering

Indian Institute of Technology Delhi

New Delhi- 110016

# Acknowledgements

First and foremost I thank God, for all and everything.

I thank my thesis advisors Prof. M Balakrishnan and Prof Anshul Kumar, for all their help, patience and guidance. They made me not only a good scholar but also a better person. Thank you.

I thank all the faculty in the Architecture group – Prof. Panda, Prof. Kolin Paul and Dr. Sarangi for reviewing my work and providing valuable insights.

I thank Rajshaker K and my friends for all the technical discussions that aided my research and deepened my understanding, as well as for all the other experiences that made life enjoyable.

I thank my wife Eden Nisku and my children: Yohannan, Kaleab and Soleyana. Without their endless love and sacrifice, I would never have been able to complete my study.

Last but not the least, I would like to thank my sisters: Etaferahu and Aynadis for supporting me throughout this journey.

**Solomon Abera**

# Abstract

Over the last couple of decades, Chip Multi-Processors (CMPs) – also called multicores, have been the leading architectural choice for computing systems ranging from battery-operated devices to high-end servers. Even though CMPs enhance performance through concurrent execution of application programs, the contention for shared resources makes the performance of individual application programs and associated energy consumption unpredictable. Applications sharing the memory hierarchy in CMP can slow down by more than 80% [1], completely degrading the performance improvement achieved by multicores.

The level of performance degradation encountered by a single application due to CMP shared resource contention depends on its own and co-runners' (applications running on other cores) memory behavior. By taking these factors into account, modeling the inter-application resource contention is essential as it gives valuable information to achieve performance and energy optimization.

In this thesis, we first propose lightweight metrics that accurately capture the potential contention among concurrently running applications on CMP. We also show and capture the time varying behavior of applications and extract them using phase-wise profiling.

With the help of these metrics, we propose and validate a methodology for optimizing performance that can be used in making contention aware scheduling decisions. In addition to this, we also propose a learning model that predicts the performance of applications with partially shared cache.

On the other hand, energy consumption is a critical parameter as it constitutes the most significant operating cost for computing clouds. Analogous to this, limited use time before the need to recharge batteries, continues to be an essential user concern in mobile devices. To optimize on power consumption, modern CMP processors are designed with Dynamic Voltage and Frequency Scaling (DVFS) support at the individual core as well as for the uncore part. This control on DVFS can provide for fine-grained control of performance and energy as

shown in this work. When applications slow down due to resource contention, typically their performance sensitivity to DVFS also reduces. Hence, the performance-energy trade-off curve of each application varies with its co-runners. In this thesis, we also model the performance-energy trade-off for applications running on a CMP platform with provision for both core and uncore DVFS. We use a learning algorithm to build the model. The proposed model is not statically tuned to support a specific DVFS policy. It builds a relationship between the DVFS steps based on the behavior of the chosen application and the co-runners represented through two simple metrics and performance. It is flexible enough to be part of any DVFS controller. We demonstrate the efficacy of our proposed model by considering various QoS policies where the user specifies the maximum permissible performance loss. The model predicts the lowest possible voltage/frequency step such that the QoS requirement is met and the energy saving is maximized.

## सार

पिछले कुछ दशकों में, चिप मल्टी-प्रोसेसर (सीएमपी) - जिसे मल्टीकोर भी कहा जाता है, बैटरी से चलने वाले उपकरणों से लेकर परिष्कृत सर्वर तक कंप्यूटिंग सिस्टम के लिए प्रमुख वास्तुशिल्प विकल्प रहा है। भले ही सीएमपी अनुप्रयोग कार्यक्रमों के समानांतर निष्पादन के माध्यम से कंप्यूटिंग गति को बढ़ाते हैं, साझा संसाधनों के लिए विवाद प्रत्येक एप्लिकेशन प्रोग्राम और संबंधित ऊर्जा खपत के प्रदर्शन को अप्रत्याशित बनाता है। सीएमपी में मेमोरी पदानुक्रम साझा करने वाले एप्लिकेशन 80% से अधिक [1] तक धीमा हो सकते हैं, मल्टीकोर द्वारा प्राप्त प्रदर्शन सुधार को पूरी तरह से खराब कर सकते हैं। सीएमपी साझा संसाधन विवाद के कारण एकल एप्लिकेशन द्वारा सामना किए गए प्रदर्शन में गिरावट का स्तर अपने और अन्य अनुप्रयोगों (अन्य कोर पर चलने वाले एप्लिकेशन) मेमोरी व्यवहार पर निर्भर करता है। इन कारकों को ध्यान में रखते हुए, अंतर-अनुप्रयोग संसाधन विवाद को मॉडलिंग करना आवश्यक है क्योंकि यह गति और ऊर्जा अनुकूलन प्राप्त करने के लिए बहुमूल्य जानकारी देता है।

इस थीसिस में, हम पहले हल्के मेट्रिक्स का प्रस्ताव करते हैं जो सीएमपी पर समवर्ती रूप से चल रहे अनुप्रयोगों के बीच संभावित विवाद को सटीक रूप से कैचर करते हैं। हम अनुप्रयोगों के समय-भिन्न व्यवहार को दिखाते और कैचर करते हैं और चरण-वार प्रोफाइलिंग का उपयोग करके उन्हें निकालते हैं। इन मेट्रिक्स की सहायता से, हम प्रदर्शन को अनुकूलित करने के लिए एक कार्यप्रणाली का प्रस्ताव और सत्यापन करते हैं जिसका उपयोग विवाद-जागरूक शेड्यूलिंग निर्णय लेने में किया जा सकता है। इसके अलावा, हम एक मशीन-लर्निंग मॉडल भी प्रस्तावित करते हैं जो आंशिक रूप से साझा किए गए कैश के साथ अनुप्रयोगों के प्रदर्शन का अनुमान लगाता है।

दूसरी ओर, ऊर्जा की खपत एक महत्वपूर्ण पैरामीटर है क्योंकि यह क्लाउड कंप्यूटिंग के लिए सबसे महत्वपूर्ण परिचालन लागत का गठन करता है। इसके अनुरूप, बैटरियों को रिचार्ज करने की आवश्यकता से पहले सीमित उपयोग समय, मोबाइल उपकरणों में एक आवश्यक उपयोगकर्ता चिंता का विषय बना हुआ है। बिजली की खपत को अनुकूलित करने के लिए, आधुनिक सीएमपी प्रोसेसर को प्रत्येक कोर के साथ-साथ अनकोर भाग के लिए डायनेमिक वोल्टेज और फ्रीक्वेंसी स्केलिंग (डीवीएफएस) समर्थन के साथ डिज़ाइन किया गया है। इस तरह के डीवीएफएस गति और ऊर्जा का बारीक नियंत्रण प्रदान कर सकते हैं, जैसा कि इस कार्य में दिखाया गया है। जब संसाधन विवाद के कारण अनुप्रयोग धीमा हो जाता है, तो यह आमतौर पर डीवीएफएस के प्रति उनकी गति संवेदनशीलता को कम कर देता है। इसलिए, प्रत्येक

एप्लिकेशन की गति और ऊर्जा ट्रेडऑफ़ वक्कर सह-निष्पादन अनुप्रयोगों के साथ भिन्न होता है। इस थीसिस में, हम कोर और अनकोर डीवीएफएस दोनों के प्रावधान के साथ सीएमपी प्लेटफॉर्म पर चलने वाले अनुप्रयोगों के लिए गति और ऊर्जा ट्रेडऑफ़ को भी मॉडल करते हैं। हम मॉडल बनाने के लिए मशीन लर्निंग एल्गोरिदम का उपयोग करते हैं। प्रस्तावित मॉडल किसी विशिष्ट डीवीएफएस नीति का समर्थन करने के लिए स्थिर रूप से ट्यून नहीं किया गया है। यह चुने हुए एप्लिकेशन के व्यवहार और दो सरल मेट्रिक्स और प्रदर्शन के माध्यम से प्रदर्शित सह-निष्पादन अनुप्रयोगों के आधार पर डीवीएफएस चरणों के बीच संबंध बनाता है। इसे किसी भी डीवीएफएस नियंत्रक का हिस्सा बनने के लिए संशोधित किया जा सकता है। हम विभिन्न क्यूओएस नीतियों पर विचार करके अपने प्रस्तावित मॉडल की प्रभावकारिता प्रदर्शित करते हैं जहां उपयोगकर्ता अधिकतम अनुमेय गति हानि को निर्दिष्ट करता है। मॉडल न्यूनतम संभव वोल्टेज/आवृत्ति कदम की भविष्यवाणी करता है जैसे कि क्यूओएस आवश्यकता पूरी हो और ऊर्जा की बचत अधिकतम हो।

Solomon abera bekele: सोलोमोन आबेरा बेकेले

# Contents

<b>Certificate</b>	<b>i</b>
<b>Acknowledgements</b>	<b>iii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Contention in CMPs and Performance-Energy Trade-off . . . . .	2
1.1.1 Source of Contention In CMPs . . . . .	2
1.1.2 Dynamic Voltage and Frequency Scaling . . . . .	3
1.1.3 Uncore Frequency Scaling . . . . .	4
1.2 Motivation and Prior work . . . . .	5
1.3 Contributions . . . . .	9
1.4 Thesis Organization . . . . .	10
<b>2 Capturing Resource Contention in CMPs</b>	<b>11</b>
2.1 Introduction . . . . .	11
2.1.1 Background and Related Work . . . . .	11
2.2 Aggressiveness . . . . .	15
2.2.1 Impact of number of co-runners on GA . . . . .	18
2.3 Sensitivity . . . . .	19

2.3.1	Bandwidth sensitivity . . . . .	23
2.4	Summary . . . . .	26
<b>3</b>	<b>PLSS: Phase-guided Locality Signature based Scheduler</b>	<b>27</b>
3.1	Introduction . . . . .	27
3.1.1	Background and Related Work . . . . .	27
3.2	PLSS Overview . . . . .	32
3.2.1	Phase Detection Mechanism . . . . .	34
3.2.2	PLSS: Algorithm . . . . .	34
3.3	Evaluation . . . . .	36
3.3.1	Evaluation Setup . . . . .	36
3.3.2	Phase Detection Validation . . . . .	36
3.3.3	Scheduling Experiments . . . . .	38
3.3.4	SDP Based S_Score [2] . . . . .	39
3.4	Summary . . . . .	42
<b>4</b>	<b>Performance-Energy Trade-off In CMPs Using Per-Core DVFS</b>	<b>43</b>
4.1	Introduction . . . . .	43
4.2	Background and Related Work . . . . .	44
4.3	Impact of Contention and CFS on Performance . . . . .	46
4.4	Overview . . . . .	49
4.4.1	Data Collection . . . . .	50
4.4.2	Model Construction . . . . .	51
4.4.3	Application replacement . . . . .	52

---

4.4.4	Application of the model . . . . .	54
4.5	Evaluation . . . . .	56
4.5.1	Experimental Setup . . . . .	56
4.5.2	Frequency Scaling in Linux . . . . .	57
4.5.3	Machine Learning Algorithms . . . . .	60
4.5.4	CFS Validation . . . . .	62
4.5.5	Simulation Based Evaluation . . . . .	68
4.6	Summary . . . . .	72
<b>5</b>	<b>Performance-Energy Trade-off Estimation for Uncore Frequency Scaling</b>	<b>73</b>
5.1	Introduction . . . . .	73
5.1.1	Background and Related Work . . . . .	74
5.2	Motivation . . . . .	76
5.3	Impact of Resource Contention and UFS on Performance . . . . .	79
5.4	Overview . . . . .	82
5.4.1	UFS Predictor Model . . . . .	83
5.5	Evaluation . . . . .	84
5.5.1	Experimental Setup . . . . .	84
5.5.2	UFS Validation . . . . .	86
5.5.3	Results Analysis . . . . .	87
5.5.4	Summary . . . . .	88
<b>6</b>	<b>Performance Prediction for Applications with Partially Shared Cache</b>	<b>91</b>
6.1	Introduction . . . . .	91

---

6.1.1	Cache Partitioning . . . . .	92
6.2	Motivation . . . . .	94
6.3	Data Collection and Model Construction . . . . .	96
6.4	Validation . . . . .	97
6.5	Summary . . . . .	98
<b>7</b>	<b>Conclusion</b>	<b>99</b>
7.1	Summary . . . . .	99
7.2	Contributions . . . . .	102
7.3	Future work . . . . .	103
	<b>Appendices</b>	<b>105</b>
<b>A</b>	<b>Background on Caches</b>	<b>107</b>
A.1	Caches . . . . .	108
A.1.1	Associativity . . . . .	108
A.1.2	Replacement Policy . . . . .	109
	<b>Bibliography</b>	<b>111</b>
	<b>List of Publications</b>	<b>125</b>
	<b>Biography</b>	<b>127</b>

# List of Figures

1.1	Simplified versions of typical a) single-core and b) multicore architectures . . .	2
1.2	Performance slowdown incurred by <code>rate</code> benchmarks from the SPEC2017 suite for their run with a particular set of seven co-runner applications on an Octa-core processor sharing 20MB L3 cache . . . . .	6
1.3	Normalized IPC of <code>xalancbmk_r</code> benchmark for its run with 18 different co-runner groups (CGs – each composed of seven benchmarks) on an Octa-core processor sharing 20MB LLC. . . . .	7
1.4	Performance sensitivity of <code>xalancbmk_r</code> benchmark to DVFS when it runs with different co-runners . . . . .	8
1.5	Thesis overview . . . . .	10
2.1	Stack distance histogram for a 8-way cache . . . . .	12
2.2	Aggressiveness score experiment . . . . .	15
2.3	Impact of contention on performance . . . . .	18
2.4	<code>Xalancbmk_r</code> 's run with co-runners constituted from different number of applications . . . . .	19
2.5	Sensitivity score experiment . . . . .	20
2.6	Impact of cache space on performance . . . . .	21
2.7	Slowdown with <code>S_Score</code> . . . . .	23
2.8	Solo-run where 8-way 8MB cache and bandwidth entirely used by $A_T$ . . . . .	24

2.9	Co-run with private 8-way 8MB LLC to $A_T$ and co-runners with shared memory bandwidth . . . . .	24
2.10	Slowdown_bw vs A_Score . . . . .	25
3.1	Quad-Core symmetric multicore with no partial sharing . . . . .	28
3.2	Quad-Core symmetric multicore with partially shared caches . . . . .	29
3.3	Overview of PLSS . . . . .	33
3.4	1bm, One Billion instruction simulation; top-to-bottom: Cache hit-rate for various cache sizes(1K to 128K), IPC, L1-D hit frequency, Number of L1 accesses(10K) and instruction distance (phase detector) . . . . .	37
3.5	mcF : one Billion instruction simulation; top-to-bottom: Cache hit-rate for various cache sizes(1K to 128K), IPC, L1-D hit frequency, Number of L1 accesses(10K) and instruction distance (phase detector) . . . . .	38
3.6	Performance degradation encountered by 12 SPEC2006 benchmarks when they are running on dual-core . . . . .	39
3.7	Slowdown encountered by 10 SPEC2006 benchmarks running on dual-core . . . . .	42
4.1	Performance response of 538.imagick_r benchmark to core frequency scaling . . . . .	47
4.2	Performance response of 523.xalancbmk_r benchmark to core frequency scaling . . . . .	48
4.3	Performance response of 549.fotonik3d_r benchmark to core frequency scaling . . . . .	49
4.4	Training and testing of the model . . . . .	51
4.5	Application replacement . . . . .	52
4.6	Modeling process of application replacement . . . . .	53
4.7	CPUFreq framework architecture [3] . . . . .	58
4.8	Impact of <code>n_estimators</code> (number of trees in the forest) on the prediction score and training time . . . . .	61

---

4.9	Impact of <code>max_features</code> on the prediction score and training time . . . . .	62
4.10	Cross-validation . . . . .	65
4.11	Validation using test set . . . . .	66
4.12	Performance-energy trade-off for different QoS Policies: Perfect Vs Random Forest Model . . . . .	69
4.13	Energy and performance loss/gain because of prediction inaccuracies . . . . .	71
5.1	Simplistic view of core and uncore in the Intel Haswell architecture . . . . .	75
5.2	Marginal processor energy savings for core DVFS . . . . .	77
5.3	Effect of UFS on performance and processor energy . . . . .	78
5.4	Performance response of 538.imagick_r benchmark to resource contention and uncore frequency scaling . . . . .	79
5.5	Performance response of 523.xalancbmk_r benchmark to resource contention and uncore frequency scaling . . . . .	80
5.6	Performance response of 549.fotonik3d_r benchmark to resource contention and uncore frequency scaling . . . . .	81
5.7	Training and testing of the model . . . . .	84
5.8	Data collection process for UFS validation . . . . .	86
5.9	Average and maximum performance degradation and energy savings observed for each of the policies . . . . .	88
6.1	Different CAT partitioning strategies . . . . .	95
6.2	Performance response of <code>xalancbmk_s</code> and <code>lbm_s</code> applications to various LLC partitioning schemes . . . . .	96
6.3	Model construction . . . . .	97
6.4	10-fold cross-validation result . . . . .	97

