

MANAGING POWER IN HETEROGENEOUS MULTICORE SYSTEMS

Dr Geoff Merrett

01-03 July 2019 | York, UK

14th Int'l Symp. Reconfigurable communication-centric Systems on Chip (ReCoSoC 2019)

PROCESSOR EVOLUTION

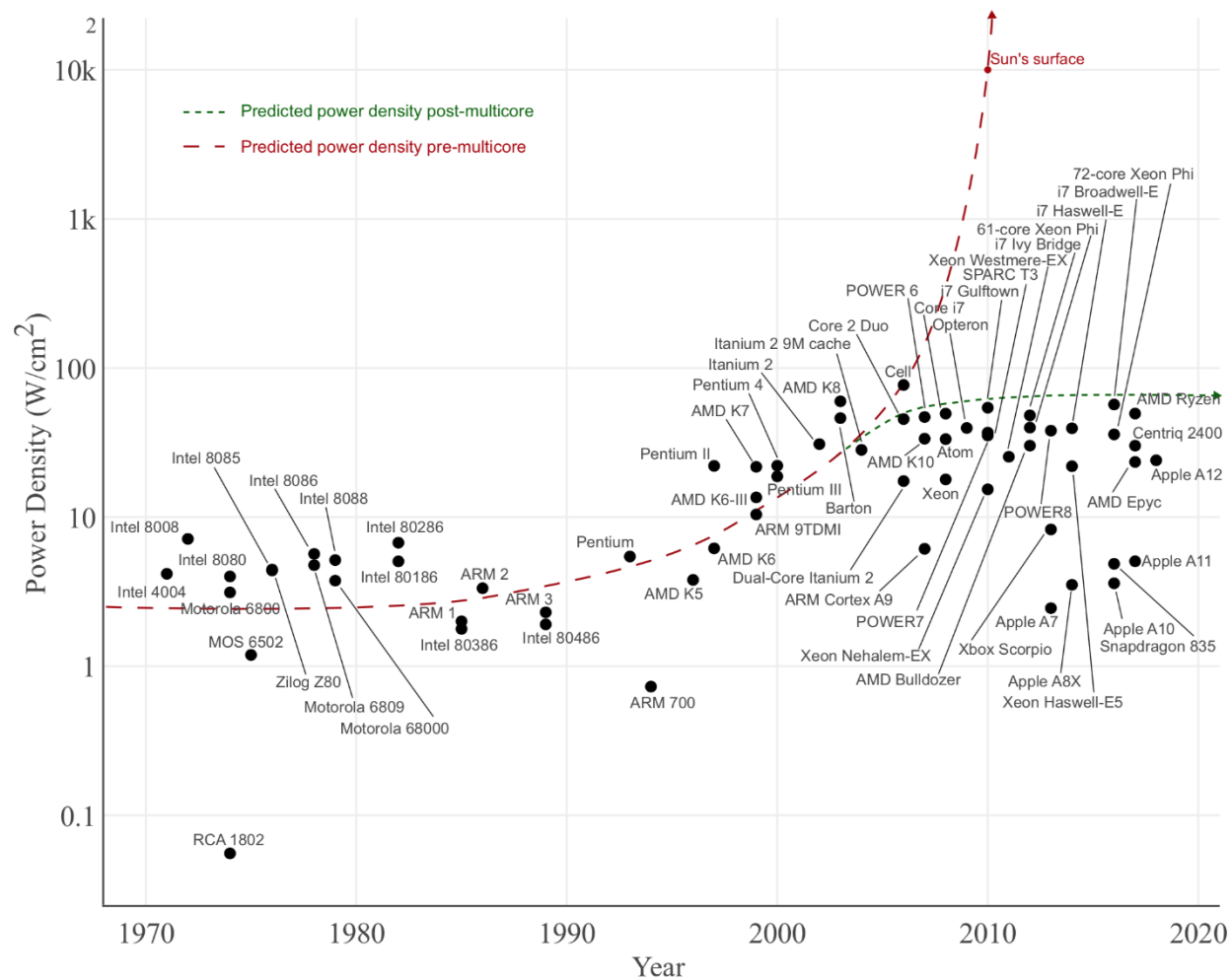
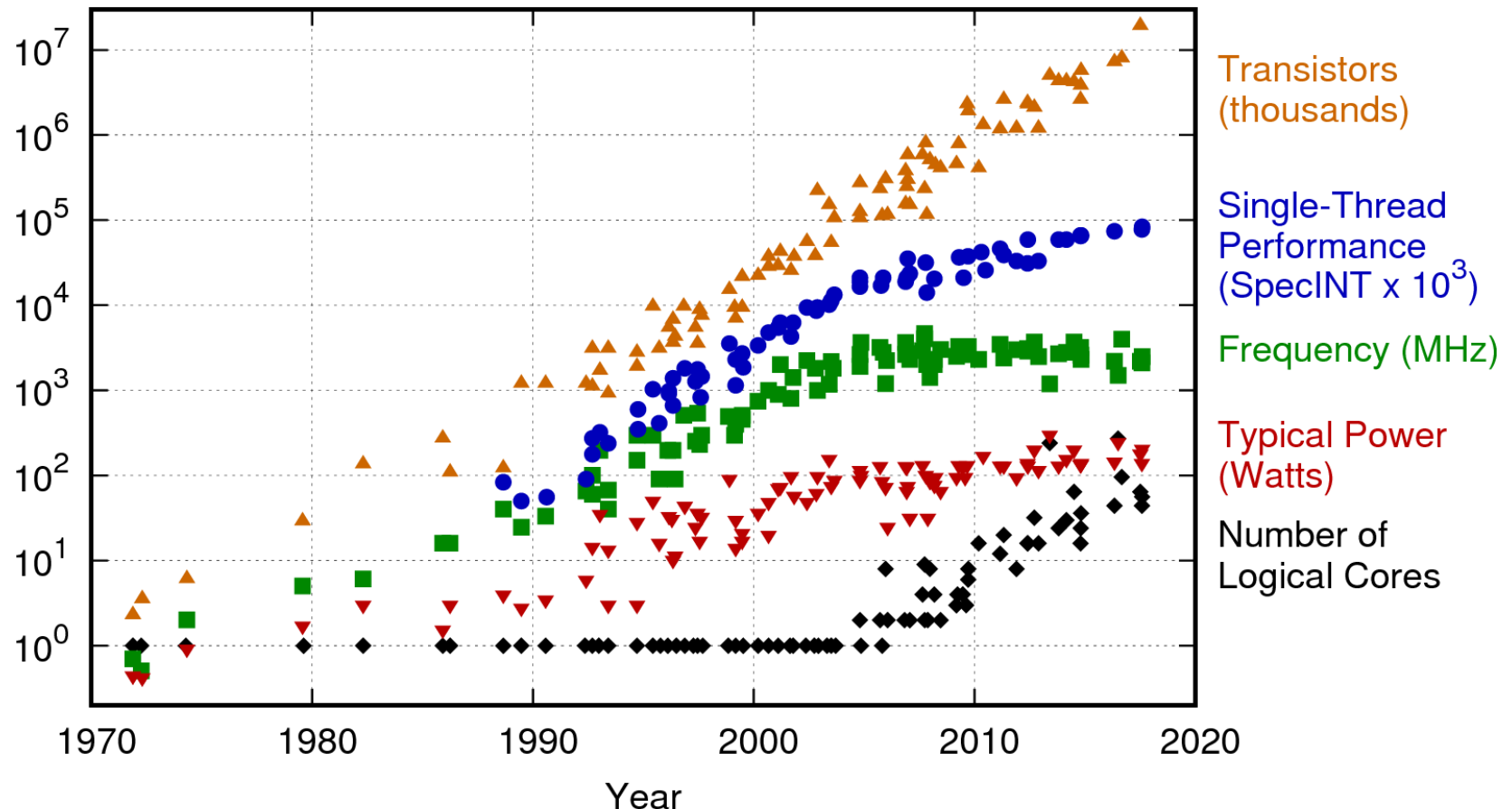


Image from Ilias Vougioukas' PhD Thesis, using data from Berkeley's extended database (Danowitz et al., 2012).

3

WE ARE MANY-CORE

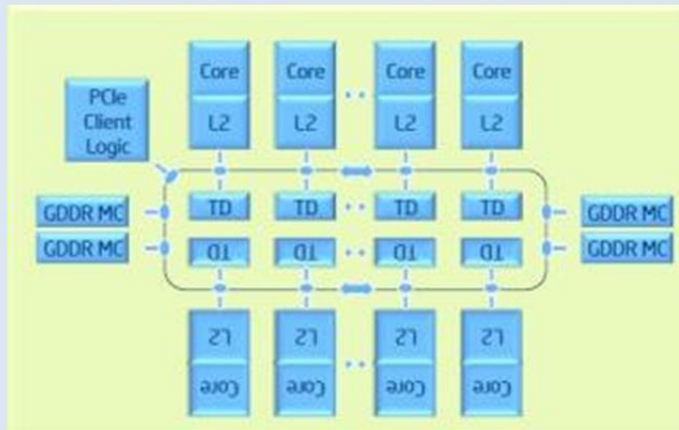
42 Years of Microprocessor Trend Data



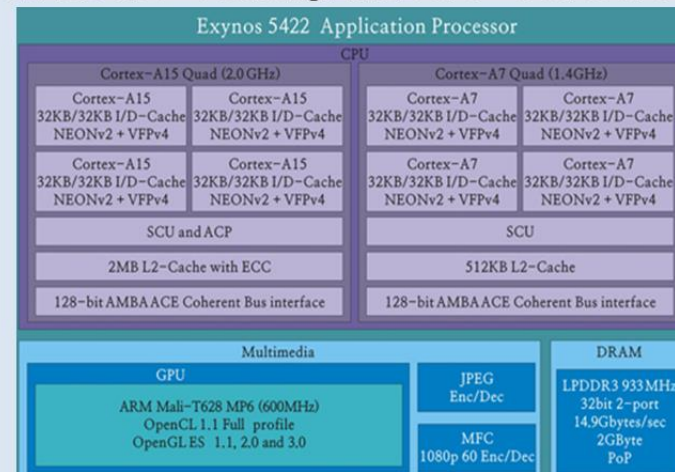
Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2017 by K. Rupp

WE ARE MANY-CORE

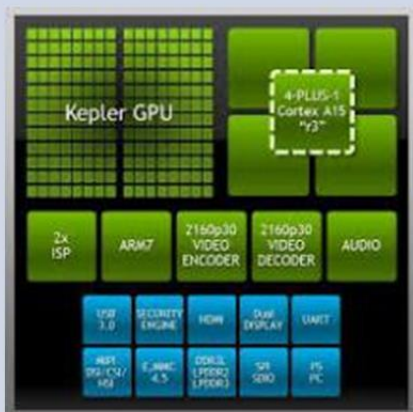
Intel Xeon Phi - Homogeneous 61 Cores



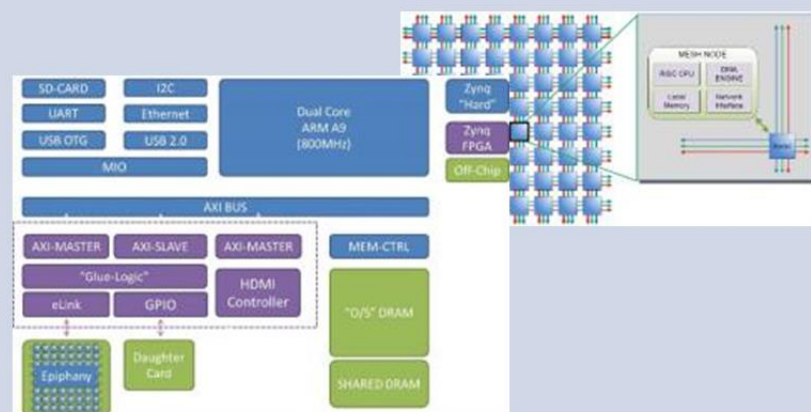
ODROID XU3 – 8 core big.LITTLE CPU + 6 cores GPU



Nvidia Jetson TK1 - Quad core CPU + 192 cores GPU

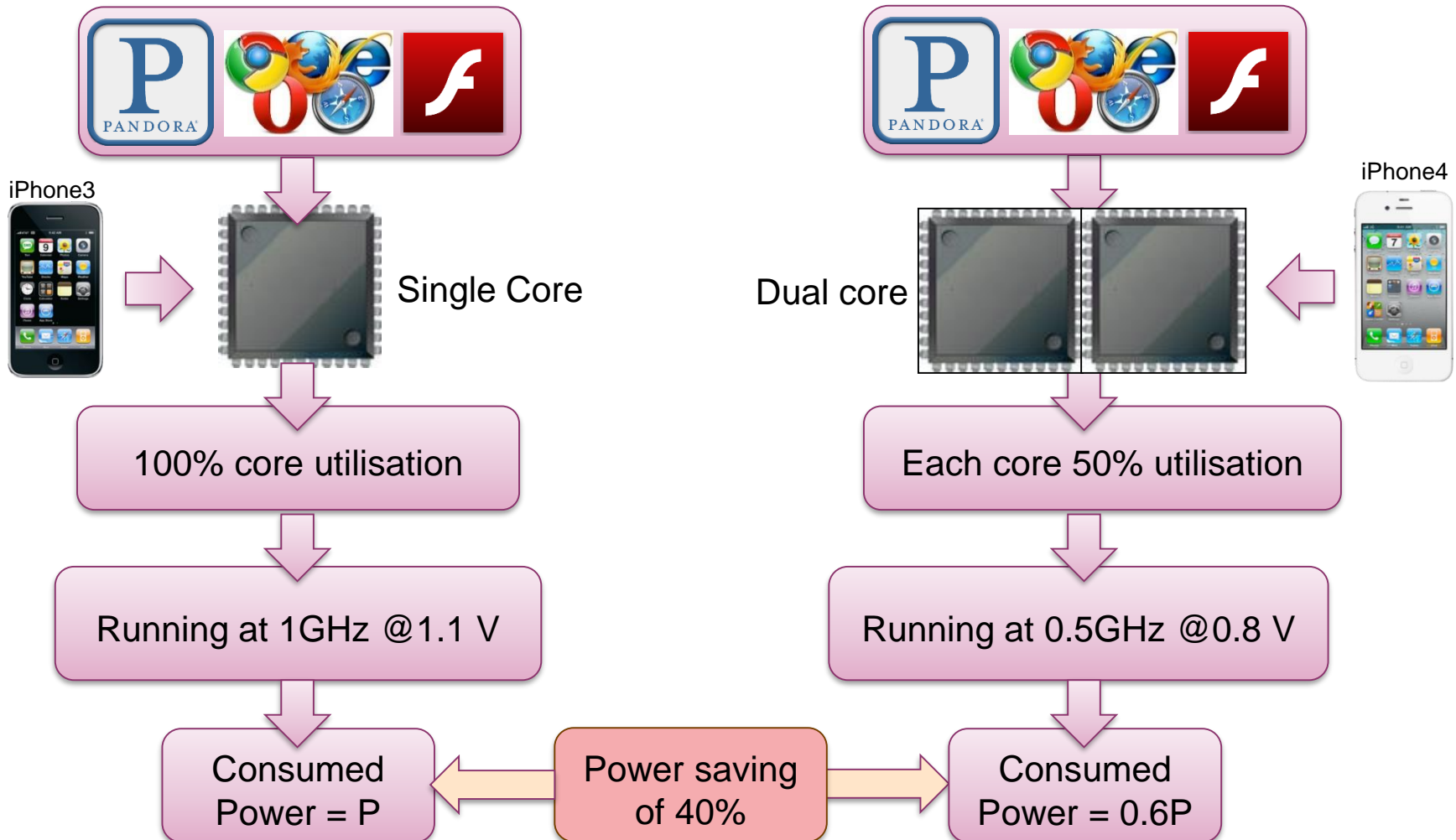


Parallela- **Dual core CPU + FPGA + 16 cores NoC**



WHY MULTI-CORE?

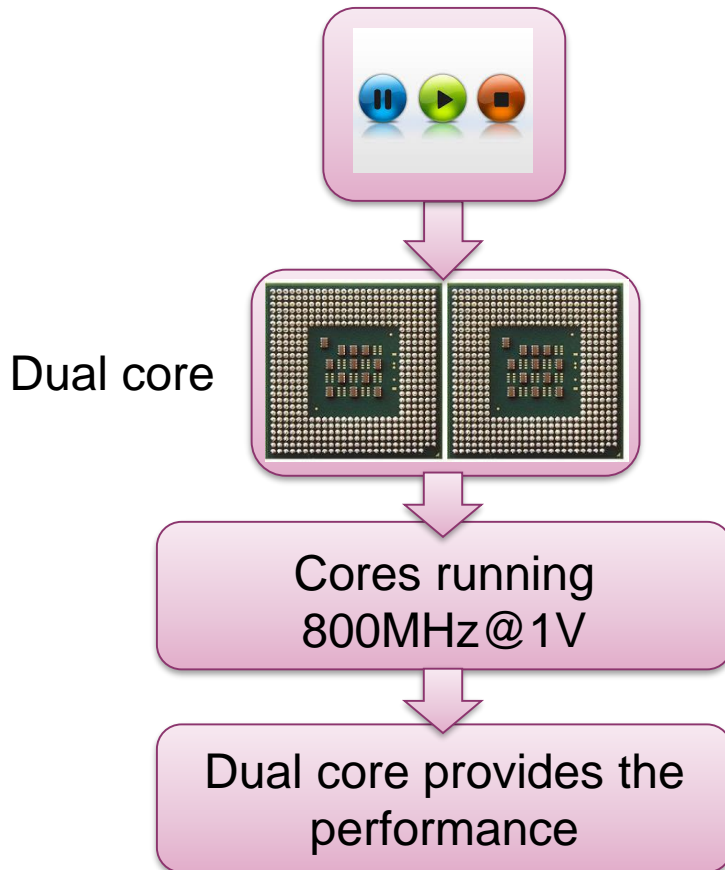
Dynamic Voltage Frequency Scaling



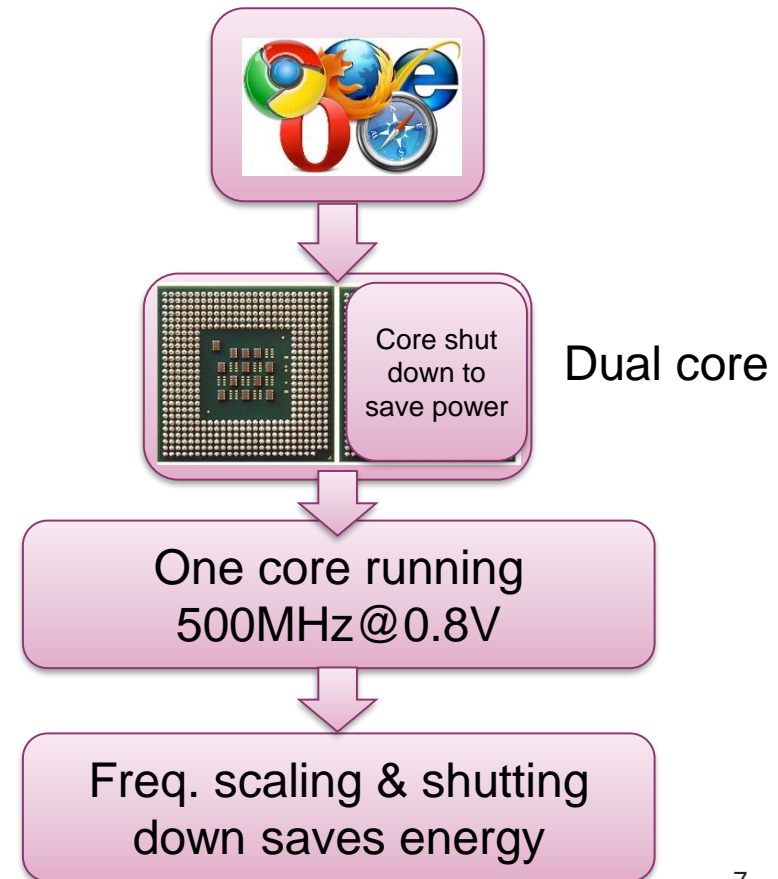
WHY MULTI-CORE?

Dynamic Power Management/Core Scaling

Multimedia (**high performance**)



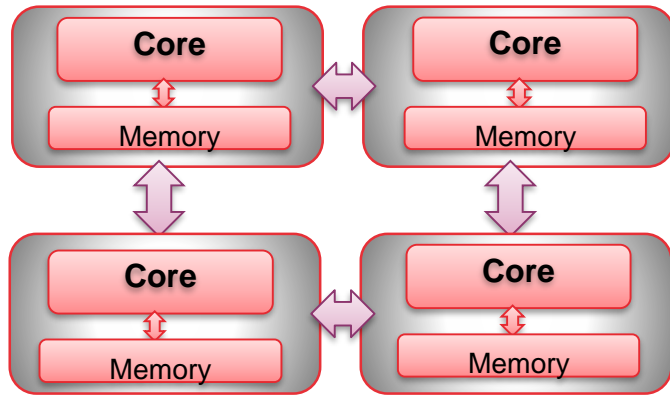
Browser (**low performance**)



WHY MULTI-CORE?

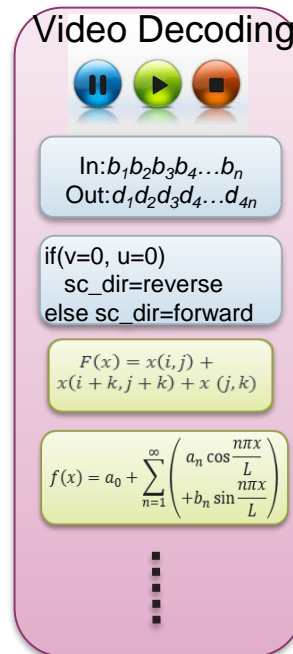
Heterogeneous Platforms

SW executing on Cores

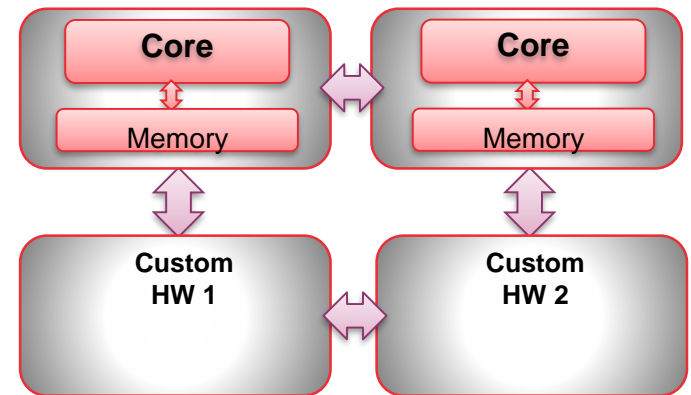


Performance: P
Energy: E (J)

Video Decoding



SW+Custom HW (ASIC)



Performance: $(P+x\%)$
Energy: $(E-x\%)$ (J)



THE PRiME PROJECT

“Enable the sustainability of **many-core scaling** by preventing the uncontrolled increase in **energy consumption** and **unreliability** through a step change in holistic design methods and **cross-layer** system optimisation.”

UNIVERSITY OF
Southampton

Imperial College
London

MANCHESTER
1824

Newcastle
University

arm

Imagination

intel

Microsoft Research

NXP

EPSRC
Engineering and Physical Sciences
Research Council

nmi
Semiconductors
to Systems

Innovate UK
Knowledge Transfer Network

www.prime-project.org

MORE COMPUTE FOR THE SAME POWER



A “*supercomputer in your hands*”, running 100s of cores with a battery lasting for a day

THE PRiME PROJECT

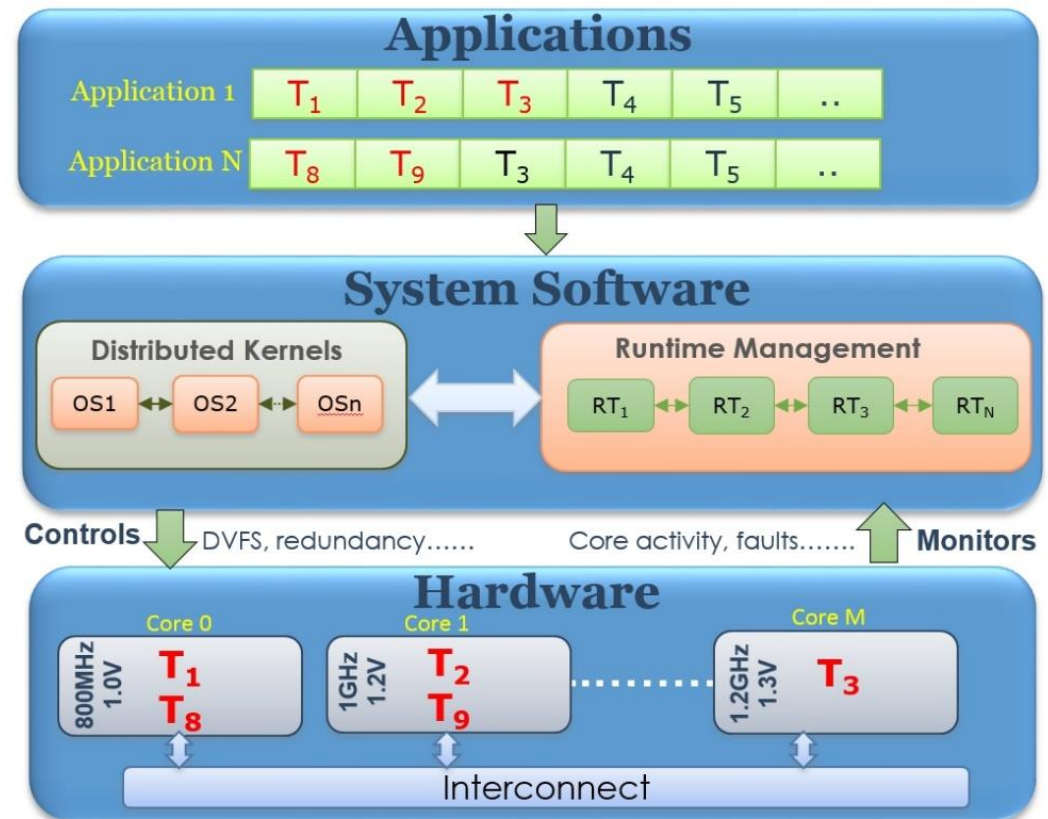
www.prime-project.org

Principles of energy-saving decision-making

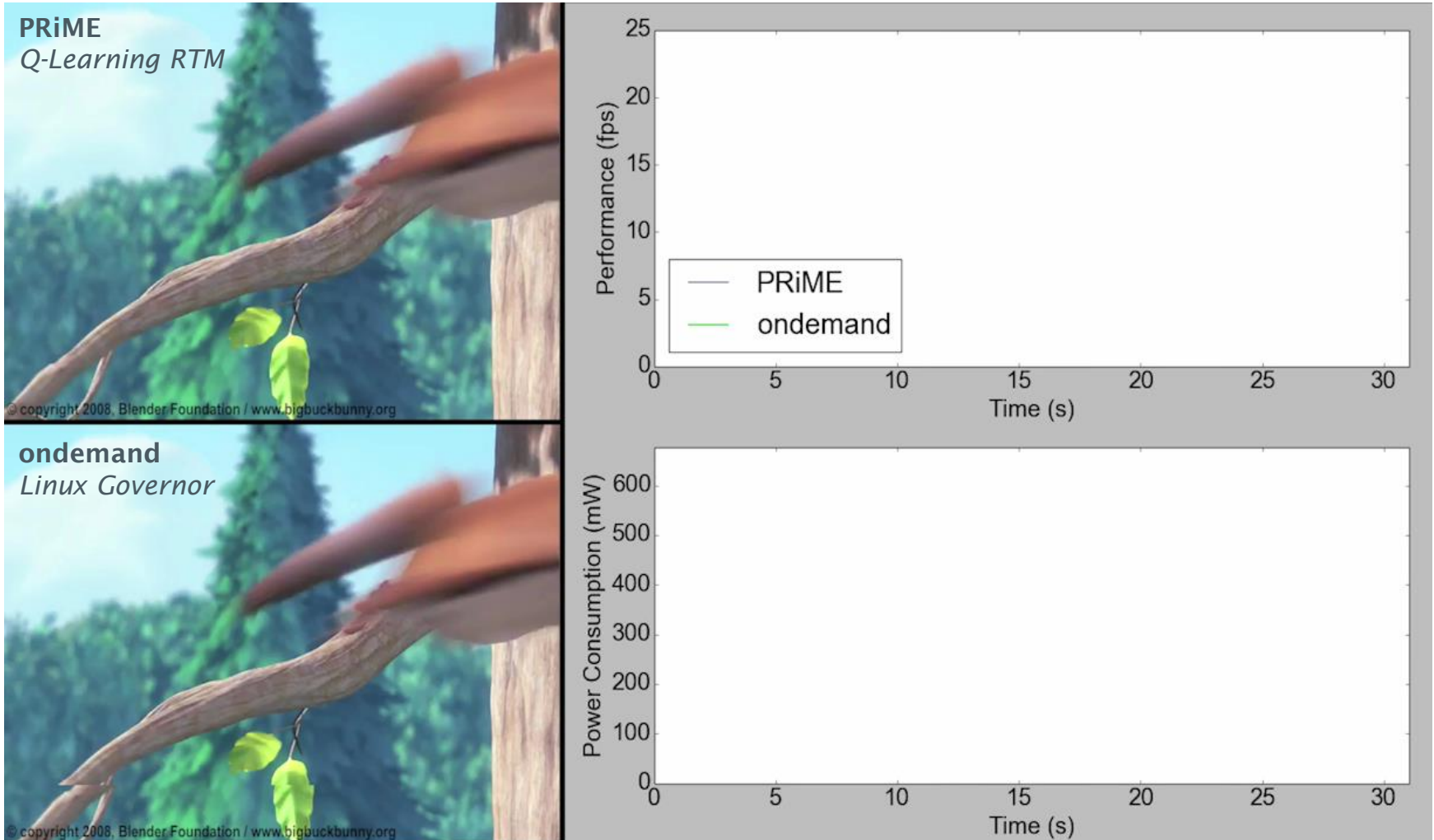
System software that makes run-time intelligent decisions to manage the system.

- **Responding:** Reacting to the *present*
- **Predicting:** Learning from the *past* informs the future

Variety of different algorithms explored



RUNTIME POWER MANAGEMENT

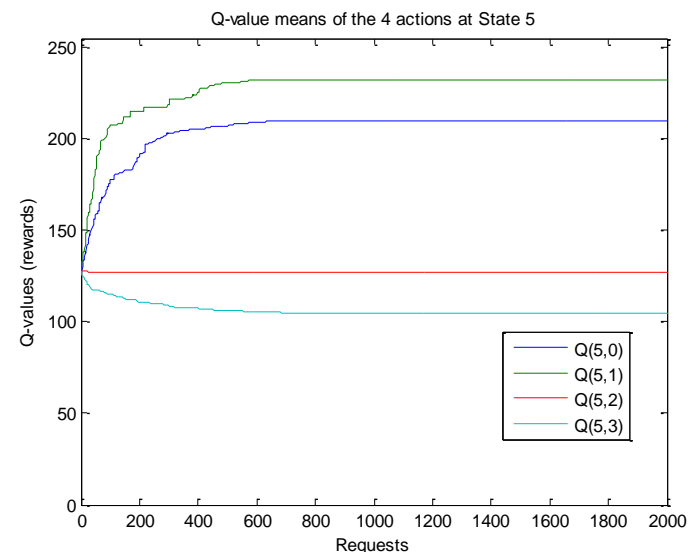
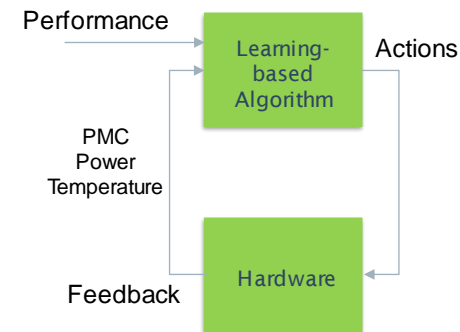


LEARNING OPTIMAL DVFS CHOICES

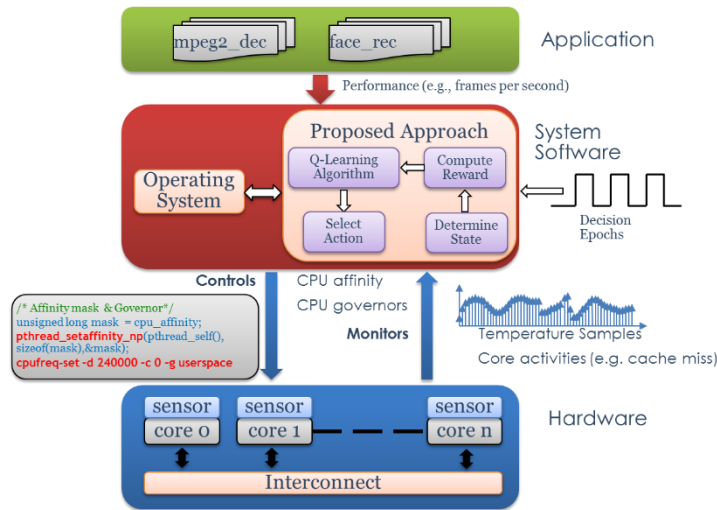
Reinforcement Learning

- Observes the system state (**workload**)
- Selects an action (**V-F pairs**)
- Changes the state (**performance**)
- Leads to a payoff (**reward/penalty**)

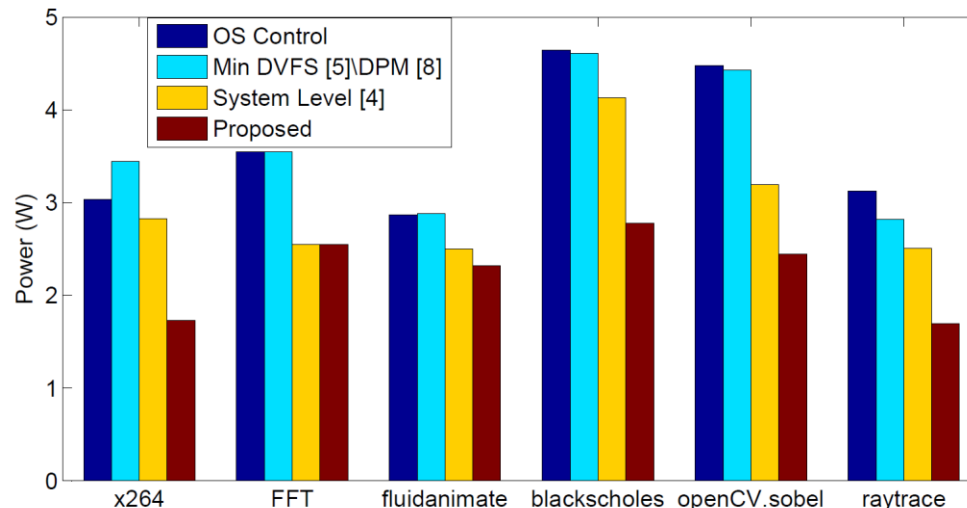
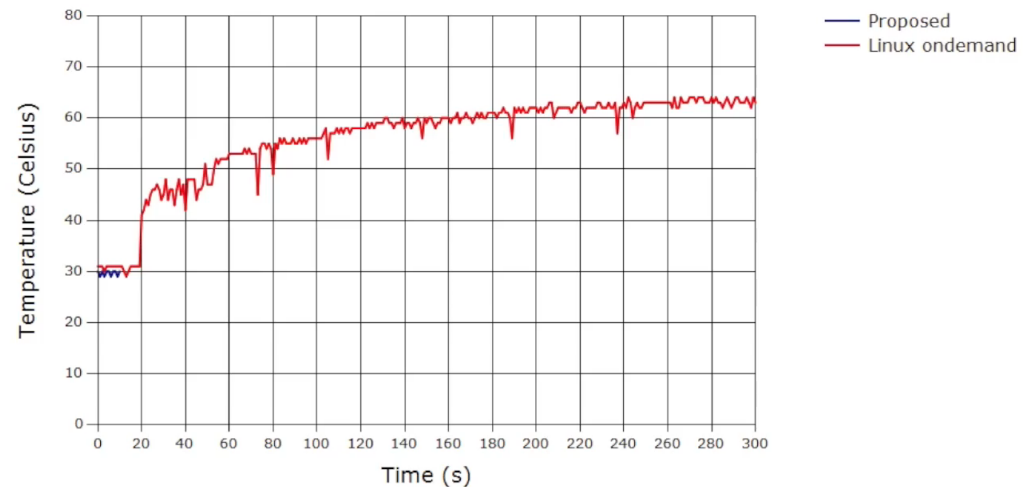
	ACTIONS (Power Modes)			
STATES (Tasks)	P0	P1	P2	P3
WD0	219	224	230	235
WD1	222	230	238	246
WD2	224	235	245	125
WD3	204	220	236	252
WD4	210	232	253	106
WD5	210	232	127	105



MANAGING THERMAL (LIFETIME) RELIABILITY



Convergence of the Reinforcement Learning Algorithm



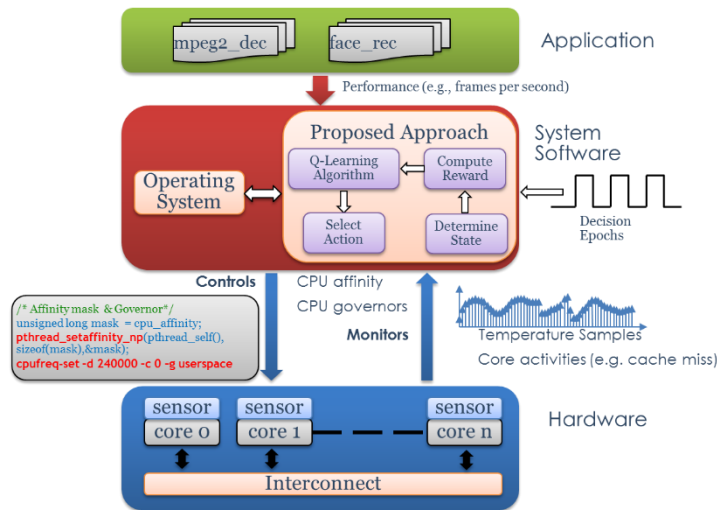
Techniques	Average Temperature	Peak Temperature
OS Controlled	76.6°C	82
System level [4]	69.9°C	79°C
Proposed	62.1°C	70°C

[4] Dhiman et al. "System-level power management using online Learning" in IEEE TCAD 2009

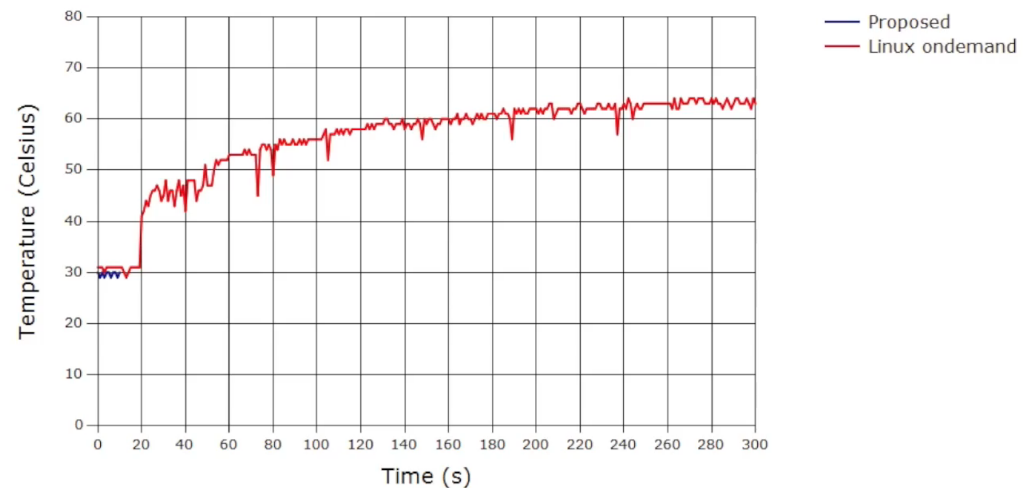
[5] Shen et al. "Achieving autonomous power management using reinforcement learning" in ACM TECS 2013

[8] Ye et al. "Learning-based power management for multicore processors via idle period manipulation" in TCAD 2014

MANAGING THERMAL (LIFETIME) RELIABILITY



Convergence of the Reinforcement Learning Algorithm



Application	Data Set	Average Temperature (Celcius)			Peak Temperature (Celcius)		
		Linux	Ge et al.	Proposed	Linux	Ge et al.	Proposed
tachyon	set 1	69.2	52.6	38.6	71.5	63	60
	set 2	50.5	44.5	43.8	57.3	56.3	52
	set 3	50.8	44.7	41.6	57.8	54.5	48.8
mpeg2_dec	clip 1	36	34	34.2	42.7	41.3	39
	clip 2	35.6	34.4	34.2	42.3	42	39.3
	clip 3	34.3	34.4	34	43	39.7	44.3

Average MTTF improvements: 5x (thermal aging); 4x (thermal cycling)

OVERVIEW

Applications

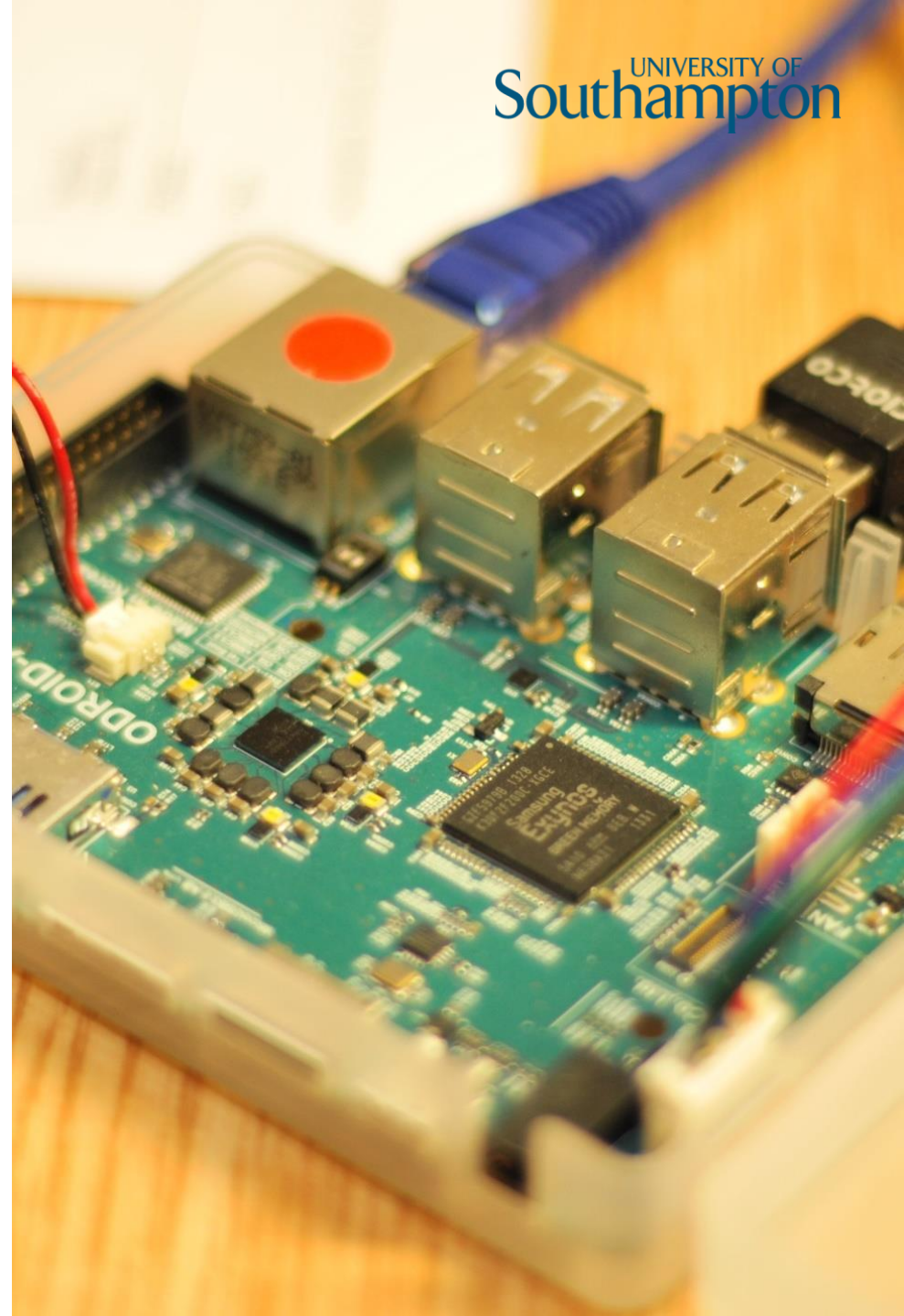
- From single > sequential > concurrent execution

Offline Characterisation

- Can we improve RTM through offline characterisation?

Towards Many-Core

- How do RTM approaches scale with number of cores?



RTMs and Application Workloads

From single > sequential > concurrent execution

QUALITY OF EXPERIENCE

- User cares about **observable performance**
 - Responsiveness, battery life, consistency, uninterrupted service
 - Doesn't really care about FLOPS, FPS, bandwidth, latency (QoS)
- Therefore, optimise for **quality of user experience** (QoE)
 - “*good-enough*” performance
 - Minimum energy usage

Bischoff, Alexander S. (2016) *User-experience-aware system optimisation for mobile systems*, University of Southampton, Electronics and Computer Science, Doctoral Thesis , 199pp.

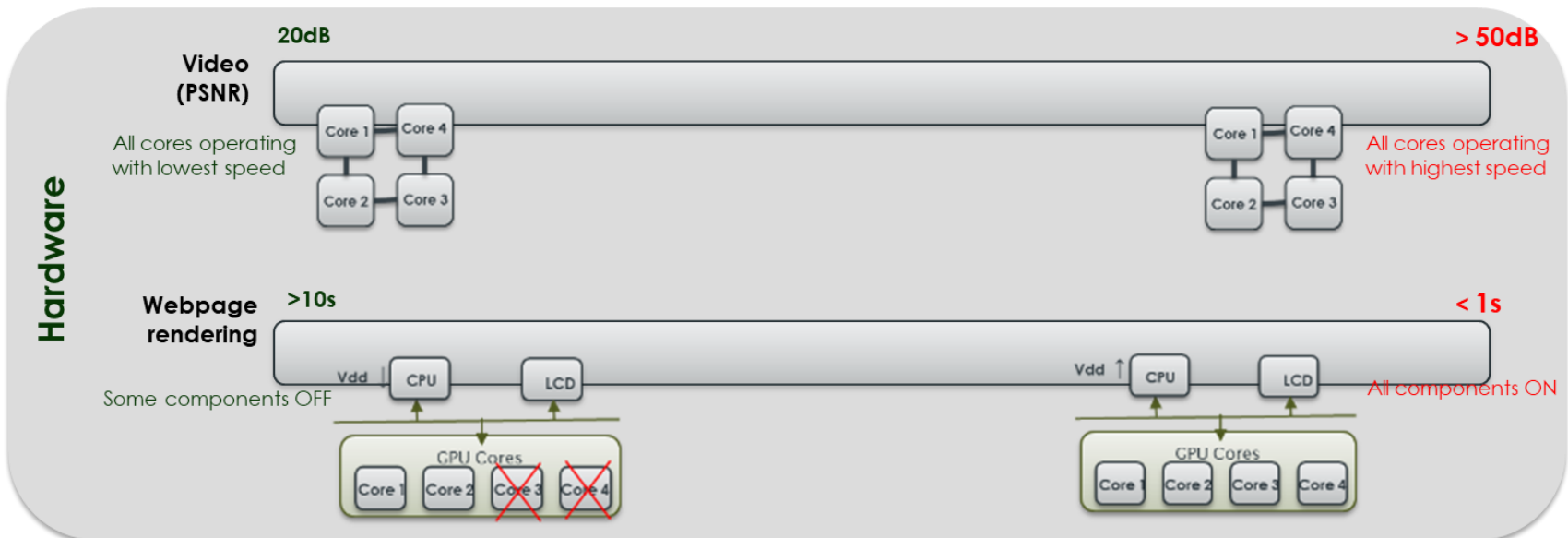
Bischoff S, Hansson A and Al-Hashimi BM. *Applying of Quality of Experience to System Optimisation*. International Workshop on Power and Timing Modeling, Optimization and Simulation (PATMOS), Germany, 2013.

QUALITY OF EXPERIENCE

Example Scenario



Runtime Management



QUALITY OF EXPERIENCE

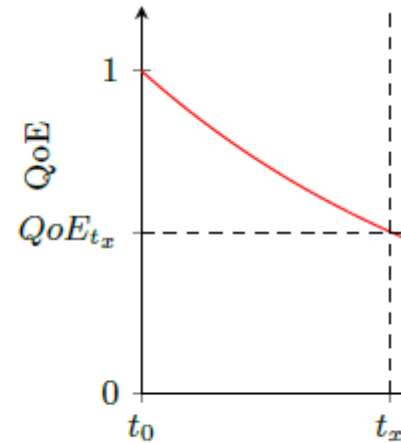
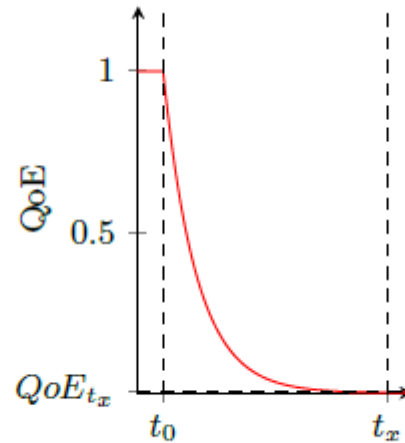
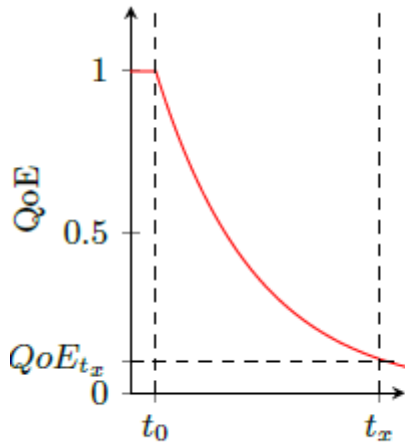
Workload Classification

Applications	Type of QoE
Audio	Throughput
Video	Throughput
Application Loading	Latency
Web Page Loading	Latency
Downloading a File	Latency
3D Gaming	Throughput
Word Processing	Latency

Delay of
<0.1 s appears instant,
1 s becomes noticeable
10 s become disruptive'

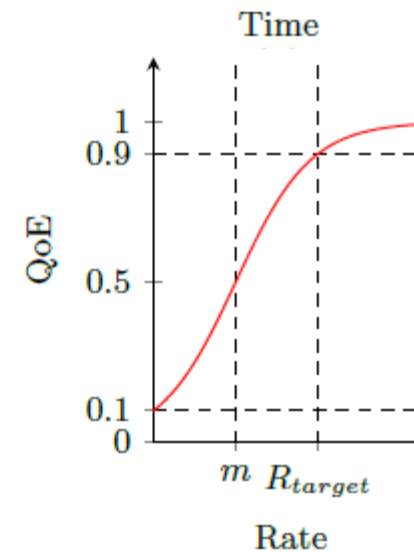
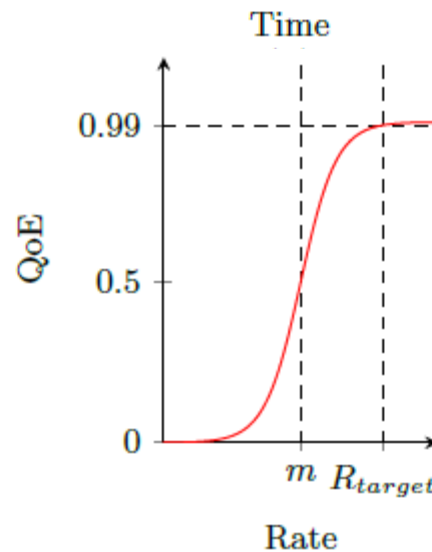
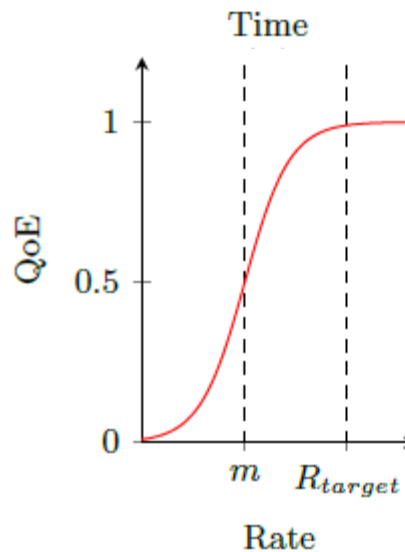
- Types of QoE:
 - Latency sensitive - complete workload in short time period
 - Throughput sensitive - complete at minimum rate

QoE CHARACTERISTICS



Latency
Sensitive

Inverse exponential



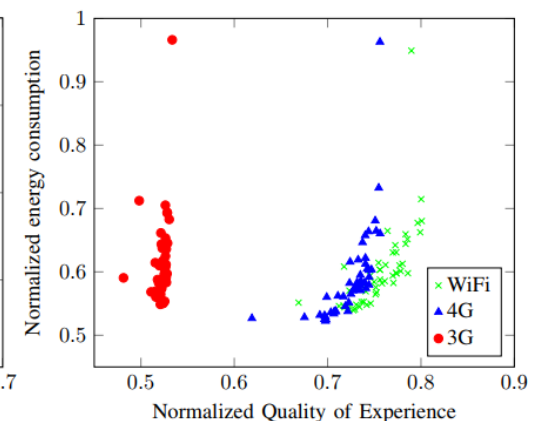
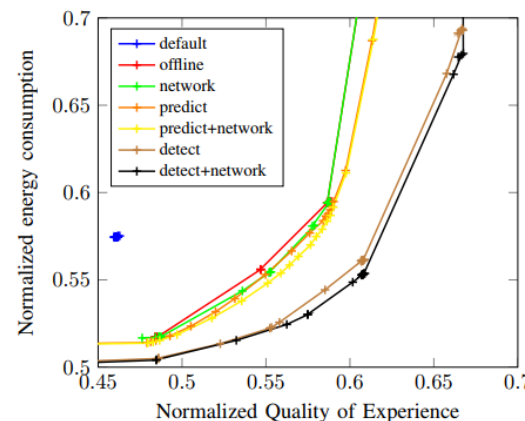
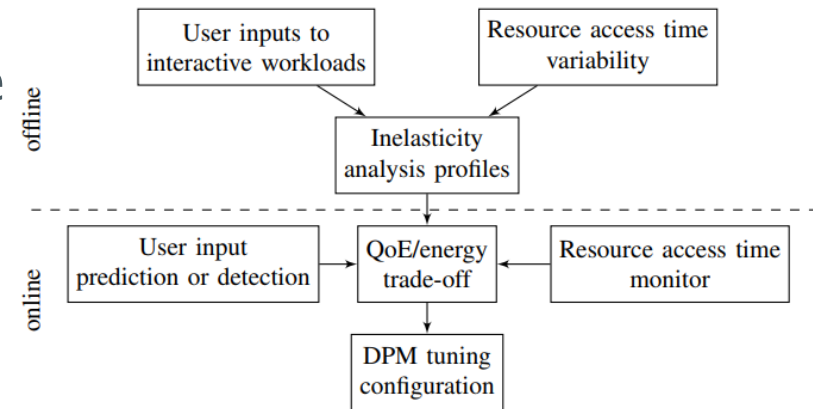
Throughput
Sensitive

Sigmoid function

TUNING DPM/RTM PARAMETERS

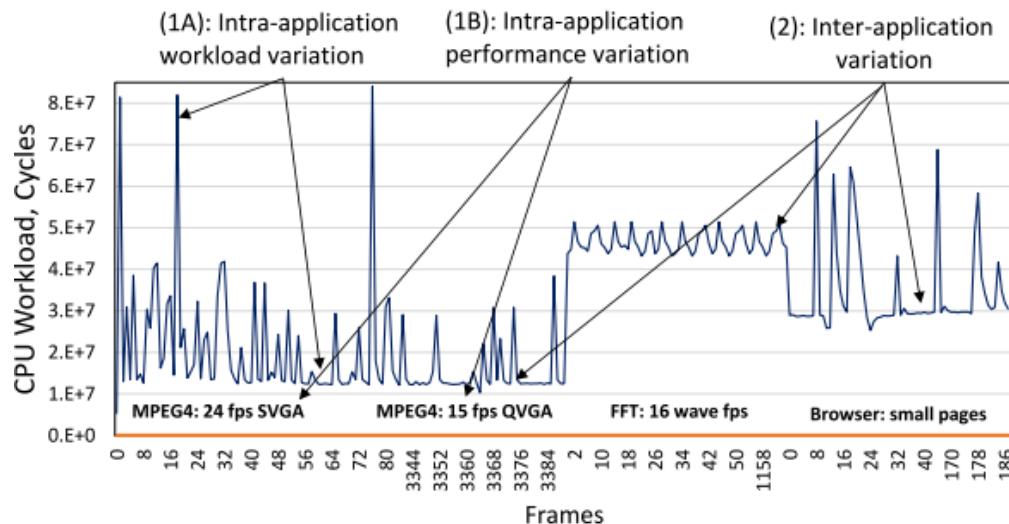
- Tune governor parameters for the executing (interactive) workload
- Account for variability in access times and user input
- Prediction/detection dependent
- Energy saving/QoE improvement compared to ‘default’, e.g.
 - 13% energy saving
 - 27% QoE improvement
 - 9% energy + 15% QoE

Exynos-5422 A15/A7, Android 6.0
Google Chrome browser workloads
Touch input emulation
Network throttling (UL, DL, RTT latency)



EXECUTING MULTIPLE APPLICATIONS

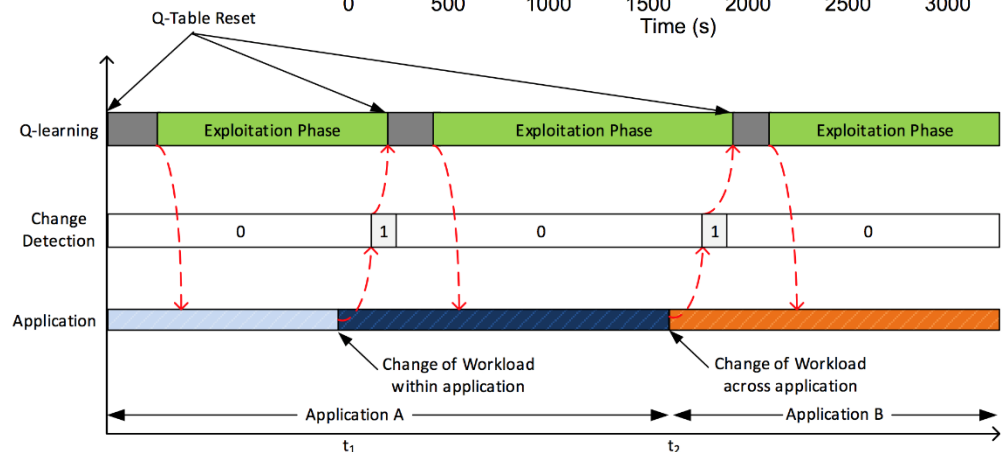
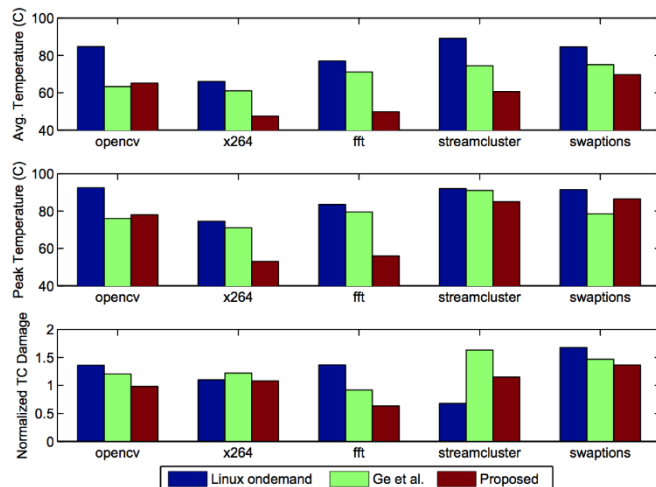
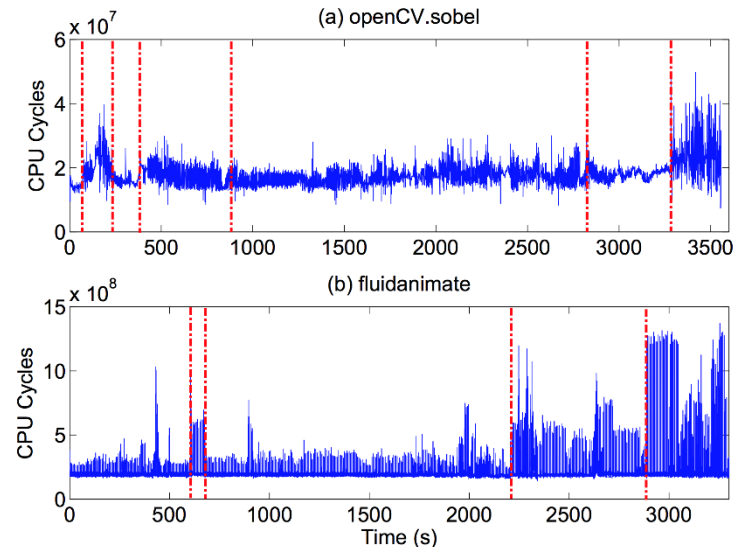
- Workload and performance variation due to:
 - Changes within an application
 - Changing applications (*sequential execution*)



- Overlapping applications (*concurrent execution*)

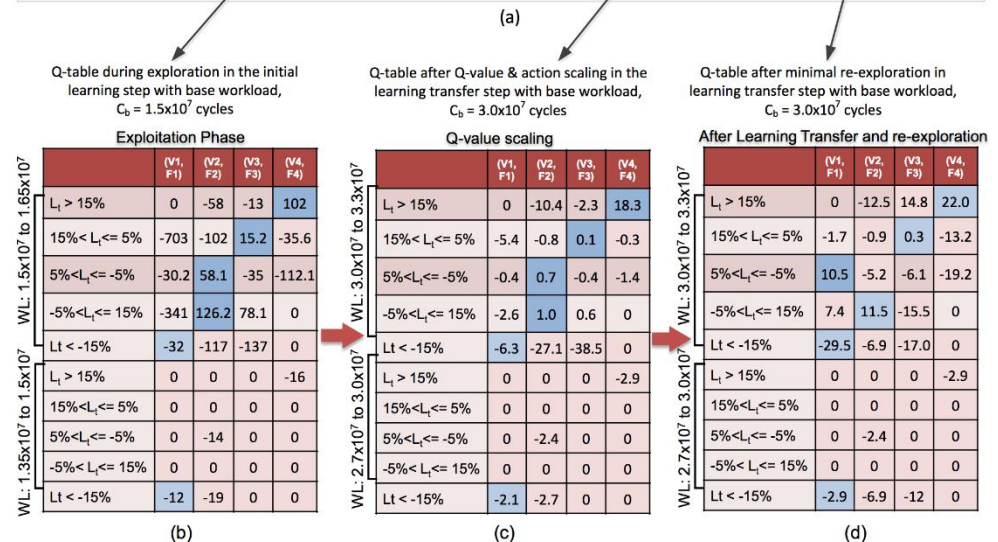
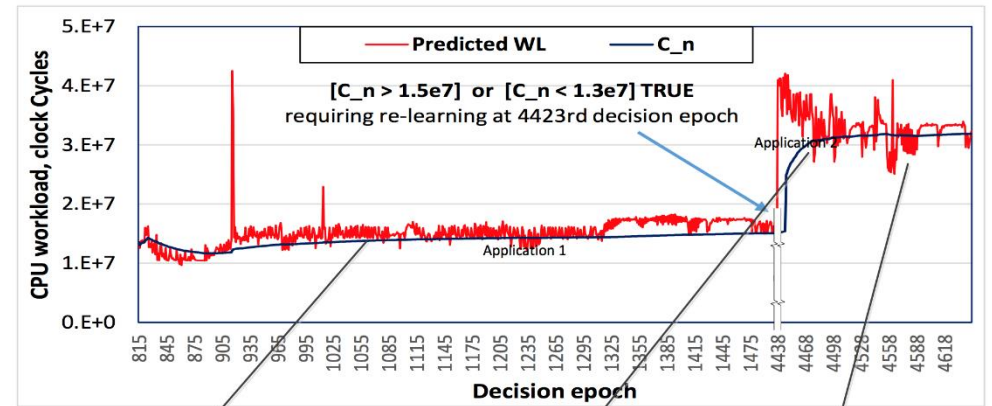
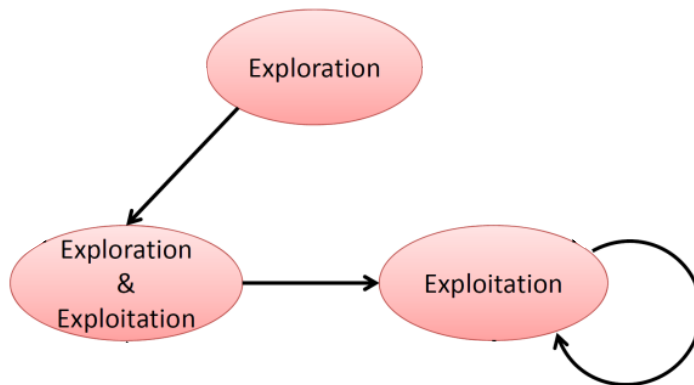
DETECTING WORKLOAD CHANGES

- Density ratio-based statistical divergence between overlapping sliding windows of CPU cycles
- Use this information to clear learning table (i.e. start afresh)



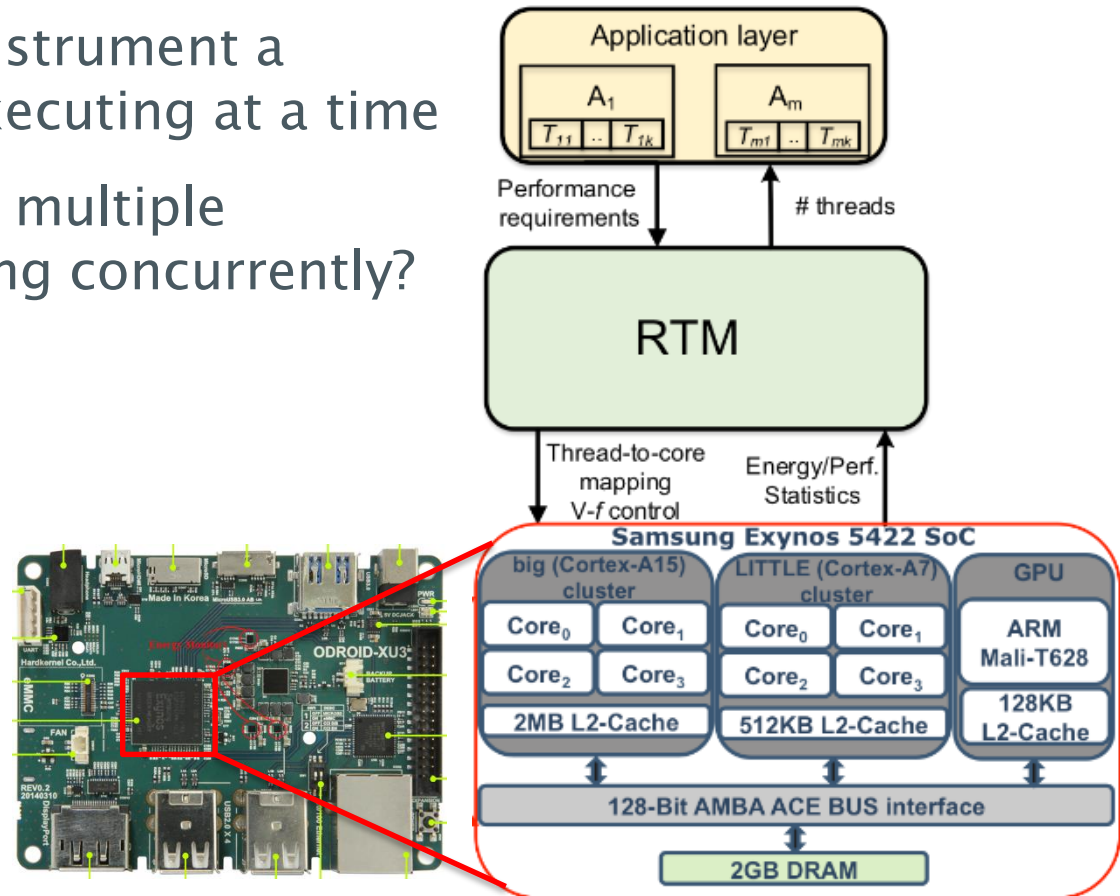
TRANSFERRING LEARNING

- Detect workload changes
- Transfer knowledge where possible
- Learn again fresh when not



RTM FOR CONCURRENT EXECUTION

- Approaches so far instrument a single application executing at a time
- How can we manage multiple applications executing concurrently?

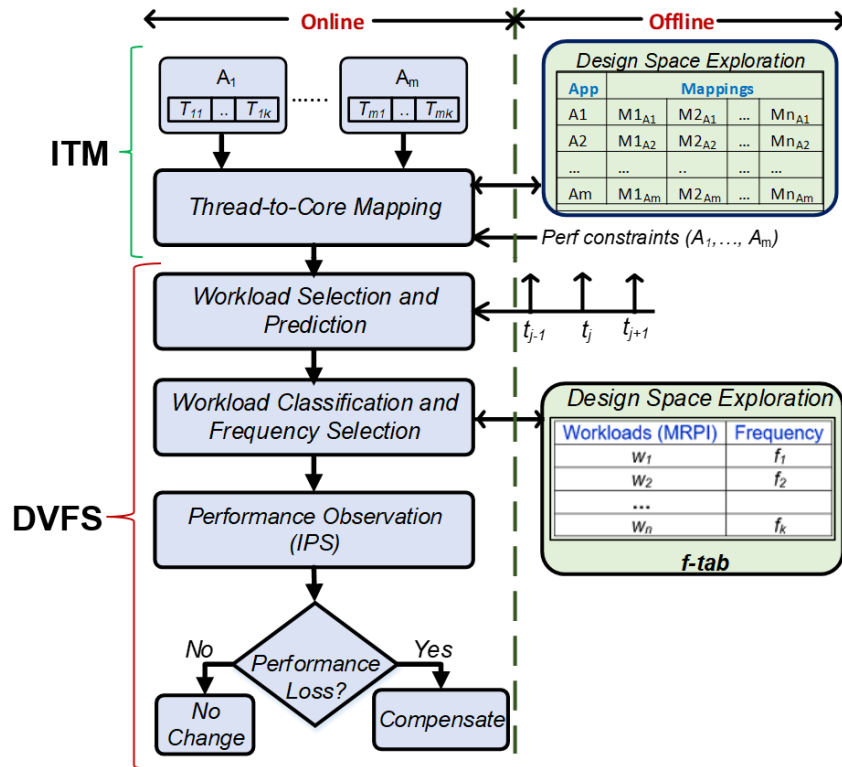


Online vs Offline

Can we improve RTM through offline characterisation?

MANAGING CONCURRENT EXECUTION

MRPI (Memory Reads Per Instruction)



- Supports concurrent execution of applications
- Inter-cluster Thread-to-core Mapping (ITM).
- MRPI informs DVFS control

Up to 33% reduction in energy compared to SoA while meeting performance requirements

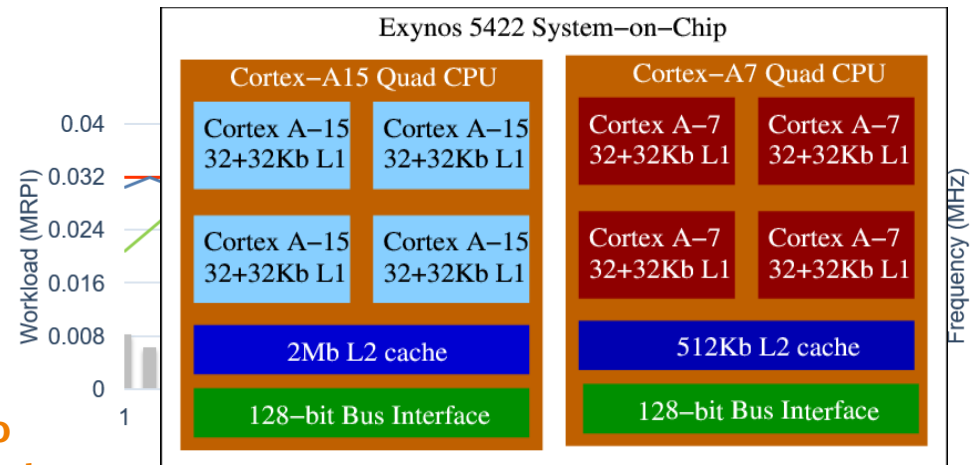
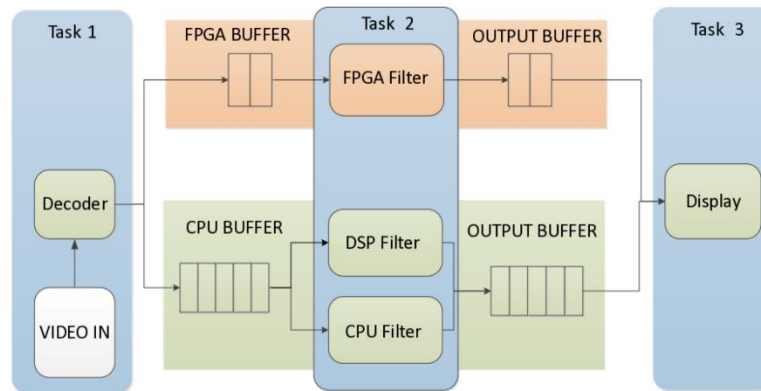


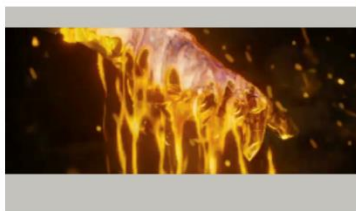
Image: Catalán et al., Performance and Energy Optimization of Matrix Multiplication on Asymmetric big.LITTLE Processors, 2015

REGRESSION-BASED RTM

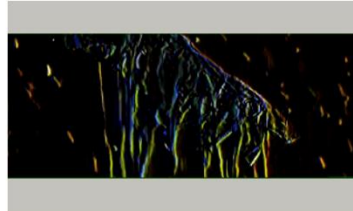
Heterogeneous Platforms



(a) Convolution filter implementation



(b) Original



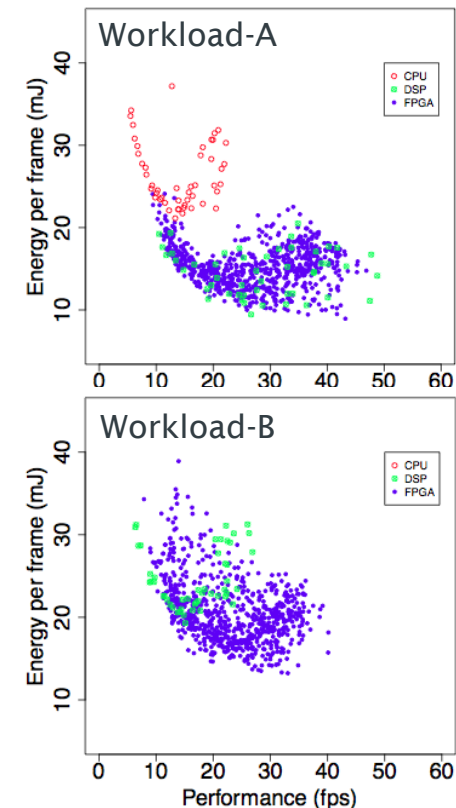
(c) Edge detected



(d) Blurred

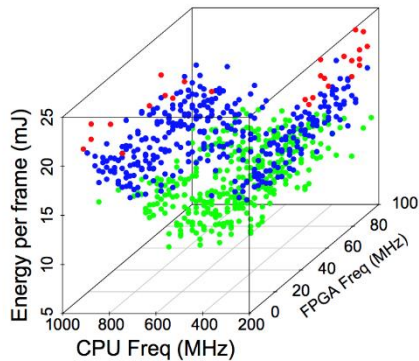
Run-time changes in:

- Performance requirements
- Application workload changes

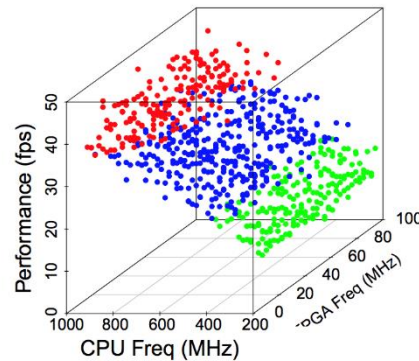


REGRESSION-BASED RTM

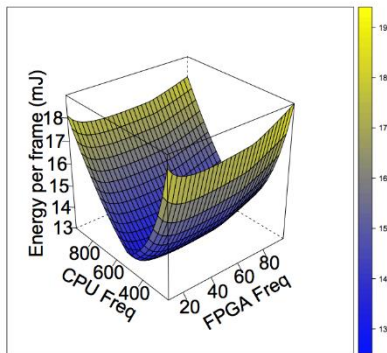
Heterogeneous Platforms



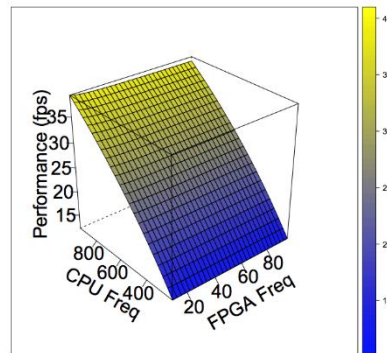
(a) FPGA measured energy



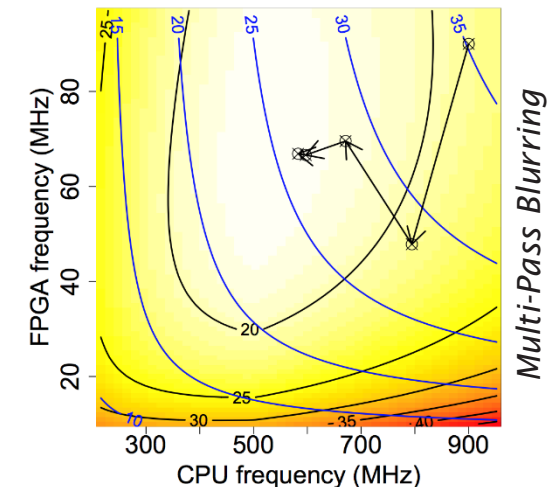
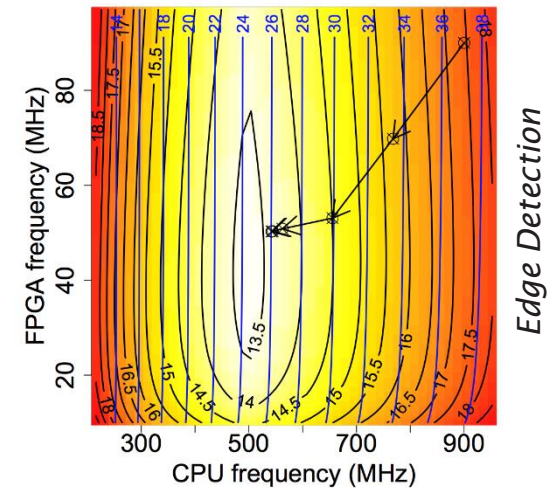
(b) FPGA measured performance

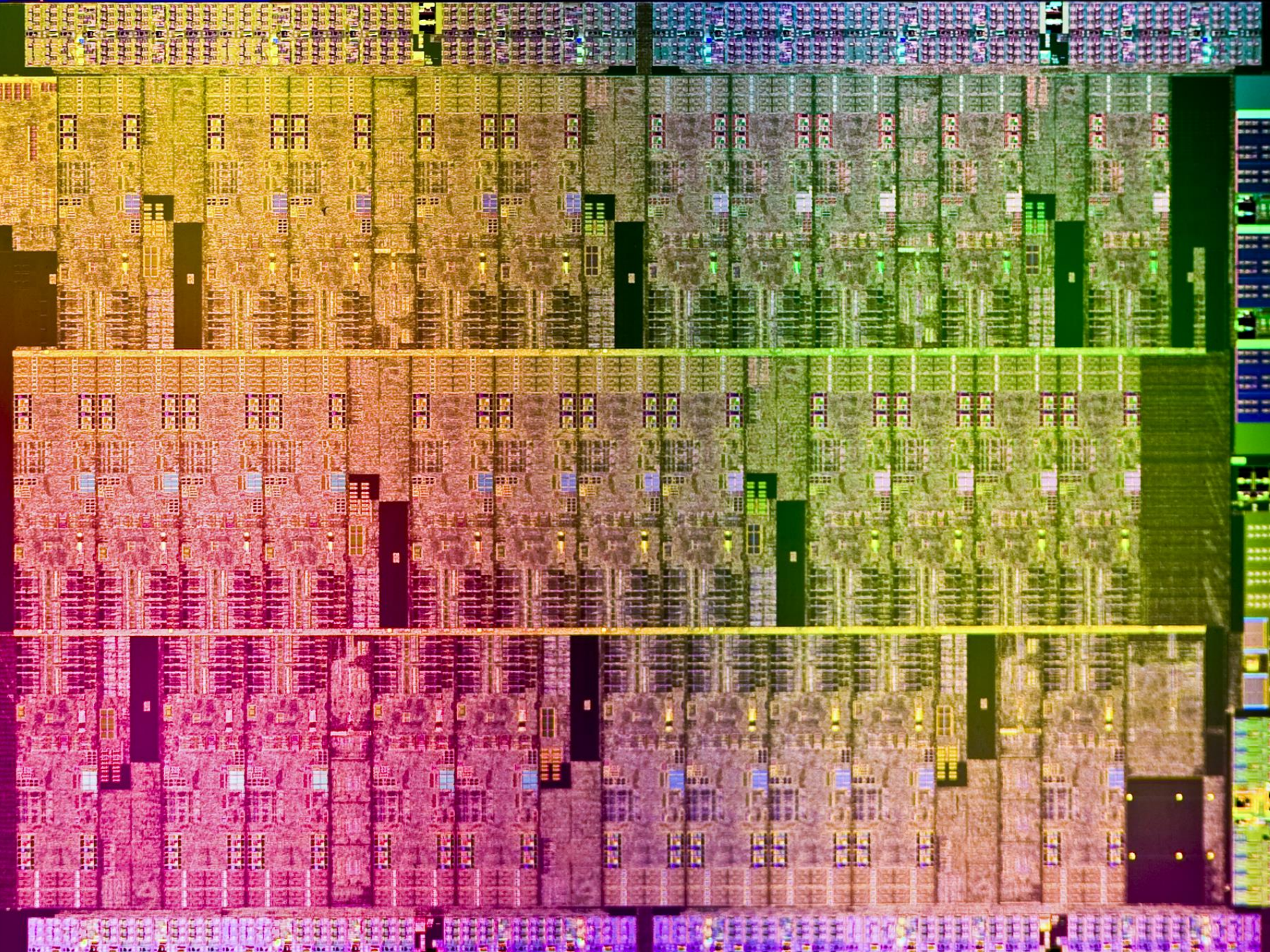


(c) FPGA modeled energy



(d) FPGA modeled performance



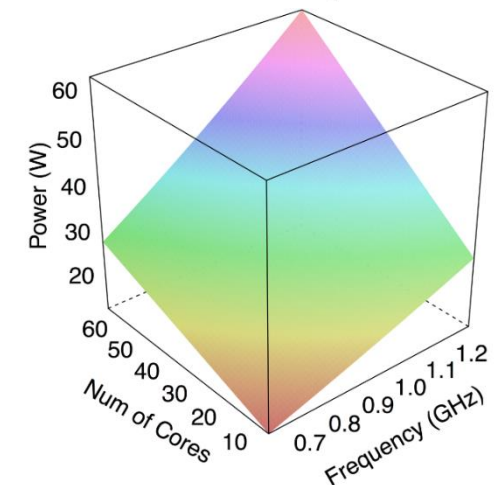
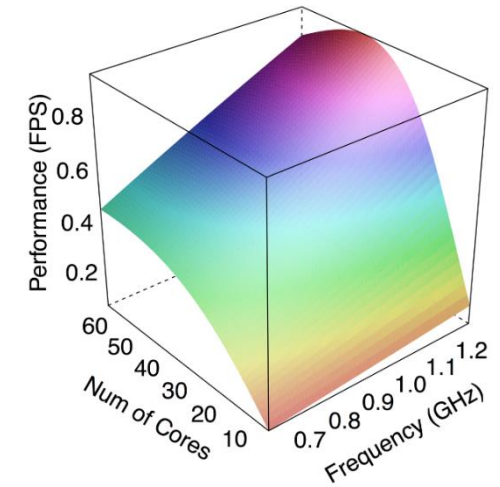
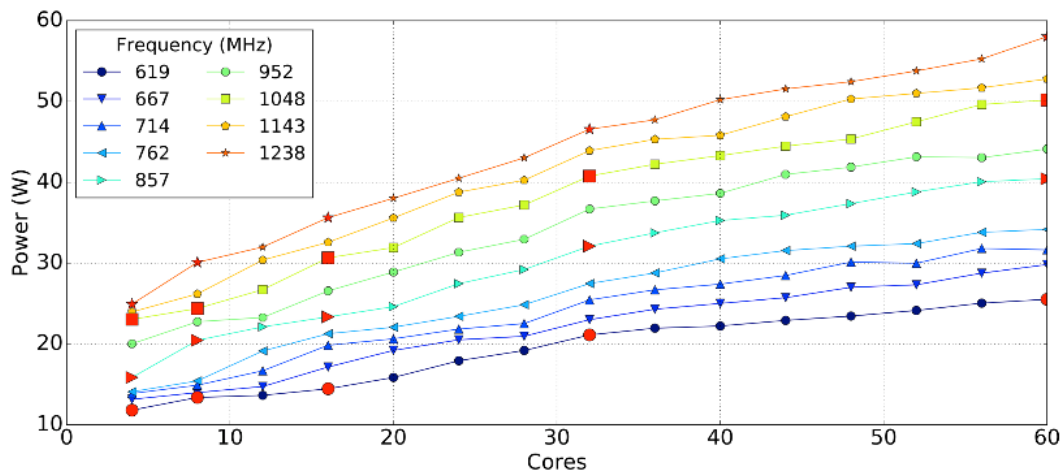
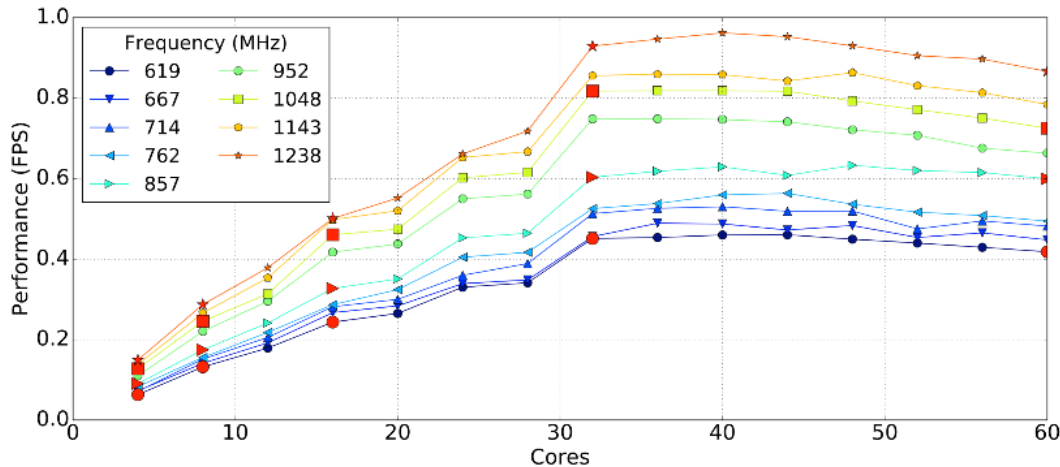


Towards Many-Core

How do RTM approaches scale
with number of cores?

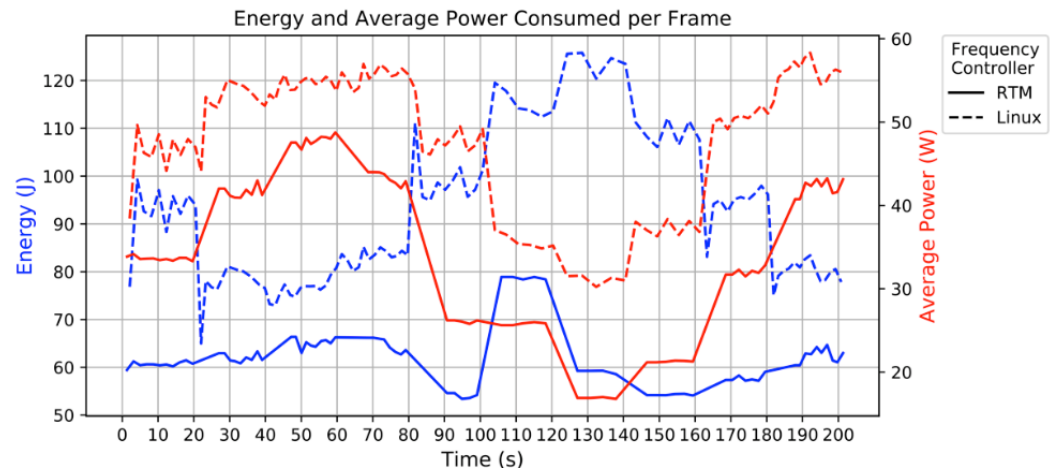
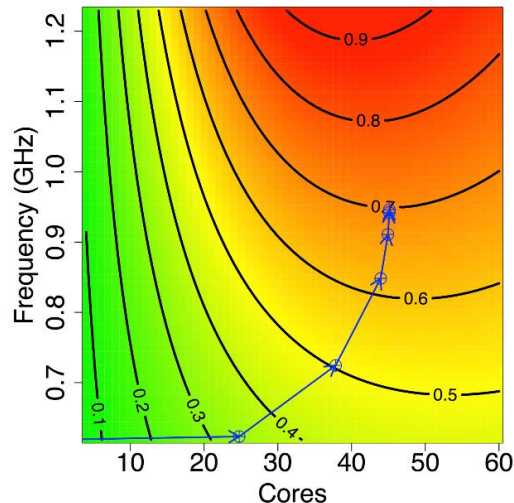
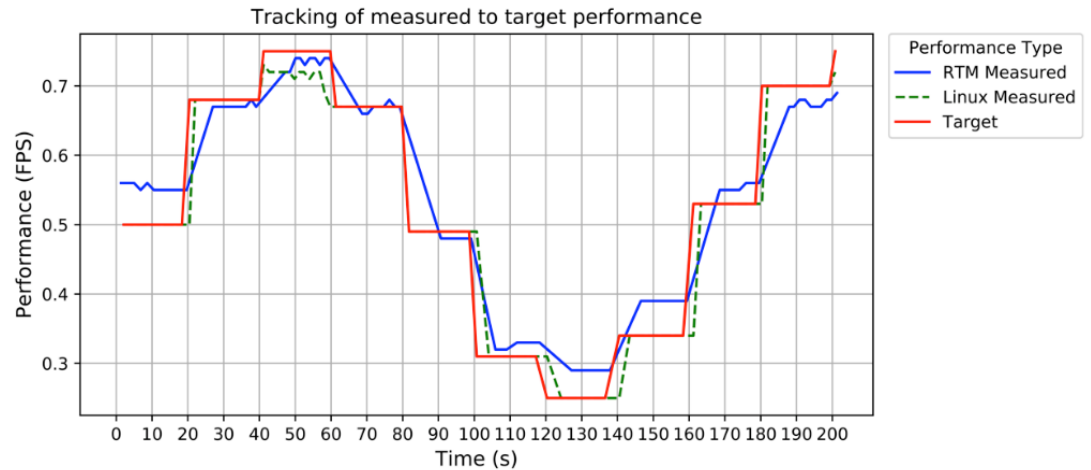
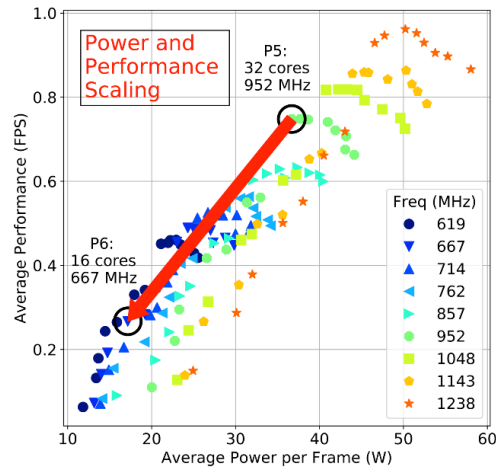
REGRESSION-BASED RTM

Model Building



REGRESSION-BASED RTM

Runtime Management



CONCURRENT RTM ON HPC SYSTEMS

- Applications targeted for HPC are usually multi-threaded
- Modern HPC often based on Non-Uniform Memory Access (NUMA) architecture
- Our Approach:
 - Platform characterized offline
 - Workload estimated based on memory-intensity, thread synchronization contention, NUMA latency
 - V - f determined using binning, while accounting for contention due to concurrent execution

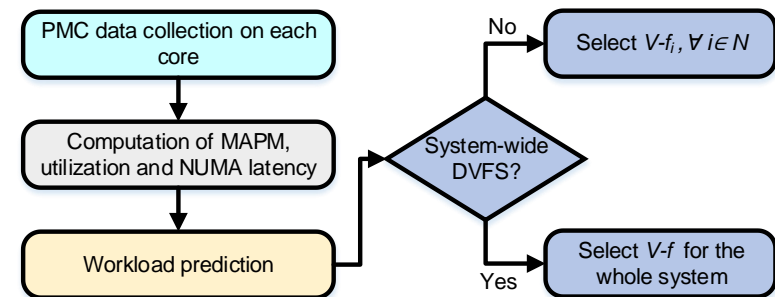
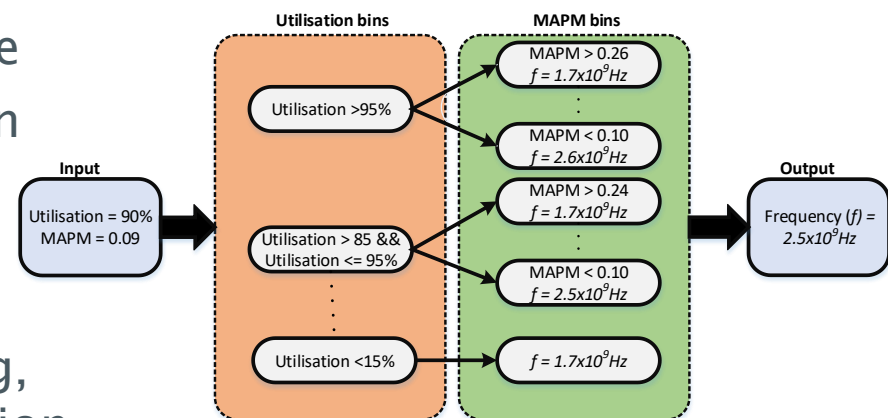
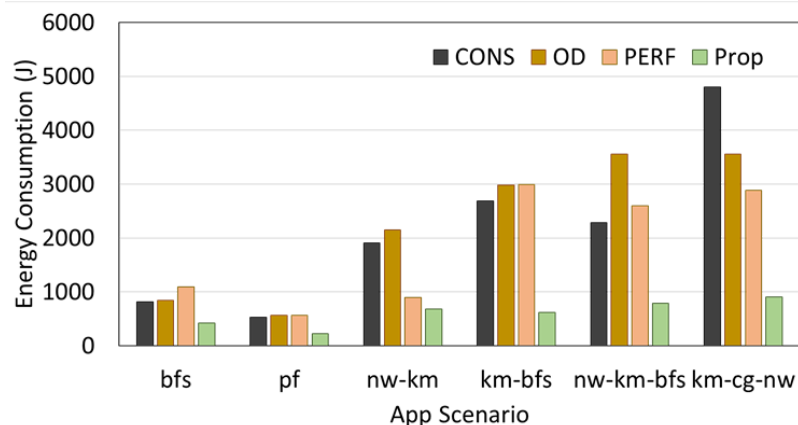


Illustration of various steps in the proposed approach

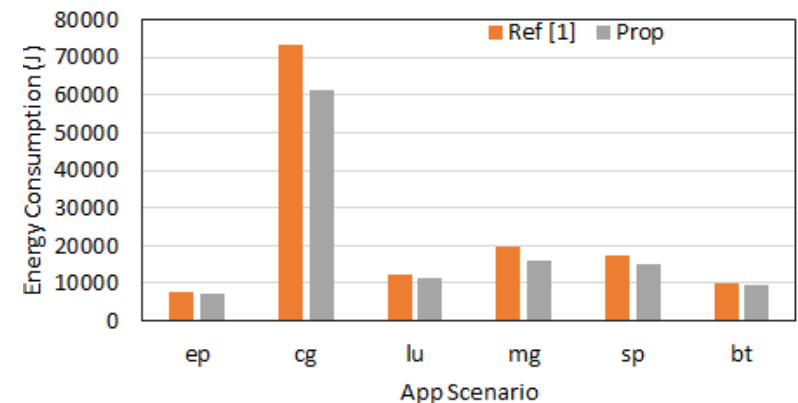


An example of V - f setting selection using binning-based approach

CONCURRENT RTM ON HPC SYSTEMS



Energy consumption of different approaches



Comparison of presented approach with Sundriyal *et al*

- Xeon E5-2630 (12 cores, 24 threads) and Xeon Phi 7620P (61 cores, 244 threads); NAS and Rodinia benchmarks
- Proposed (Prop) approach achieves energy savings of up to 81% (Xeon) and 61% (Phi) compared to Linux's governors
- Outperforms Sundriyal *et al.* by 10% in energy efficiency and 3.7% in performance

OPEN SOURCE TOOLS

MEASURING/MODELLING POWER

www.powmon.ecs.soton.ac.uk and www.gemstone.ecs.soton.ac.uk

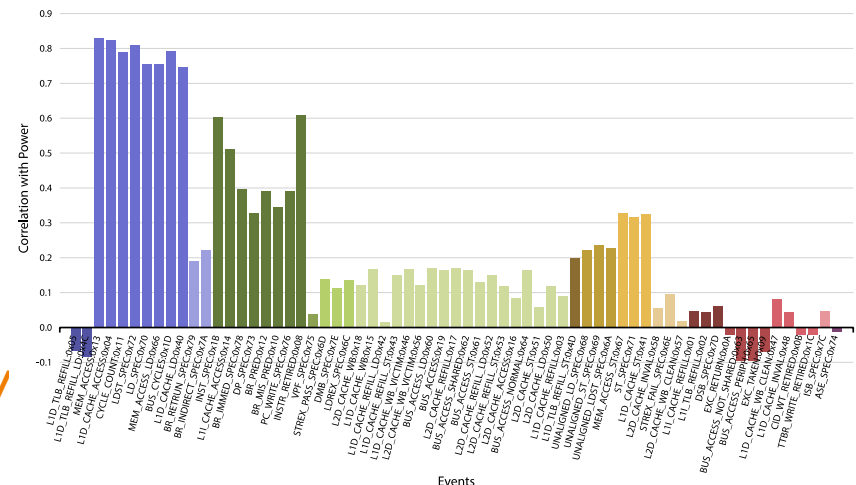
Why Power Estimation?

- Few platforms allow hardware measurement of power consumption
- RTMs need to make decisions based on real-time ‘measurements’
- Offline DPM strategy evaluation/design-space exploration

PowMon: A *Stable*, “Top-Down” Approach to Power Modelling

PMCs (Performance Monitoring Counters)

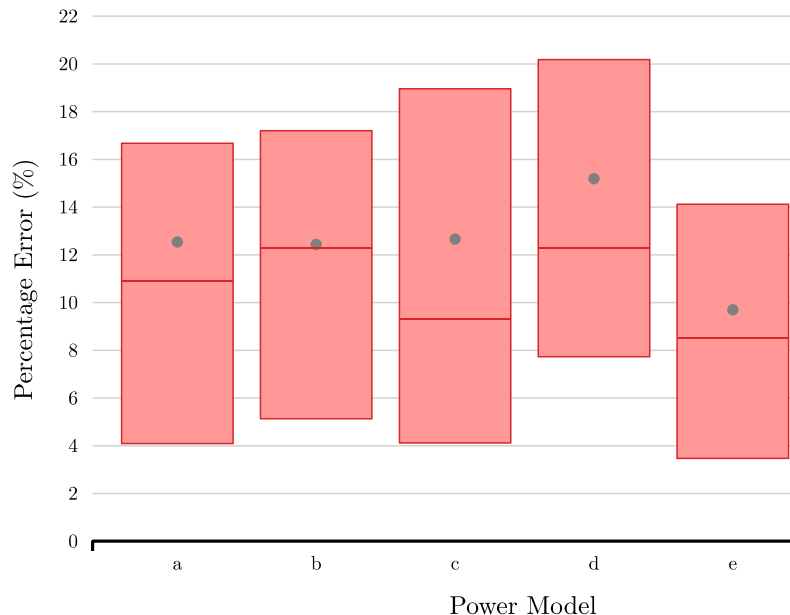
- On several platforms, low overhead, many events available
- ...but a small number (e.g. 4-6) can be monitored **simultaneously**



POWMON: STABLE POWER MODELLING

www.powmon.ecs.soton.ac.uk and www.gemstone.ecs.soton.ac.uk

Our **stable** approach achieves a low average error and narrow error distribution compared to existing techniques.



Training: Small set of 20 workloads

Testing: Full set of 60 workloads

- [a] M. Pricopi, T. S. Muthukaruppan, V. Venkataramani, T. Mitra, and S. Vishin, "Power-performance modeling on asymmetric multi-cores," CASES '13.
- [b] M. Walker et al., "Run-time power estimation for mobile and embedded asymmetric multi-core cpus," HIPEAC Workshop Energy Efficiency with Hetero. Comp. 2015
- [c] S. K. Rethinagiri et al., "System-level power estimation tool for embedded processor based platforms," RAPIDO '14. New York, 2014.
- [d], [e] R. Rodrigues et al, "A study on the use of performance counters to estimate power in microprocessors," IEEE TCAS II, vol. 60, no. 12, pp. 882-886, Dec 2013.

M. J. Walker et al., "Accurate and Stable Run-Time Power Modeling for Mobile and Embedded CPUs," in *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 36, no. 1, pp. 106-119, Jan. 2017.

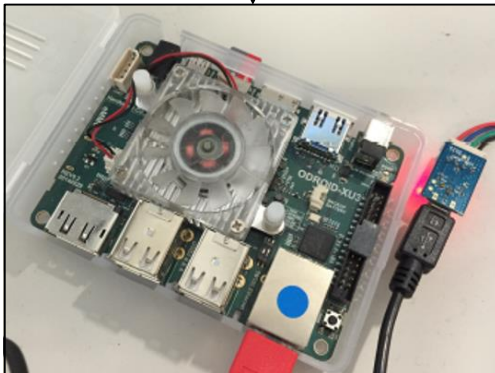
POWMON: METHODOLOGY

www.powmon.ecs.soton.ac.uk

1. Run workloads

@ different DVFS levels

39 workloads used: MiBench, LMBench, Roy Longbottom, ParMiBench and ALPBench

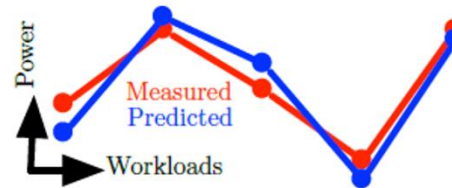


ODROID-XU3

Exynos-5422
4x Cortex-A7
4x Cortex-A5

5. Validate

- K-fold cross validation
- $R^2 : > 0.99$
- Error: 2.8 – 3.7%



4. Build Model

- OLS multiple regression
- Considers collinearity and heteroscedasticity
- “sensible” equation

2. Record

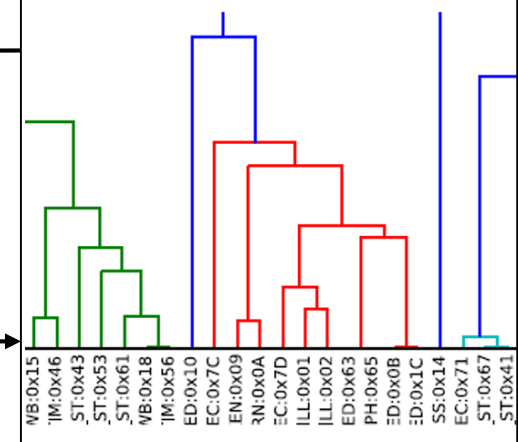
- PMCs
- Power, Voltage, Temperature, etc.

6. Uses

- OS Run-time management
- Reference for research
- gem5 add-on

3. Choose PMCs

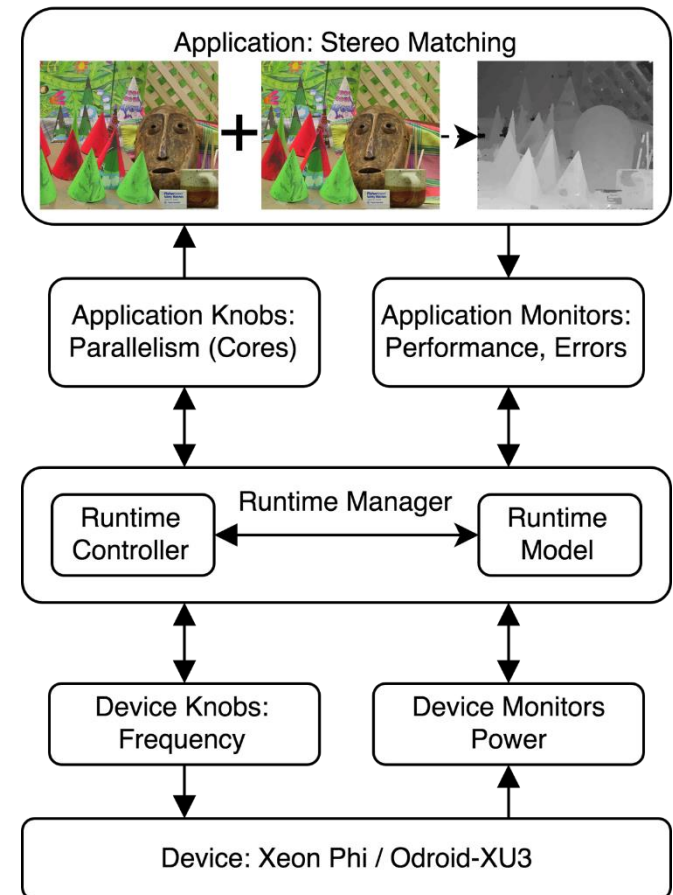
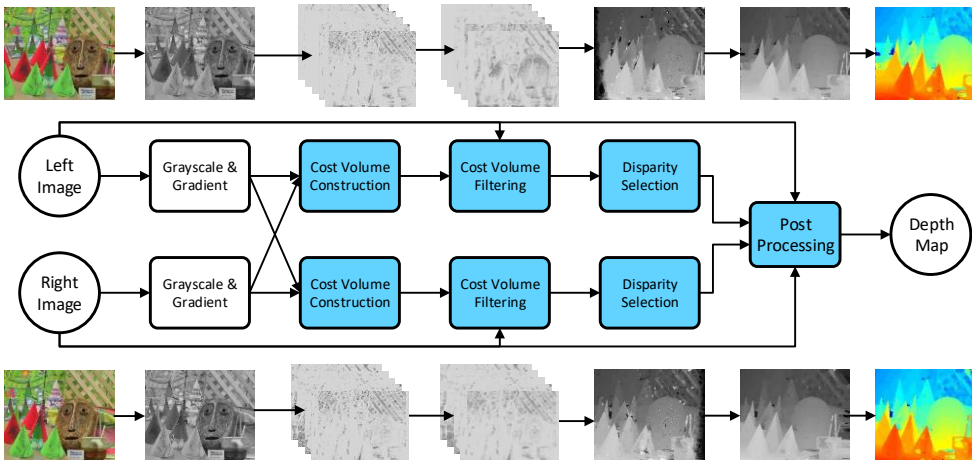
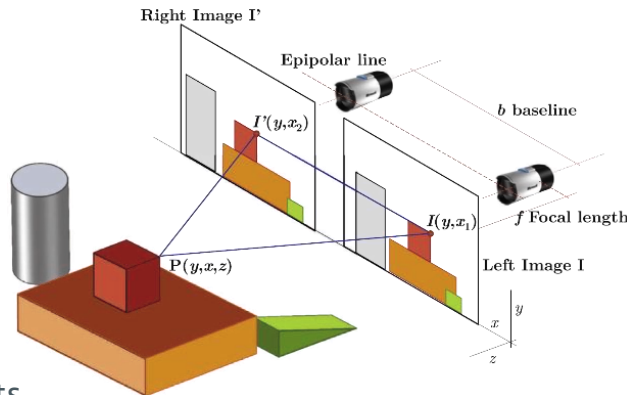
Hierarchical cluster analysis,
Correlation matrix analysis,
Exhaustive search, etc.



STEREO MATCHING APPLICATION

<http://github.com/PRiME-project/PRiMEStereoMatch>

- Processes still images, video or a camera feed
- OpenCL supported
- Includes test datasets



IMPROVING USABILITY/COMPARISON

The PRiME RTM Framework

App-RTM interaction

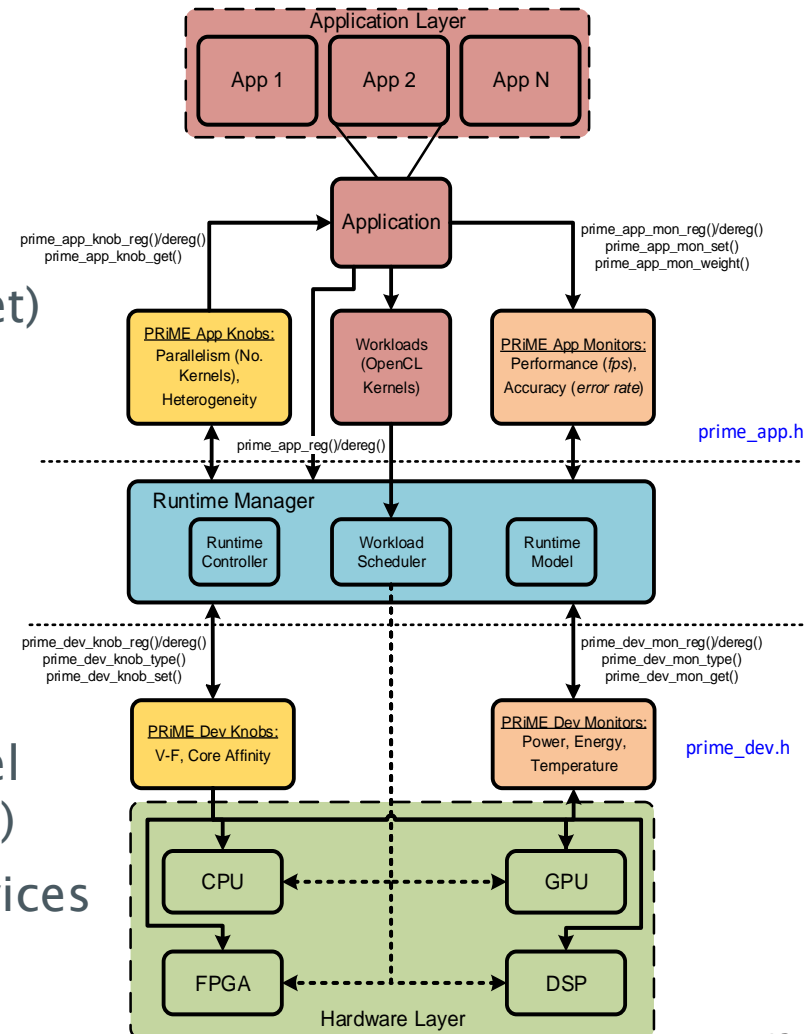
- Communication of controls/monitors (e.g. register/deregister)
- Update controls/monitors (e.g. set/get)

RTM-Device interaction

- Communication of controls/monitors (e.g. register/deregister)
- Update controls/monitors

OpenCL Workload Management

- Tuneable parameters (e.g. thread level parallelism – kernels, device selection)
- Scheduling of kernels to compute devices



THE PRiME RTM FRAMEWORK

Experimental Evaluation

Application-RTM interaction

- 3 Applications

Application	Name	Const.	Space	Allowed/target values
Jacobi	Iterations	knob	disc	$\mathbb{N} \in [1, \infty)$
	Data type	knob	disc	{float, double}
	Device type	knob	disc	{CPU, GPU/FPGA}
	Throughput	mon	cont	$\mathbb{R} \in [10, \infty)$
	Error	mon	cont	$\mathbb{R} \in (-\infty, 1e^{-12}]$
Video decoder	Throughput	mon	cont	$\mathbb{R} \in [25, \infty)$
Whetstone	Threads	knob	disc	$\mathbb{N} \in [1, \infty)$
	Throughput	mon	cont	$\mathbb{R} \in [2.5, \infty)$

Run-Time Managers (RTMs)

- Q-Learning
- MRPI
- Exhaustive characterisation

RTM-Device interaction

- 2 Platforms

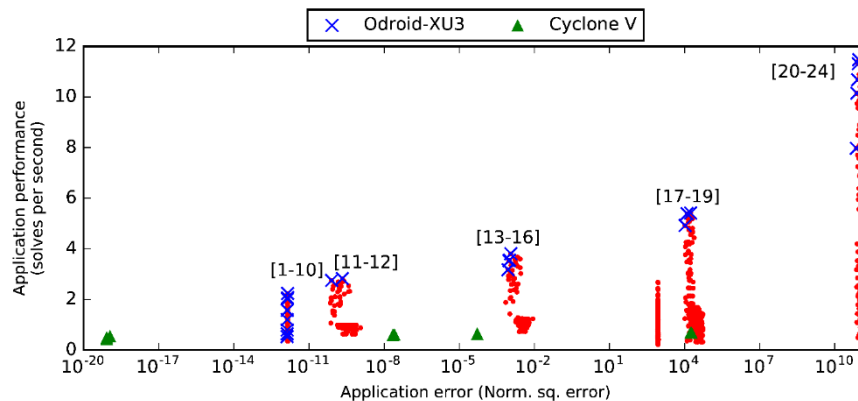
Plat.	Const.	Space	Type	For	No.
Odroid-XU3	knob	disc	GOVERNER	A7 cluster	1
		disc	GOVERNER	A15 cluster	1
		disc	FREQ	A7 cluster	1
		disc	FREQ	A15 cluster	1
		disc	FREQ_EN	GPU DVFS	1
		disc	FREQ	GPU	1
		disc	PMC_CTRL	A7 cores	16
		disc	PMC_CTRL	A15 cores	24
	mon	cont	POW	Clusters, RAM, GPU, SoC	5
		cont	TEMP	A15 cores	4
		cont	TEMP	GPU	1
		disc	CYCLE	A7 cores	4
		disc	CYCLE	A15 cores	4
		disc	PMC	A7 cores	16
Cyclone V	knob	disc	PMC	A15 cores	24
		disc	PMC	A15 cores	24
	mon	cont	VOLT	A9 cluster, peripherals	4
		cont	VOLT	FPGA, peripherals	3
		cont	POW	A9 cluster, peripherals	5
		cont	POW	FPGA, peripherals	4
		cont	POW	SoC	1

PRiME RTM FRAMEWORK

Experimental Evaluation: Exhaustive Characterisation

Characterisation (DSE)

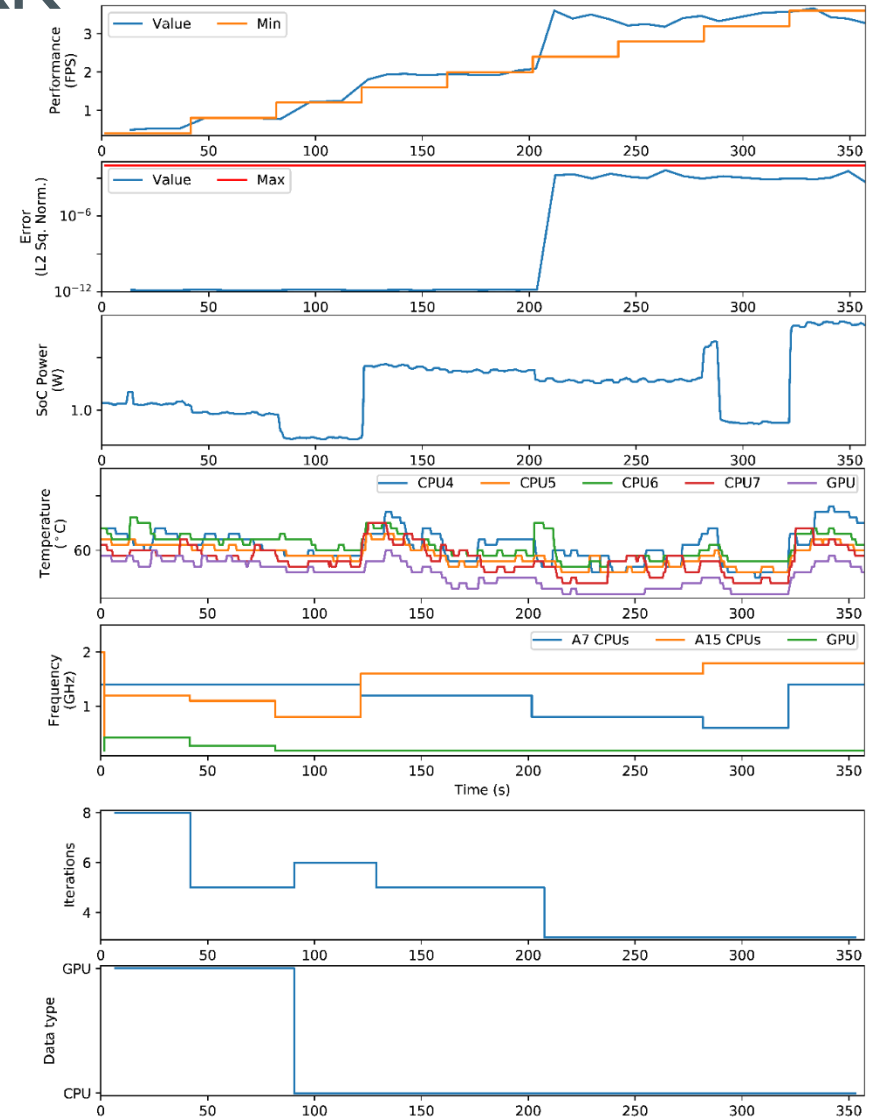
- Jacobi application
- 2 platforms



Run-Time Management

- User requirements (FPS/error)
- RTM excludes points
- Optimises remainder for power

<http://www.prime-project.org/>



Many-Core Computing

Hardware and software

Computing has moved away from a focus on performance-centric serial computation, instead towards energy-efficient parallel computation. This provides continued performance increases without increasing clock frequencies, and overcomes the thermal and power limitations of the dark-silicon era. As the number of parallel cores increases, we transition into the many-core computing era. There is considerable interest in developing methods, tools, architectures and applications to support many-core computing.

The primary aim of this edited book is to provide a timely and coherent account of the recent advances in many-core computing research. Starting with programming models, operating systems and their applications; the authors present runtime management techniques, followed by system modelling, verification and testing methods, and architectures and systems. The book ends with some examples of innovative applications.

About the Editors

Bashir M. Al-Hashimi is an international leader in the theory and practice of energy-efficient computing, undertaking fundamental and experimental research on hardware-software co-design for energy efficiency and hardware reliability. He is the Dean of the Faculty of Engineering and Physical Sciences, University of Southampton, UK, and founder and co-director of the Arm-ECS Research Centre. He was appointed Commander of the Order of the British Empire (CBE) in the 2018 Queen's Birthday Honours for services to computer engineering and to industry, elected Fellow of the Royal Academy of Engineering in 2013, and Fellow of the Institute of Electrical and Electronics Engineers in 2009. He has received a number of international awards, graduated 40 PhD students, and has published 400 technical papers and five research books.

Geoff V. Merrett is an Associate Professor at the School of Electronics and Computer Science at the University of Southampton, where he is head of the Centre for Internet of Things and Pervasive Systems. He is internationally known for his research into the system-level energy management of mobile and self-powered embedded systems, and he has published around 170 journal and conference papers in these areas. He was a co-investigator on the £5.6M PRIME Programme Grant on runtime power and reliability management of many-core systems, where he led the applications and cross-layer interaction theme. He was General Chair of the European Workshop on Microelectronics Education in 2016, is an Associate Editor for IET CDT, and is a member of the IET and IEEE.

1598 978-1-78561-582-5



9 781785 615825

Computing
Materials, Circuits & Devices

The Institution of Engineering and Technology • www.theiet.org
978-1-78561-582-5



Many-Core Computing
Hardware and software



The Institution of
Engineering and Technology

Many-Core Computing

Hardware and software

Edited by

Bashir M. Al-Hashimi and Geoff V. Merrett

Edited by Al-Hashimi and Merrett



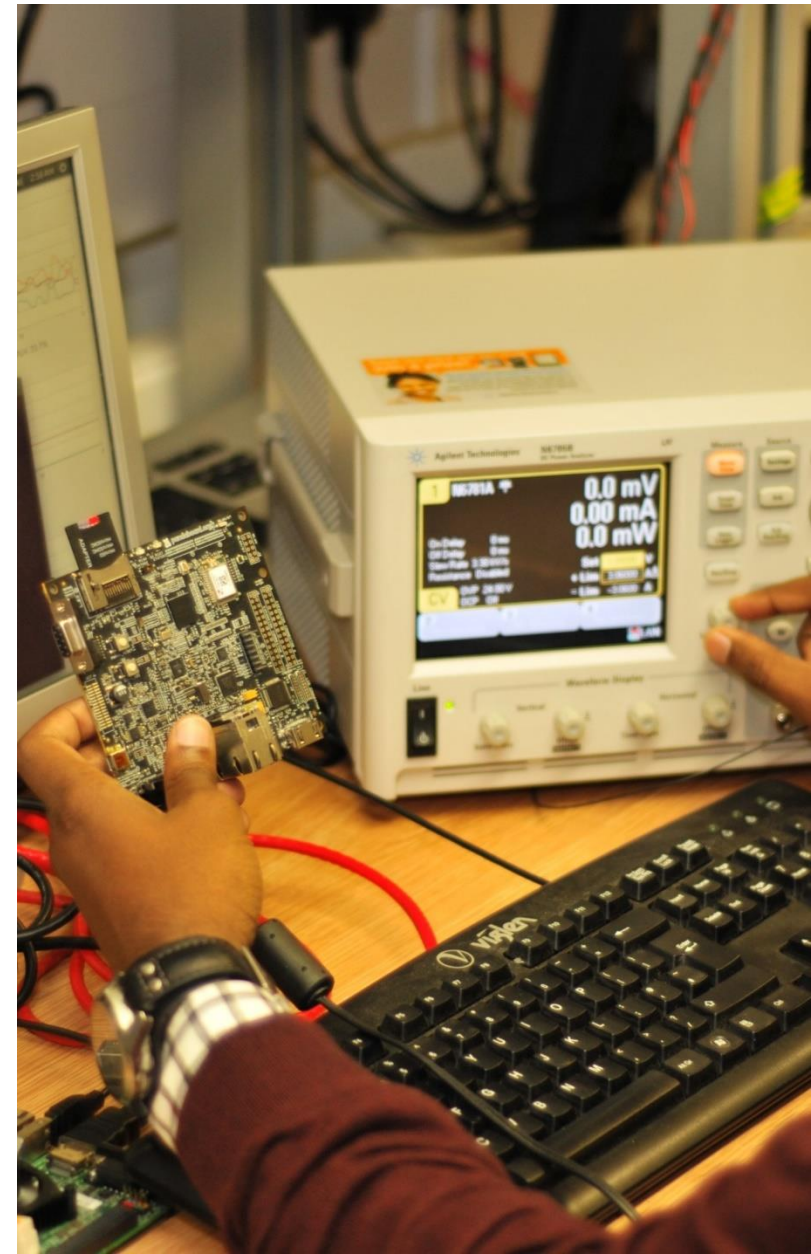
CONCLUSIONS

Runtime Power Management

- Single > multiple > concurrent applications
- Online *vs* offline+online approaches
- >> Number of cores
- Homogeneous *vs* Heterogeneous platforms

Tools and Support www.prime-project.org

- PowMon power estimation
www.powmon.ecs.soton.ac.uk
www.gemstone.ecs.soton.ac.uk
- PRiME RTM Framework
github.com/PRiME-project/PRiME-Framework
- PRiMEStereoMatch application
github.com/PRiME-project/PRiMEStereoMatch





Any Questions?

UNIVERSITY OF
Southampton

Dr Geoff V Merrett

Associate Professor | Head of Centre

Centre for IoT and Pervasive Systems

Tel: +44 (0)23 8059 2775

Email: gvm@ecs.soton.ac.uk | www.geoffmerrett.co.uk

Highfield Campus, Southampton, SO17 1BJ UK