# Adaptive Resource Management through Self-Awareness[+]

## Nikil D. Dutt

**Center for Embedded and Cyber-Physical Systems (CECS)**

**University of California, Irvine**

**dutt@uci.edu**

**http://www.ics.uci.edu/~dutt**

**https://duttgroup.ics.uci.edu**

[+] **Joint work with Tiago Mück, Bryan Donyanavard, Kasra Moazzemi, Amir Rahmani, Santanu Sarma, Biswadip Maity**

**Dutt Research Group**

# Self-Awareness ?

## Self-awareness

From Wikipedia, the free encyclopedia

*Not to be confused with Self-concept, Self-consciousness, Self-perception, or Self image.*

**This article has multiple issues.** Please help **improve**  [hide]
**it** or discuss these issues on the **talk page**.

- This article **may require cleanup** to meet Wikipedia's **quality standards**. *(March 2009)*
- This article **needs attention from an expert on the subject**. *(May 2009)*

**Self-awareness** is the capacity for introspection and the ability to recognize oneself as an individual separate from the environment and other individuals.

**Contents** [hide]

The mirror test is a simple measure of self-awareness.

Co                                                                                                         #2

# Computational Self-* Properties

- **Self-Awareness** [Hinchey2006]: System is aware of its *self states* and *behaviors*

- **Context-Awareness** [Parashar 2005] : System is aware of *context – i.e., its operational environment*

---

- *Self-configuring ->* capability of reconfiguring automatically
- *Self-healing* [Robertson2005] *-> self-diagnosing and self-repairing*
- *Self-optimizing-> capability of self-tuning* or *Self-adjusting*
- *Self-protecting ->* capability of detecting dangerous outcomes (e.g. security breaches) and recovering from their effects
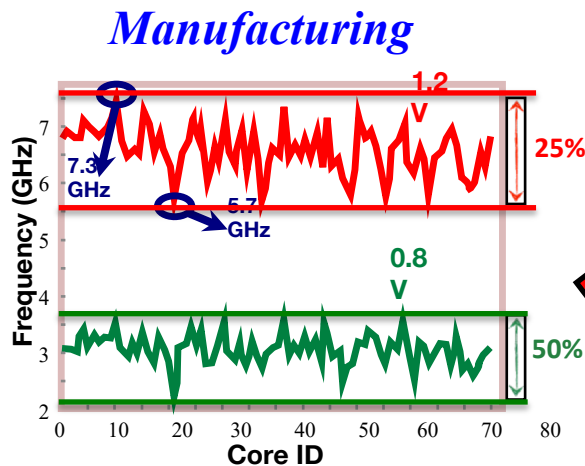
# Outline

- Computational Self-Awareness

- Why Self-Aware Chips?

- Cross-Layer Sensing & Actuation

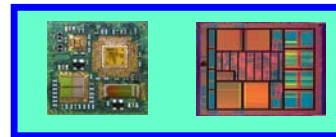- Towards Self-Aware Chips

- Supervisory Control & Coordination

# Why On-Chip Self-Awareness (1)?

https://duttgroup.ics.uci.edu

[Source: NSF Variability Expedition Project]

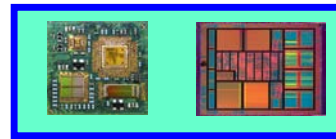# Why On-Chip Self-Awareness (1)?

## *Variability-induced challenges*

*Manufacturing*

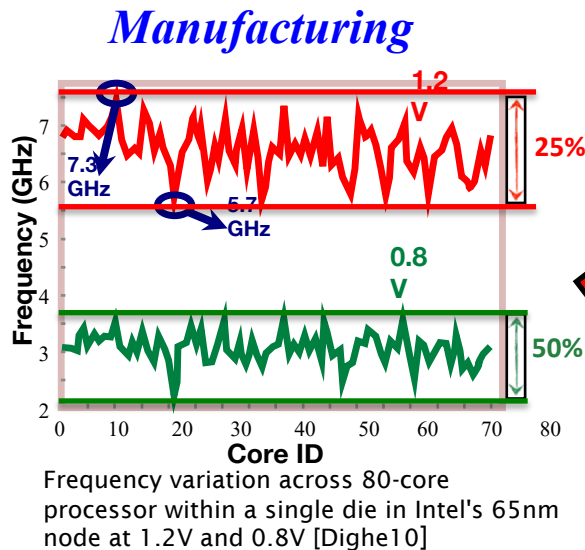

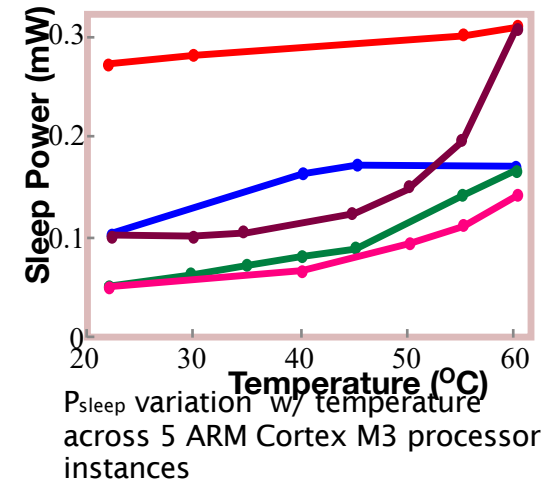Frequency variation across 80-core processor within a single die in Intel's 65nm node at 1.2V and 0.8V [Dighe10]

https://duttgroup.ics.uci.edu

[Source: NSF Variability Expedition Project]

# Why On-Chip Self-Awareness (1)?

## *Variability-induced challenges*

*Environment*



$P_{sleep}$ variation w/ temperature across 5 ARM Cortex M3 processor instances

*Manufacturing*



Frequency variation across 80-core processor within a single die in Intel's 65nm node at 1.2V and 0.8V [Dighe10]

https://duttgroup.ics.uci.edu

[Source: NSF Variability Expedition Project]
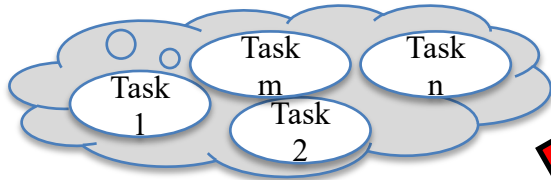
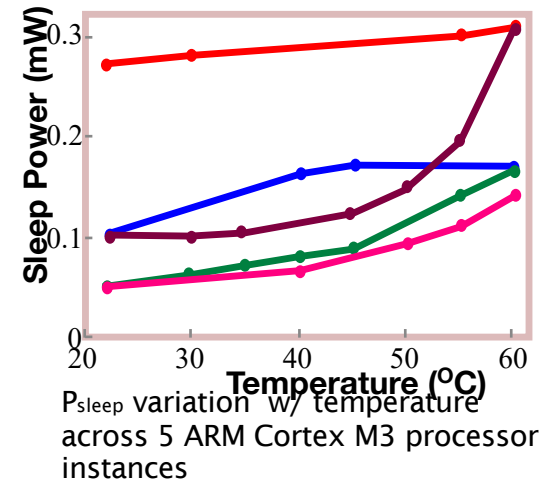# Why On-Chip Self-Awareness (1)?

## *Variability-induced challenges*

**Applications:** *varying compute, memory, communication*



Task 1

Task m

Task 2

Task n

**Environment**



$P_{sleep}$ variation w/ temperature across 5 ARM Cortex M3 processor instances

**Manufacturing**



1.2 V

25%

7.3 GHz

5.7 GHz

0.8 V

50%

Frequency variation across 80-core processor within a single die in Intel's 65nm node at 1.2V and 0.8V [Dighe10]

## *Triple Whammy!*

https://duttgroup.ics.uci.edu    [Source: NSF Variability Expedition Project]
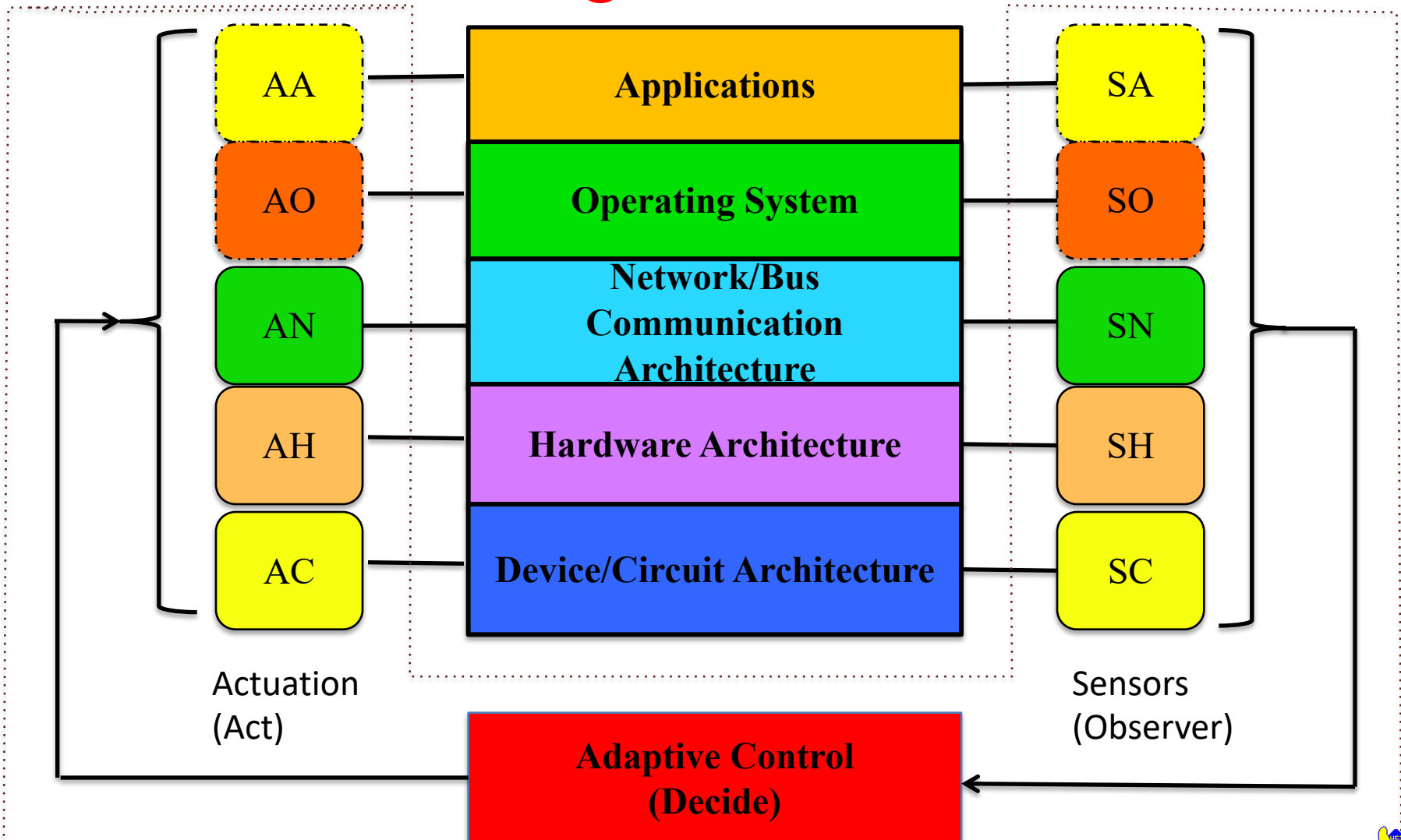
# Why On-Chip Self-Awareness (2) ?

- Chips must adapt to:
  - Performance, Power, Resilience, Security,….
- Provide Guarantees
- Dynamically manage multi-dimensional trade-offs
  - Performance, Power/Energy, Thermal,…..
  - QoS, TDP, Wear-out, ….

**Exploit Computational Self-Awareness**

# Outline

- Computational Self-Awareness

- Why Self-Aware Chips?

- Cross-Layer Sensing & Actuation

- Towards Self-Aware Chips

- Supervisory Control & Coordination

# Cross-Layer Physical/Virtual Sensing & Actuation

https://duttgroup.ics.uci.edu

# Examples of Virtual Sensors and Actuators Across Layers of CPSoC

| Layers | Virtual/Physical Sensors | Virtual/Physical Actuators |
|---|---|---|
| **Application** | Execution Time, Workload Power, Energy, | Loop perforation Algorithmic Choice |
| **Operating System** | System Utilization Peripheral States | Task Allocation, Scheduling, Migration, Duty Cycling |
| **Network/Bus Communication** | Bandwidth; Packet/Flit status; Channel Status, Congestion, Latency | Adaptive Routing Dynamic Bandwidth Allocation Ch. no and direction |
| **Hardware Architecture** | Cache misses, Miss rate; access rate; IPC, Throughput, ILP/MLP, Core asymmetry | Cache Sizing; Reconfiguration, Resource Provision Static/Dynamic Redundancy |
| **Circuit/Device** | Circuit Delay, Aging, leakage Temperature, oxide breakdown | DVFS, DFS, DVS ABB, Clock and Power-gating |

# Outline

- Computational Self-Awareness

- Why Self-Aware Chips?

- Cross-Layer Sensing & Actuation

- Towards Self-Aware Chips

- Supervisory Control & Coordination

# Self-Reflection & Introspection



- Ability to create a *self-model (introspect)*

- Ability to model their own body/structure (usually known *self-modeling*)

- Ability to model their own *behavior*

- *Metacognition capacity*: *'models one's own thinking'*, *'think about thinking'*

- *System with **two/multiple minds:** one **being modeled** and other **doing modeling***

https://duttgroup.ics.uci.edu          #14

# Reflex vs Reflect

**Reflexive, Reactive**



- **Actions driven solely on external feedback**
  - **E.g., our autonomic nervous systems**

# Reflex vs Reflect

**Reflexive, Reactive**

**Reflection, Introspection**



- **Actions driven solely on external feedback**
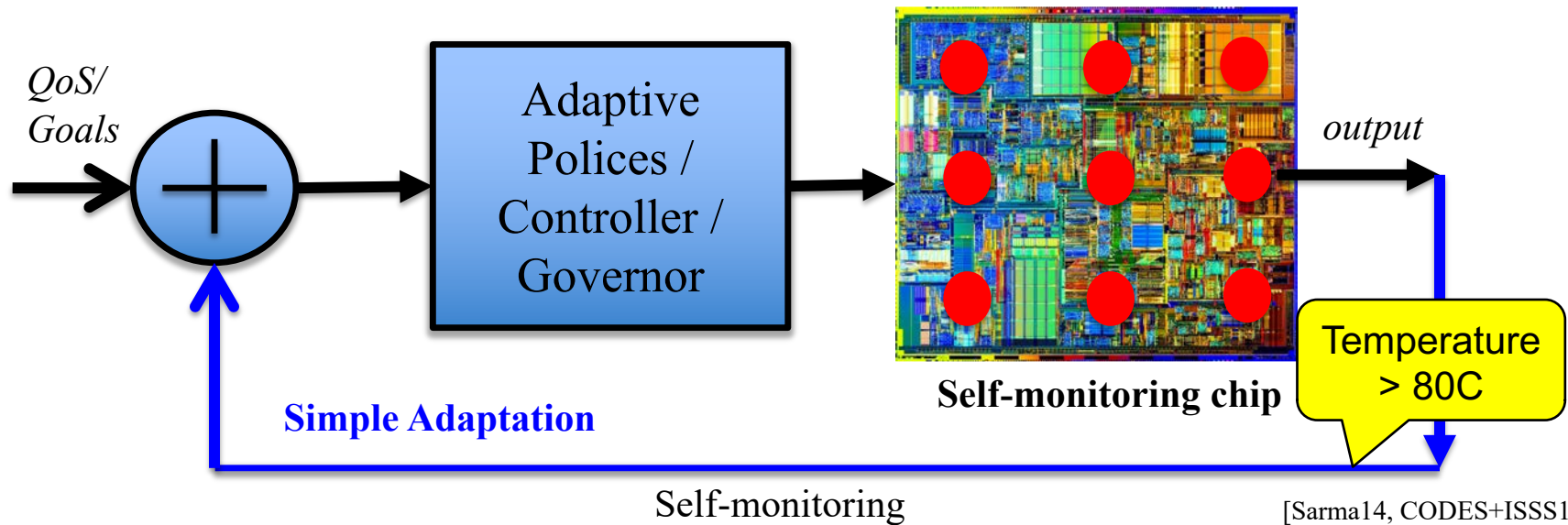  - *E.g., our autonomic nervous systems*

- **Consider past and future outcomes**
  - *E.g., planning, strategies, policies, …*
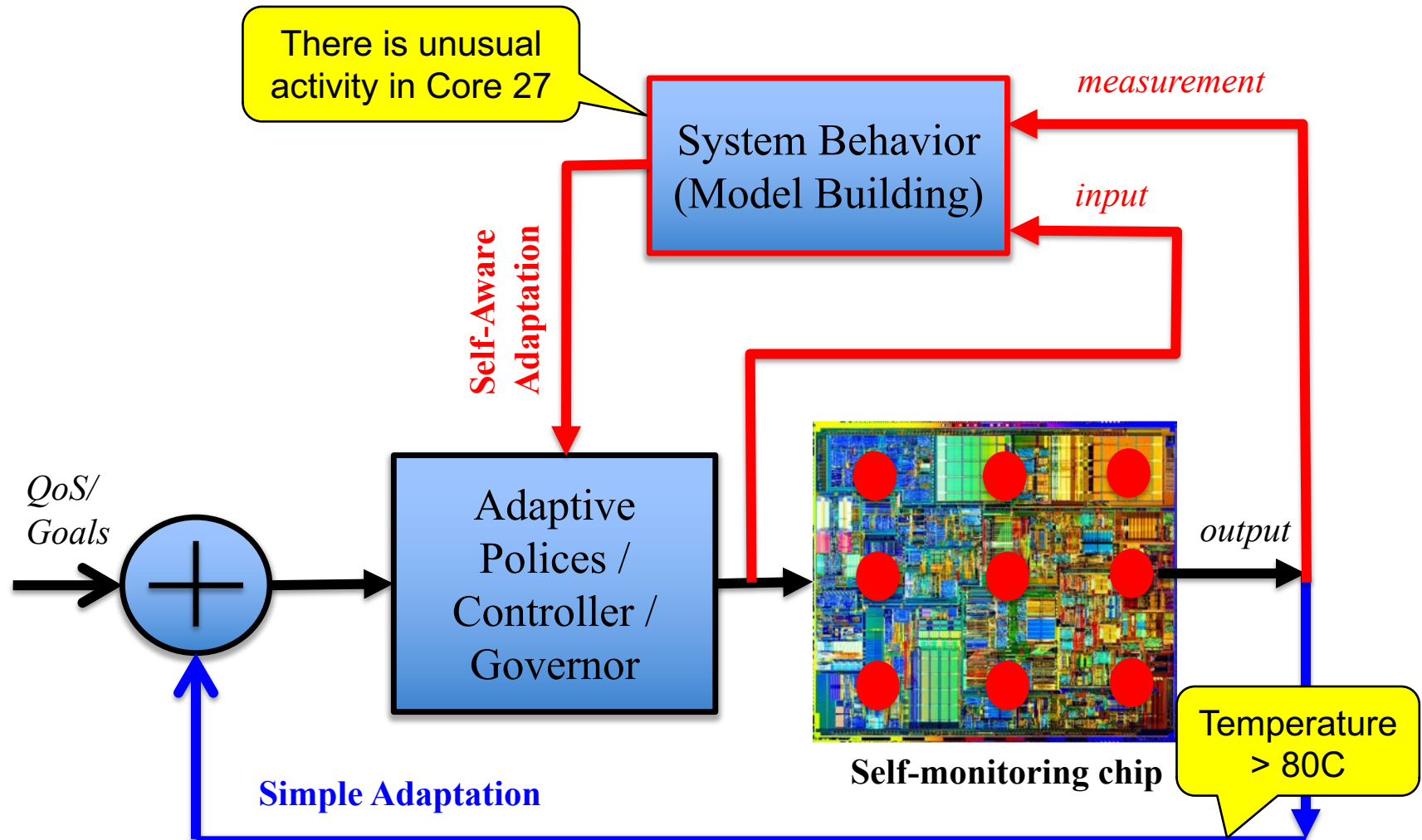
# Towards Self-Aware Chips: What we do now

**Reflexive, Reactive**



QoS/ Goals → ⊕ → Adaptive Polices / Controller / Governor → output

**Self-monitoring chip**

Temperature > 80C

**Simple Adaptation**

Self-monitoring

[Sarma14, CODES+ISSS14]

# Towards Self-aware chips
## *Beyond simple reactive models*



There is unusual activity in Core 27

System Behavior
(Model Building)

*measurement*

*input*

**Self-Aware Adaptation**

*QoS/ Goals*

Adaptive
Polices /
Controller /
Governor

*output*

Self-monitoring chip

Temperature > 80C

**Simple Adaptation**

Self-monitoring and **Self-modeling**

[Sarma14, CODES+ISSS14]

# Towards Self-aware chips
## Beyond simple reactive models



**Reflection, Introspection**

*measurement*

System Behavior
(Model Building)

*input*

**Self-Aware Adaptation**

*QoS/ Goals*

Adaptive Polices / Controller / Governor

*output*

**Self-monitoring chip**

**Simple Adaptation**

Self-monitoring and **Self-modeling**

[Sarma14, CODES+ISSS14]

# Today: "Reflexive" Resource Management
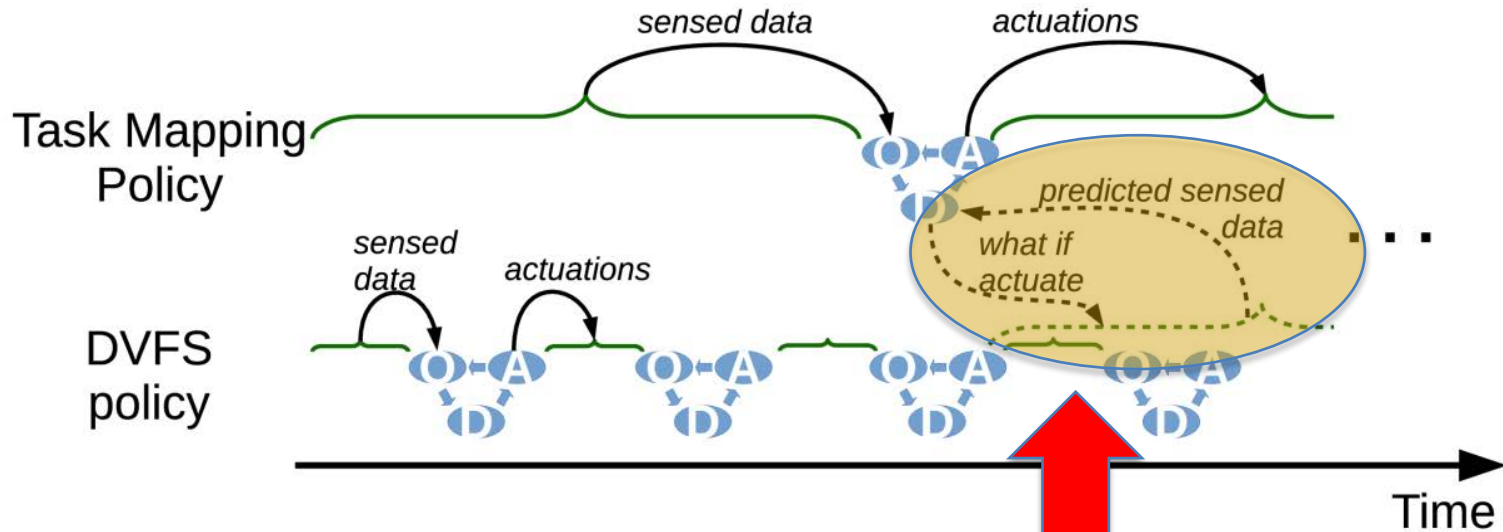
- Dynamic Voltage/Frequency scaling (DVFS)

Scale core frequency

**DVFS policy**

Core load increased

- Observe-Decide-Adapt approaches

*sensed data*  *actuations*  . . .

O-A  O-A  O-A  O-A
D    D    D    D

Time

# Relfe**X**ive vs Refle**CT**ive Resource Management



- *Relfe**X**ive ODA*:  decisions taken based on
  - **past** observations (purely reflexive)   OR
  - predictions made from **past** observations

# RefleXive vs RefleCTive Resource Management



- *RefleXive ODA*: decisions taken based on
  - **past** observations (purely reflexive)   OR
  - predictions made from **past** observations

- **RefleCTive approach**: considers **future** events that could happen in the next iteration of the ODA loop

# Adaptive Resource Management

- Use concept of **reflection**
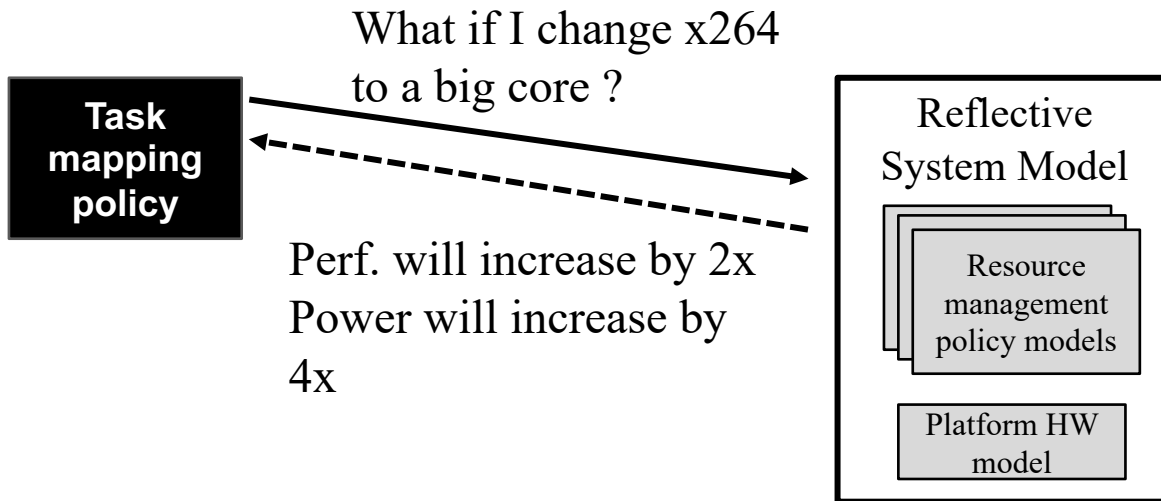  - **Reflection**: *change your actions based on both external feedback and **introspection** (i.e., self-assessment)*

# Adaptive Resource Management

- Use concept of **reflection**

  - Reflection: *change your actions based on both external feedback and introspection (i.e., self-assessment)*

- **Reflective resource management combines:**

  - Current system state assessed from sensing information (e.g., readings from performance counters, power sensors, etc.)

  - Models to predict the behavior of other system components before performing an action

# MARS: Our coordination approach

- Coordination though **reflective** resource management
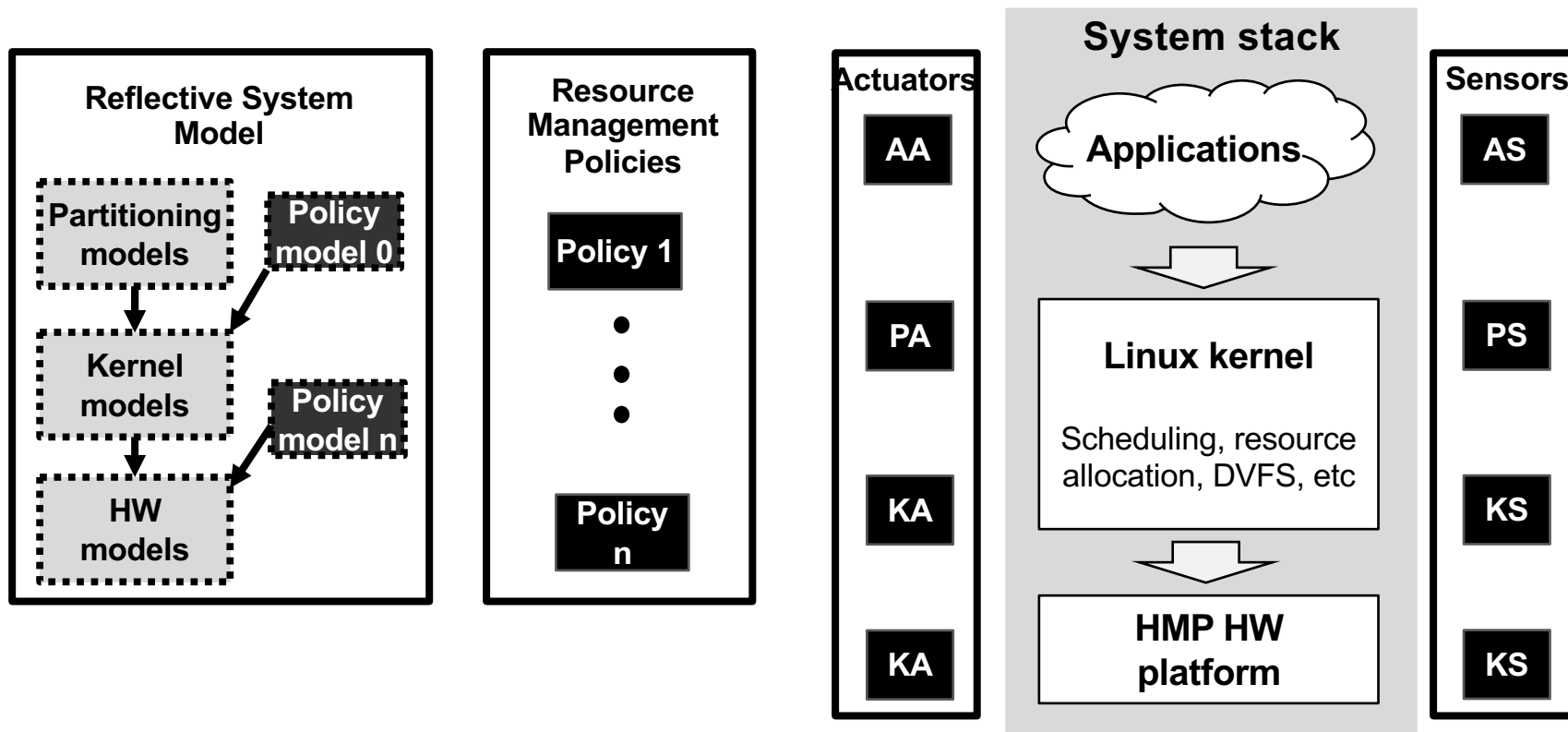  - **MARS: M**iddleware for **A**daptive **R**eflective **S**ystems

What if I change x264
to a big core ?

**Task mapping policy**

Perf. will increase by 2x
Power will increase by 4x

Reflective System Model

Resource management policy models

Platform HW model

# Do we have room for reflection ?

- Systems actuations happen at different timescales

| Task mapping | DVFS | Scheduling |
|---|---|---|
| Which core executes a task? (100's – 10's ms) | What's the next frequency? (10's ms – 10's us) | Should I preempt a task ? (<1 us) |

← **Coarser grained actuation** → **Finer grained actuation**

- Some actuations happen quickly with little room for reasoning

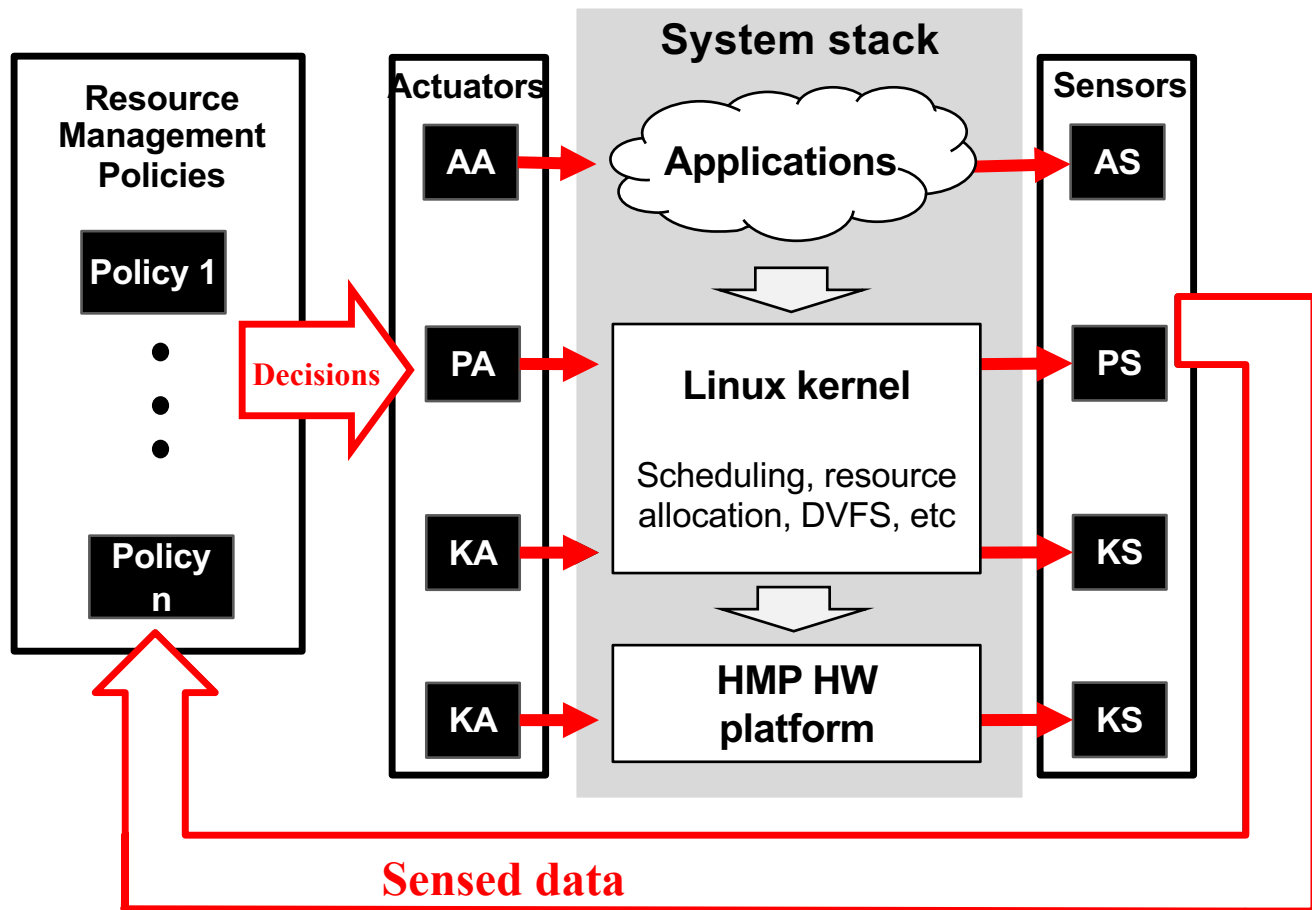- Other actuations can occur on larger timescales
  - **Task mapping, Wear-leveling (for aging)….**

# MARS middleware for reflective resource management

**Reflective System Model**

- Partitioning models
- Policy model 0
- Kernel models
- Policy model n
- HW models

**Resource Management Policies**

- Policy 1
- •
- •
- •
- Policy n

**Actuators**

- AA
- PA
- KA
- KA

**System stack**

Applications

⬇

**Linux kernel**

Scheduling, resource allocation, DVFS, etc

⬇

**HMP HW platform**

**Sensors**

- AS
- PS
- KS
- KS

# MARS middleware for reflective resource management

# MARS middleware for reflective resource management

# MARS middleware for reflective resource management
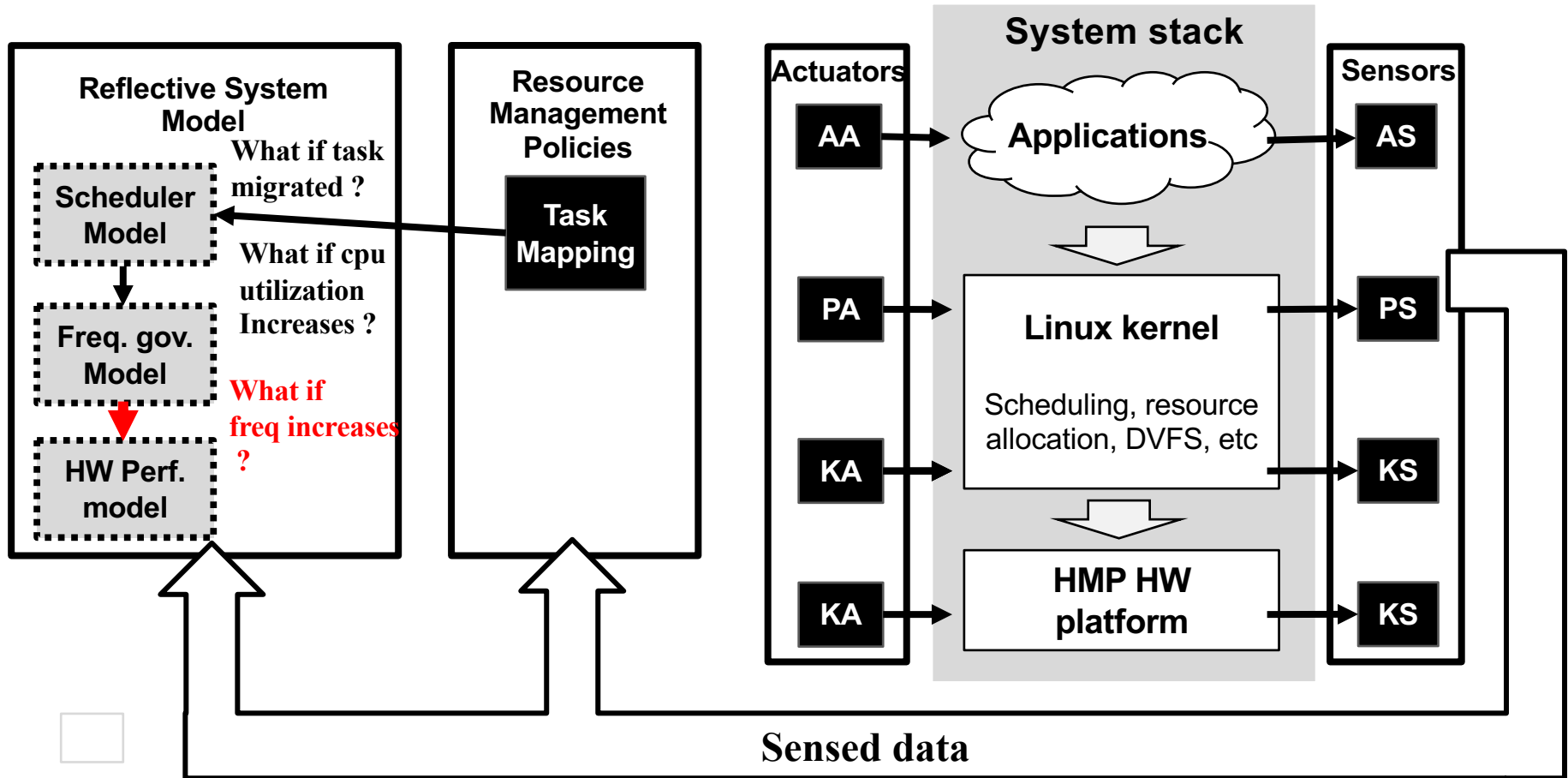
# Example: Reflective Task Mapping



**Reflective System Model**

**Resource Management Policies**

**Task Mapping**

**System stack**

**Actuators**

**Sensors**

AA → Applications → AS

PA → Linux kernel → PS

Scheduling, resource allocation, DVFS, etc

KA → KS

KA → HMP HW platform → KS

**Sensed data**

Reflective middleware

https://duttgroup.ics.uci.edu

10/10/18

# Example: Reflective Task Mapping



**Reflective System Model**

**Scheduler Model**

*What if task migrated ?*

**Resource Management Policies**

**Task Mapping**

**System stack**

**Actuators**

AA

PA

KA

KA

**Applications**

**Linux kernel**

Scheduling, resource allocation, DVFS, etc

**HMP HW platform**

**Sensors**

AS

PS

KS

KS

**Sensed data**

# Example: Reflective Task Mapping

**Reflective System Model**

- Scheduler Model
- Freq. gov. Model

What if task migrated ?

**What if cpu utilization Increases ?**

**Resource Management Policies**

- Task Mapping

**System stack**

**Actuators**

- AA
- PA
- KA
- KA

Applications

Linux kernel

Scheduling, resource allocation, DVFS, etc

HMP HW platform

**Sensors**

- AS
- PS
- KS
- KS

**Sensed data**

# Example: Reflective Task Mapping



**Reflective System Model**

- Scheduler Model
- Freq. gov. Model
- HW Perf. model

What if task migrated ?

What if cpu utilization Increases ?

**What if freq increases ?**

**Resource Management Policies**

Task Mapping

**Actuators**

AA
PA
KA
KA

**System stack**

Applications

Linux kernel

Scheduling, resource allocation, DVFS, etc

HMP HW platform

**Sensors**

AS
PS
KS
KS

**Sensed data**

# Example: Reflective Task Mapping



**Reflective System Model**

Scheduler Model

What if task migrated ?

Freq. gov. Model

What cpu utilization Increases ?

HW Perf. model

What if freq increases ?

Task perf is X

**Resource Management Policies**

Task Mapping

**System stack**

Actuators

AA

Applications

PA

Linux kernel

Scheduling, resource allocation, DVFS, etc

KA

KA

HMP HW platform

Sensors

AS

PS

KS

KS

**Sensed data**

# Example: Reflective Task Mapping



**Reflective System Model**

Scheduler Model

Freq. gov. Model

HW Perf. model

**What if task migrated ?**

**What cpu utilization Increases ?**

**What if freq increases ?**

**Task perf is X**

**Resource Management Policies**

Task Mapping

**New task-to-core assignment**

**Actuators**

AA

PA

KA

KA

**System stack**

Applications

Linux kernel

Scheduling, resource allocation, DVFS, etc

HMP HW platform

**Sensors**

AS

PS

KS

KS

**Sensed data**

# SPARTA improvements

- 8-core big.LITTLE Exynos SoC
  - 4x big
  - 4x LITTLE
- Workload mixes (4 tasks each)
  - Mibench
  - x264 (Parsec)
- SPARTA vs Linux's GTS
- **Avg. improvements of 16% in energy efficiency without performance degradation**

**Donyanavard, B., Mück, T., Sarma, S., & Dutt, N.,** *SPARTA: Runtime Task Allocation for Energy Efficient Heterogeneous Many-cores.* **CODES+ISSS '16**

# MARS: Middleware for Adaptive Reflective Computer Systems

- Framework and tools for developing reflective resource/power management policies

  - Use models to predict system behavior

  - Enable easy adaptation to runtime changes

  - Case studies show promise

MARS framework is open source

**https://github.com/duttresearchgroup/MARS**

# Outline

- Computational Self-Awareness

- Why Self-Aware Chips?

- Cross-Layer Sensing & Actuation

- Towards Self-Aware Chips

- **Supervisory Control & Coordination**

# Goals and Autonomy

**Goal**



- **Single, straightforward objective**
  - **E.g., hit the pin**

# Goals and Autonomy

**Goal**



**Model Imperfection**



**+**

- **Single, straightforward objective**
  - **E.g., hit the pin**

- **What happens when we introduce unpredictability?**
  - **E.g., balls with different sizes, shapes weights; uneven or damaged surfaces**

# Goals and Autonomy

## Supervision



- **Constrain behavior so we are always headed toward the goal**
  - **E.g., bumpers**

# Goals and Autonomy

## Supervision



- **Constrain behavior so we are always headed toward the goal**
  - **E.g., bumpers**
- **Bonus: what about when we have more complex or multiple goals?**



https://duttgroup.ics.uci.edu

# SPECTR: Our Supervisory Approach

- Autonomy and robustness through supervisory control

*Rahmani, A. M., **Donyanavard, B., Mück, T.,** Moazzemi, K., Jantsch, A., Mutlu, O., & **Dutt, N.,** *SPECTR: Formal Supervisory Control and Coordination for Many-core Systems Resource Management.* **ASPLOS '18**
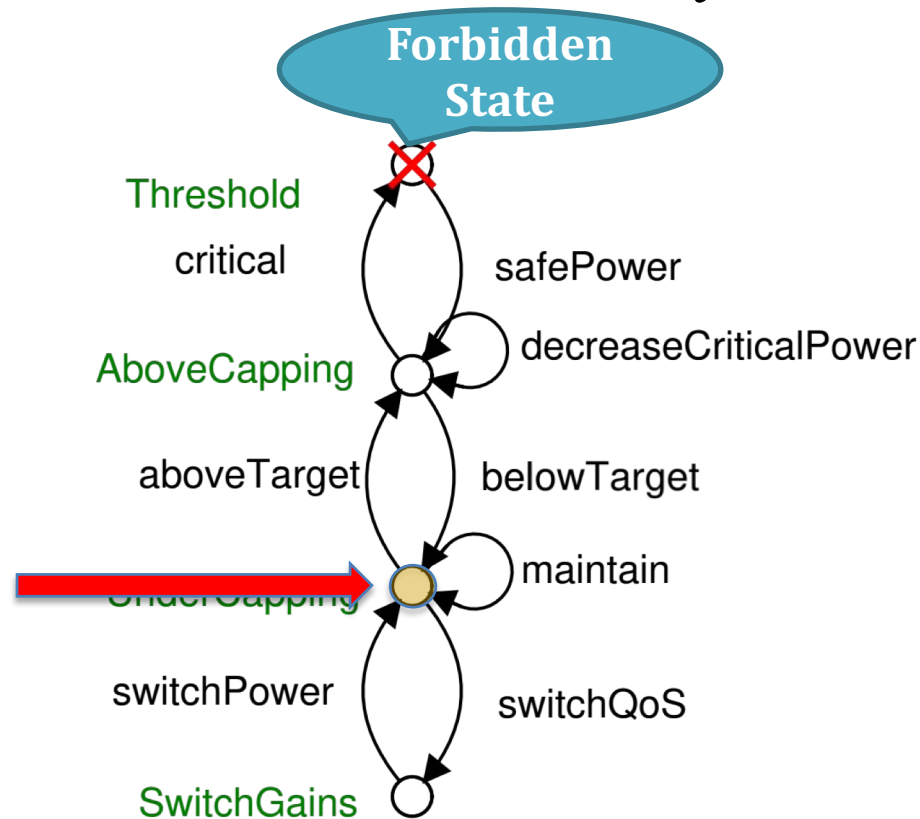
# SPECTR: Our Supervisory Approach

- Autonomy and robustness through supervisory control



**Low-level controllers satisfy objective**

*Rahmani, A. M., **Donyanavard, B.**, **Mück, T.**, Moazzemi, K., Jantsch, A., Mutlu, O., & **Dutt, N.**, *SPECTR: Formal Supervisory Control and Coordination for Many-core Systems Resource Management.* **ASPLOS '18**

# SPECTR: Our Supervisory Approach

- Autonomy and robustness through supervisory control

**Supervisor bounds behavior of controllers, manages goal**



**Low-level controllers satisfy objective**

*Rahmani, A. M., **Donyanavard, B.**, **Mück, T.**, Moazzemi, K., Jantsch, A., Mutlu, O., & **Dutt, N.**, *SPECTR: Formal Supervisory Control and Coordination for Many-core Systems Resource Management*. **ASPLOS '18**

# SPECTR: Our Supervisory Approach

- Autonomy and robustness through supervisory control
- Case Study*

**Supervisor bounds behavior of controllers, manages goal**



**Low-level controllers satisfy objective**



**ODROID-XU3 platform contains an Exynos 5422 Octa-core SoC**

*Rahmani, A. M., **Donyanavard, B.**, **Mück, T.**, Moazzemi, K., Jantsch, A., Mutlu, O., & **Dutt, N.**, *SPECTR: Formal Supervisory Control and Coordination for Many-core Systems Resource Management*. **ASPLOS '18**

# SPECTR: Our Supervisory Approach

- Autonomy and robustness through supervisory control
- Case Study*

**Supervisor bounds behavior of controllers, manages goal**



**ODROID-XU3 platform contains an Exynos 5422 Octa-core SoC**

**Low-level controllers satisfy objective**

*Rahmani, A. M., **Donyanavard, B.**, **Mück, T.**, Moazzemi, K., Jantsch, A., Mutlu, O., & **Dutt, N.**, *SPECTR: Formal Supervisory Control and Coordination for Many-core Systems Resource Management.* ASPLOS '18

# Example: Power Capping

- Specify desired behavior via accepted and forbidden states to restrict the behavior of the system

# Example: Power Capping

- Specify desired behavior via accepted and forbidden states to restrict the behavior of the system



**If power is in safe region, prioritize QoS**

# Example: Power Capping

- Specify desired behavior via accepted and forbidden states to restrict the behavior of the system

**If power exceeds threshold,**

# Example: Power Capping

- Specify desired behavior via accepted and forbidden states to restrict the behavior of the system



**If power exceeds threshold, reduce power**

# Example: Power Capping

- Specify desired behavior via accepted and forbidden states to restrict the behavior of the system

**If power exceeds threshold, reduce power...until it lowers again**

# Example: Power Capping

- Specify desired behavior via accepted and forbidden states to restrict the behavior of the system



**If power exceeds threshold, reduce power...until it lowers again**

# SPECTR Demonstration

*QoS Task: x264*

# SPECTR Demonstration



QoS Application: 🌐 📞 📹
Non-QoS Applications: 🖨 🗐 🛠

| A15 | A15 | A15 | A15 |
Big cores cluster

| A7 | A7 | A7 | A7 |
Little cores cluster

*QoS Task: x264*

**Safe Phase**: QoS app only
**SPECTR** satisfies FPS with
minimum power



Threshold
critical — safePower
AboveCapping — decreaseCriticalPower
aboveTarget — bellowTarget
UnderCapping — maintain
switchPower — switchQoS
SwitchGains

# SPECTR Demonstration

QoS Application: 🌐 📞 📹  Non-QoS Applications: 🖨 📄 🔧

| A15 | A15 | A15 | A15 |
|-----|-----|-----|-----|

**Big cores cluster**

| A7 | A7 | A7 | A7 |
|----|----|----|----|

**Little cores cluster**

*QoS Task: x264*

**Emergency Phase: TDP reduced in response to thermal event SPECTR satisfies FPS and power**

# SPECTR Demonstration



QoS Application
Non-QoS Applications

| A15 | A15 | A15 | A15 |
Big cores cluster

| A7 | A7 | A7 | A7 |
Little cores cluster

*QoS Task: x264*

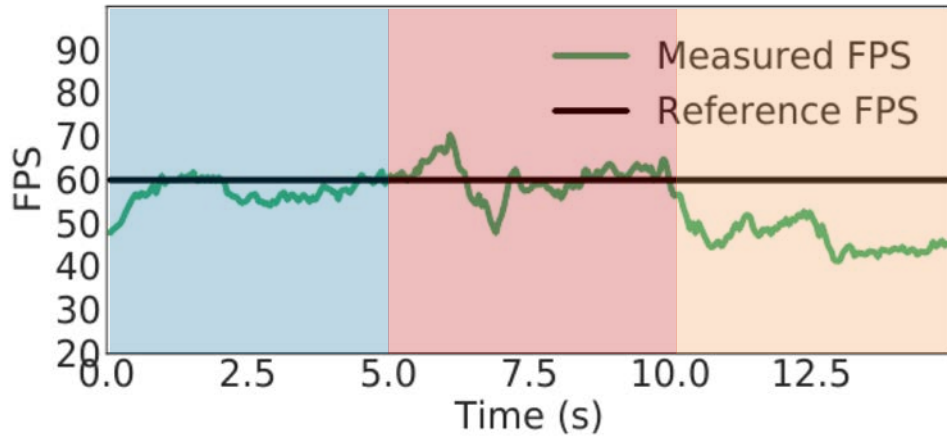**Disturbance Phase: TDP returned to normal, background tasks introduced SPECTR prioritizes power capping**



Threshold
critical — safePower
AboveCapping — decreaseCriticalPower
aboveTarget — bellowTarget
UnderCapping — maintain
switchPower — switchQoS
SwitchGains

# SPECTR Demonstration



*QoS Task: x264*



**SPECTR meets FPS target when possible, while honoring power cap**
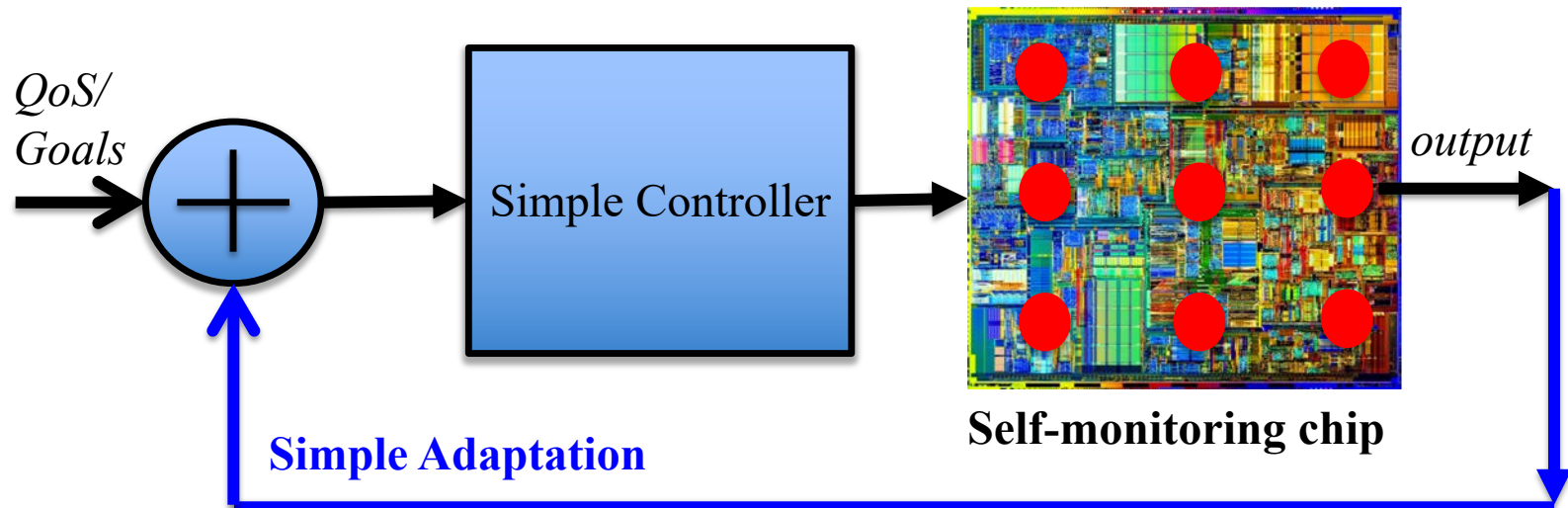
# Outline

- Computational Self-Awareness

- Why Self-Aware Chips?

- Cross-Layer Sensing & Actuation

- Towards Self-Aware Chips

- Supervisory Control & Coordination

- Wrap-up

# Key Take-Away 1:
# Cross-Layer Physical/Virtual Sensing & Actuation
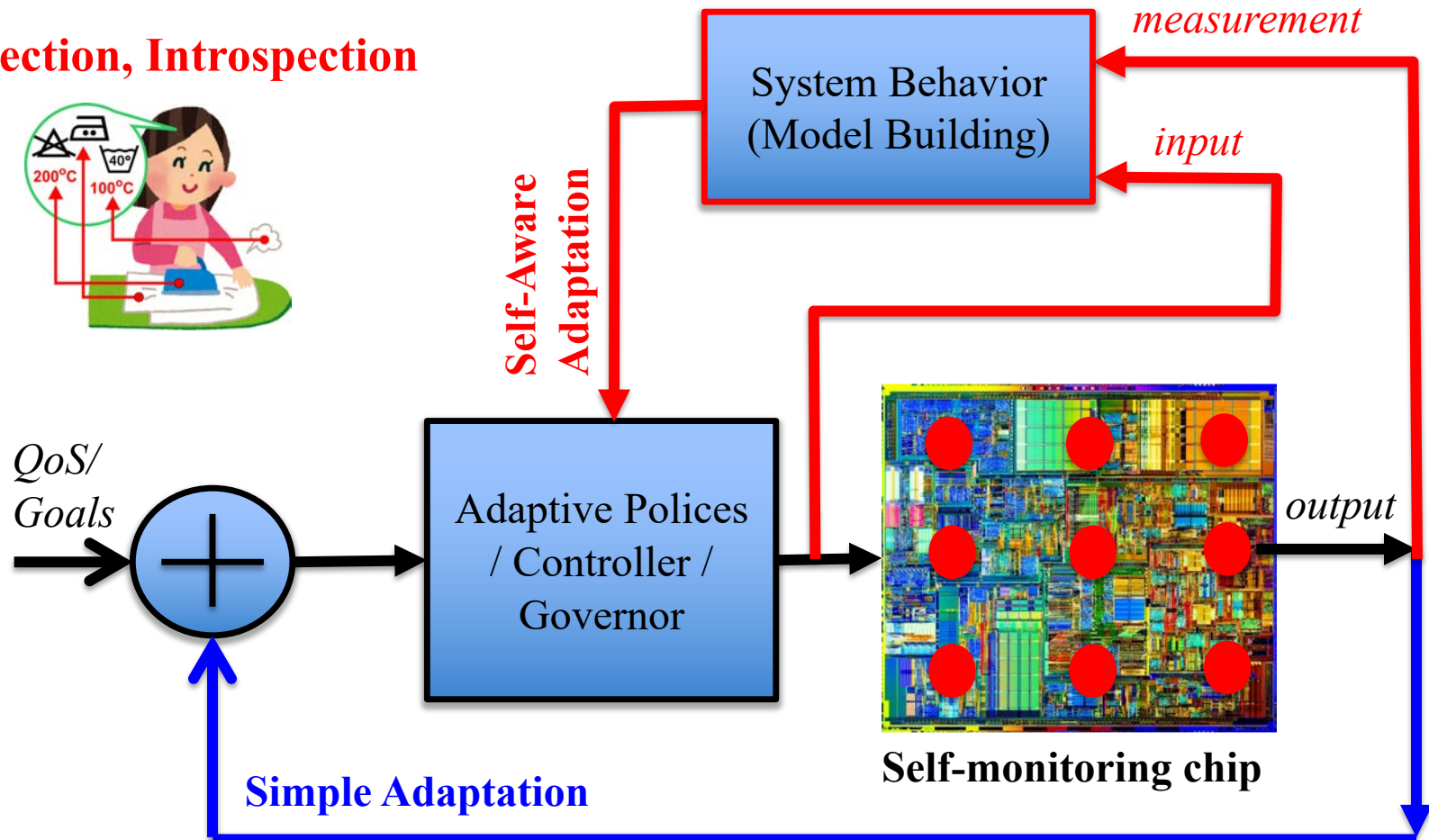
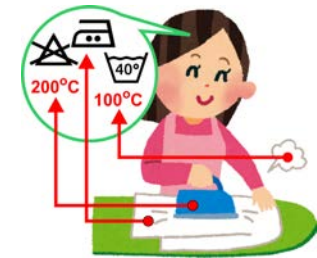# *From today's chips*

**Reflexive, Reactive**



QoS/ Goals → ⊕ → Simple Controller → Self-monitoring chip → *output*

**Simple Adaptation**

**Self-monitoring chip**

Self-monitoring and **simple adaptation**

CECS

# Key Take-Away 2:
## *Towards on-chip self-awareness*

**Reflection, Introspection**



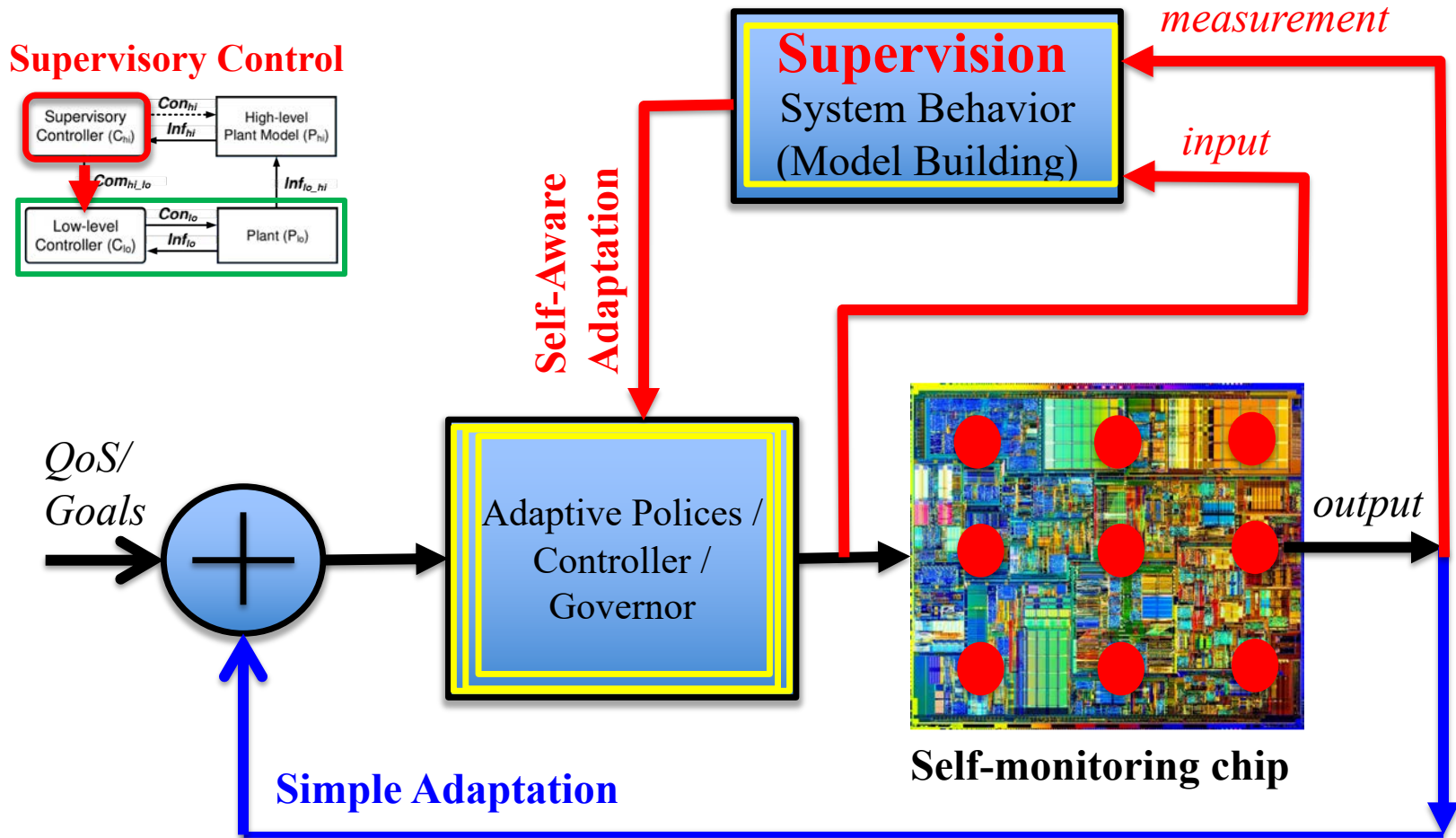**System Behavior
(Model Building)**

*measurement*

*input*

**Self-Aware Adaptation**

*QoS/ Goals*

**Adaptive Polices / Controller / Governor**

*output*

**Self-monitoring chip**

**Simple Adaptation**

Self-monitoring and **Self-modeling**        **[Sarma14, CODES+ISSS14]**

CECS

# Key Take-Away 3:
## *Supervisory Control & Coordination*

**Supervisory Control**



**Supervision**
System Behavior
(Model Building)

*measurement*

*input*

**Self-Aware Adaptation**

QoS/
Goals

Adaptive Polices /
Controller /
Governor

*output*

**Self-monitoring chip**

**Simple Adaptation**

**Supervisory Control & Coordination**    **[Rahmani18, ASPLOS 2018]**

Special Issue on Self-Awareness in Systems on Chip 2017

# IEEE
# Design&Test

NOVEMBER/DECEMBER 2017



ALEXANDER TORRES 2017

## Special Issue on Self-Awareness in Systems on Chip 2017

· Self-Awareness in Systems on Chip—A Survey
· Health Management for Self-Aware SoCs Based on IEEE 1687 Infrastructure
· KOCL: Power Self-Awareness for Arbitrary FPGA-SoC-Accelerated OpenCL Applications
· A Self-Aware Architecture for PVT Compensation and Power Nap in Near-Threshold Processors
· Self-Adaptive Timing Repair

CEDA — IEEE Council on Electronic Design Automation

CAS — IEEE CIRCUITS AND SYSTEMS SOCIETY

SSCS

tttc

◆ IEEE

# Contents

CECS

# SelPhyS 2019 and TCPS Special Issue



Call for Papers: Special Issue on Self-Awareness in Resource Constrained Cyber-Physical Systems

Inspired by biological examples, self-awareness has become a hot research topic in a variety of disciplines and its applicability has been explored in various application domains. The topic owes its attractiveness to its promise to facilitate highly resilient, adaptive and outstandingly efficient behaviors. Thus, self-awareness holds the promise to promote dependability in all types of smart gadgets and artificial agents in the interconnected world of future.

However, the challenges raised by these new promising features are also significant, not le... because they have a profound impact on the way we design, validate and test incorporating awareness. If a system smartly adapts to changing needs and environment, how do we vali... functionality at design time? How do we specify the correct functionality in the first place? W... the relevant trade-offs? How can we quantify uncertainties and variabilities in a meaningful... deal with them in the design process? These are only some of the pressing questions that h... addressed before these new features can be exploited.

The ACM Transactions on Cyber-Physical Systems seeks original manuscripts for a special i... "Self-Awareness in Resource Constrained Cyber-Physical Systems" which will cover recent development on methods, architecture, design, validation and application of resource-cons... cyber-physical systems that exhibit a degree of self-awareness.

**Submission Guidelines:**

Authors should submit their journal version at Manuscript Central adhering to the formatting instructions on the TCPS Web page, and indicate that you are submitting to the Special Issue on Self-Awareness in Resource Constrained Cyber-Physical Systems" on the first page and in the field "Author's Cover Letter:" in Manuscript Central). For additional questions, please send an email to any of the guest editors: p.lewis@aston.ac.uk, axel.jantsch@tuwien.ac.at, dutt@uci.edu.

**Submission Guidelines:**

Submission deadline:   7 September, 2018
Notification of First Round: 7 December, 2018
Submission of Revision: 8 February, 2019
Final Notification: 12 April 2019
Final Paper Due: 23 May 2019

**Guest Editors Contacts:**

Peter Lewis, p.lewis@aston.ac.uk
Axel Jantsch, axel.jantsch@tuwien.ac.at
Nikil Dutt, dutt@uci.edu

All ACM Journals | See Full Journal Index

# Ongoing Efforts

- More heterogeneity (CPU+GPU+DSP+NPU+FPGA+.....)
    - Reconfigure workloads at runtime to freely migrate between resources
    - Complex predictive models

- Distributed management
    - Propagating sensing info across non-coherent processing units

- Non-compute resources
    - Memory and I/O

# Ongoing Challenges

- ## Self-trained models
    - Add feedback for error correction
    - Challenging for models that are non-linear and/or based on heuristics

- ## Machine learning
    - Replacement for analytical/heuristic-based models ?
    - Unsupervised machine learning to mine sensing data and find patterns for optimizing policies or creating new ones

- ## Policy supervisors
    - Provide formal or stronger guarantees

# Acknowledgements

- Dutt Research Group
  - **Tiago Mück**, Amir Rahmani, Santanu Sarma, Majid Shoushtari, Bryan Donyanavard, Kasra Moazzemi, Roger Hsieh, Jurngyu Park, Hossein Tajik, Biswadip Maity

- Collaborating Faculty:
  - UCI: Fadi Kurdahi
  - Austria:  TU Wien: Prof. Axel Jantsch
  - Germany: TUM: Prof. Andreas Herkersdorf, TUB: Prof. Rolf Ernst

- NSF Information Processing Factory (IPF) project

CECS

# Questions?



**Dutt Research Group:   http://duttgroup.ics.uci.edu/**