# Learning-based Runtime Management of Energy-Efficient and Reliable Many-Core Systems

Rishad A. Shafik, Anup K. Das, Sheng Yang, Geoff V. Merrett & Bashir M. Al-Hashimi
School of ECS, University of Southampton, SO17 1BJ, UK, *e-mail: {ras1n09,akd1g13,gvm,bmah}@ecs.soton.ac.uk*

## I. INTRODUCTION

Silicon technology scaling has enabled the fabrication of many interconnected cores on a single chip for current and future generations of computing systems. The emergence of such systems has facilitated computing performance at unprecedented levels with application parallelization and architectural support. However, higher device-level integration and faster speed in these systems have rendered exponentially increased power density and energy consumption [1]. Due to such high power density, these systems are increasingly being confronted with emerging reliability challenges. A major challenge is significantly reduced lifetime of these systems due to higher operating temperatures and their variations, which accelerate the device wearout mechanisms through electromigration, dielectric breakdown, etc. Indeed, many-core computing with reduced energy and improved reliability, while dealing with the conflicting performance trade-offs, is highly challenging [3].

This paper highlights and demonstrates our research works to date, which address the above challenges through intelligent runtime management algorithms. The algorithms are implemented through cross-layer interactions between the three layers: application, runtime and hardware, forming one of our core themes of working together. The annotated application tasks communicate the performance, energy or reliability requirements to the runtime. With such requirements, the runtime exercises the hardware through various control knobs and gets the feedback of these controls through the performance monitors. The aim is to learn the best possible hardware controls during runtime to achieve energy-efficiency and improved reliability, while meeting the specified application requirements.

## II. LEARNING-BASED APPROACHES AND RESULTS

The paper will specifically report the methodology and experimental results of the following two runtime management approaches. The first approach [5] implements a reinforcement learning-based adaptive runtime thermal management through thread-to-core affinity and processor DVFS controls. The aim is to reduce different thermal emergencies (peak, average and thermal cycles) and extend the lifetime reliability, while meeting the application-specified performance requirements. The experimental results of this approach show that significant lifetime improvement (up to 7x) can be achieved compared to the existing methods (see Fig. 1). The second approach [6] demonstrates the use of programming model based runtime management of energy-efficient and reliable many-core systems. The approach uses hierarchical learning of dynamic concurrency throttling (DCT) and DVFS controls of processor cores in many-core systems. The experimental results show that up to 21% energy can be reduced compared to the existing methods, while

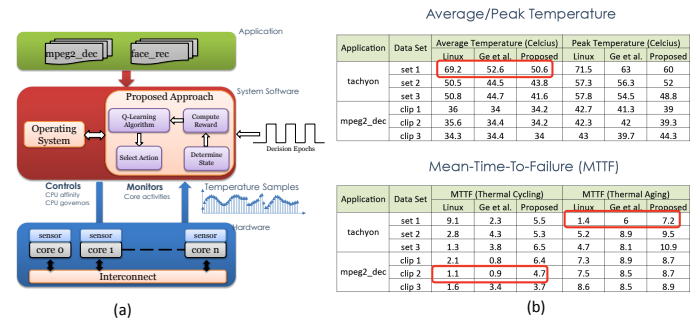lifetime reliability can also be improved significantly (see Fig. 2).



Fig. 1. (a) Learning-based cross-layer approach for reliability improvement of multi-core systems through thread affinity and DVFS controls, and (b) experimental results demonstrating improvement in average/peak temperatures and lifetime reliability using ALPbench applications
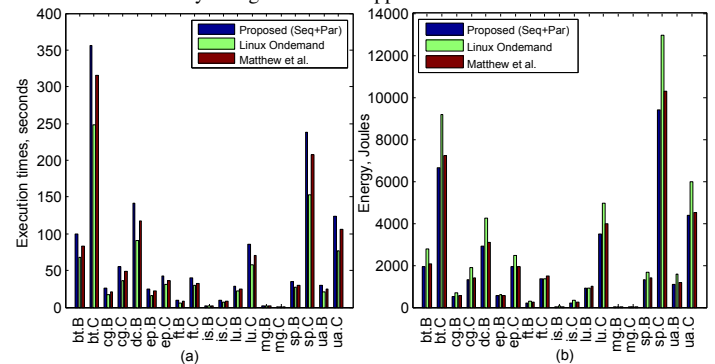


Fig. 2. Comparative evaluation of (a) performance (execution times, in seconds) and (b) energy consumptions (in Joules) of the OpenMP-based runtime energy minimization approach through dynamic concurrency throttling (DCT) and DVFS controls using NAS parallel benchmark applications

The validation of the approaches through experimental results showed that learning-based approaches effectively minimize energy and improve lifetime reliability when different layers work together interactively and intelligently, while also meeting the ever changing application performance requirements.

## REFERENCES

[1] H. Esmaeilzadeh *et al.*. Dark silicon and the end of multicore scaling. *38th ISCA*, 2011.
[2] D. Brodowski and N. Golde. Linux CPUFreq–CPUFreq governors. *Linux Kernel*. [Online]: http://www.mjmwired.net/kernel/Documentation/cpu-freq/governors.txt
[3] Y. Ge and Q. Qiu. Dynamic Thermal Management for Multimedia Applications Using Machine Learning. In DAC, 2011.
[4] C-M. Matthew *et al.*. Prediction models for multi-dimensional power-performance optimization on many cores. *ICPAC*, 2008.
[5] A.K. Das *et al.* Reinforcement learning-based inter-and intra-application thermal optimization for lifetime improvement of multicore systems. in *DAC*, pp.1–6, June, 2014.
[6] R.A. Shafik *et al.* Adaptive Energy Minimization of OpenMP Parallel Applications on Many-Core Systems. (under review).