

Theme 2

High Integrity Run-time Management and Optimisation

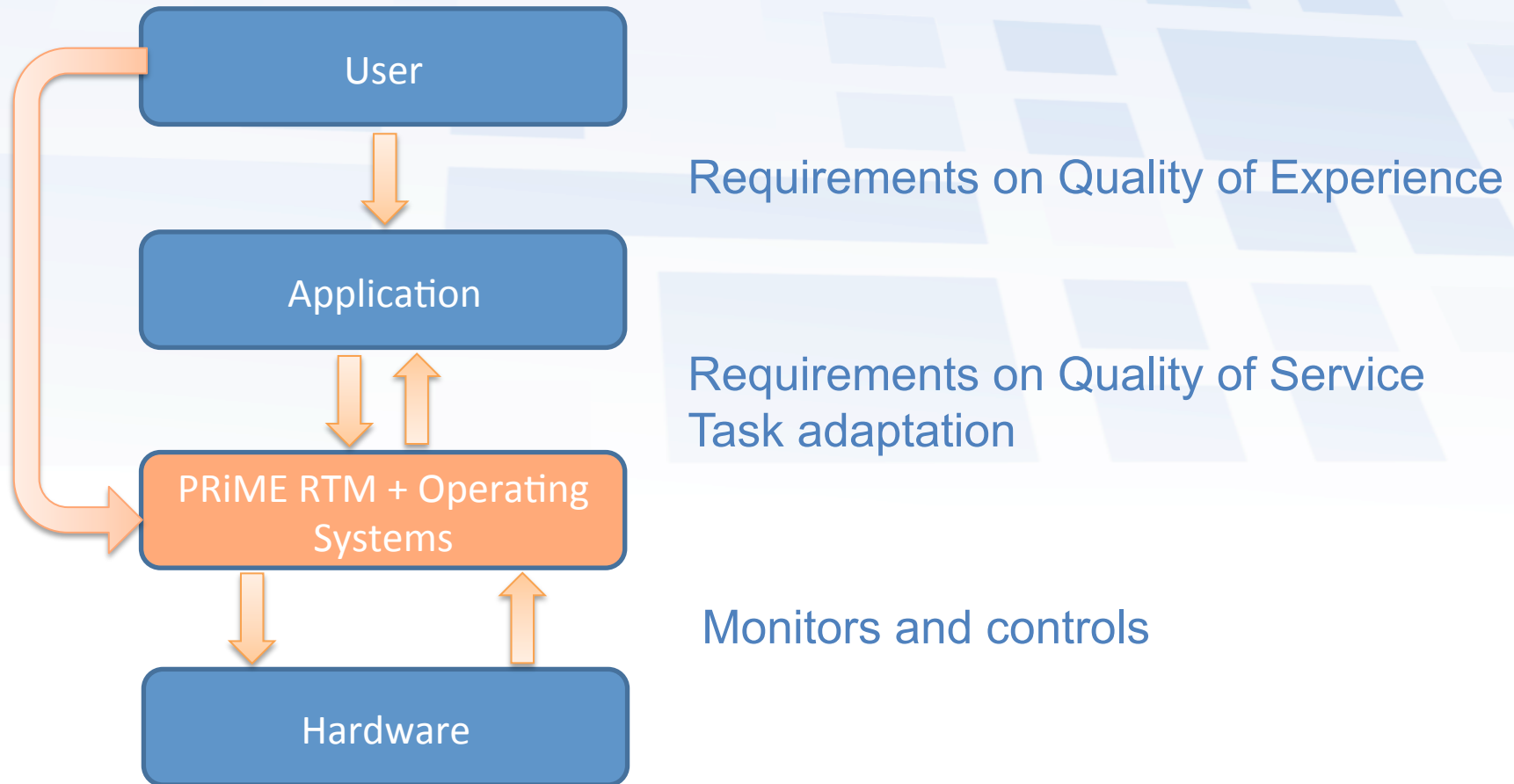
Michael Butler, Geoff Merrett, University of Southampton
PRiME Industry Day
London, 30 June 2015



Theme 2 Aim

- Theme 2 is developing sound and effective *design principles, algorithms and mechanisms* for *runtime management (RTM)* of power and faults to achieve a balance between performance, energy efficiency and resilience.
- Key challenges:
 - address a diverse range of many-core architectures
 - configurable for different optimisation strategies
 - support distinct usage scenarios with various degrees of criticality
 - do not compromise performance, energy efficiency or reliability.

PRiME RTM and its environment



Where are the opportunities?

- Energy saving
 - Dynamic Frequency and Voltage Control (DVFS)
 - Core shutdown
 - General purpose versus specialised cores
 - Reducing communications overhead
- Reliability management
 - Reduced thermal cycling for lifetime reliability
 - Fault tolerance
 - Reduced quality of experience/service
 - Frame rate, resolution, accuracy,...

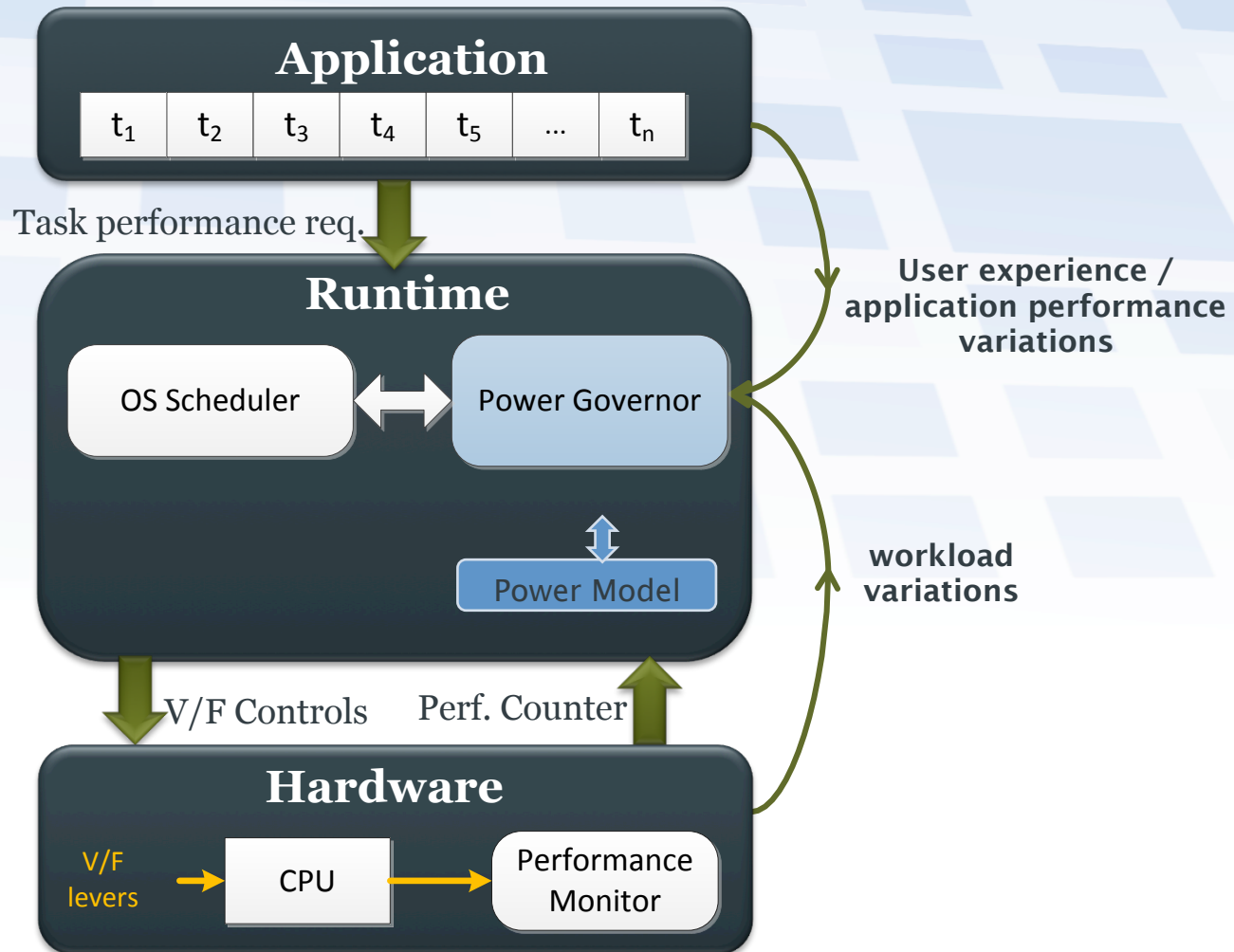
Highest Frequency (Performance)



Lowest Frequency (PowerSave)

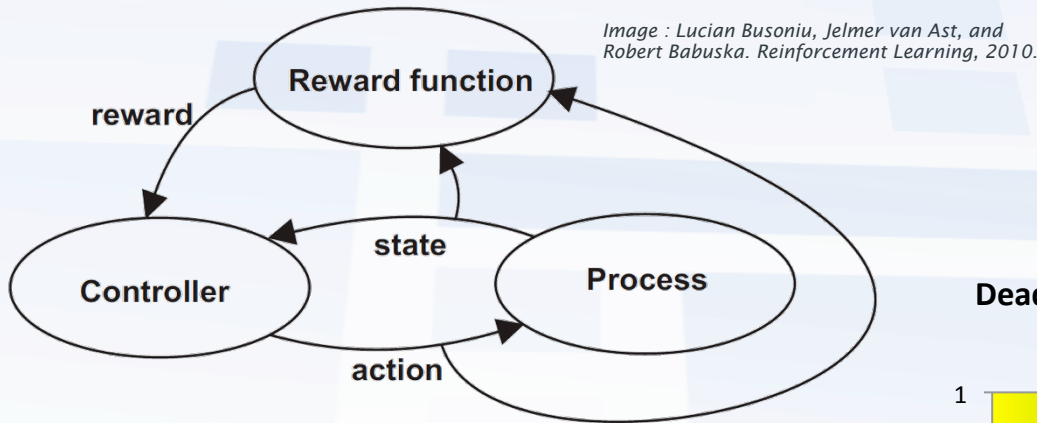


Run-time Power Management

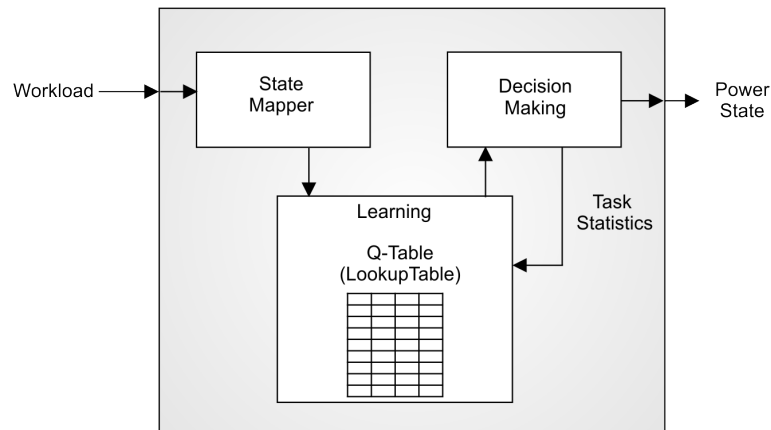


Learning the best DVFS mode

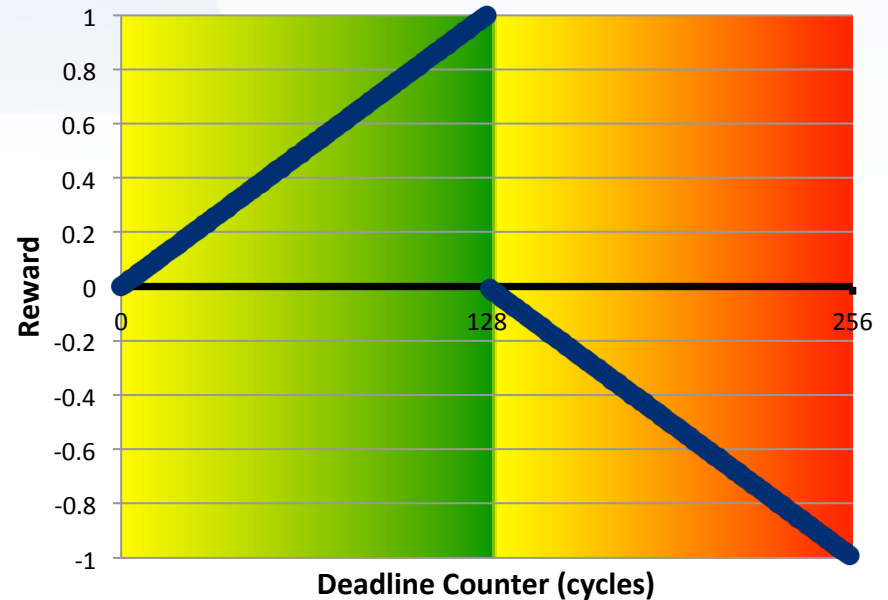
- Reinforcement Learning



- For Power Management

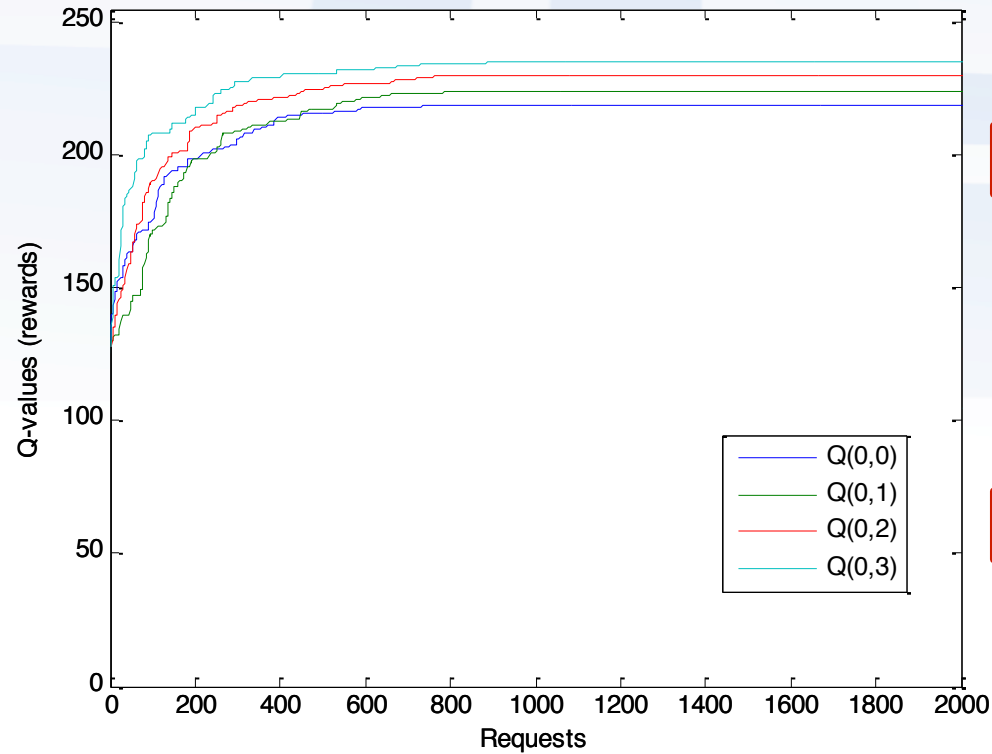


**Deadline Counter and its effect on the Reward
(Deadline = 128 cycles)**



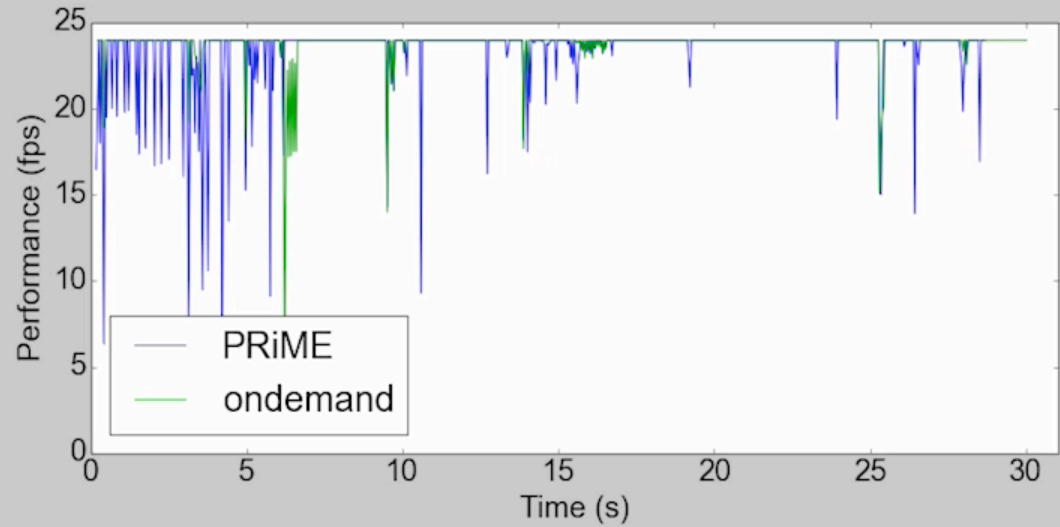
Q Table (and convergence)

Q-value means of the 4 actions at State 6

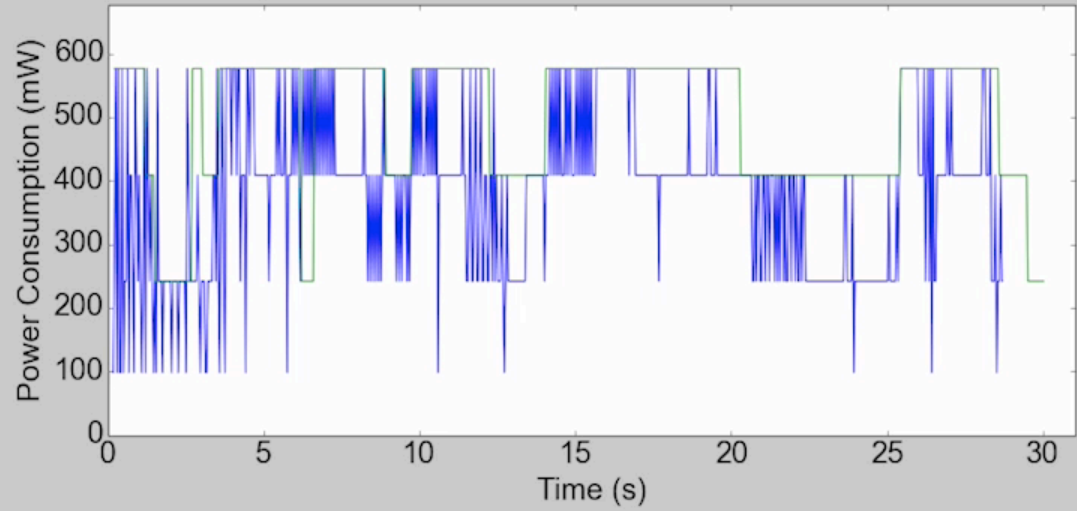


	ACTIONS (Power Modes)			
STATES (Tasks)	P0	P1	P2	P3
WD0	0	0	0	0
WD1	0	0	0	0
WD2	0	0	0	0
WD3	0	0	0	0
WD4	0	0	0	0
WD5	0	0	0	0
WD6	0	0	0	0
WD7	0	0	0	0
WD8	0	0	0	0

PRiME



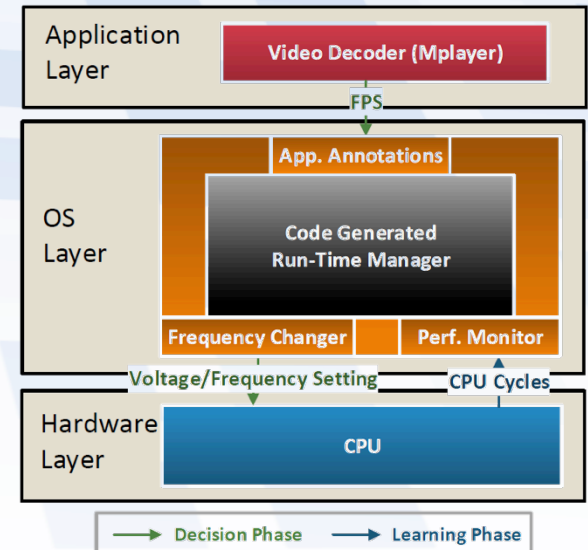
ondemand



(c) copyright 2008, Blender Foundation / www.bigbuckbunny.org

RTM Modelling/Verification

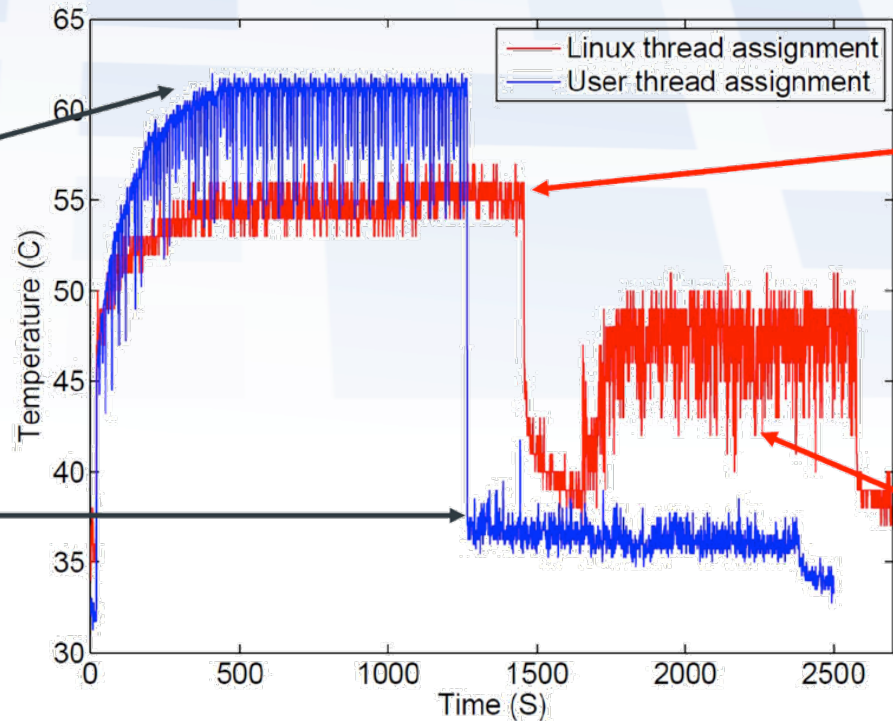
- Model-based verification and code generation
 - Event-B formal modelling of DVFS RTM
 - Code generation from Event-B model of DVFS RTM for media decoder application running on embedded boards
 - Modelica modelling and co-simulation of temperature/power in 4 core system



Lifetime-Aware Power Management

Face recognition
Thermal cycles:
High
Average Temperature:
High
Reliability Problem:
All

Mpeg encoding
Thermal cycles:
Low
Average Temperature:
Low
Reliability Problem:
Less

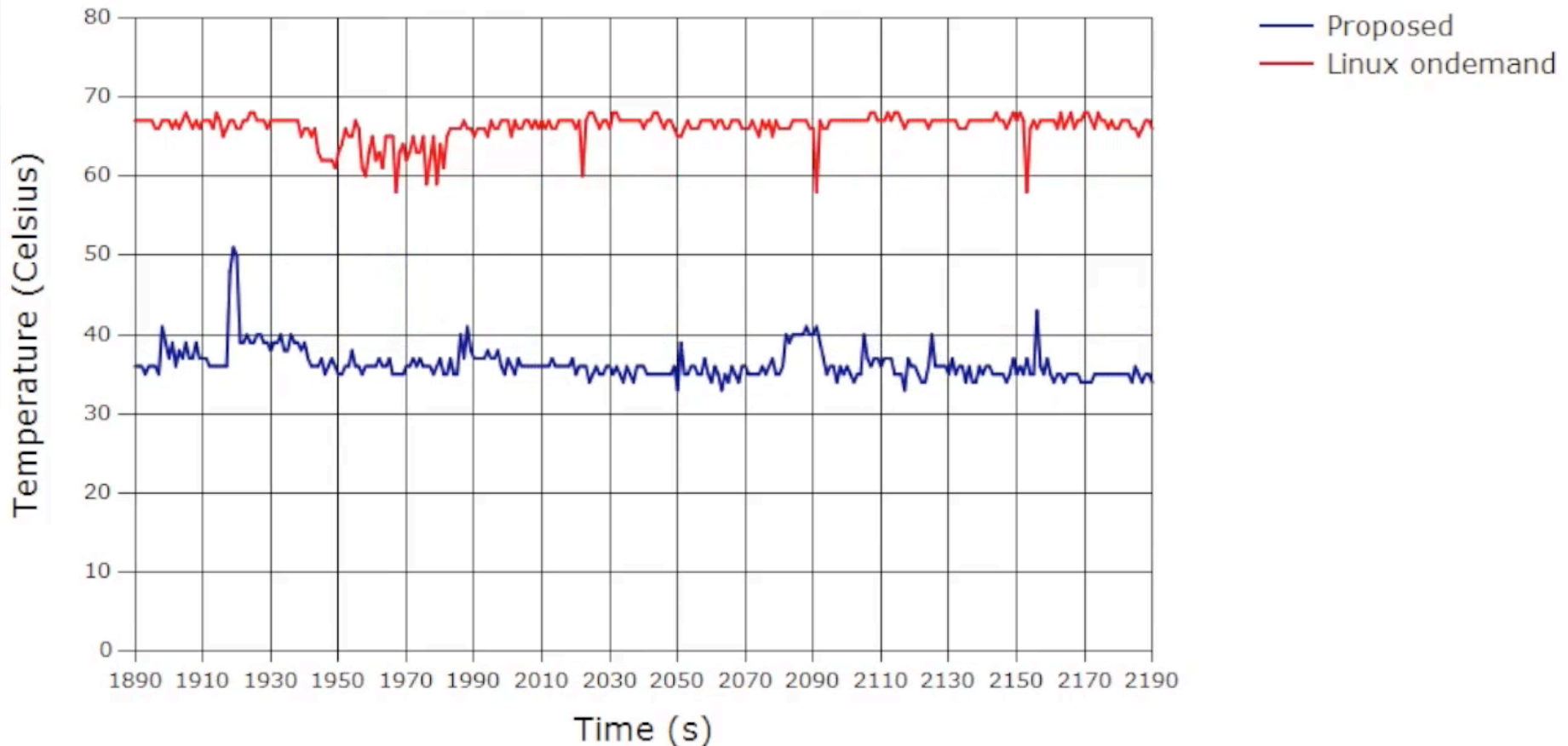


Face recognition
Thermal cycles:
Less
Average Temperature:
High
Reliability Problem:
NBTI, EM

Mpeg encoding
Thermal cycles:
High
Average Temperature:
Low
Reliability Problem:
Thermal Stress → TDDB

Stages of the Learning Algorithm

Convergence of the Reinforcement Learning Algorithm



PRiME RTM Algorithms

Power and reliability management of multi-core using machine learning approaches



1. Energy-reliability joint optimisation on ARM multi-core controlling DVFS and thread affinity
 - for energy optimisation and lifetime reliability
2. Power minimisation on heterogeneous architecture (CPU, FPGA, DSP)
3. Centralised DFVS / Shutdown on multi-core Xeon Phi
4. Decentralised learning-based runtime energy minimisation using multi-agent approach running on Xeon-Phi
5. Runtime power modelling for embedded processors

RTM Framework Specification

- A number of application-specific runtime prototypes have been developed for experimentation and demonstration
- PRiME will *not* develop a single runtime product
 - This would require fixing on an architecture and OS thus limiting the applicability of PRiME results
- To enable working together and reuse, we identified need to reach a common understanding of
 - Purpose,
 - Architecture, and
 - Interfaces to runtime

RTM Framework Specification

- Common structure for defining *purpose* of runtime features (requirements):
 - Structured around system layers
 - Has made it easier to see commonalities and differences
- Architecture
 - Continuum from centralised to fully decentralised
 - Relationship to Operating System structure
 - Shared memory vs message-passing
- Interfaces with
 - users, applications, controls, monitors, communications, OS

What's next?

Decentralised many-core power/reliability management

- Experiment with a wider range of applications for current multi-core platforms, including concurrency.
- Target RTM for more complex heterogeneous architectures
- Identify and exploit additional hardware monitors
- Improved models of energy usage in memory and communications
- Develop techniques for reducing energy consumption in memory and communications
- Apply code generation to a wider range of architectures
- Maintain runtime specification up-to-date



Thank you!

Any Questions?

Dr Geoff Merrett
Associate Professor



PRiME Applications and Demonstrators Lead

University of Southampton
Highfield Campus, Southampton, SO17 1BJ UK

Tel: +44 (0)23 8059 2775

Email: gvm@ecs.soton.ac.uk

Web: www.geoffmerrett.co.uk