

# NATIONAL UNIVERSITY OF SINGAPORE

School of Computing

## PH.D DEFENCE - PUBLIC SEMINAR

**Title:**           **Energy-Time Performance of Heterogeneous Computing Systems: Models and Analysis**

**Speaker:**       Ms Lavanya Ramapantulu

**Date/Time:**   17 August 2016, Wednesday, 10:00 AM to 11:30 AM

**Venue:**           Executive Classroom, COM2-04-02

**Supervisor :**   Dr Teo Yong Meng, Associate Professor, School of Computing

### Abstract:

While heterogeneity is increasingly becoming the norm in most computing platforms today, one of the key challenges is to determine the set of energy-time efficient system configurations among the large system configuration space to execute a parallel application. This large configuration space offers a new opportunity to improve the match between parallel application demands and system resources to achieve efficient energy-time performance. This thesis presents an approach to address this challenge using a measurement-driven analytical model that determines both time and energy efficient system configurations. Based on our taxonomy of heterogeneous computing systems, we first propose a core analytical model for a baseline heterogeneous system representing inter-node heterogeneity and consisting of brawny and wimpy nodes. The proposed core model is scalable for different types of heterogeneity and is formulated using parametric values obtained from baseline measurements of the application for better accuracy. The key novelties of our approach include modeling both inter and intra-node resource overlaps and resource contention.

Among heterogeneous systems, intra-node heterogeneous systems with Vector Processing Units (VPUs) are increasingly being adopted in the Top500 supercomputers as they offer accelerated performance gains. Secondly, the impact of heterogeneity in parallel programs leads to the wider adoption of hybrid programming models for scientific applications. Hybrid programming models are gaining traction as they exploit system resources and parallelism at both inter- and intra-node levels. The scalability of our proposed core model is shown by extending it to both intra-node heterogeneous system and hybrid programs. Key model extensions include (i) inter- and intra-core contentions for VPUs in a Many Integrated Core (MIC) architecture system, and (ii) inter- and intra-node communication for hybrid programs.

With the advent of heterogeneity at both program and system level, it is non-trivial for

application developers or users to choose an energy and time optimal configuration from the large configuration space. The proposed core model and its extensions are applied to determine energy-time efficient system configurations for inter-node heterogeneous system, intra-node heterogeneity with VPUs and hybrid programs. In determining these efficient system configurations, we exposed a number of insights. Firstly, there are multiple Pareto-optimal "sweet-spot" configurations that can be approximated using a distinct energy-deadline Pareto-frontier. These configurations facilitate energy-time trade-offs such as to minimize energy used for a given execution-time deadline and/or to minimize execution time for a given energy budget. Our analysis shows that for inter-node heterogeneous clusters and hybrid programs, energy savings of up to 75% can be achieved by selecting Pareto-optimal configurations as opposed to non-optimal configurations. Furthermore, among the Pareto-optimal configurations hybrid programs reduces the energy used by up to 65% at the expense of 18% increase in execution time.

With the explosion of the configuration space, we show that the Pareto-frontier can be analytically established using the node performance-to-power ratios (PPR). We show that the Pareto-frontier can be further improved by replacing low PPR nodes with higher PPR nodes using our power substitution ratio. Additionally, our energy proportionality analysis reveals that inter-node heterogeneous clusters enable the scaling of the energy-proportionality wall by exposing sub-linear energy-proportional configurations. To further optimize the Pareto-frontier, we introduce a new metric called useful computation ratio (UCR) to quantify the degree of resource contentions and communication overheads in an execution. Lastly, we show how UCR and Pareto-optimal configurations can be used in conjunction by system designers to gain further insights into system resource imbalances, and how application developers can further fine-tune their hybrid programs.