

DEPARTMENT OF COMPUTER SCIENCE

PhD Degree Oral Presentation

PhD Candidate:	Mr. Qiang WANG
Date	August 7, 2020 (Friday)
Time:	10:30 am – 12:30 pm (35 mins presentation and 15 mins Q & A)
Venue:	Zoom ID: 970 1107 1417 (The password and direct link will only be provided to registrants)
Registration:	https://bit.ly/sem-zm (Deadline: 5:00pm, 6 August 2020)

Performance and Power Modeling of GPU Systems with Dynamic Voltage and Frequency Scaling

Abstract

To address the ever-increasing demand for computing capacities, more and more heterogeneous systems have been designed to use both general-purpose and special purpose processors. The huge energy consumption of them raises new environmental concerns and challenges. Besides performance, energy efficiency is another key factor to be considered by system designers and consumers. In particular, contemporary graphics processing units (GPUs) support dynamic voltage and frequency scaling (DVFS) to balance computational performance and energy consumption. However, accurate and straightforward performance and power estimation for a given GPU kernel under different frequency settings is still lacking for real hardware, which is essential to determine the best frequency configuration for energy saving. In this thesis, we investigate how to improve the energy efficiency of GPU systems by accurately modeling the effects of GPU DVFS on the target GPU kernel.

First, we present a benchmark suite EPPMiner for evaluating the performance, power, and energy of different heterogeneous systems. EPPMiner consists of 16 benchmark programs that cover a broad range of application domains, and it shows a great variety in the intensity of utilizing the processors. We have implemented a prototype of EPPMiner that supports OpenMP, CUDA, and OpenCL, and demonstrated its usage by three showcases. The showcases justify that GPUs provide much better energy efficiency than other types of computing systems, and especially illustrate the effectiveness of GPU Dynamic Voltage and Frequency Scaling (DVFS) on the energy efficiency of GPU applications.

Second, we reveal a fine-grained analytical model to estimate the execution time of GPU kernels with both core and memory frequency scaling. Compared to the cycle-level simulators, which are too slow to apply on real hardware, our model only needs one-off micro-benchmarks to extract a set of hardware parameters and kernel performance counters without any source code analysis.

Third, we design a cross-benchmarking suite, which simulates kernels with a wide range of instruction distributions. The synthetic kernels generated by this suite can be used for model pre-training or as supplementary training samples. We then build machine learning models to predict the execution time and runtime power of a GPU kernel under different voltage and frequency settings.

At last, we establish a new DDL job scheduling framework which organizes DDL jobs as Directed Acyclic Graphs (DAGs) and considers communication contention between nodes. We then propose an efficient job placement algorithm, LWF- κ , to balance the GPU utilization and consolidate the allocated GPUs for each job. When scheduling the communication tasks, we propose Ada-SRSF for the DDL job scheduling problem to address the communication contention issue.

***** ALL INTERESTED ARE WELCOME *****