

**SCHEDULING AND CHARACTERIZATION OF
COMPUTER VISIONWORKLOADS ON
HETEROGENEOUS SYSTEMS**

DIKSHA MOOLCHANDANI



**AMAR NATH AND SHASHI KHOSLA SCHOOL OF
INFORMATION TECHNOLOGY
INDIAN INSTITUTE OF TECHNOLOGY DELHI**

MARCH 2022

© Indian Institute of Technology Delhi (IITD), New Delhi, 2022

**SCHEDULING AND CHARACTERIZATION OF
COMPUTER VISION WORKLOADS ON
HETEROGENEOUS SYSTEMS**

by

DIKSHA MOOLCHANDANI

**Amar Nath and Shashi Khosla School of Information
Technology**

Submitted

**in fulfilment of the requirements of the degree of Doctor of Philosophy
to the**



**INDIAN INSTITUTE OF TECHNOLOGY DELHI
MARCH 2022**

Certificate

This is to certify that the thesis titled **SCHEDULING AND CHARACTERIZATION OF COMPUTER VISION WORKLOADS ON HETEROGENEOUS SYSTEMS** being submitted by **Ms. DIKSHA MOOLCHANDANI** for the award of **Doctor of Philosophy** in Information Technology is a record of bonafide work carried out by her under my guidance and supervision at the Amar Nath and Shashi Khosla School of Information Technology, 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.

Smruti Ranjan Sarangi
Professor
Dept. of Computer Science & Engg.
Indian Institute of Technology Delhi.

Anshul Kumar
Emeritus Professor
Dept. of Computer Science & Engg.
Indian Institute of Technology Delhi.

Acknowledgements

When I started counting my blessings, my whole life turned around.

- Willie Nelson

No duty is more urgent than giving thanks.

- James Allen

When eating fruit, remember the one who planted the tree.

- Vietnamese Proverb

My Ph.D. experience was a roller coaster ride with lots of ups and downs. Today when I look back, it was every bit worth it. It was a mind-blowing transformation of my personality, thought process, and nature. I believe there are more people instrumental in my growth than I remember. Nevertheless, I will make the best effort to acknowledge most of them.

Firstly, I am grateful to God. I had always believed that not even a leaf moves without His will, and whatever happens is always for our good (we always realize this a few years later). I experienced it during my Ph.D. when after a lot of difficult phases, I got **Prof. Smruti Sarangi** as my advisor. I thank him for taking me on from another lab. This thesis is a result of his firm belief in me. He chose not to give up on me and worked relentlessly to bring out the best in me. His invaluable pieces of advice not only shaped me as a researcher but also as a person. I thank him for guiding me when I was utterly lost and letting me explore when I was ready. He did his best to make me an independent researcher and made me identify my strengths. I thank him for making me understand that patience is the key to overcoming all the challenges. Now, when I look through my journey, I realize that His plan was perfect, and I could not have asked for a better advisor than Prof. Sarangi.

I also thank my co-advisor, **Prof. Anshul Kumar**. I took up a Ph.D. because he made me believe that I was capable of innovating. I thank him for believing in me even when I had not started research. I also thank **Prof. Kolin Paul** who was the first person in IIT to have believed in me and admitted me as a project associate in his lab, which enabled the cascading set of events that made me pursue a Ph.D. I thank my student research committee members: **Prof.**

Preeti R. Panda, Prof. M. Balakrishnan, and Prof. Shouribrata Chatterjee for providing valuable feedback on my research. I would also like to thank **Prof. José F. Martínez** for collaborating with me on my first paper and providing me the pointers to think critically. I thank **Prof. Manuel Mazzara** and **Prof. Ilya Afanasyev** for mentoring and hosting me during my internship at Innopolis University. I express my deepest gratitude to **Prof. Naveen Garg** and **Prof. Chetan Arora** for providing valuable feedback on my papers.

I express my heartfelt gratitude to Ms. Aditi, Ms. Vandana, Ms. Manju, Mr. Suresh, and Mr. Rajesh for their speedy assistance in all administrative matters. I started my research with a senior student, Rajshekar Kalayappan. I am thankful to him for teaching me a lot about computer architecture and architectural simulators. He had been instrumental in giving me a kick-start in research. I would also like to thank my best friend, Ashish Kumar. Despite his busy schedule, he always lent a listening ear to my problems and provided sound advice both professionally and personally. I thank him for his jokes and for always being there for me during my struggling graduate days. One of my very first friends in IIT was Sapna Bhardwaj. We first met by chance and then fortunately formed a team for the distributed computing course. I learned a lot from her during the course projects. Her sharp skills in learning any new computing language and debugging codes helped me in honing my skills. I am indebted to Sapna for being a great friend and teammate.

I found many more great friends in Priyanka Singla, Hameedah Sultan, Chrystle Lobo, Sandeep Kumar, Omais Shafi, Ismi Abidi, Shubhankar Singh, Janib Bashir, Nivedita Shrivastava, Shailja Pandey, Garima Modi, Chahat Bansal, and Ankita Raj. I am thankful to Priyanka for the thoughtful discussions that we had during our coffee breaks and for attending interesting courses with me. I am grateful to Hameedah for being ever-ready to talk to me during the pandemic and providing valuable pieces of advice at different points in my Ph.D. I thank Sandeep for motivating me throughout, from our first collaboration till I got a job, and making me believe in myself. I am grateful to Nivedita for being a great friend, teammate and staying awake with me countless nights so that we could finish the project even when her part was complete. I thank Janib and Shubhankar for having insightful discussions with me and guiding me about the plans after Ph.D. I am grateful to Shailja, Omais, Ismi, and Garima for keeping a regular tab on me during the pandemic. I thank Garima for taking care of me and bringing me homemade lunch so that I don't have to go to mess in the scorching sun. I was blessed to have Chahat, Ankita, and Chrystle in my life for cracking me up every time we met.

Apart from my lab friends, I found a great friend in my roommate, Dilpreet Kaur. We met in a computer architecture course in the first semester and decided to stay together from thereon. She is the best roommate I could have asked for, always lively and cheerful, kind-hearted, and

always ready to help me out in any situation. I would also like to thank all my interns for believing in me and providing me with learning opportunities. I especially thank Sudhanshu and Kishore for helping me with the parts of the papers. I also thank my friends from Innopolis University: Geesara Prathap, Swati Megha, Vivek Kumar, Jasna Jankovic, Maria Naumcheva, Inna Baskakova, and Andrey Sadovykh. They made my internship a memorable experience.

I would also like to thank my cousins in Delhi: Pavan, Preeti, Krunal, and Vraj. Their liveliness re-energized me. I thank DST for funding my research and IRD for releasing my stipend on time. This list is incomplete without mentioning the security guards of the Himadri hostel, Kailash hostel, and SIT building. All thanks to their efforts that I was able to walk through the campus without fear of dogs. I express my heartfelt thanks to the warden, caretaker, and mess staff for catering to my requests. I also thank the cleaning staff of SIT building and Himadri hostel for maintaining proper hygiene (Cleanliness is next to Godliness).

No words will ever be enough to thank my parents: Dr. Kishore Moolchandani and Mrs. Sonia Moolchandani, for their immense support and understanding throughout my life, especially through my graduate years. I can never compensate for their endless sacrifices in this lifetime. I thank them for instilling in me the value of education, time, discipline, perseverance, and hard work. I also thank my sister, Gunja Moolchandani. All through these years, she had been my box of happiness, lucky charm, my mood booster, my stressbuster, and my biggest inspiration. I cannot thank her enough for the sacrifices she made during the pandemic to make my life easier.

Diksha Moolchandani

Abstract

The coming decade will be the era of smart, intelligent systems; many of these will have computer vision algorithms at their core. Such systems include self-driving cars, UAVs, robots, mobile nodes with edge computing support, and AR/VR systems. Hence, accelerating computer vision workloads is a very important problem as of today and has begun to attract a considerable amount of attention in both academia and industry. For building such systems we need to focus on both the hardware and software aspects.

As of 2021, the community appears to have significantly migrated towards CNN (convolutional neural network) based algorithms that run on GPUs and custom accelerators including FPGAs. A significant proportion of state-of-the-art systems comprises highly optimized CNN-based algorithms running on bespoke accelerators and GPUs. This trend is justified because such design choices have proven their mettle against traditional algorithms that primarily relied on multicore processors for parallelization.

However, we argue that this trend is showing signs of saturation, does not hold all the time, and may not prevail in the future. We systematically analyze, deconstruct, and in some cases debunk such unipolar notions in this thesis. We, instead, offer a more nuanced argument and envision a system that is an amalgam of all kinds of hardware platforms and software algorithms that together provide the best performance or power efficiency – as evinced in our experiments. We first start with showing that for traditional vision algorithms, GPUs are not the best choice. We show that as the number of concurrent applications (conventional vision algorithms) running on GPUs increases, the performance drops due to increased contention at the shared resources. In contrast, the performance of a bag of these concurrent applications on multicore processors scales well with an increase in the number of concurrent applications due to well-developed contention management policies in multicores and the inherent nature of conventional vision algorithms. This work is extended to create an accurate predictor that predicts the performance (91% accuracy) and power (96% accuracy) of a set of concurrently running vision applications on a GPU given their CPU execution characteristics and the execution time/power of individual

workloads on a GPU.

This work establishes multicore CPUs as legitimate systems for executing conventional vision applications. Next, we generalize this idea and consider a system with diverse computing elements that include small cores, large cores, possibly GPUs and accelerators, and then propose a scheduling algorithm. Given the large number of choices and attendant trade-offs, existing approaches find it very hard to compute optimal thread execution schedules for such systems. We developed a fast game theory-based technique that gives us near-optimal schedules very quickly without sacrificing fairness. The approach yields a performance benefit of 17% over the nearest competing algorithm.

Subsequently, we survey and characterize the CNN-based accelerators for vision systems. We found that the energy efficiency of the accelerators designed for CNNs relies on the amount of data reuse, the sparsity in the data, and the ineffectuality in the bits used for representing the data. Then, we choose the best CNN-based stereo algorithm and the best traditional computational geometry-based algorithm and compare their depth estimation accuracies head-to-head for a self-driving task. The results are *similar*, there is no clear winner. However, if we start considering inclement weather such as rain or fog, then the results are mixed. In roughly a third of the cases, the CNN-based solution is much better and the reverse holds true for roughly another third of the cases. For the rest, the accuracies are similar. We thus design an ensemble predictor that predicts the best algorithm for a given scene. This helps us maximize the overall accuracy and reduce the error in depth estimation substantially. We achieve up to an 18% (6.25% average) improvement in accuracy over the best state-of-the-art algorithm with a negligible area overhead (0.003 mm^2) while respecting all real-time constraints. Our algorithm also yields a confidence of prediction that can be utilized by higher decision-making layers. Such experiments prove that CNNs are not always the best; they have robustness issues and don't generalize very well as compared to their counterparts that use computational geometry too.

By this point we establish that in a real system we need both CPUs as well as GPUs and accelerators; we also need CNN-based and conventional computational geometry-based algorithms to create an accurate and robust vision system. Next, we consider other knobs that characterize cyber-physical systems such as UAVs, which use vision algorithms for their navigation or as standalone applications. Here, apart from the right choice of the computer vision algorithm, there are many more parameters that are important such as collision-free path planning, power efficiency, and flight duration. We show that when we include all these parameters with the vision system, optimizing the full set of parameters for the entire system becomes a very complex problem. It is large, multi-dimensional, non-convex, and not amenable to solution in real-time. We use game theory again and replace this intractable optimization problem with a simpler

yet approximate game theory equivalent: a set of players with different payoffs and strategies. Now, a game theory solver (Gambit [1]) can very quickly (10 – 100× faster than solving the optimization problem) obtain approximate yet good-quality solutions to our optimization problem.

संक्षेप

आने वाला दशक स्मार्ट, इंटेलिजेंट सिस्टम का युग होगा; इनमें से कई के मूल में कंप्यूटर विज्ञान एल्गोरिदम होंगे। इस तरह के सिस्टम में सेल्फ-ड्राइविंग कार, यूएवी, रोबोट, एज कंप्यूटिंग सपोर्ट वाले मोबाइल नोड्स और एआर / वीआर सिस्टम शामिल हैं। इसलिए, कंप्यूटर विज्ञान वर्कलोड में तेजी लाना आज की एक बहुत ही महत्वपूर्ण समस्या है और इसने अकादमिक और उद्योग दोनों में काफी ध्यान आकर्षित करना शुरू कर दिया है। ऐसी प्रणालियों के निर्माण के लिए हमें हार्डवेयर और सॉफ्टवेयर दोनों पहलुओं पर ध्यान देने की आवश्यकता है।

२०२१ तक, समुदाय काफी हद तक सीएनएन (कॉवोल्यूशनल न्यूरल नेटवर्क) आधारित एल्गोरिदम की ओर माइग्रेट हो गया है जो जीपीयू और कस्टम त्वरक जैसे एफपीजीए पर चलते हैं। अत्याधुनिक प्रणालियों के एक महत्वपूर्ण अनुपात में अत्यधिक अनुकूलित सीएनएन-आधारित एल्गोरिदम शामिल हैं जो बीस्पोक त्वरक और जीपीयू पर चल रहे हैं। यह प्रवृत्ति उचित है क्योंकि इस तरह के डिजाइन विकल्पों ने पारंपरिक एल्गोरिदम के खिलाफ अपनी योग्यता साबित कर दी है जो मुख्य रूप से समानांतरकरण के लिए मल्टीकोर प्रोसेसर पर निर्भर थे।

हालांकि, हम तर्क देते हैं कि यह प्रवृत्ति संतृप्ति के संकेत दिखा रही है, हर समय नहीं रहती है, और भविष्य में प्रबल नहीं हो सकती है। हम इस थीसिस में व्यवस्थित रूप से विश्लेषण, विघटन, और कुछ मामलों में ऐसी एकध्रुवीय धारणाओं को खारिज करते हैं। इसके बजाय, हम एक अधिक सूक्ष्म तर्क प्रस्तुत करते हैं और एक ऐसी प्रणाली की कल्पना करते हैं जो सभी प्रकार के हार्डवेयर प्लेटफॉर्म और सॉफ्टवेयर एल्गोरिदम का एक मिश्रण है जो एक साथ सर्वोत्तम प्रदर्शन या शक्ति दक्षता प्रदान करते हैं - जैसा कि हमारे प्रयोगों में दिखाया गया है। हम सबसे पहले यह दिखाना शुरू करते हैं कि पारंपरिक विज्ञान एल्गोरिदम के लिए, जीपीयू सबसे अच्छा विकल्प नहीं है। हम दिखाते हैं कि जैसे-जैसे जीपीयू पर चलने वाले समवर्ती अनुप्रयोगों (पारंपरिक दृष्टि एल्गोरिदम) की संख्या बढ़ती है, साझा संसाधनों के बढ़ते विवाद के कारण प्रदर्शन गिर जाता है। इसके विपरीत, मल्टीकोर प्रोसेसर पर इन समवर्ती अनुप्रयोगों के एक बैग का प्रदर्शन समवर्ती अनुप्रयोगों की संख्या में वृद्धि के साथ अच्छी तरह से बढ़ता है। इसका कारण मल्टीकोर में अच्छी तरह से विकसित विवाद प्रबंधन नीतियाँ और पारंपरिक विज्ञान एल्गोरिदम की अंतर्निहित प्रकृति है। इस कार्य में हम एक सटीक भविष्यवक्ता बनाते हैं जो एक जीपीयू पर समवर्ती रूप से चल रहे विज्ञान अनुप्रयोगों के एक सेट के प्रदर्शन (९१% सटीकता) और शक्ति (९६% सटीकता) की भविष्यवाणी करता है, जो कि उनके सी पी यू निष्पादन विशेषताओं और व्यक्तिगत कार्यभार के जीपीयू पर निष्पादन समय पर आधारित है।

यह कार्य पारंपरिक विज्ञान अनुप्रयोगों को क्रियान्वित करने के लिए वैध सिस्टम के रूप में मल्टीकोर सीपीयू को स्थापित करता है। इसके बाद, हम इस विचार को सामान्य बनाते हैं और विविध कंप्यूटिंग तत्वों वाली एक प्रणाली पर विचार करते हैं जिसमें छोटे कोर, बड़े कोर, संभवतः जीपीयू और त्वरक शामिल होते हैं, और फिर एक शेड्यूलिंग एल्गोरिदम का प्रस्ताव करते हैं। बड़ी संख्या में विकल्पों और परिचर ट्रेड-ऑफ को देखते हुए, मौजूदा दृष्टिकोणों को ऐसी प्रणालियों के लिए इष्टतम थ्रेड निष्पादन शेड्यूल की गणना करना बहुत कठिन लगता है। हमने एक तेज़ गेम थ्योरी-आधारित तकनीक विकसित की है जो हमें निष्पक्षता का त्याग किए बिना बहुत जल्दी-इष्टतम कार्यक्रम देती है। हमारा दृष्टिकोण १७% से अधिक का प्रदर्शन लाभ देता है निकटतम प्रतिस्पर्धी एल्गोरिथ्म पर।

इसके बाद, हम दृष्टि प्रणालियों के लिए सीएनएन-आधारित त्वरक का सर्वेक्षण और लक्षण वर्णन करते हैं। हमने पाया कि सीएनएन के लिए डिज़ाइन किए गए त्वरक की ऊर्जा दक्षता निर्भर करती है निम्नलिखित बातों पर – डेटा के पुनः उपयोग की मात्रा, डेटा में विरलता, और डेटा का प्रतिनिधित्व करने के लिए उपयोग किए जाने वाले बिट्स में अप्रभावीता। फिर, हम सबसे अच्छा सीएनएन-आधारित स्टीरियो एल्गोरिदम और सबसे अच्छा पारंपरिक कम्प्यूटेशनल ज्यामिति-आधारित एल्गोरिदम चुनते हैं और सेल्फ-ड्राइविंग कार्य के लिए उनकी गहराई अनुमान सटीकता की तुलना करते हैं। परिणाम समान हैं, कोई स्पष्ट विजेता नहीं है। हालांकि, अगर हम बारिश या कोहरे जैसे खराब मौसम पर विचार करना शुरू करते हैं, तो परिणाम मिले-जुले होते हैं। मोटे तौर पर एक तिहाई मामलों में सीएनएन-आधारित समाधान बहुत बेहतर होता है और लगभग एक तिहाई मामलों के लिए रिवर्स सही होता है। बाकी के लिए, सटीकता समान हैं। इस प्रकार हम एक समूह भविष्यवक्ता डिज़ाइन करते हैं जो किसी दिए गए दृश्य के लिए सर्वश्रेष्ठ एल्गोरिथ्म की भविष्यवाणी करता है। यह हमें समग्र सटीकता को अधिकतम करने में मदद करता है और गहराई से आकलन में त्रुटि को काफी हद तक कम करता है। हम सभी वास्तविक समय की बाधाओं का सम्मान करते हुए एक नगण्य क्षेत्र ओवरहेड (0.003 मिमी²) के साथ सर्वश्रेष्ठ अत्याधुनिक एल्गोरिदम पर सटीकता में 18% (6.25% औसत) सुधार प्राप्त करते हैं। हमारा एल्गोरिदम भविष्यवाणी का विश्वास भी पैदा करता है जिसका उपयोग उच्च निर्णय लेने वाली परतों द्वारा किया जा सकता है। ऐसे प्रयोग साबित करते हैं कि सीएनएन हमेशा सर्वश्रेष्ठ नहीं होते हैं; उनके पास मजबूती के मुद्दे हैं और वे अपने समकक्षों की तुलना में बहुत अच्छी तरह से सामान्यीकरण नहीं करते हैं जो कम्प्यूटेशनल ज्यामिति का भी उपयोग करते हैं।

इस बिंदु तक हम यह स्थापित करते हैं कि एक वास्तविक प्रणाली में हमें सीपीयू के साथ-साथ जीपीयू और त्वरक दोनों की आवश्यकता होती है; हमें सीएनएन-आधारित और पारंपरिक कम्प्यूटेशनल ज्यामिति-आधारित एल्गोरिदम की भी आवश्यकता है एक सटीक और मजबूत विज्ञान प्रणाली बनाने के लिए। इसके बाद, हम अन्य नॉक्स पर विचार करते हैं जो यूएवी जैसे साइबर-भौतिक प्रणालियों की विशेषता रखते हैं, जो अपने नेविगेशन के लिए विज्ञान एल्गोरिदम का उपयोग करते हैं। कंप्यूटर विज्ञान एल्गोरिथ्म के सही विकल्प के अलावा, कई और पैरामीटर हैं जो महत्वपूर्ण हैं जैसे कि टकराव मुक्त पथ योजना, बिजली दक्षता और उड़ान अवधि। हम दिखाते हैं कि जब हम इन सभी मापदंडों को विज्ञान सिस्टम के साथ शामिल करते हैं, तो पूरे सिस्टम के लिए मापदंडों के पूर्ण सेट को अनुकूलित करना एक बहुत ही जटिल समस्या बन जाती है। यह कई जगहों पर बड़ा, बहु-आयामी, गैर-उत्तल है, और वास्तविक समय में समाधान के लिए उत्तरदायी नहीं है। हम फिर से गेम थ्योरी का उपयोग करते हैं और इस कठिन अनुकूलन समस्या को एक सरल लेकिन अनुमानित गेम थ्योरी समकक्ष के साथ प्रतिस्थापित करते हैं: विभिन्न भुगतान और रणनीतियों वाले खिलाड़ियों का एक सेट। अब, एक गेम थ्योरी सॉल्वर बहुत जल्दी (ऑप्टिमाइज़ेशन समस्या को हल करने की तुलना में 10-100x तेज़) हमारी ऑप्टिमाइज़ेशन समस्या के अनुमानित लेकिन अच्छी गुणवत्ता वाले समाधान प्राप्त कर सकते हैं।

Contents

Certificate	i
Acknowledgements	iii
Abstract	vii
1 Introduction	1
2 Performance and Power Prediction for Concurrent Execution on GPUs	9
2.1 Introduction	9
2.1.1 Contributions and Organization of the Chapter	11
2.2 Background	13
2.2.1 Multi-application Concurrent Execution	13
2.2.2 ML-based Prediction Models	14
2.3 Related Work	17
2.3.1 Performance Models	17
2.3.2 Power Models	18
2.3.3 Challenges: Power and Performance Prediction for Concurrent Execution	19
2.4 Motivation	21

2.4.1	Overview of the Benchmarks	21
2.4.2	Experimental Setup	21
2.4.3	Performance Variation of the Workloads	22
2.4.4	Energy Variation of the Workloads	23
2.5	Methodology	25
2.5.1	Defining the Features	25
2.5.2	Creating the Data Points	30
2.5.3	Collection of Features	31
2.5.4	Predictor Models	31
2.6	Performance Predictor: Results and Analysis	33
2.6.1	Cross-validation	33
2.6.2	Comparison of Different Feature Combinations	34
2.6.3	Sensitivity of the Prediction Error	34
2.6.4	Analysis of the Decision Paths	35
2.7	Power Predictor: Results and Analysis	37
2.7.1	Experimental Setup	37
2.7.2	Cross-Validation	38
2.7.3	Comparison of Different Feature Combinations	39
2.7.4	Testing the Model: Random Split	42
2.8	Generalization to a Higher Concurrency	44
2.9	Conclusion	46
3	VisSched: An Auction-based Scheduler for Vision Workloads on Heterogeneous Processors	47

3.1	Introduction	47
3.1.1	Scope of our Work	48
3.1.2	Contributions	48
3.2	Auction-Based Scheduling	50
3.2.1	Overview	50
3.2.2	Nash Equilibria	50
3.2.3	Auction Theoretic Principles	51
3.3	Related Work	53
3.3.1	Scheduling on Heterogeneous Cores	53
3.3.2	Market Mechanisms for Scheduling of Heterogeneous Tasks	53
3.3.3	Auction Theory for Job-shop and Network Scheduling	54
3.4	Characterization of Workloads	56
3.4.1	Experimental Setup	56
3.4.2	Characterization of the MEVBench Suite	57
3.5	Design	59
3.6	Implementation of VisSched	60
3.6.1	Mapping the Set of Cores to Phases: Steps ❶, ❷ and ❸	60
3.6.2	Predicting the Phase of the Next Interval: Step ❹	61
3.6.3	Auctioning process: : Steps ❺ and ❻	61
3.6.4	Auctioning process: The working	64
3.6.5	Starvation	66
3.7	Evaluation of VisSched	68
3.7.1	Performance and ED^2	68

3.7.2	Fairness	71
3.7.3	Setting the Parameters	71
3.7.4	Other Statistics	72
3.7.5	Justification for the Mix of Cores	74
3.7.6	Sensitivity to the Big and Small Core Counts	75
3.7.7	Extending VisSched to include a LACore accelerator	75
3.7.8	Generality of the Scheme	77
3.7.9	Scalability of the Scheme	77
3.8	Conclusion	78
4	Accelerating CNN Inference on ASICs: A Survey	79
4.1	Introduction	79
4.1.1	Scope of the Survey	81
4.1.2	Organization	82
4.2	Background	84
4.2.1	Overview of a Convolutional Neural Network (CNN)	84
4.2.2	Running Example: Convolution Layer	88
4.2.3	Reference Architecture of a CNN Accelerator	89
4.3	Taxonomy	92
4.4	Reduction in Computation Time	93
4.4.1	Exploiting the Inherent Parallelism in CNNs	93
4.4.2	Pattern-based Computation Reduction in Convolution	96
4.4.3	Removal of Ineffectual Computations	98
4.4.4	Prediction-driven Computation Reduction	100

4.4.5	Improving the PE Utilization	102
4.5	Reduction of the Memory Access Time	105
4.5.1	Data Reuse: Temporal Reuse	105
4.5.2	Data Reuse: Spatial Reuse	108
4.5.3	Eliminate Loading of Ineffectual Data	110
4.5.4	Miscellaneous Techniques	115
4.6	Reduction in the Memory Footprint	117
4.6.1	Data and Weight Compression	117
4.6.2	Reduction of Ineffectual Bits in the Binary Representation	119
4.6.3	Reduction of Precision	123
4.7	Industrial Designs of CNN Accelerators	125
4.7.1	SIMD-based design	125
4.7.2	VLIW + SIMD	125
4.7.3	MIMD paradigm	126
4.7.4	VLIW with Systolic Computation	126
4.8	Discussion	128
4.8.1	Accuracy of the Competing Architectures	130
4.8.2	Comparison of the schemes	132
4.9	Challenges in ASIC Design of CNN Accelerators	135
4.10	Conclusion	137
4.11	Future Directions	138
5	PredStereo: An Accurate Real-time Stereo Vision System	141
5.1	Introduction	141

5.1.1	Contributions	143
5.2	Background and Related Work	145
5.2.1	Depth Estimation	145
5.2.2	Related Work	145
5.3	Characterization of Stereo Vision Workloads	146
5.3.1	CNN-based Workloads	146
5.3.2	Traditional Workloads	148
5.4	Design of the Predictor	150
5.4.1	Data Augmentation	150
5.4.2	The Selection Predictor	151
5.4.3	Sensitivity Analysis of the Predictor	154
5.4.4	The Confidence Predictor	154
5.5	Architecture	156
5.5.1	Hardware Predictor	156
5.6	Evaluation of PredStereo	159
5.6.1	Setup	159
5.6.2	Synthesis Results	160
5.6.3	Performance Comparison	160
5.6.4	Accuracy Comparison	162
5.6.5	Analysis of the Disparity Estimation Error in a Self-Driving Scenario	162
5.6.6	Qualitative Evaluation	164
5.6.7	Novelty vis-a-vis Related Work	165
5.7	Conclusion	166

6	Game Theory-based Parameter Tuning for Energy-efficient Path Planning on Modern UAVs	167
6.1	Introduction	167
6.2	Background	171
6.2.1	Navigation in UAVs	171
6.2.2	Game Theory Preliminaries	172
6.2.3	Relating Optimization Problems to Game Theory	173
6.3	Modeling the Energy Consumption of UAVs	175
6.3.1	Experimental Setup	175
6.3.2	Data Collection	175
6.3.3	Data Pre-processing	176
6.3.4	Quantification and Modeling	177
6.3.5	Accuracy Comparison of Different Power Models	182
6.4	Experimental Setup	184
6.4.1	Overview	184
6.4.2	Setup for Sensitivity Analyses	185
6.4.3	Gazebo and Rviz for UAV Simulation	186
6.4.4	Creation of Virtual Worlds	186
6.4.5	IPOPT and AMPL	187
6.4.6	Setup for Measuring the Power and Performance of Vision Algorithms	187
6.5	Formulation of the Optimization Problem	189
6.5.1	Collection of Data Points	189
6.5.2	Curve Fitting: Hover Time and Path Length	190

6.5.3	Formulation of the Optimization Problem	191
6.5.4	Sensitivity to the Power Envelope (P_{drone})	193
6.6	Game Theory	195
6.6.1	Sensitivity Analyses	195
6.6.2	Game Setup	196
6.6.3	Choosing the Hyper-parameters	198
6.7	Results	200
6.7.1	Comparison of Optimization-based and Game Theory-based Approaches for RRT*	200
6.7.2	Performance Comparison of Solvers on BeagleBone Black for RRT*+Vision	201
6.7.3	Results from the Game Theory-based Approach for RRT*+Vision . . .	202
6.7.4	Comparison of Game Theory and Optimization-based Approaches for RRT*+Vision: Gambit vs IPOPT	203
6.8	Related Work	205
6.8.1	Parameter Tuning	205
6.8.2	Power Modeling	205
6.9	Conclusion	207
7	Conclusion and Future Work	209
7.1	Contributions	209
7.2	Future Work	212
	Bibliography	213
	List of Publications	243

Biography

245

List of Figures

1.1	An overview of the work done	4
2.1	(a) CPU performance (normalized), (b) GPU performance (normalized), and (c) GPU/CPU performance	22
2.2	(a) CPU energy (normalized), (b) GPU energy (normalized), and (c) GPU/CPU energy	23
2.3	Black box diagram of the predictor	25
2.4	Clustering of workloads	29
2.5	Silhouette score (a score above 0.7 means that the clusters are tightly bound) .	29
2.6	MAPE (%) for leave-one-out cross validation for a concurrency of 2	33
2.7	Comparison of different feature combinations for a concurrency of 2	33
2.8	MAPE (%) for different feature combinations for a concurrency of 2	34
2.9	Percentage of the test data points containing a feature in their decision path for a concurrency of 2	36
2.10	Radar plot of the frequency of each feature in the decision making for different test points for a concurrency of 2	36
2.11	Violin Plots showing the error distribution of the 20 combinations for 6 ML models for each benchmark	39
2.12	Strip Plot showing the MAPE for 20 combinations of features using LOOCV .	39

2.13 Heatmaps of different combinations and ML models (results show the MAPE (%) error)	40
2.14 Strip plot showing the MAPE for 20 combinations of features	43
2.15 Split violin plot showing the error distribution with LOOCV and random data split validation	43
2.16 Performance prediction for 3 concurrent applications	44
2.17 Power prediction for 3 concurrent applications	44
3.1 Clustering of all the phases across the benchmarks	57
3.2 (a) Silhouette plot of the clusters, (b) Phase-wise correlation of <i>Sift</i> and <i>Surf</i>	58
3.3 Overview of the design	59
3.4 Hardware Scheduler	60
3.5 A row of the PT (pattern table)	65
3.6 Method of generating the PT	65
3.7 Perf. comparison: 20 threads on 16 cores	68
3.8 Perf. comparison for a bag-of-tasks:20 threads on 16 cores	70
3.9 Normalized ED^2 (w.r.t <i>SLICC</i>) for a bag-of-tasks:20 threads on 16 cores	70
3.10 Fairness for the bag-of-task workloads	71
3.11 Search space for 2 bidders	73
3.12 Constrained search space for 4 bidders	73
3.13 Energy vs slowdown for different combinations of fast (fc) and slow cores (sc) (normalized to 16fc)	74
3.14 Sensitivity of VisSched w.r.t. big and small core counts	75
3.15 Perf. comparison of VisSched with state-of-the-art (accelerator included in all the schemes)	76

3.16 Performance of Cortexsuite [31] and self-driving applications: 20 threads on 16 cores	76
3.17 Performance comparison: 64 threads on 32 cores	77
4.1 Evolution of CNNs in recent years (data taken from [126])	80
4.2 A convolution layer and a pooling layer in a CNN	86
4.3 Algorithm for a tiled convolution layer (adapted from Figure 5 in [158], names of variables have been changed)	89
4.4 (a) Reference dataflow architecture for CNNs (adapted from [140, 159], and (b) Normalized energy (w.r.t ALU) for accessing data from the memory hierarchy.	90
4.5 Different types of architectures: (a) 1D Systolic, (b) 2D systolic, (c) 1D array, and (d) 2D matrix (adapted from [160]). DLY is the delay.	91
4.6 Taxonomy of CNN accelerators	92
4.7 Venn diagram of hardware optimizations	92
4.8 Factorization of the dot product (adapted from [175])	96
4.9 Redundant data detection circuit	98
4.10 Kernel decomposition scheme for the single input and multiple output case	100
4.11 (a) Prediction based (adapted from Song et al. [183]), and (b) SNAPEA (adapted from Akhlaghi et al. [184])	101
4.12 Weight lookaside and Weight lookahead (adapted from [185])	103
4.13 Dataflow architectures for (a) Output stationary, (b) Weight stationary, (c) No local reuse , and (d) Row stationary (adapted from [140])	106
4.14 Step Indexing (adapted from [199])	111
4.15 ZFNaf encoding format (adapted from [202])	112
4.16 (a) Working of Cnvlutin and Cnvlutin2 (adapted from [203])	113
4.17 NZVL encoding scheme (adapted from [205])	114

4.18	Weight storage architecture	118
4.19	Ineffectual bits in the precision	118
4.20	(a) Parallel multiplier unit, (b) Stripes with equivalent throughput, and (c) Pragmatic with equivalent throughput (adapted from [219])	120
4.21	Concept of LUT based PE (LBPE) (adapted from [224])	122
5.1	Error difference of SGM and Highres for images in KITTI 2015 dataset	143
5.2	Error difference of SGM and Highres for augmented KITTI dataset	143
5.3	Example of data augmentation on KITTI dataset. Annotations show the algorithm that achieves better accuracy on the frame	151
5.4	Accuracy of different classifiers	153
5.5	Frequency of features in the decision making process	153
5.6	Sensitivity of the predictor to brightness	154
5.7	Sensitivity of the predictor to the percentage of dark regions	154
5.8	Architecture of the SoC with ASIC	156
5.9	Architecture of the hardware predictor.	157
5.10	Sensitivity of the total time to the image size	160
5.11	Execution time per frame	161
5.12	Reduction in the safety buffer time	164
5.13	Results of disparity estimation on KITTI-2015 training images. The first column shows the left image of the stereo image pair. The second and the third columns show the disparity map obtained by Highres [9] and SGM [247], respectively.	164
5.14	Visualization of some features of KITTI-2015 training images. The first column shows the left image of the stereo image pair. The second and the third columns show the dark regions (in white) and the histogram of the grayscale image, respectively.	165

6.1	Drones used in our experiments	175
6.2	Power module and Pixhawk flight controller	176
6.3	Phases in a flight of the <i>20in4</i> drone	176
6.4	Controlled flight	177
6.5	Random flight	177
6.6	Hovering power	178
6.7	Power consumption for vertically upward motion of the <i>20in4</i> and <i>20in8</i> drones	179
6.8	Power consumption for vertically downward motion of the <i>20in4</i> and <i>20in8</i> drones	180
6.9	Power consumption for the horizontal motion of the <i>20in4</i> and <i>20in8</i> drones . .	181
6.10	Mean absolute error (%) of ML models for estimating the drone power	182
6.11	Overview of our approach	186
6.12	RMSE comparison of different models for World 1	190
6.13	RMSE comparison of different models for World 2	190
6.14	RMSE comparison of different models for World 3	190
6.15	Sensitivity of path length to P_{drone} for the <i>20in4</i> drone	193
6.16	Sensitivity of path length to P_{drone} for the <i>20in8</i> drone	193
6.17	Search space for the optimization objective keeping resolution=0.4, and obstacle avoidance distance=0.1.	194
6.18	Search space for the optimization objective keeping resolution=0.7, and obstacle avoidance distance=0.1.	194
6.19	No. of obstacles v/s hover time	195
6.20	No. of obstacles v/s path length	195
6.21	Resolution of the path v/s hover time	195

6.22	3D Occupancy map and the planned path for World 1	201
6.23	3D Occupancy map for World 2 (tree density 0.1 trees/m^2)	201
6.24	3D Occupancy map for World 3 (tree density 0.2 trees/m^2)	201
6.25	Path length for the <i>20in4</i> drone	202
6.26	Path length for the <i>20in8</i> drone	202
6.27	Hover time for the <i>20in4</i> drone	202
6.28	Hover time for the <i>20in8</i> drone	202
6.29	Path lengths for the <i>20in4</i> drone	203
6.30	Path lengths for the <i>20in8</i> drone	203

List of Tables

2.1	Comparison with related work	18
2.2	Benchmarks (derived from MEVBench)	21
2.3	Baseline system	21
2.4	Features used in the performance estimation model	26
2.5	Features used in the power estimation model	27
2.6	Configurations of ML models	30
2.7	Combination of features and MAPE (%) of the regressors	38
3.1	Details of the baseline system (source [101])	56
3.2	Heuristics for phase prediction	61
3.3	Terms and definitions	62
3.4	Description of the competing techniques	69
3.5	Cache statistics for the bag-of-task workloads	70
3.6	Area overheads of the hardware structures	72
3.7	Phase prediction, and the auction process for 200M instruction run	72
3.8	Linear algebra kernels in the benchmarks	76
4.1	Comparison of FPGAs and ASICs (adapted from the observations in [130, 134])	81

4.2	Glossary of loop iterators in a tiled CNN	93
4.3	Loop variables targeted for parallelization	93
4.4	Types of parallelism exploited by different architectural implementations	94
4.5	Types of reuse exploited by different proposals	108
4.6	Types of optimizations exploited by different proposals	129
4.7	Accuracy loss of several recently proposed techniques	131
4.8	Comparison of dataflows (data obtained by running Timeloop [195])	133
4.9	Performance comparison of the proposals w.r.t. [234] (taken from original papers)	134
5.1	3-pixel error and execution times for CNN-based stereo algorithms	146
5.2	Execution time per frame, energy consumption, and error for CNN-based stereo algorithms	147
5.3	Error comparison	148
5.4	Features and their description	152
5.5	Configurations of the predictors	152
5.6	Accuracy of the classifiers for different scenarios	153
5.7	Mean Absolute Error of the confidence predictor for different scenarios	155
5.8	Overheads of the hardware structures	159
5.9	Schemes and their disparity estimation methods	161
5.10	Comparison of the error in the disparity estimation of different schemes in several weather conditions	162
6.1	Tunable parameters of the RRT* algorithm	185
6.2	Baseline system	186

6.3	Details of the system (source [74])	186
6.4	Benchmarks (derived from MEVBench)	188
6.5	Comparison of optimization-based and game theory-based approaches	200
6.6	Comparison of planning time (sec) for random and best configurations	201
6.7	Best platform	203

