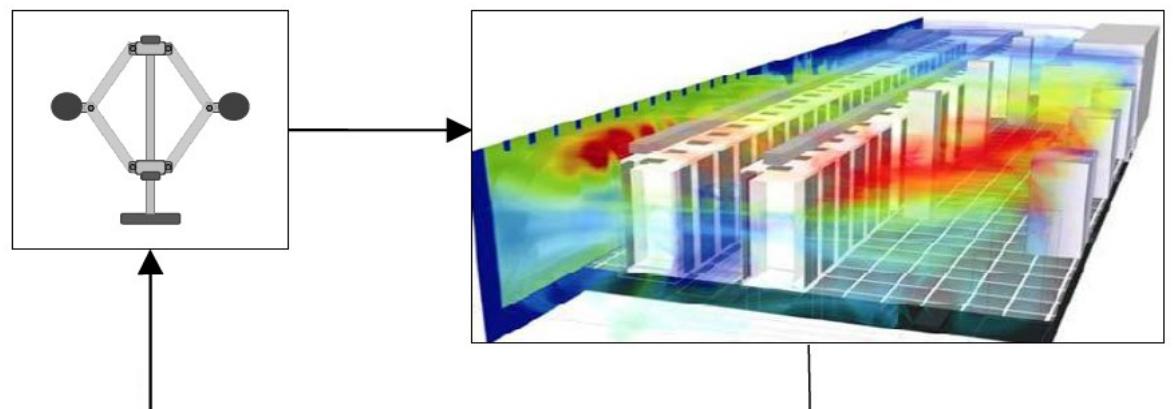


Models and Control Strategies for Data Center Energy Efficiency

Ph.D. Defense

Luca Parolini

Dept. of Electrical and Computer Engineering
Carnegie Mellon University



Feb. 27th 2012

Data center examples



■ Facebook's data center in North Carolina, US

- \$ 450 million project
- ~28,000 m² (300,000 ft²)
- Operated by 35 - 45 full-time employees

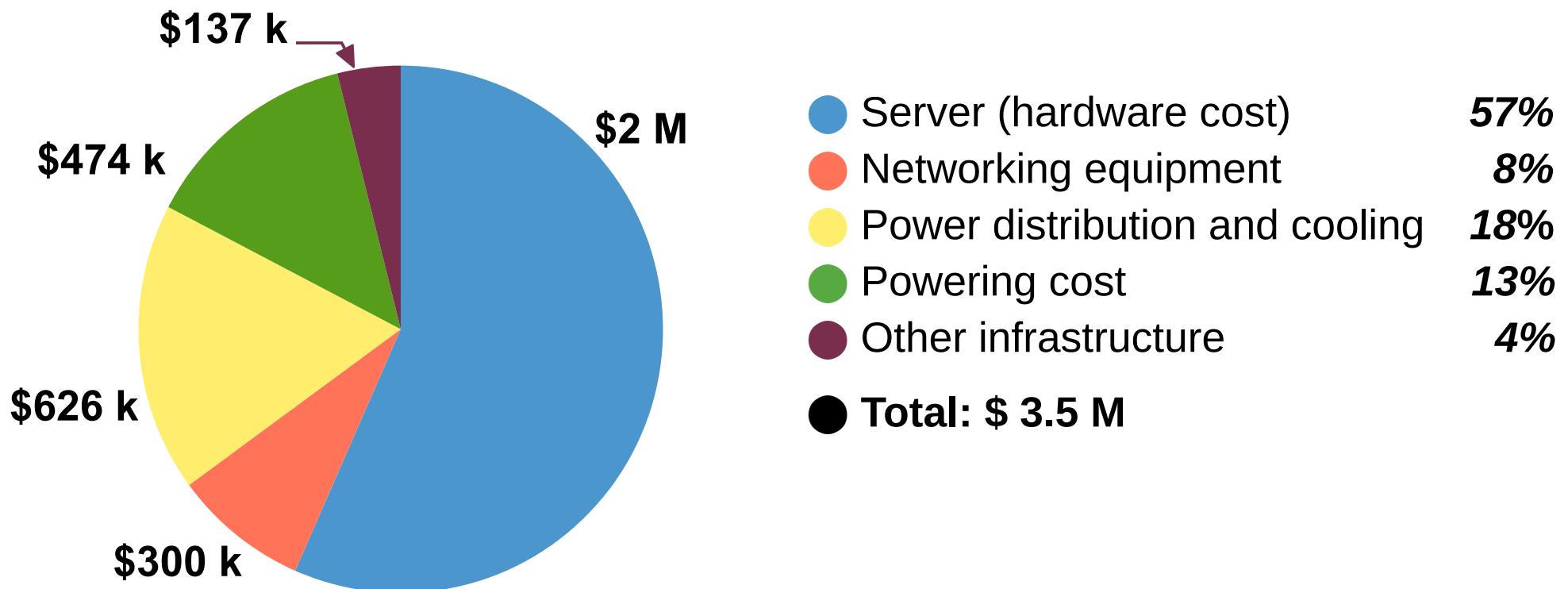
■ Racks

- Contain 42 (1U) servers in a rack
- 1U server: 480mm x 800mm x 44mm

	<i>Idle power</i>	<i>Peak power</i>
Server	200 W	350 W
Rack	8.4 kW	~15 kW

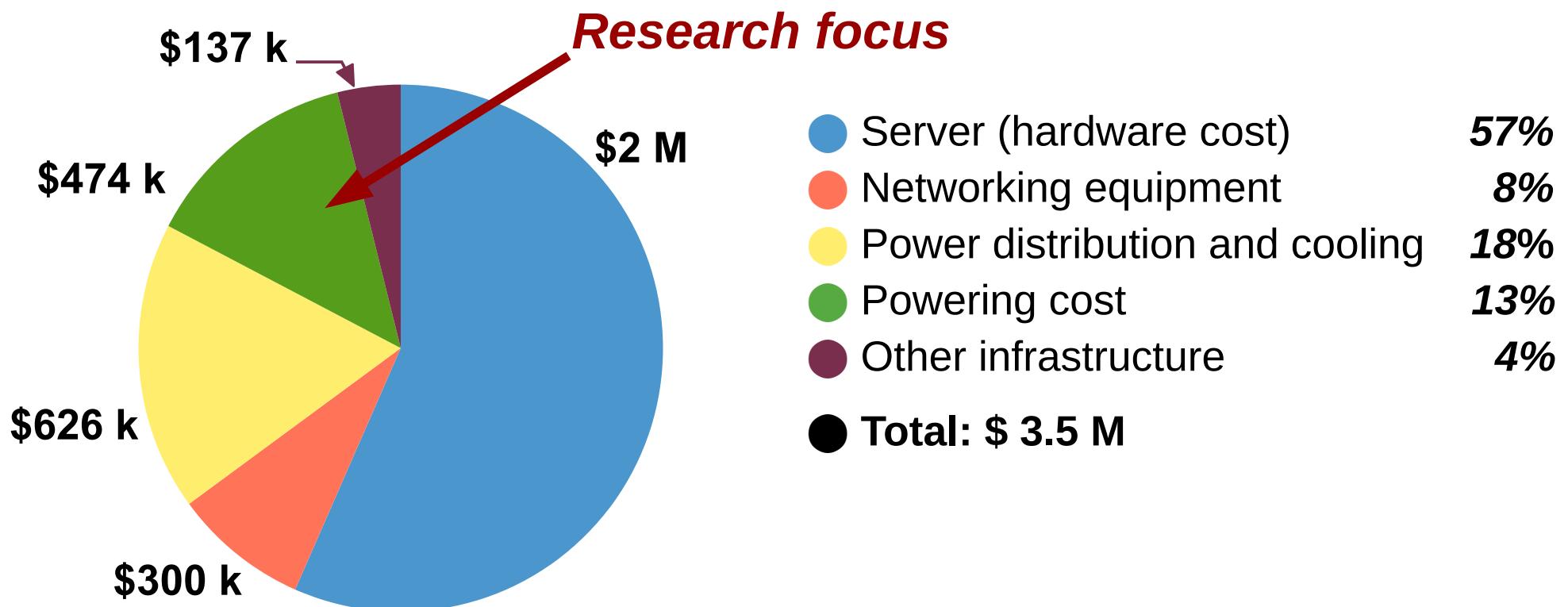
Monthly operating cost

- Large-scale facility, 50k servers
 - Facility cost amortized over 10 years
 - Server cost amortized over 3 years
 - Servers consume 70% of total power consumption

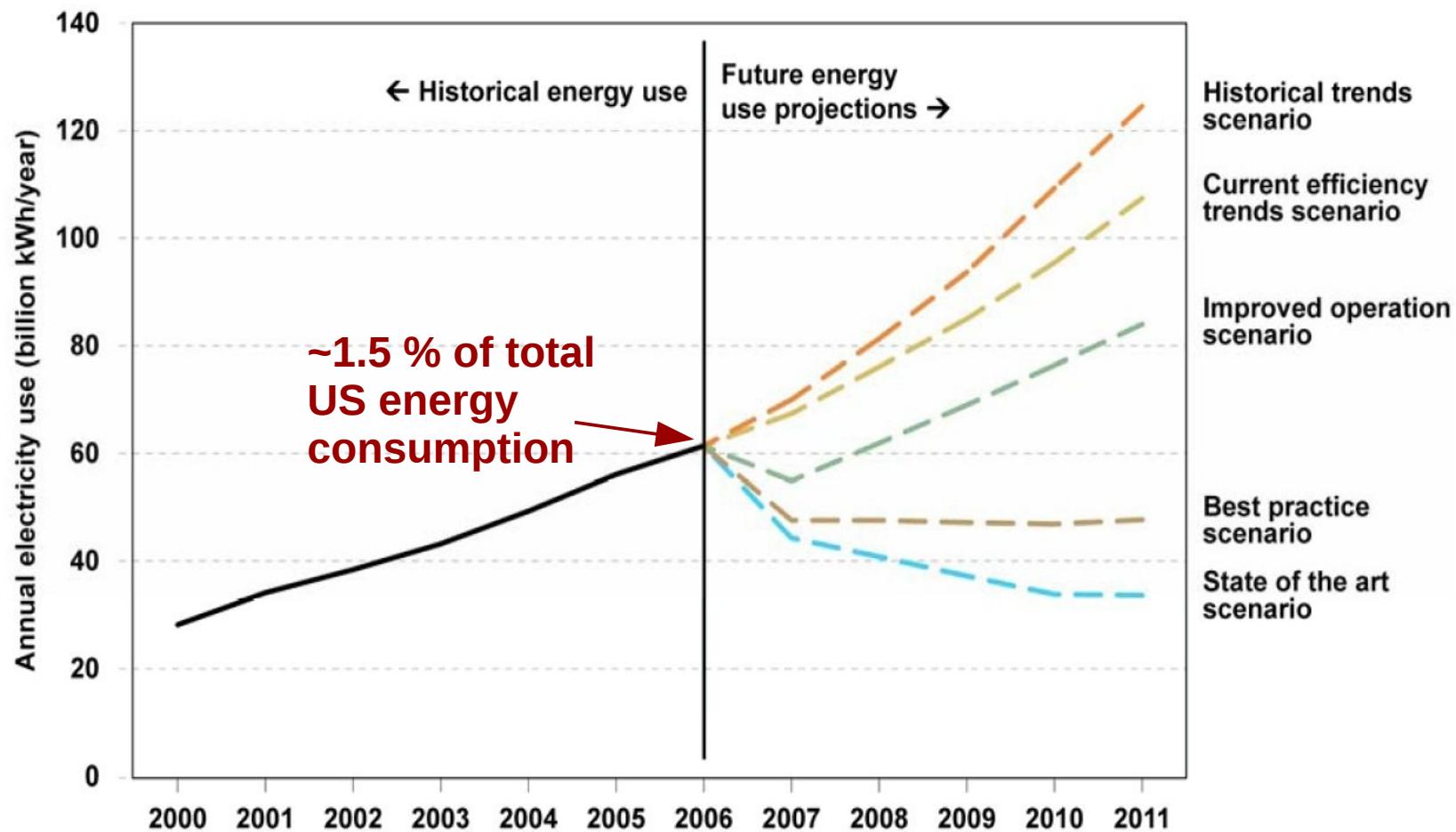


Monthly operating cost

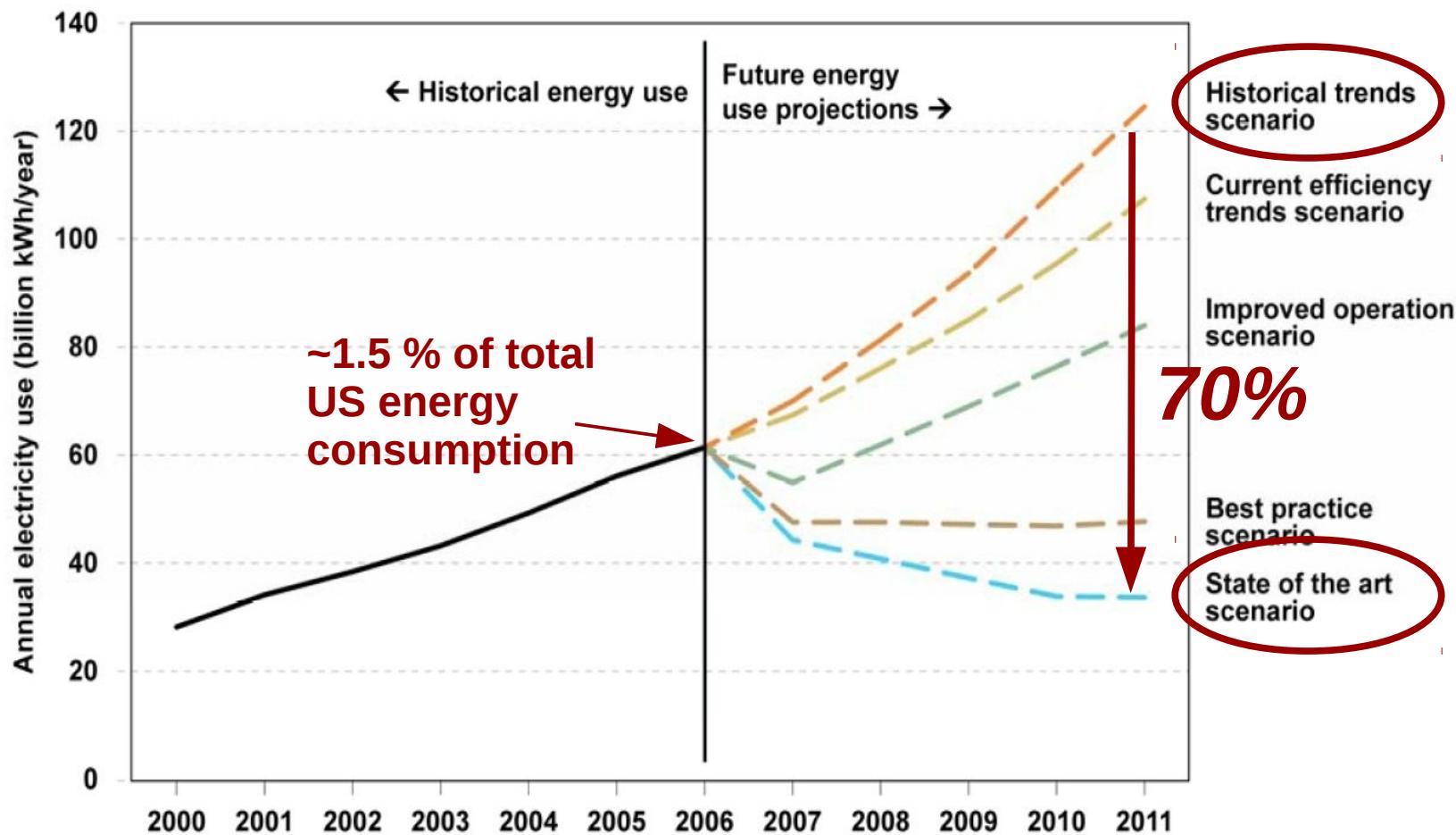
- Large-scale facility, 50k servers
 - Facility cost amortized over 10 years
 - Server cost amortized over 3 years
 - Servers consume 70% of total power consumption



Data center electricity consumption in US

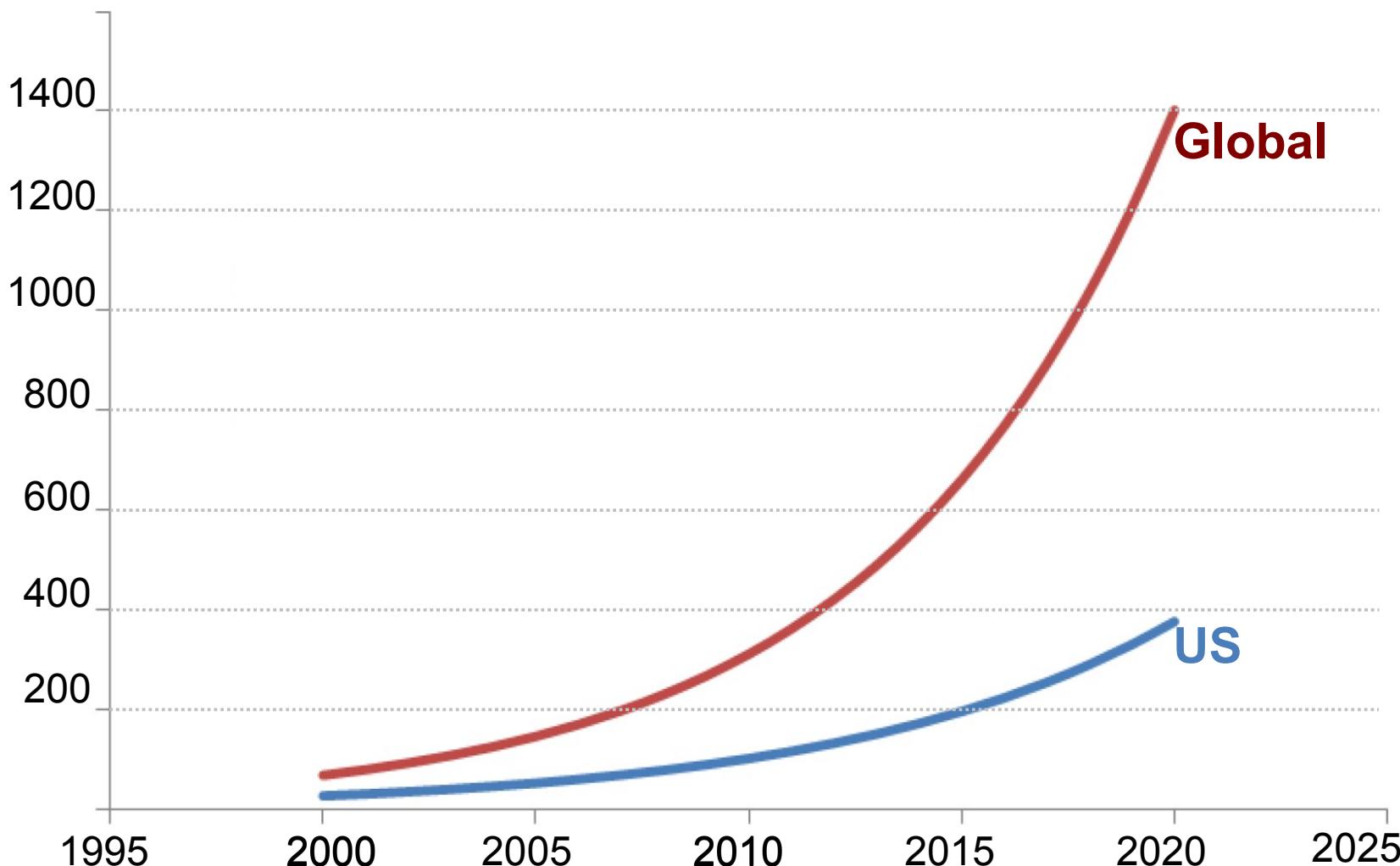


Data center electricity consumption in US



Data center electricity consumption

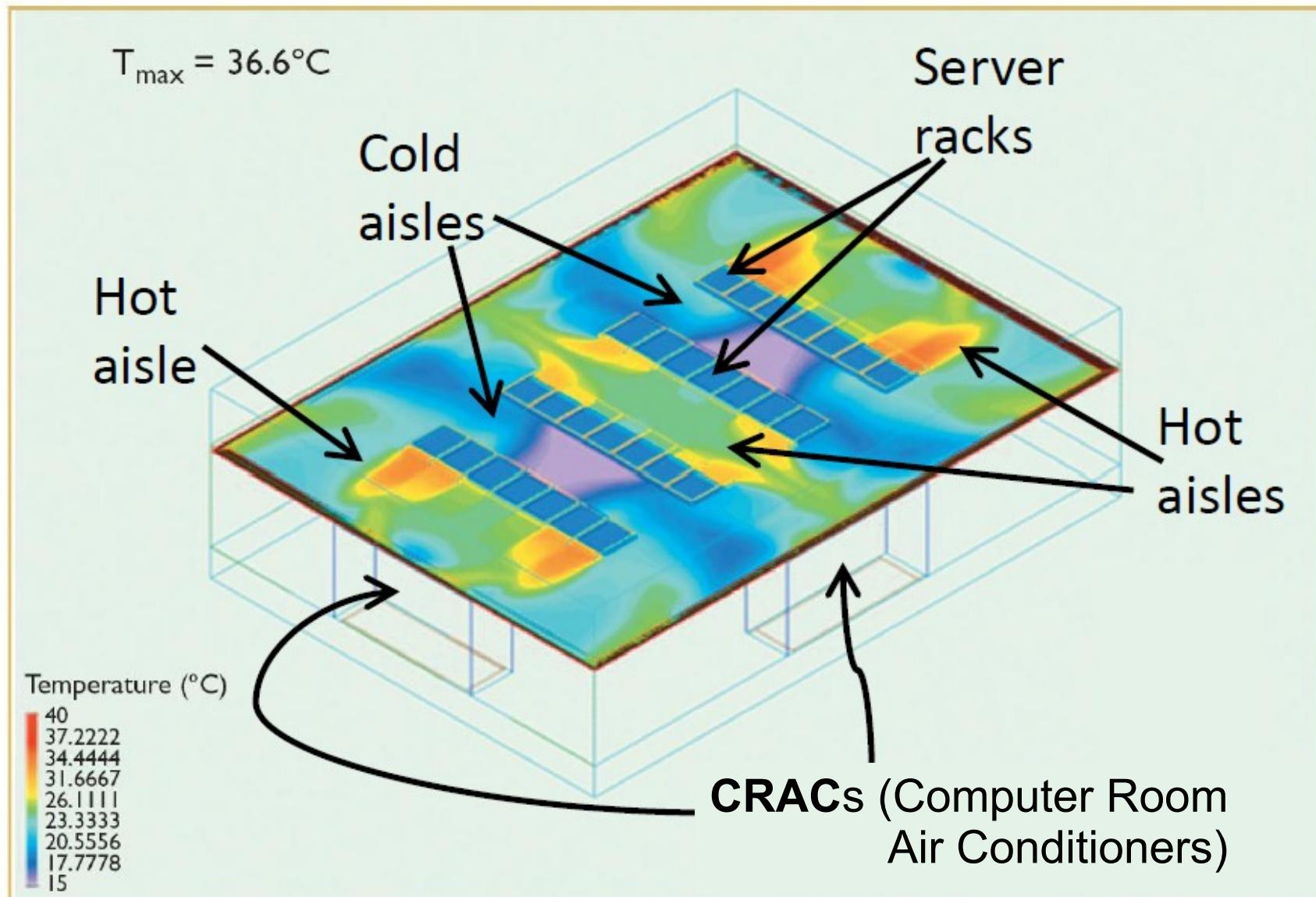
(Billion kWh / year)



Thermal constraints

- Without proper cooling, chip temperatures would exceed safe operation limits
 - *Almost complete transformation of electrical energy into heat*
- Chip temperature generally unobservable
- Industrial approach
 - Bound server *inlet air temperature*

Temperature distribution



Contributions of the dissertation

- **A modeling framework for data centers**
 - Facilitates the development of cyber-physical models of data centers
- **A control strategy**
 - Takes advantage of the modularity typically found in data centers
- **Multiple control-oriented models**
 - Make it possible to formulate and to solve, in real-time, optimal control problems
- **Multiple control algorithms**
 - Representative abstractions of different approaches to the data center control
- **A collection of MATLAB routines to simulate the evolution of a data center (Data center simulator)**
 - Makes it possible to analyze the performance of multiple controllers in a variety of scenarios

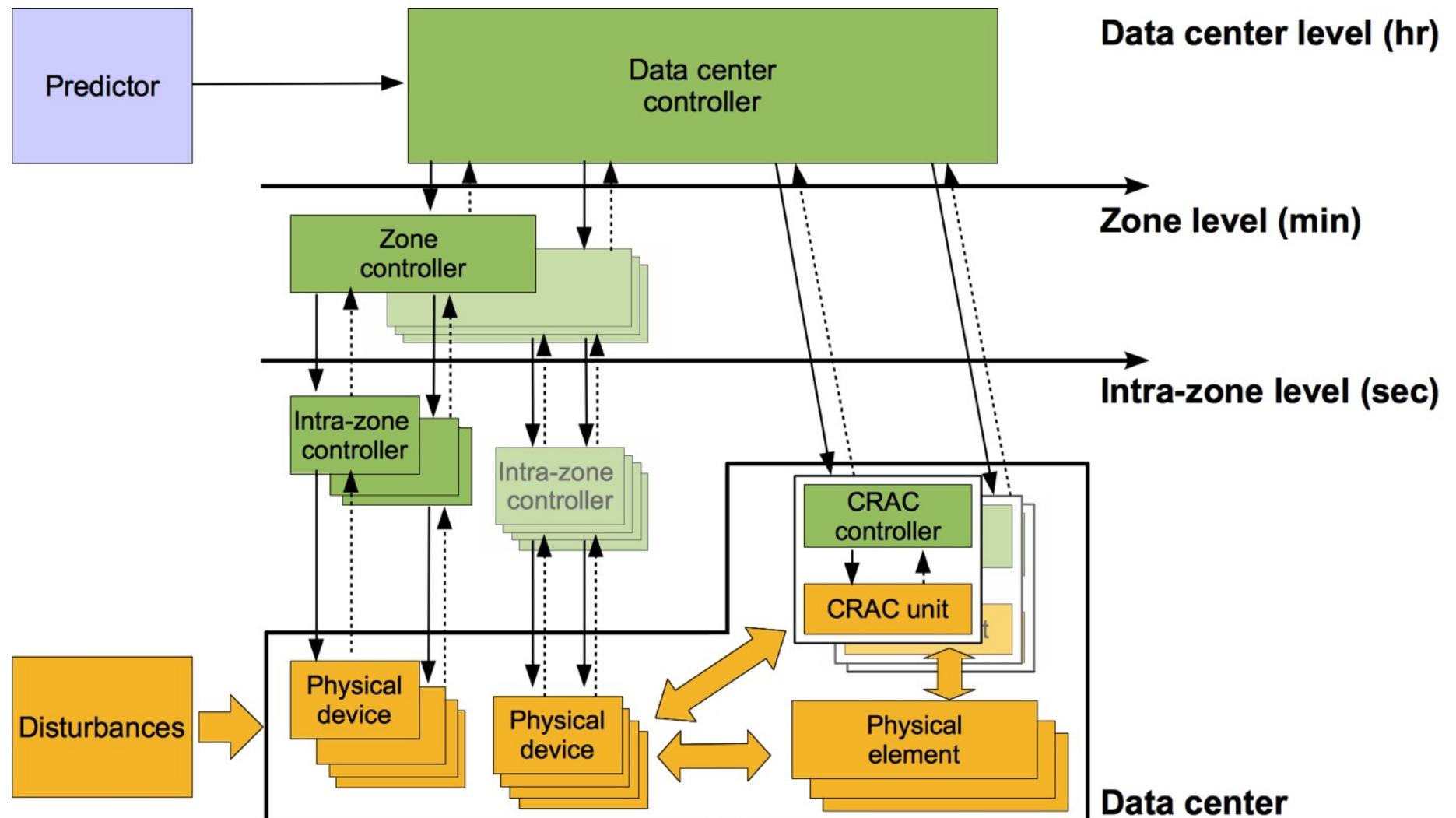
Contributions of the dissertation

- **A modeling framework for data centers**
 - Facilitates the development of cyber-physical models of data centers
- ***A control strategy***
 - Takes advantage of the modularity typically found in data centers
- ***Multiple control-oriented models***
 - Make it possible to formulate and to solve, in real-time, optimal control problems
- ***Multiple control algorithms***
 - Representative abstractions of different approaches to the data center control
- **A collection of MATLAB routines to simulate the evolution of a data center (Data center simulator)**
 - Makes it possible to analyze the performance of multiple controllers in a variety of scenarios

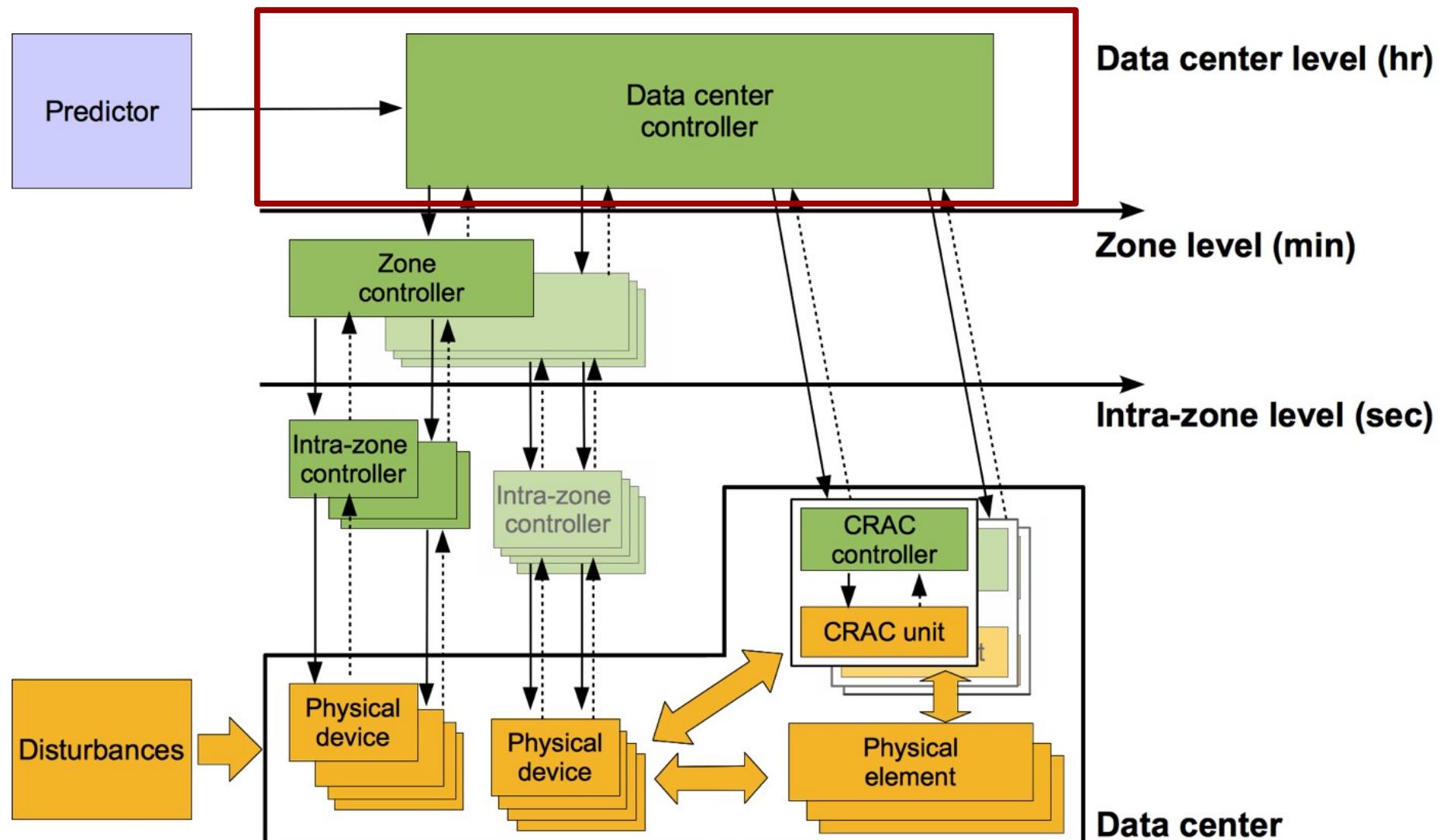
Outline

- Introduction
- Proposed control strategy
- Performance analysis for constant job arrival rate
- Interaction with the smart-grid
- Zone-level control
- Conclusion and future work

Hierarchical control strategy



Hierarchical control strategy



Data-center-level controller

- Focuses on processes in the hours time scale
 - Considers only the mean values of the predicted processes
- Groups servers into zones
 - *Power consumption of a zone is proportional to the amount of jobs executed*

Data-center-level controller

- Focuses on processes in the hours time scale
 - Considers only the mean values of the predicted processes
- Groups servers into zones
 - *Power consumption of a zone is proportional to the amount of jobs executed*
- Defines set-points for
 - Scheduling of jobs among the zones
 - Migration of jobs among the zones
 - Rate at which jobs have to be executed in every zone
 - Reference temperature of CRAC units
- Considers
 - Computational and thermal dynamics
 - Nonlinear efficiency of the CRAC units, i.e. coefficient of performance (COP)

Proposed data-center-level controllers

■ Baseline controller

- Open-loop controller
- Set-points set for the worst-case scenario
- Lower bound of the performance of other control approaches

Proposed data-center-level controllers

■ Baseline controller

- Open-loop controller
- Set-points set for the worst-case scenario
- Lower bound of the performance of other control approaches

■ Uncoordinated controller

- Based on model-predictive-control (MPC)
- Neglects the coupling between the cyber and the physical parts
- Represents the control approaches typically found in modern data centers

Proposed data-center-level controllers

■ Baseline controller

- Open-loop controller
- Set-points set for the worst-case scenario
- Lower bound of the performance of other control approaches

■ Uncoordinated controller

- Based on model-predictive-control (MPC)
- Neglects the coupling between the cyber and the physical parts
- Represents the control approaches typically found in modern data centers

■ Coordinated controller

- Based on MPC
- Considers the coupling between the cyber and the physical parts
- Takes full advantage of a cyber-physical model of a data center

Proposed data-center-level controllers

■ Baseline controller

- Open-loop controller
- Set-points set for the worst-case scenario
- Lower bound of the performance of other control approaches

■ Uncoordinated controller

- Based on model-predictive-control (MPC)
- Neglects the coupling between the cyber and the physical parts
- Represents the control approaches typically found in modern data centers

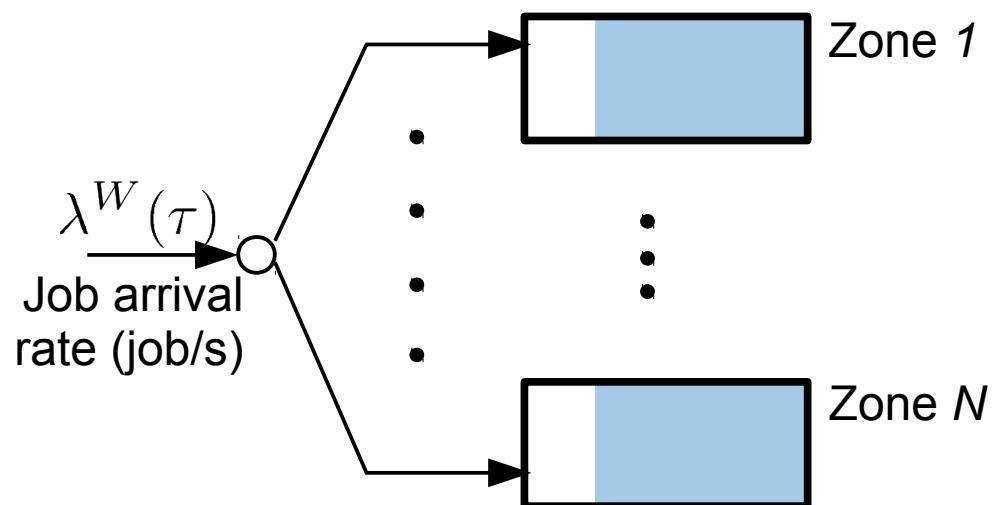
■ Coordinated controller

- Based on MPC
- Considers the coupling between the cyber and the physical parts
- Takes full advantage of a cyber-physical model of a data center

■ *The different controllers help us to understand the impact of considering the cyber-physical nature of data centers*

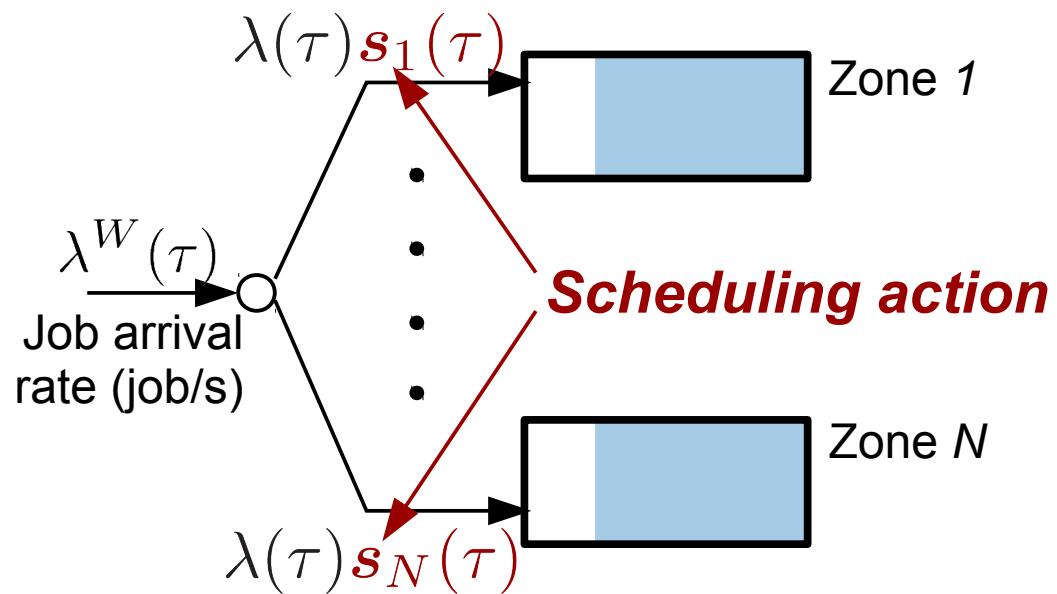
Computational network

- Describes the process of data exchange within the data center and between the data center and external users
- Based on a fluid approximation of the job execution and arrival processes
 - First-order approximation of a queuing system



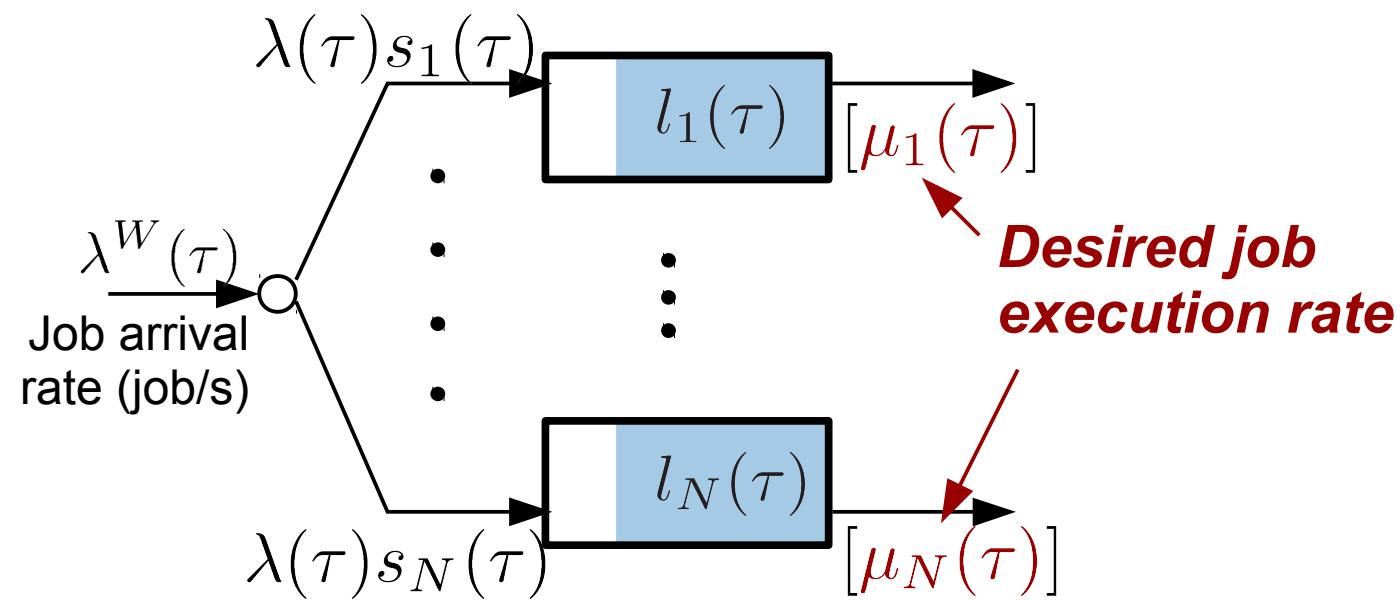
Computational network

- Describes the process of data exchange within the data center and between the data center and external users
- Based on a fluid approximation of the job execution and arrival processes
 - First-order approximation of a queuing system



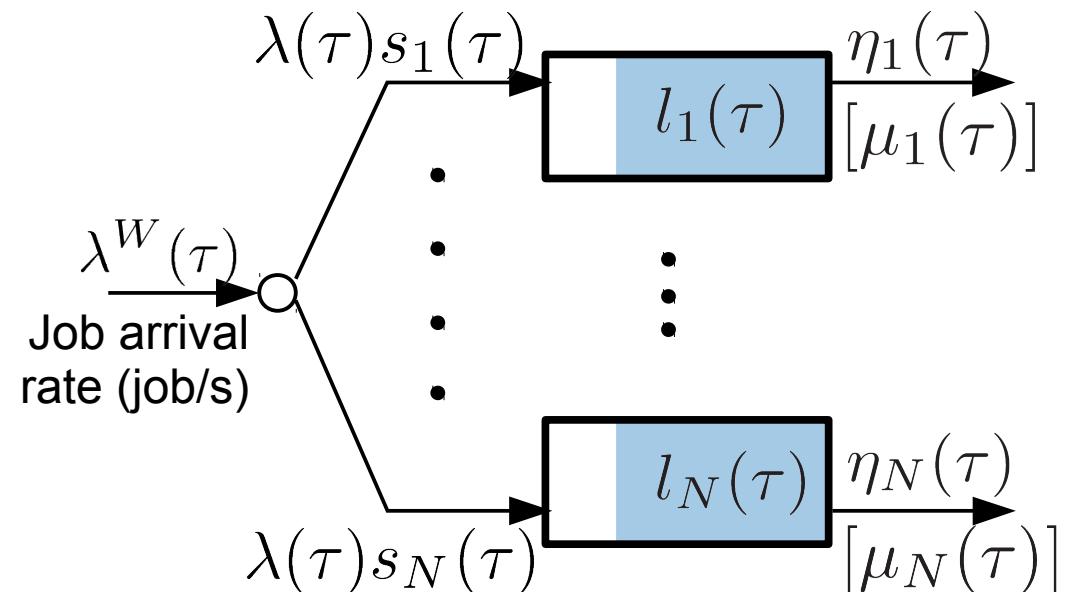
Computational network

- Describes the process of data exchange within the data center and between the data center and external users
- Based on a fluid approximation of the job execution and arrival processes
 - First-order approximation of a queuing system



Computational network

- Describes the process of data exchange within the data center and between the data center and external users
- Based on a fluid approximation of the job execution and arrival processes
 - First-order approximation of a queuing system



$$a_1(\tau) = \lambda(\tau)s_1(\tau)$$

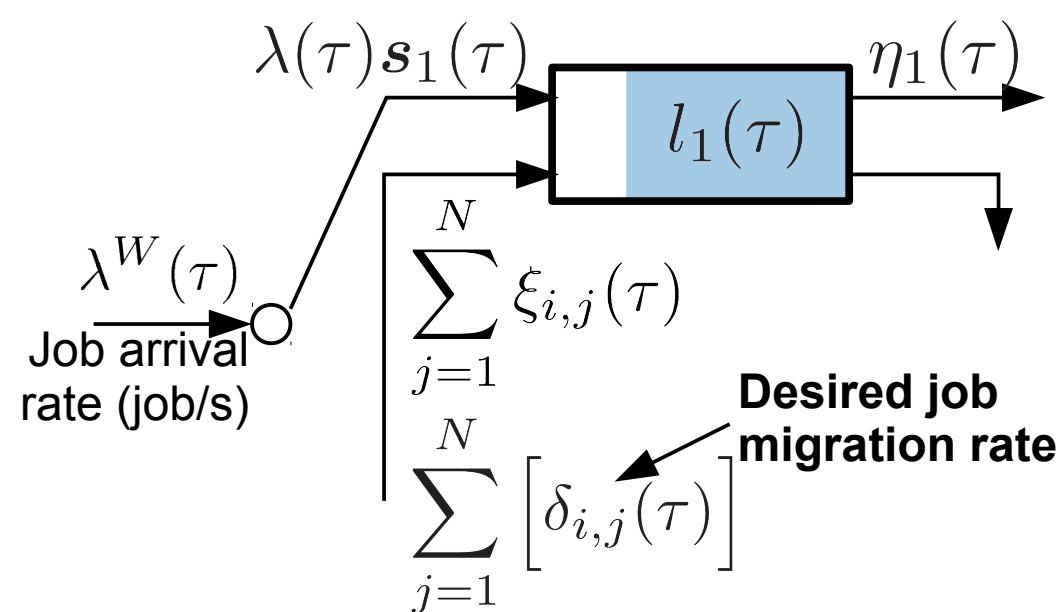
$$d_1(\tau) = \eta_1(\tau)$$

$$\dot{l}_1(\tau) = a_1(\tau) - d_1(\tau)$$

$$\eta_1(\tau) = \begin{cases} \mu_1(\tau) & \text{if } l_1(\tau) > 0 \\ a_1(\tau) & \text{or } a_1(\tau) > \mu_1(\tau) \\ a_1(\tau) & \text{otherwise} \end{cases}$$

Computational network

- Describes the process of data exchange within the data center and between the data center and external users
- Based on a fluid approximation of the job execution and arrival processes
 - First-order approximation of a queuing system



$$a_1(\tau) = \lambda(\tau)s_1(\tau) + \sum_{j=1}^N \xi_{1,j}(\tau)$$

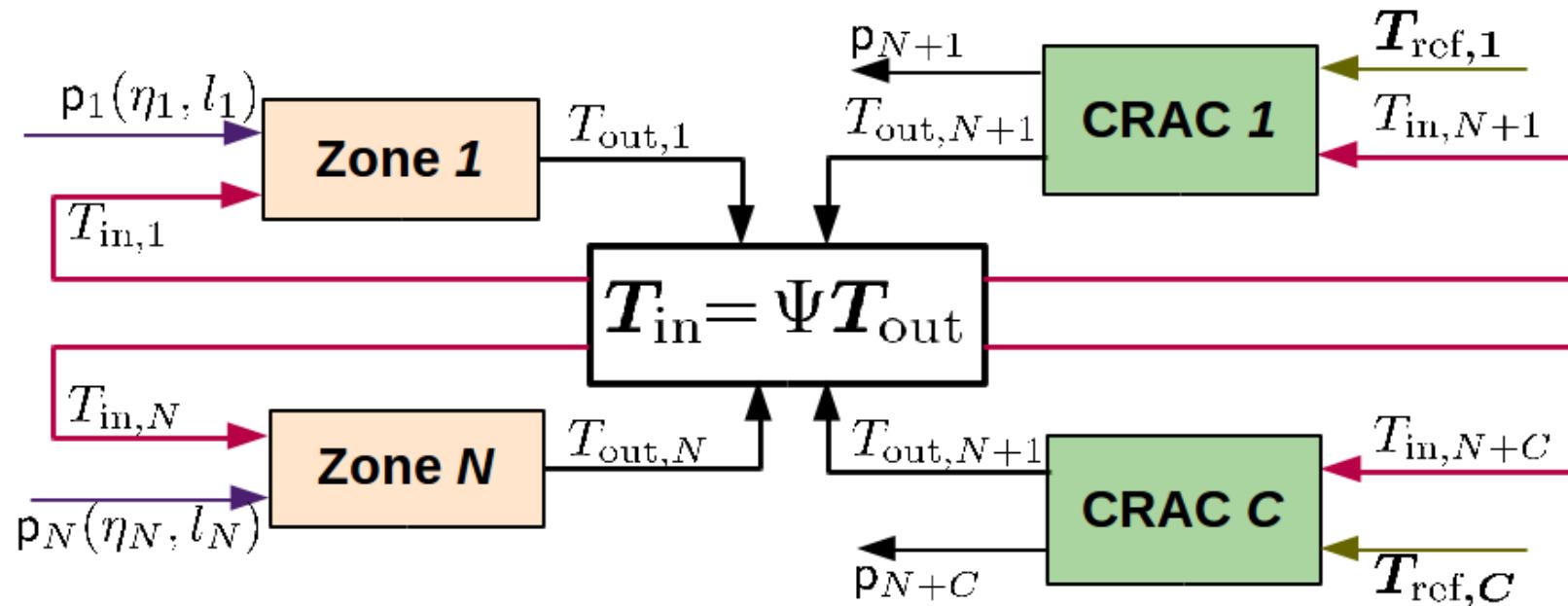
$$d_1(\tau) = \eta_1(\tau) + \sum_{j=1}^N \xi_{j,1}(\tau)$$

$$\dot{l}_1(\tau) = a_1(\tau) - d_1(\tau)$$

$$\eta_1(\tau) = \begin{cases} \mu_1(\tau) & \text{if } l_1(\tau) > 0 \\ a_1(\tau) & \text{otherwise} \end{cases}$$

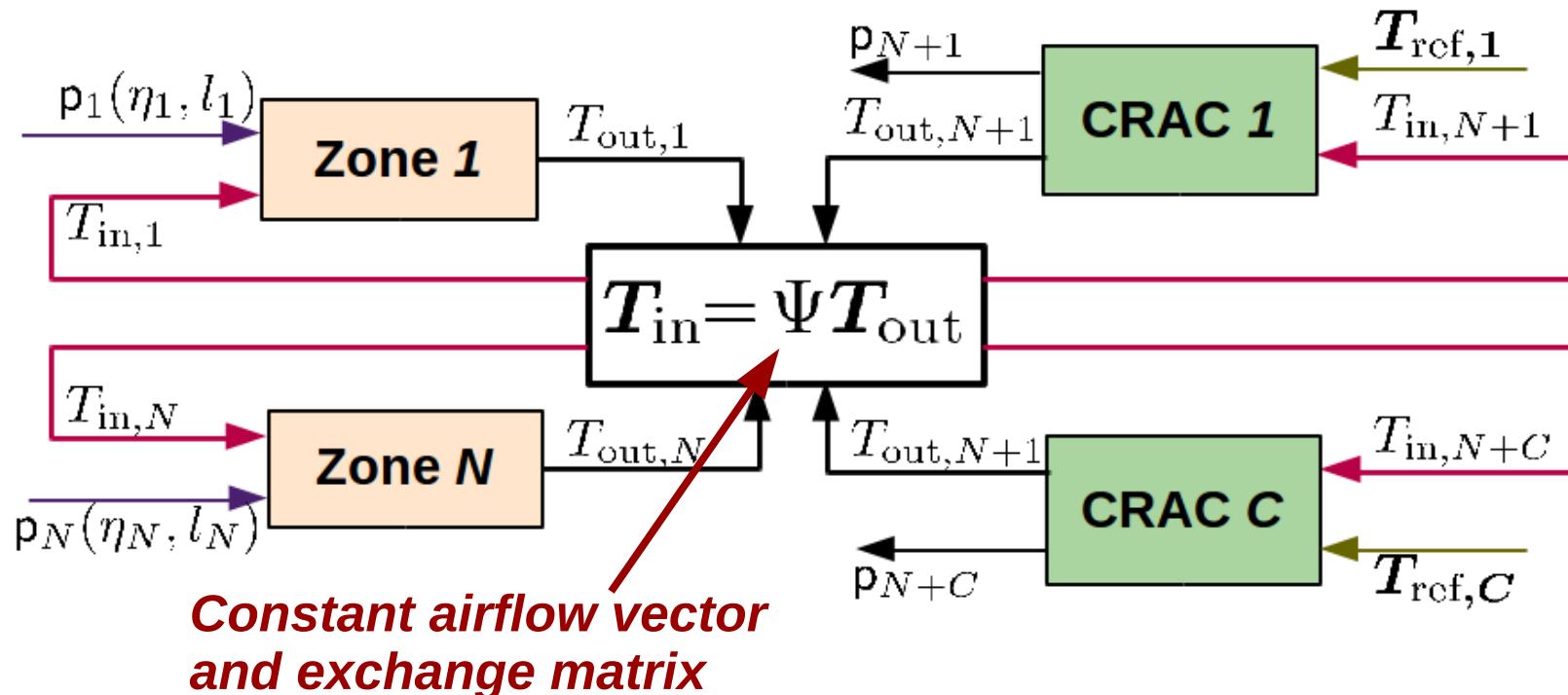
Thermal network

- Describes the process of energy exchange within the data center and between the data center and the external environment



Thermal network

- Describes the process of energy exchange within the data center and between the data center and the external environment



- Inlet temperature constraint
- CRAC power consumption

$$T_{in}(\tau) = \Psi T_{out}(\tau) \leq \bar{T}_{in}$$

$$p_i(t) = \frac{\dot{Q}_i(t)}{COP_i(T_{out,i}(t))}$$

Heat removed rate (W)

Variables at step k

- No environment nodes considered in this table

		Variables	
Input	Controllable	Job scheduling	$s(k)$
		Desired job execution rate	$\mu(k)$
		Desired job migration rate	$\delta(k)$
		CRAC unit reference temperature	$T_{ref}(k)$
	Uncontrollable	Job arrival rate	$\lambda^w(k)$
Output	Job departure rates from every zone		$\eta(k)$
	Power consumption of CRAC nodes		$p_c(k)$
	Input temperatures of zones and CRAC nodes		$T_{in}(k)$
	Zone power consumption		$p_N(k)$
State	Number of jobs in every zones		$I(k)$
	Output temperatures of zones and CRAC nodes		$T_{out}(k)$

Baseline & uncoordinated controllers

■ Baseline controller

$$\mu(k) = \bar{\mu} \quad \delta(k) = 0 \quad s(k) = 1 \frac{1}{N} \quad T_{\text{ref}}(k) = \underline{T}_{\text{ref}}$$

Baseline & uncoordinated controllers

■ Baseline controller

$$\mu(k) = \bar{\mu} \quad \delta(k) = 0 \quad s(k) = 1 \frac{1}{N} \quad T_{\text{ref}}(k) = \underline{T}_{\text{ref}}$$

■ Uncoordinated controller

$$\min_{\mathcal{M}, \mathcal{S}, \mathcal{D}} \sum_{h=k}^{k+\mathcal{T}-1} \mathbf{1}^T \hat{\mathbf{p}}_{\mathcal{N}}(h|k)$$

Minimize expected zone power consumption

s.t. $\hat{l}(k|k) = l(k)$

for all $h = k, \dots, k + \mathcal{T} - 1$

computational dynamics

QoS constraints

$0 \leq \hat{\mu}(h|k) \leq \bar{\mu}, \quad \hat{\delta}(h|k) = 0$

$0 \leq \hat{s}(h|k) \leq 1, \quad B_s \hat{s}(h|k) = 1$

- Computational dynamics
- Quality of service (QoS) constraints
- Control constraints

$$\mathcal{D} = \left\{ \hat{\delta}(k|k), \dots, \hat{\delta}(k+\mathcal{T}-1|k) \right\} \quad \mathcal{T}_{\text{ref}} = \left\{ \hat{T}_{\text{ref}}(k|k), \dots, \hat{T}_{\text{ref}}(k+\mathcal{T}-1|k) \right\}$$

$$\mathcal{M} = \left\{ \hat{\mu}(k|k), \dots, \hat{\mu}(k+\mathcal{T}-1|k) \right\} \quad \mathcal{S} = \left\{ \hat{s}(k|k), \dots, \hat{s}(k+\mathcal{T}-1|k) \right\}$$

Baseline & uncoordinated controllers

■ Baseline controller

$$\mu(k) = \bar{\mu} \quad \delta(k) = 0 \quad s(k) = 1 \frac{1}{N} \quad T_{\text{ref}}(k) = \underline{T}_{\text{ref}}$$

■ Uncoordinated controller

$$\min_{\mathcal{M}, \mathcal{S}, \mathcal{D}} \sum_{h=k}^{k+\mathcal{T}-1} \mathbf{1}^T \hat{\mathbf{p}}_{\mathcal{N}}(h|k)$$

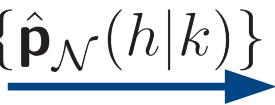
s.t. $\hat{l}(k|k) = l(k)$
for all $h = k, \dots, k + \mathcal{T} - 1$

computational dynamics

QoS constraints

$$0 \leq \hat{\mu}(h|k) \leq \bar{\mu}, \quad \hat{\delta}(h|k) = 0$$

$$0 \leq \hat{s}(h|k) \leq 1, \quad B_s \hat{s}(h|k) = 1$$



$$\{\hat{\mathbf{p}}_{\mathcal{N}}(h|k)\}$$

$$\min_{\mathcal{T}_{\text{ref}}} \sum_{h=k}^{k+\mathcal{T}-1} \mathbf{1}^T \hat{\mathbf{p}}_{\mathcal{C}}(h|k)$$

s.t. $\hat{T}_{\text{out}}(k|k) = T_{\text{out}}(k)$
for all $h = k, \dots, k + \mathcal{T} - 1$

thermal dynamics

$$\begin{aligned} \underline{T}_{\text{ref}} \leq \hat{T}_{\text{ref}}(h|k) &\leq \bar{T}_{\text{ref}} \\ \hat{T}_{\text{in}}(h+1|k) &\leq \bar{T}_{\text{in}} \end{aligned}$$

$$\mathcal{D} = \left\{ \hat{\delta}(k|k), \dots, \hat{\delta}(k+\mathcal{T}-1|k) \right\}$$

$$\mathcal{M} = \left\{ \hat{\mu}(k|k), \dots, \hat{\mu}(k+\mathcal{T}-1|k) \right\}$$

$$\mathcal{T}_{\text{ref}} = \left\{ \hat{T}_{\text{ref}}(k|k), \dots, \hat{T}_{\text{ref}}(k+\mathcal{T}-1|k) \right\}$$

$$\mathcal{S} = \left\{ \hat{s}(k|k), \dots, \hat{s}(k+\mathcal{T}-1|k) \right\}$$

Coordinated controller

- Considers the computational and the thermal dynamics in the same optimization problem

$$\min_{\mathcal{M}, \mathcal{S}, \mathcal{D}, \mathcal{T}_{ref}} \left(\sum_{h=k}^{k+\mathcal{T}-1} \mathbf{1}^T \hat{\mathbf{p}}_{\mathcal{N}}(h|k) + \mathbf{1}^T \hat{\mathbf{p}}_{\mathcal{C}}(h|k) \right)$$

← **Minimize expected zone and CRAC power consumption**

s.t. $\hat{\mathbf{l}}(k|k) = \mathbf{l}(k), \quad \hat{\mathbf{T}}_{out}(k|k) = \mathbf{T}_{out}(k)$
for all $h = k, \dots, k + \mathcal{T} - 1$
computational dynamics,
thermal dynamics,
QoS constraints,

} ← **Considers both thermal and computational dynamics**

$$\begin{aligned} \mathbf{0} \leq \hat{\boldsymbol{\mu}}(h|k) \leq \overline{\boldsymbol{\mu}}, \quad \hat{\boldsymbol{\delta}}(h|k) = \mathbf{0} \\ \mathbf{0} \leq \hat{\mathbf{s}}(h|k) \leq \mathbf{1}, \quad B_s \hat{\mathbf{s}}(h|k) = \mathbf{1}, \\ \underline{\mathbf{T}}_{ref} \leq \hat{\mathbf{T}}_{ref}(h|k) \leq \overline{\mathbf{T}}_{ref}, \quad \hat{\mathbf{T}}_{in}(h+1|k) \leq \overline{\mathbf{T}}_{in}, \\ \hat{\mathbf{p}}(h|k) = A_\alpha \hat{\boldsymbol{\eta}}(h|k) + B_\beta \hat{\mathbf{l}}(h|k) \end{aligned}$$

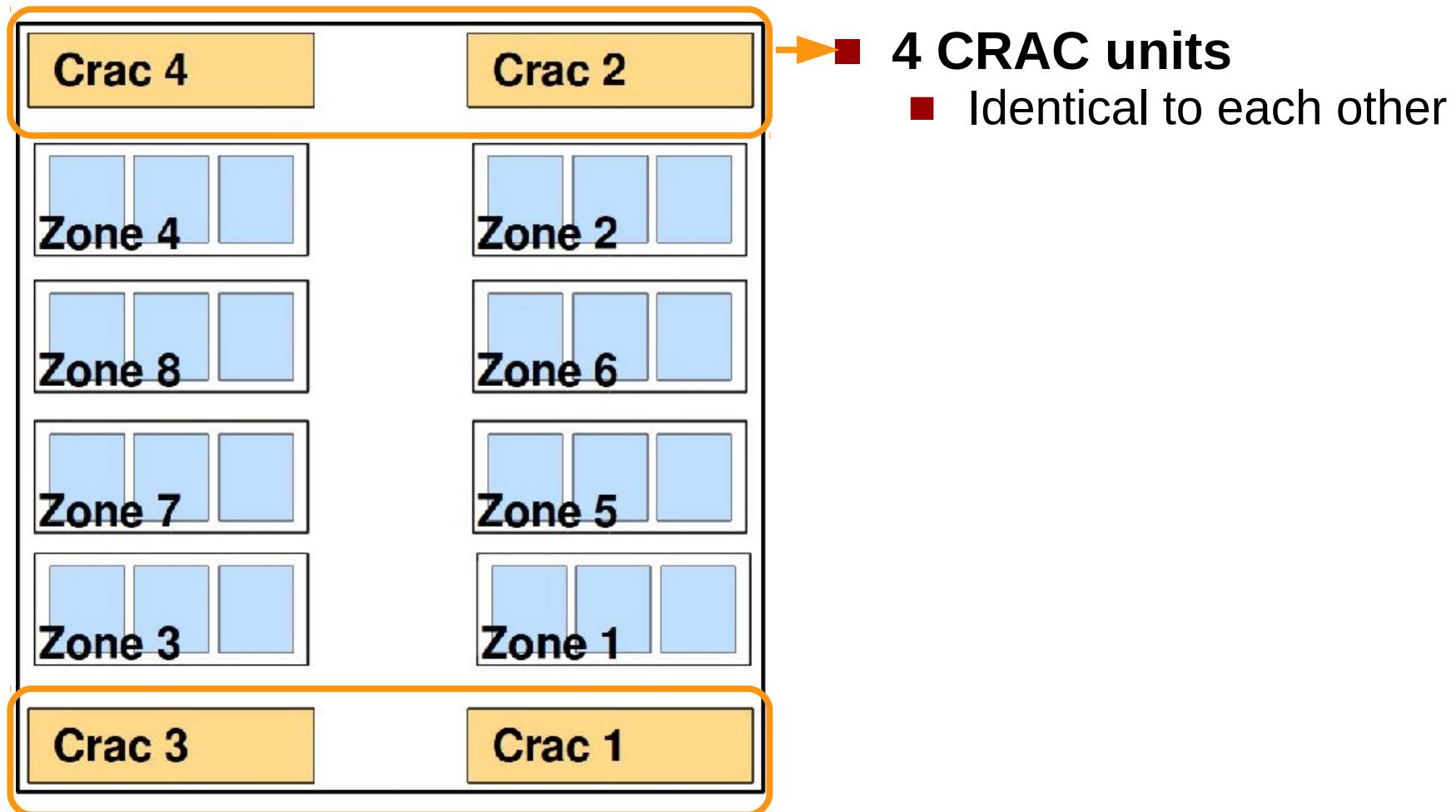
← **Thermal-computational coupling**

Outline

- **Introduction**
- **Proposed control strategy**
- **Performance analysis for constant job arrival rate**
- **Interaction with the smart-grid**
- **Zone-level control**
- **Conclusion and future work**

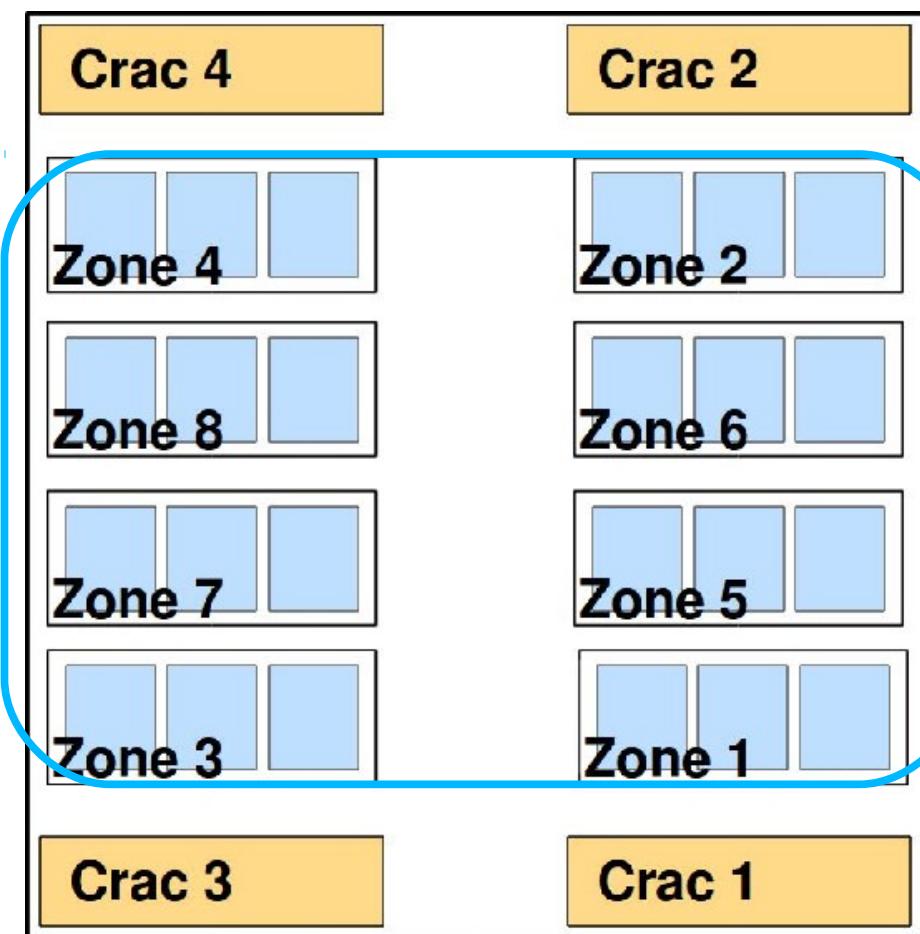
Data center layout

- Highlights the effects of choosing energy efficient servers or efficient cooled servers



Data center layout

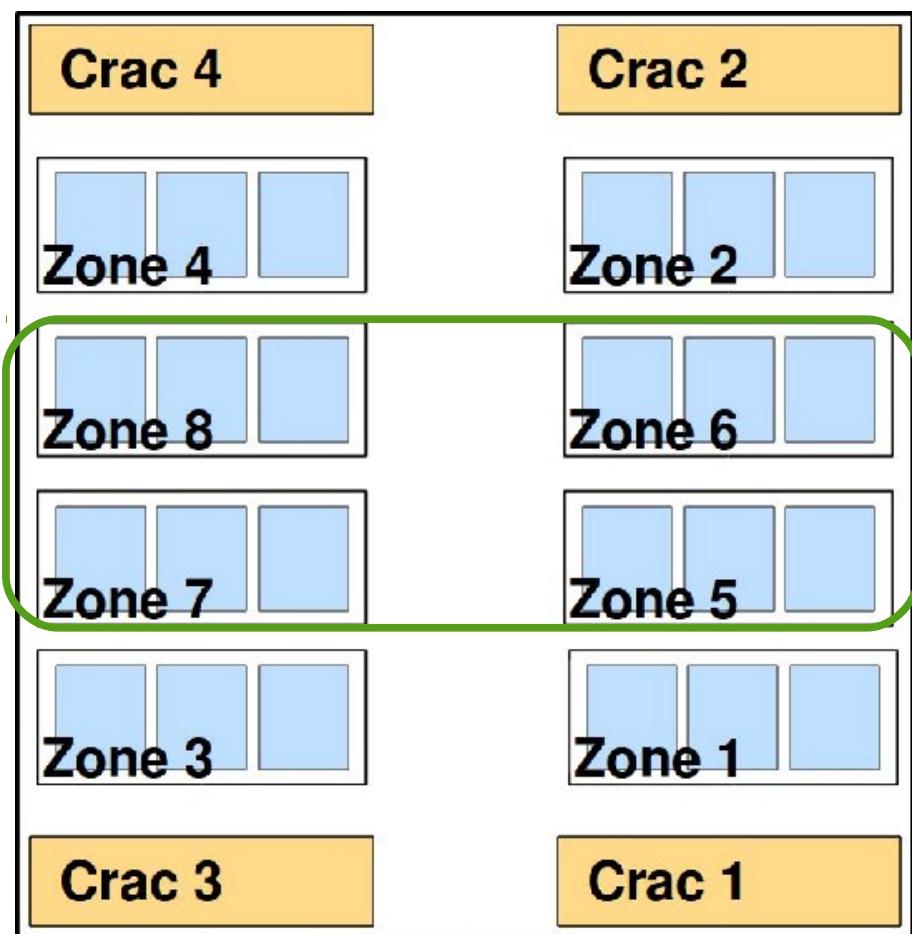
- Highlights the effects of choosing energy efficient servers or efficient cooled servers



- 4 CRAC units
 - Identical to each other
- 8 Zones
 - 3 Racks each (126 servers per zone)

Data center layout

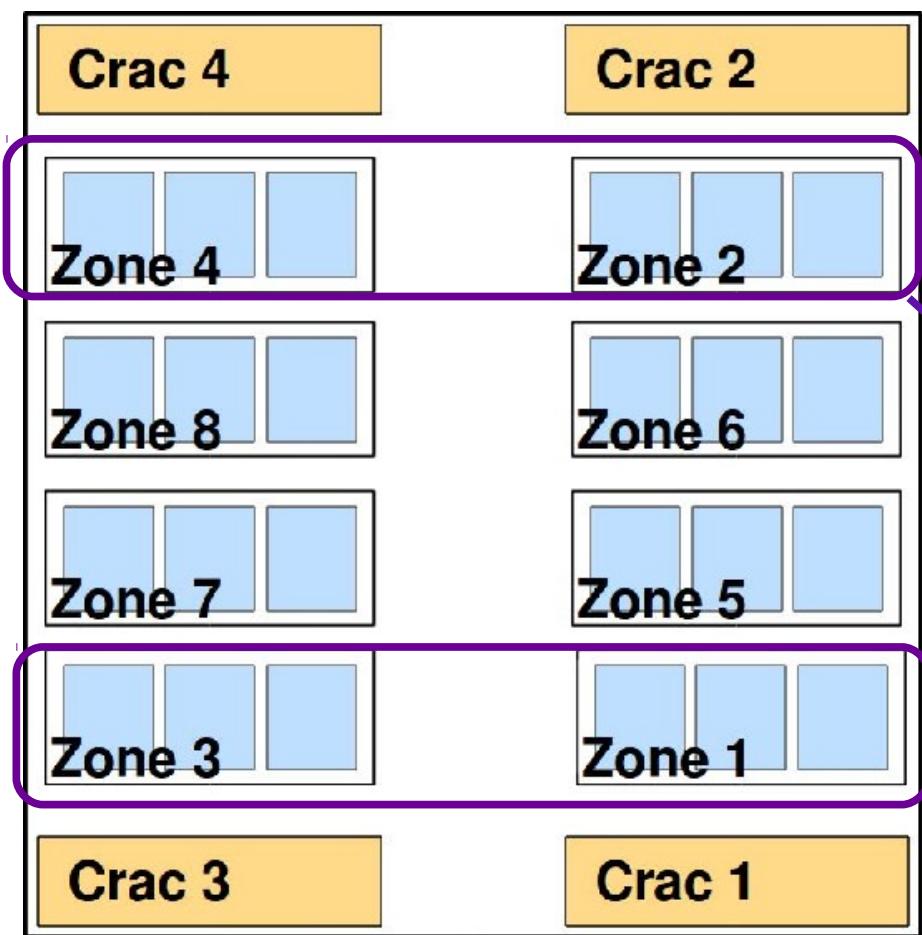
- Highlights the effects of choosing energy efficient servers or efficient cooled servers



- 4 CRAC units
 - Identical to each other
- 8 Zones
 - 3 Racks each (126 servers per zone)
 - Energy efficient servers

Data center layout

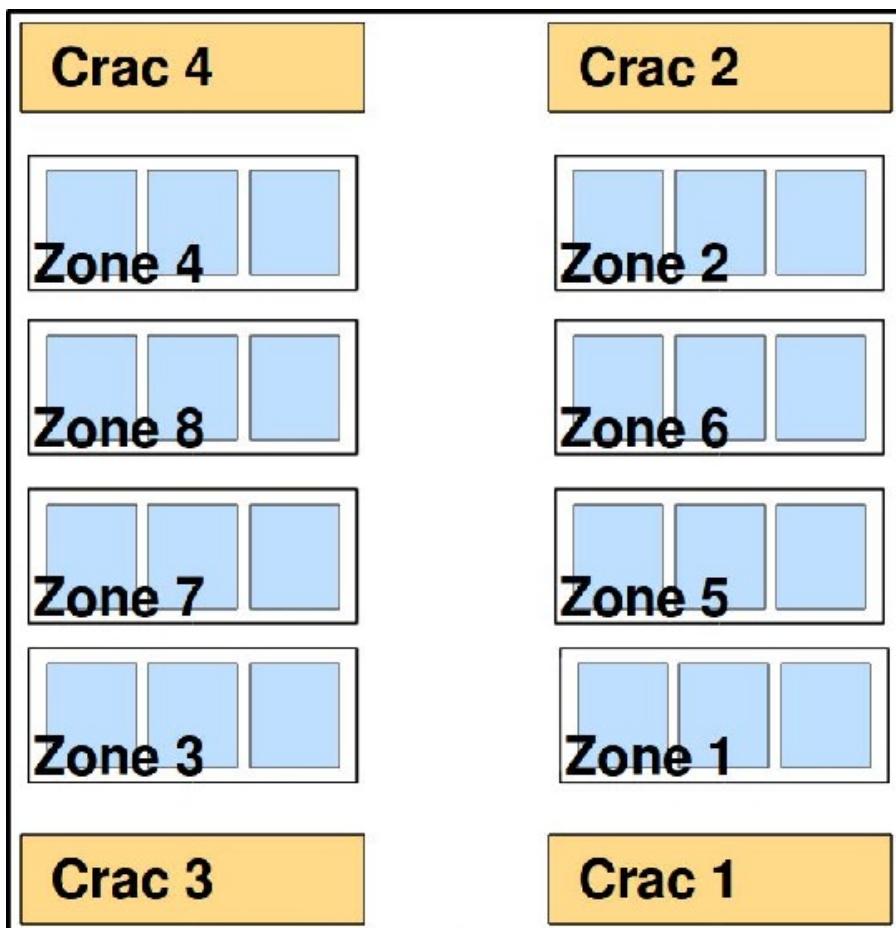
- Highlights the effects of choosing energy efficient servers or efficient cooled servers



- 4 CRAC units
 - Identical to each other
- 8 Zones
 - 3 Racks each (126 servers per zone)
 - Energy efficient servers
 - Efficiently cooled servers

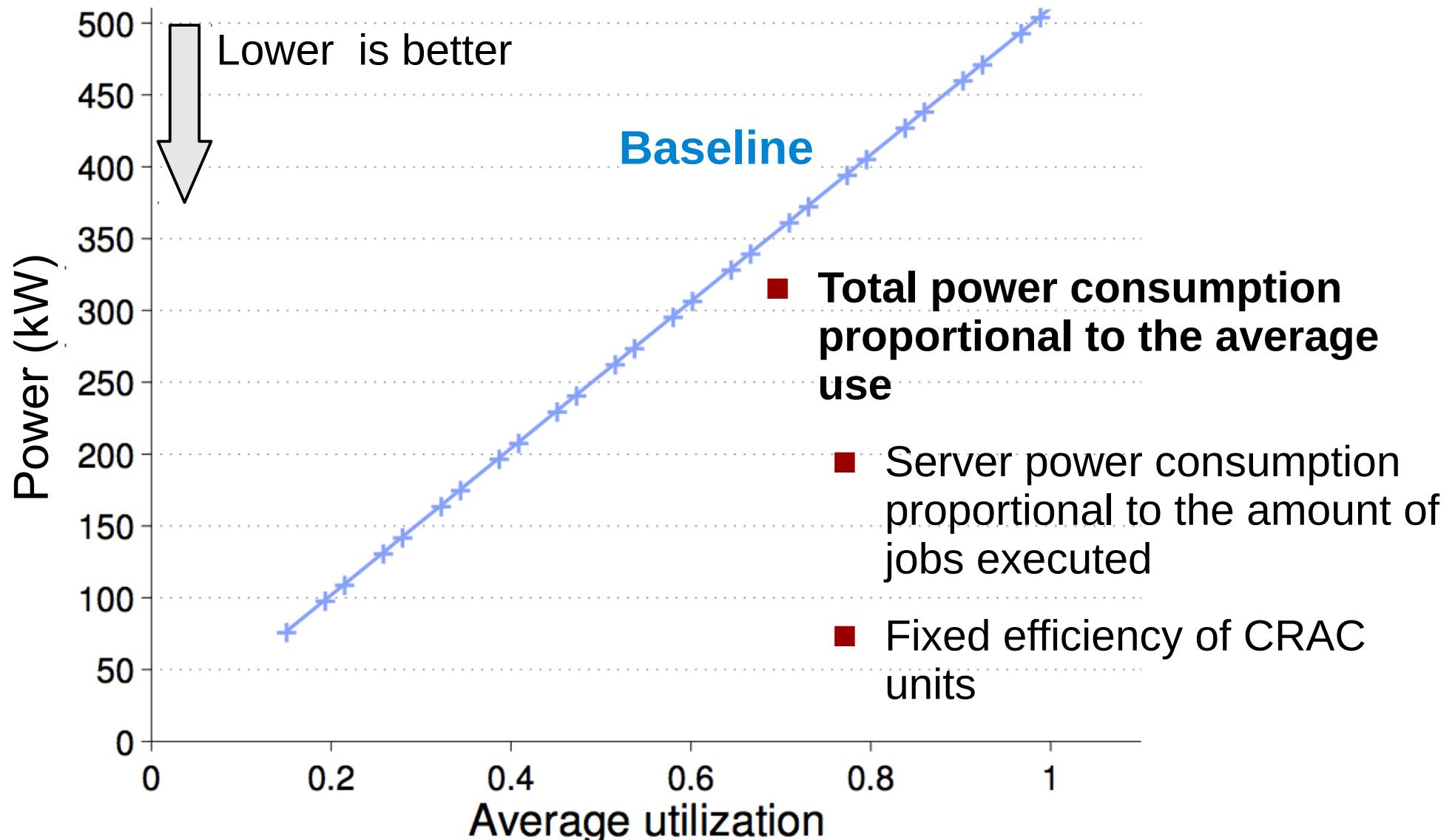
Data center layout

- Highlights the effects of choosing energy efficient servers or efficient cooled servers

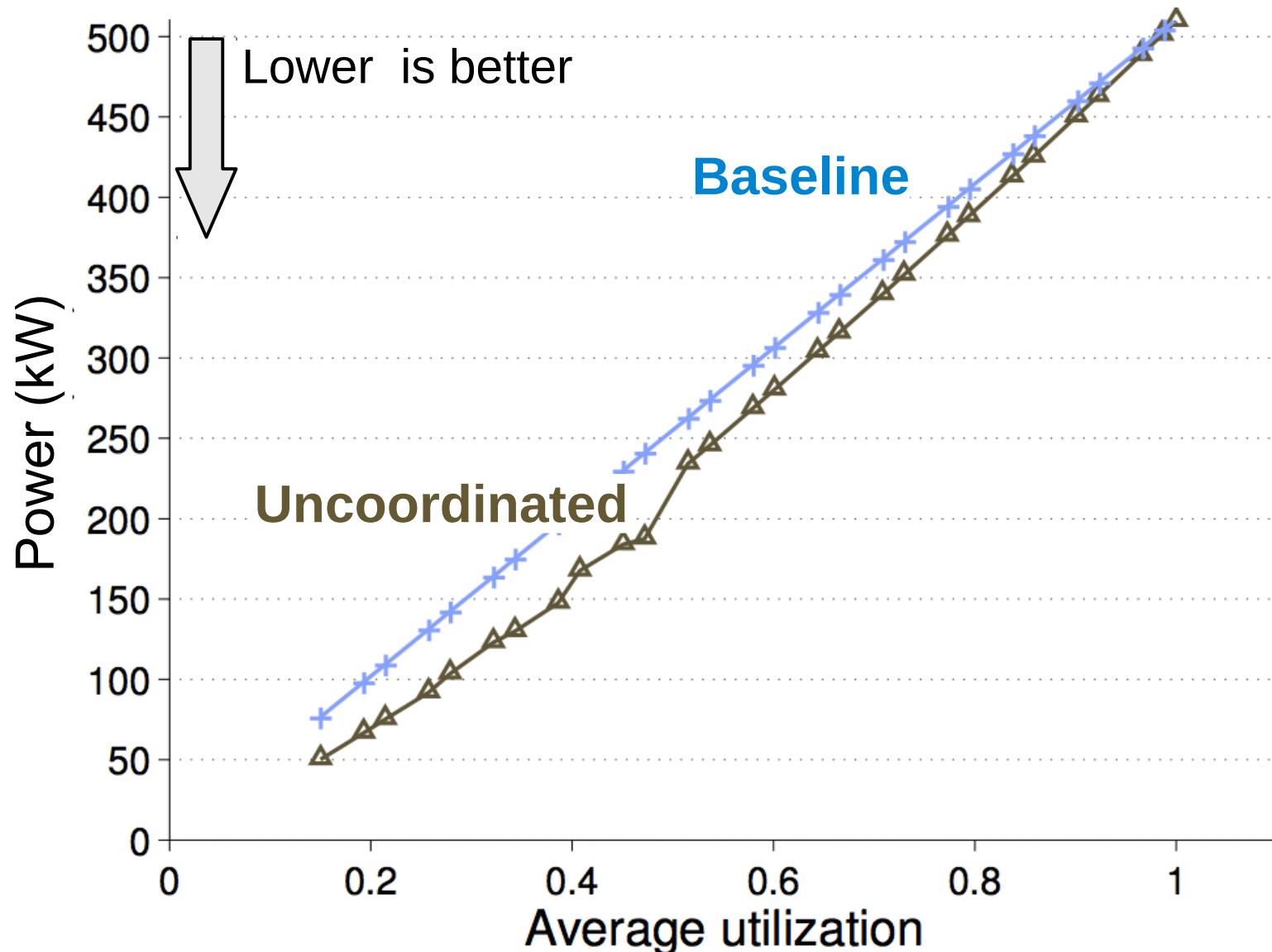


- 4 CRAC units
 - Identical to each other
- 8 Zones
 - 3 Racks each (126 servers per zone)
 - Energy efficient servers
 - Efficiently cooled servers
- Optimal control problems
 - Written in TomSym
 - Solver: KNITRO 7.0

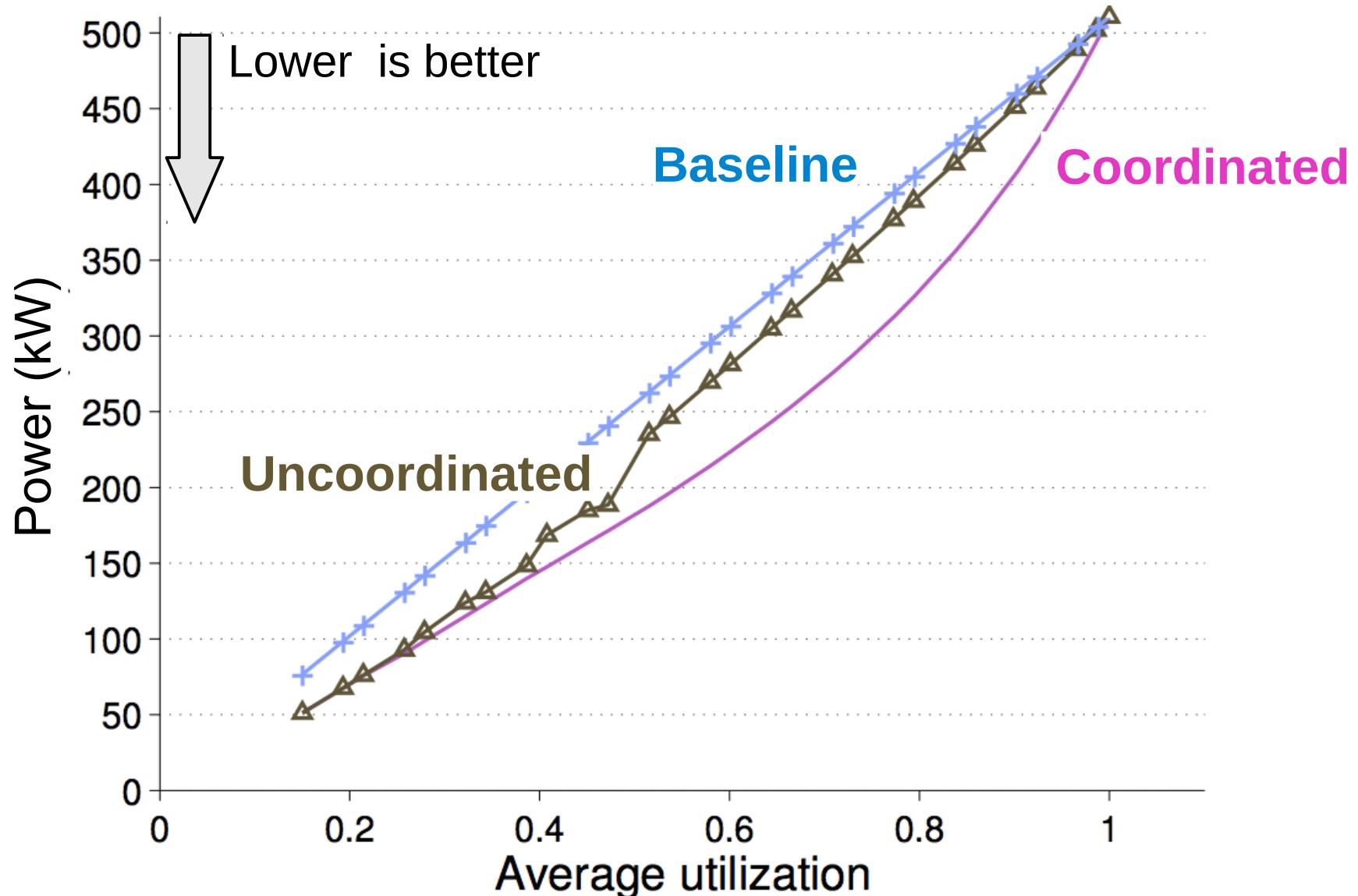
Total power consumption



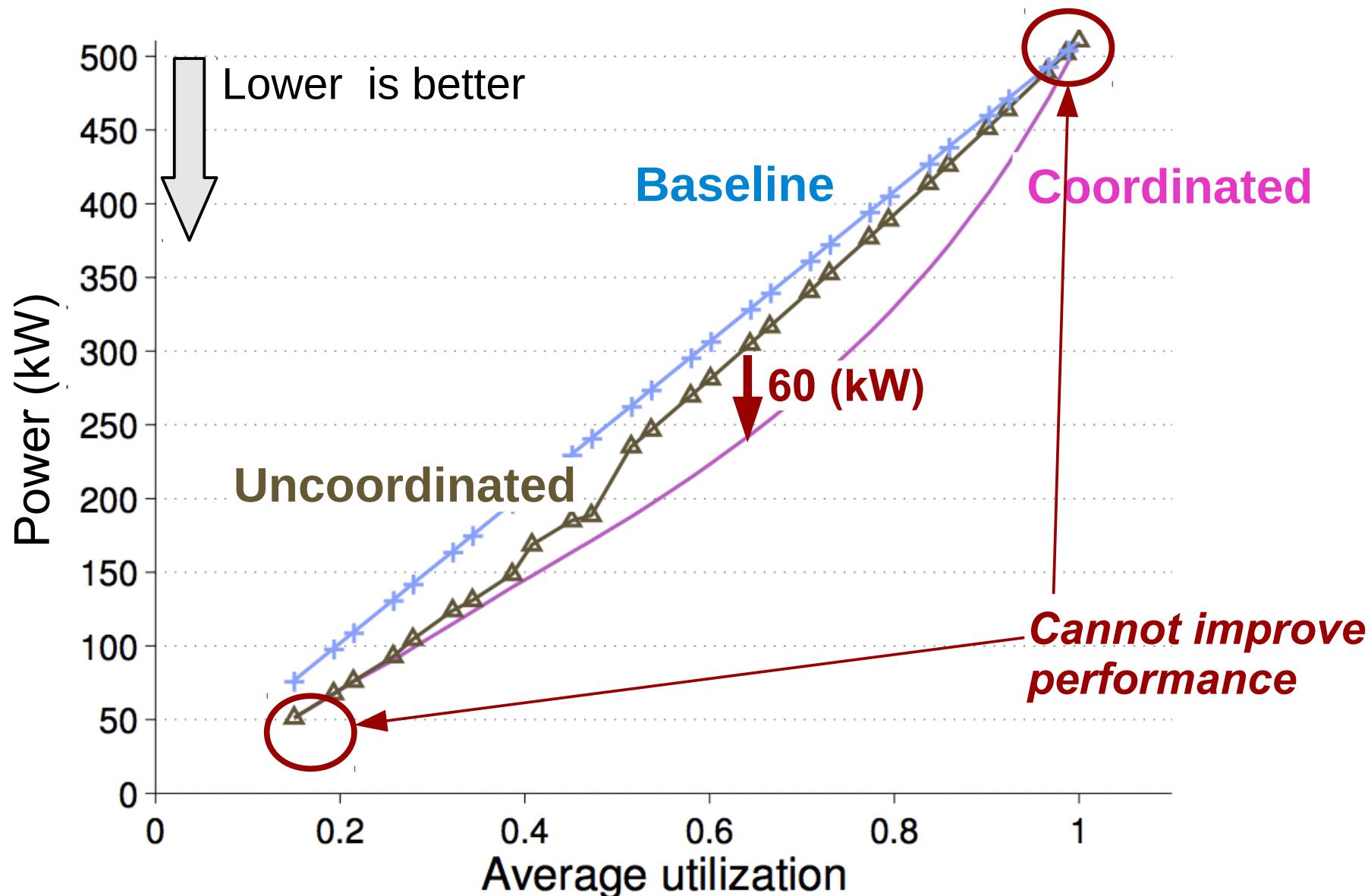
Total power consumption



Total power consumption

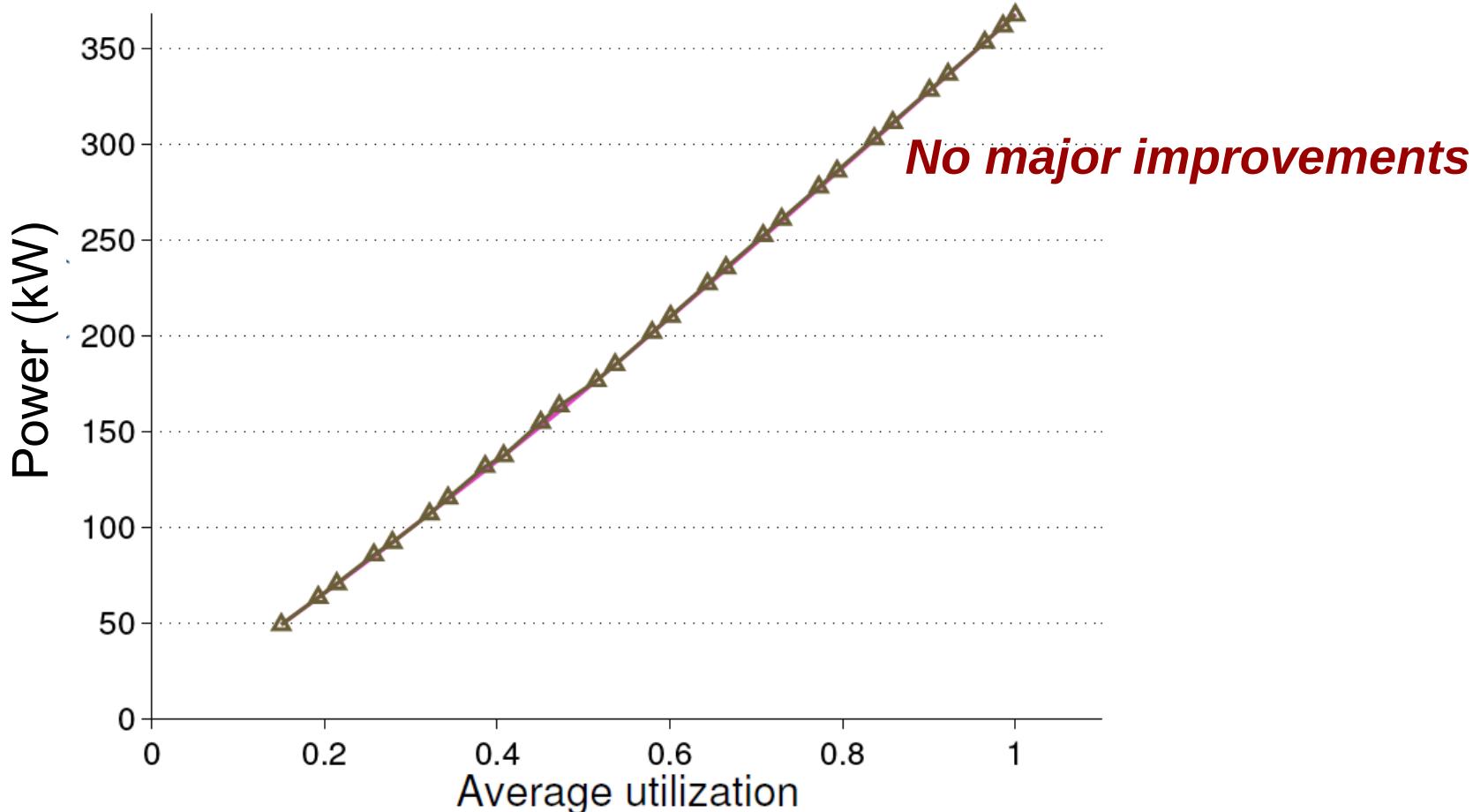


Total power consumption



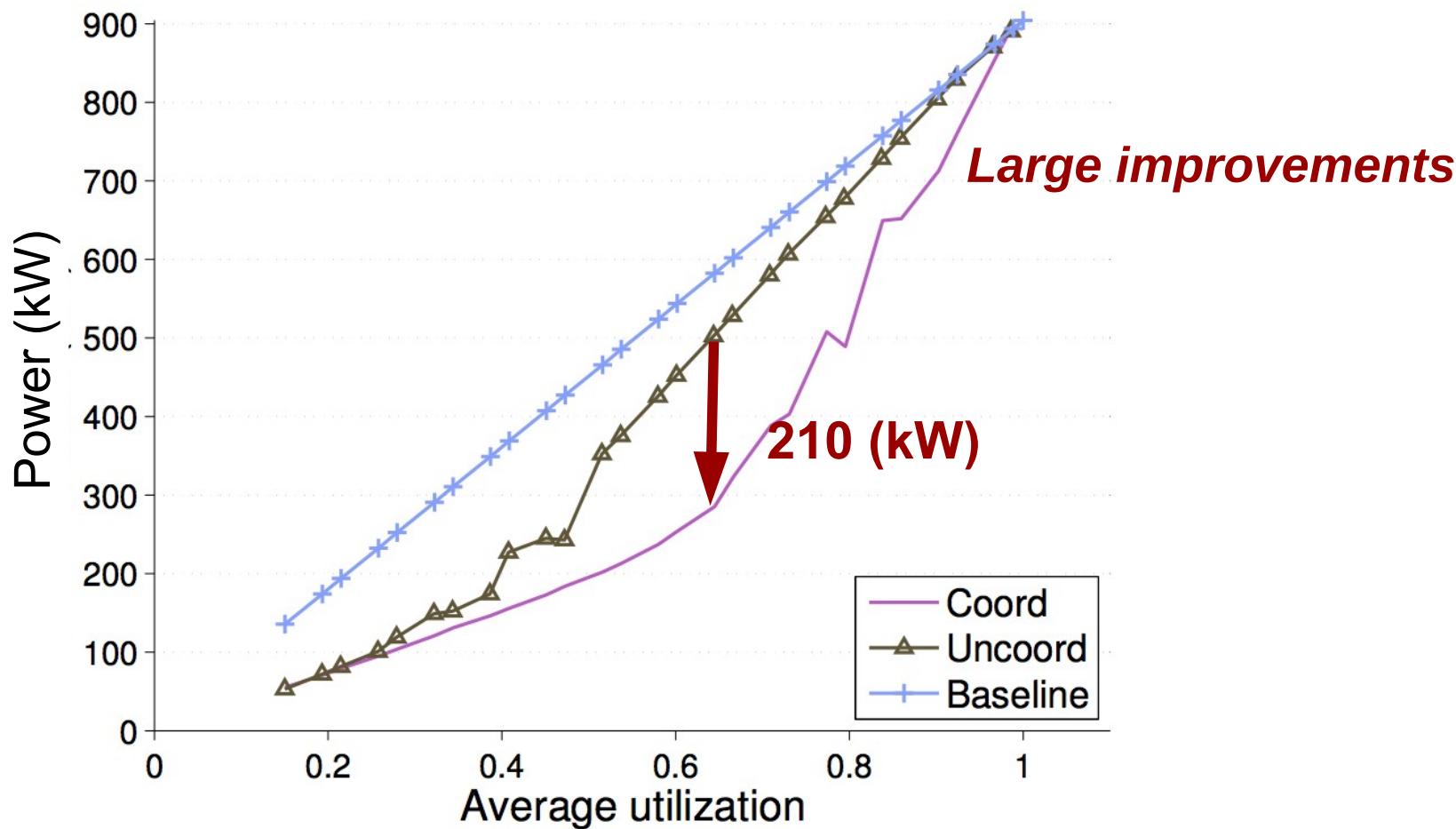
Total power consumption – case 2

- How do the controllers perform when all of the zones are efficiently cooled?



Total power consumption – case 3

- How do the controllers perform when large variability exists among the zone cooling efficiency?



Cyber-Physical index

- Given a data center
 - How much energy can we save by using a coordinated controller, with respect to an uncoordinated controller?

Cyber-Physical index

- Given a data center
 - How much energy can we save by using a coordinated controller, with respect to an uncoordinated controller?
- Cyber-Physical index (CPI), values in [0,1]
 - When CPI is close to 1, then a coordinated approach is advisable
 - When CPI is close to 0, then an uncoordinated approach tends to be as efficient as a coordinated approach
- CPI is function of the relative sensitivity of the zones

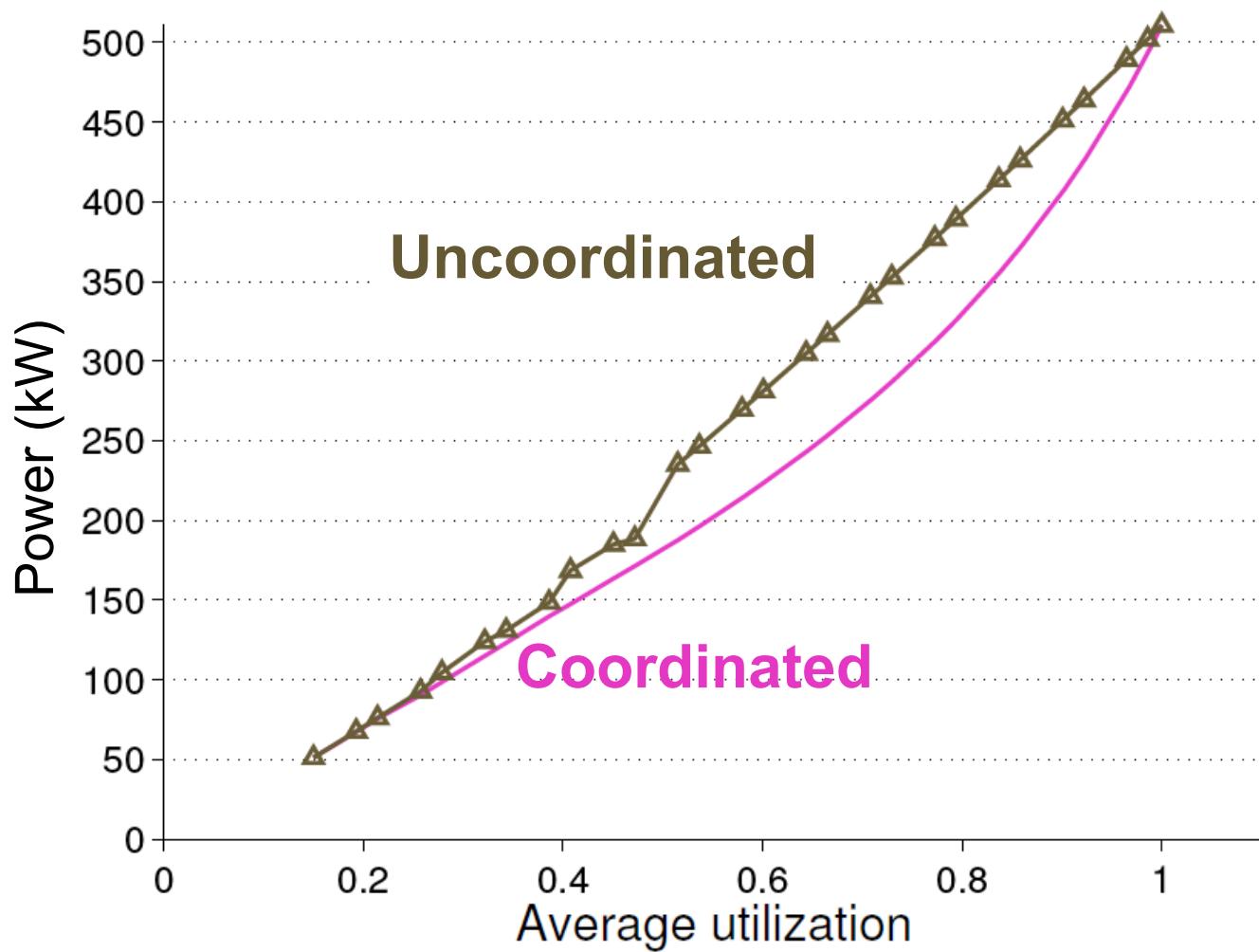
$$CPI = k \text{ std} \left(\begin{bmatrix} S_1 & \dots & S_N \end{bmatrix} \right)$$

$$S_i = \left\| \frac{\partial T_{\text{in},i}}{\partial \mathbf{T}_{\text{ref}}} \Big|_{(\mathbf{l}, \mathbf{T}_{\text{out}})} \right\|_2 / \left\| \frac{\partial T_{\text{in},i}}{\partial \mathbf{z}} \Big|_{(\mathbf{l}, \mathbf{T}_{\text{out}})} \right\|_2$$

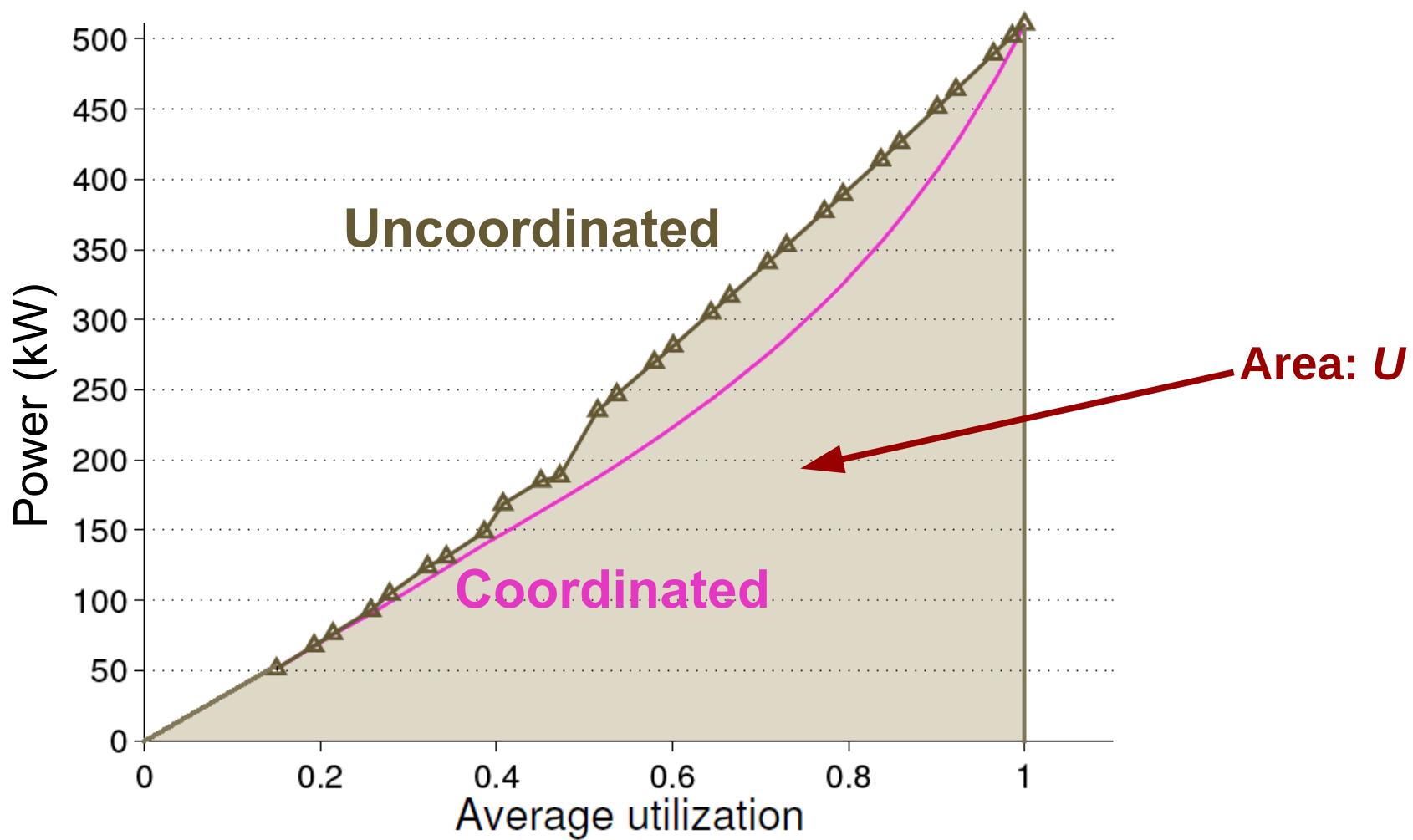
*Relative sensitivity
of the i^{th} zone*

$$\mathbf{z} = [\mathbf{T}_{\text{ref}}^T \boldsymbol{\eta}^T]^T$$

Relative efficiency

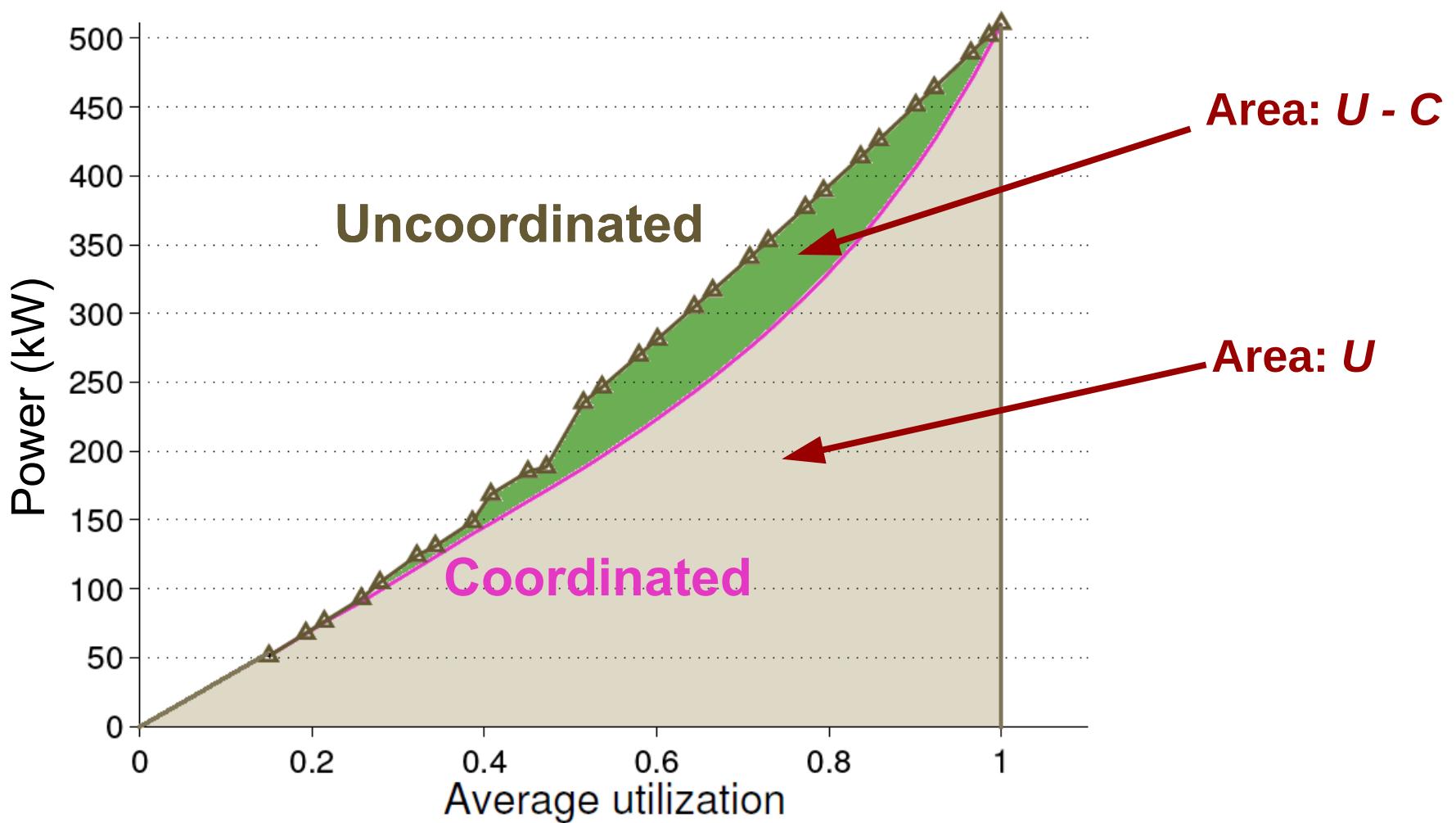


Relative efficiency



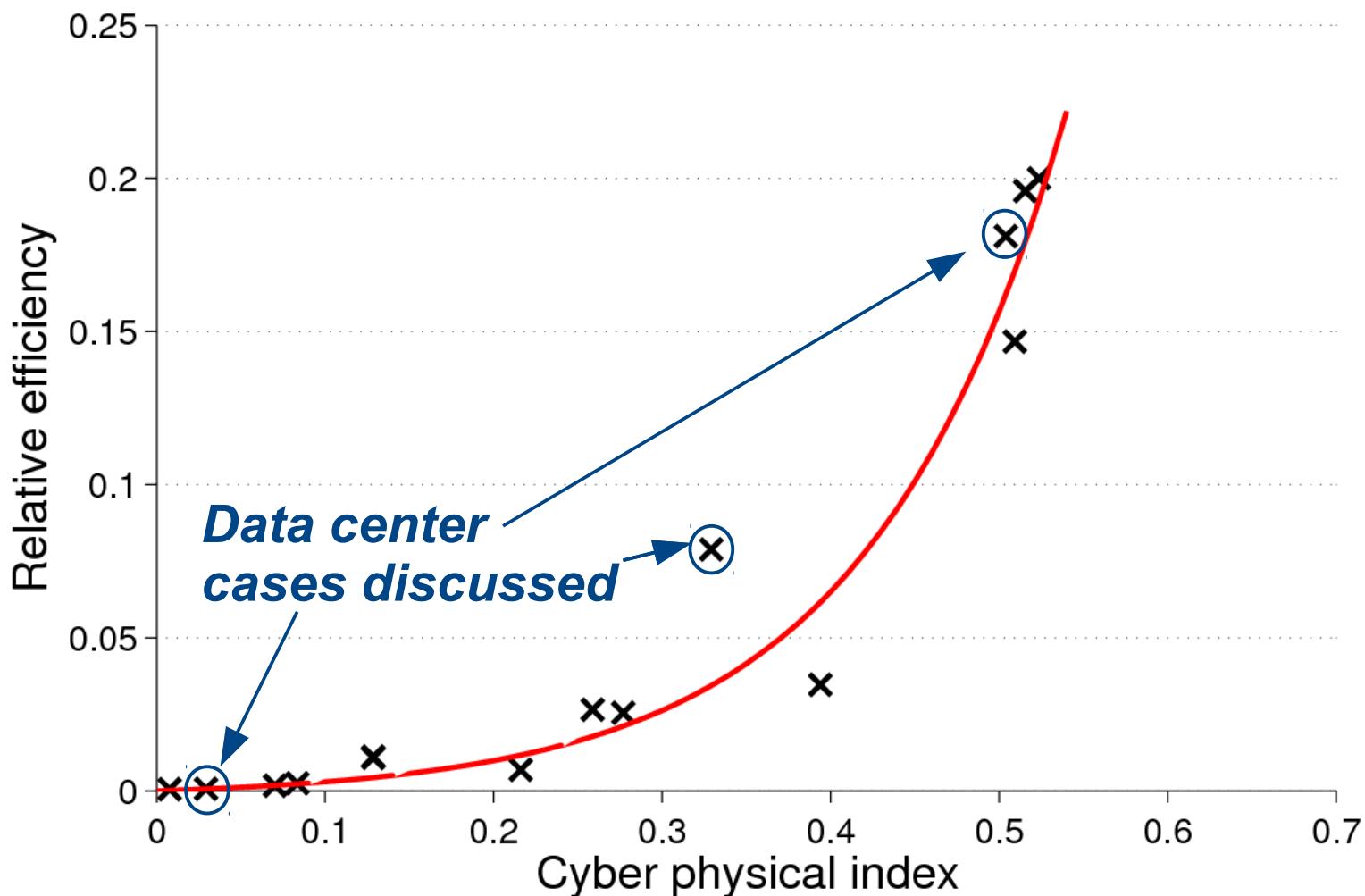
Relative efficiency

- Relative efficiency = $\frac{U - C}{U}$



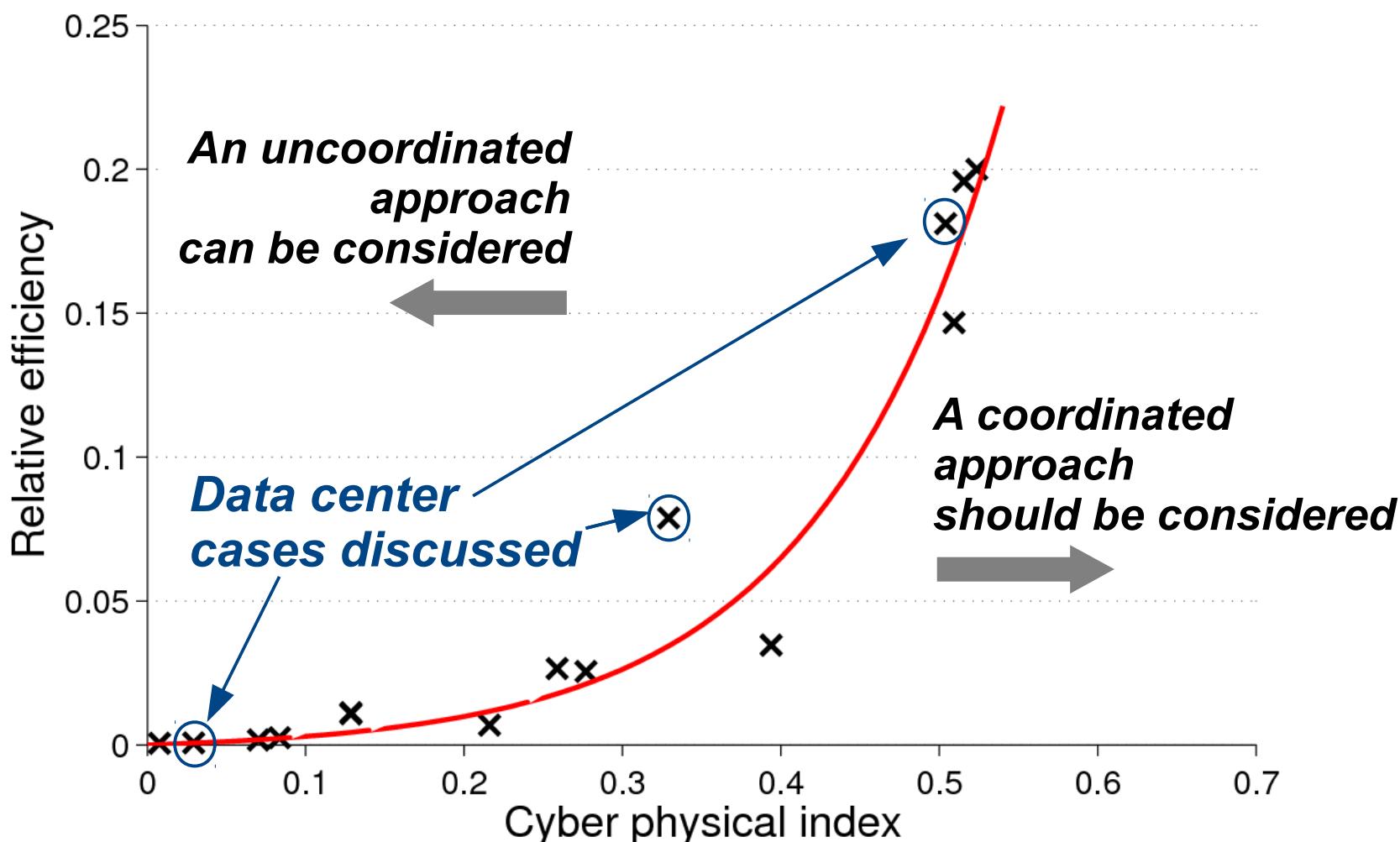
CPI and relative efficiency

- Every point represents a different data center



CPI and relative efficiency

- Every point represents a different data center



Outline

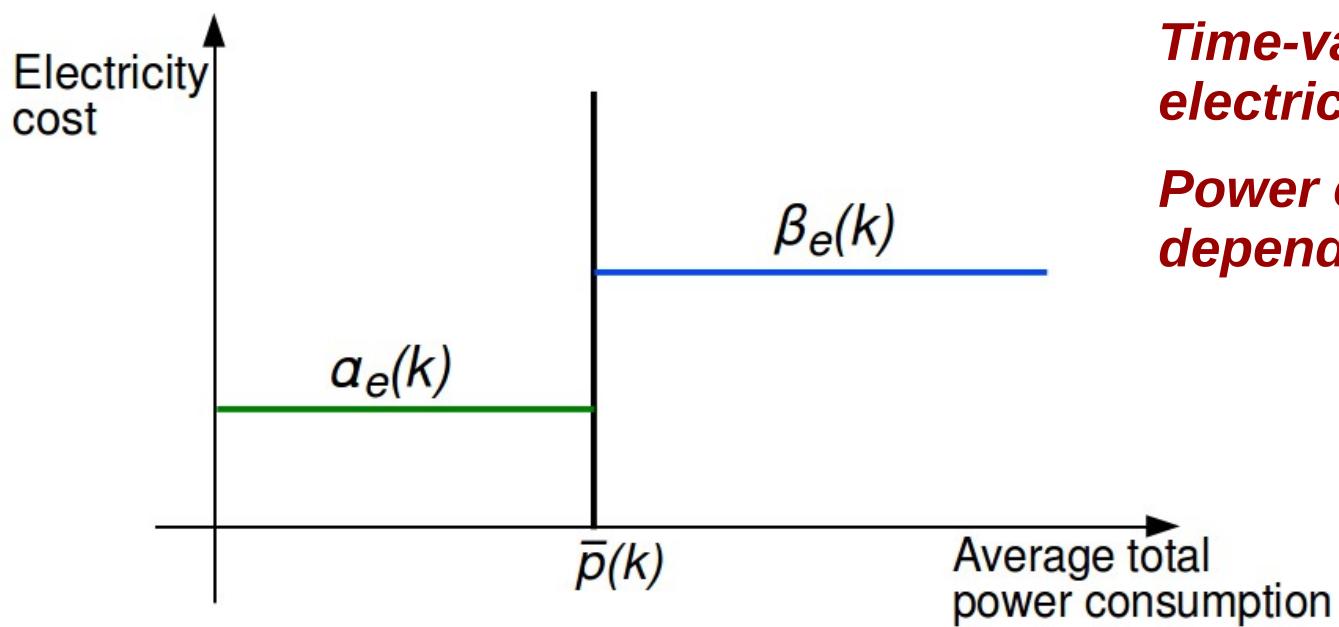
- **Introduction**
- **Proposed control strategy**
- **Performance analysis for constant job arrival rate**
- **Interaction with the smart-grid**
- **Zone-level control**
- **Conclusion and future work**

Interaction with the smart-grid

- **Goal of the data-center-level controller**
 - Minimize run-time operating cost
- **Run-time operating cost**
 - Cost of powering the data center
 - Defined by the service level agreement (SLA) with the power-grid, SLA_G
 - Revenue induced by executing the jobs with a certain QoS
 - Defined by the SLA with the users, SLA_U
 - Cost of migrating jobs among the zones

Interaction with the smart-grid

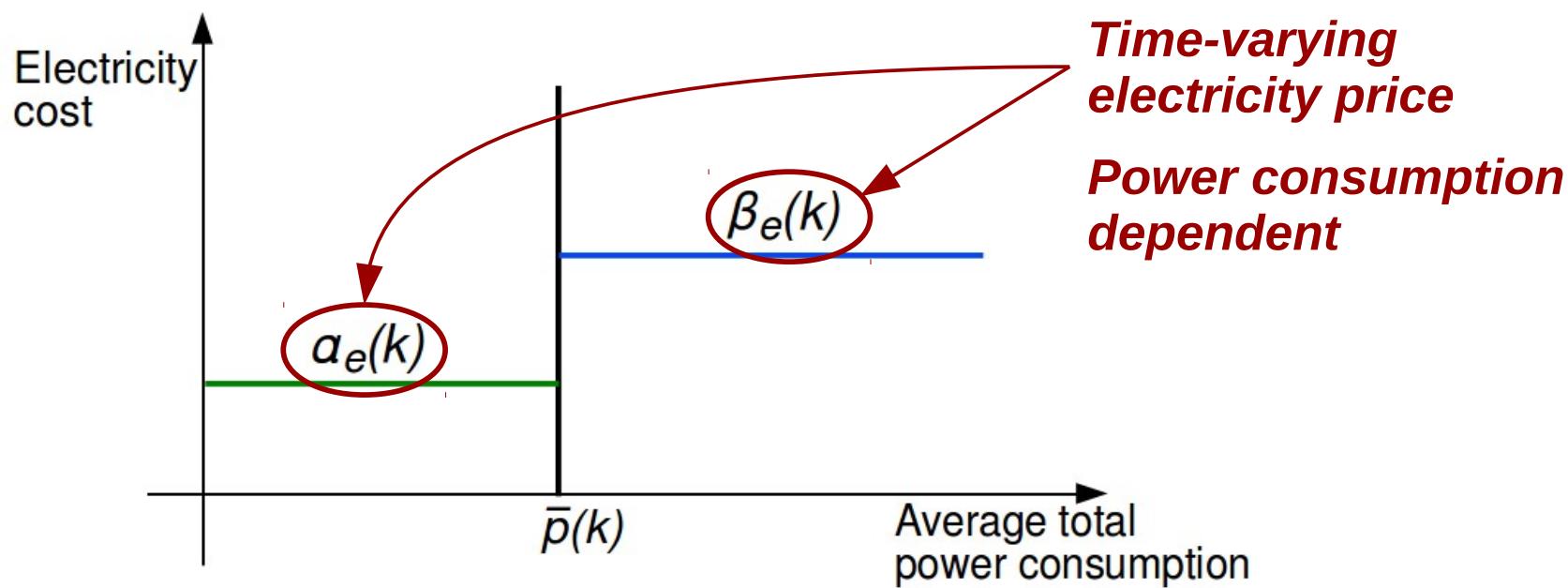
- Goal of the data-center-level controller
 - Minimize run-time operating cost
- Run-time operating cost
 - Cost of powering the data center
 - Defined by the service level agreement (SLA) with the power-grid, SLA_G



*Time-varying
electricity price*
*Power consumption
dependent*

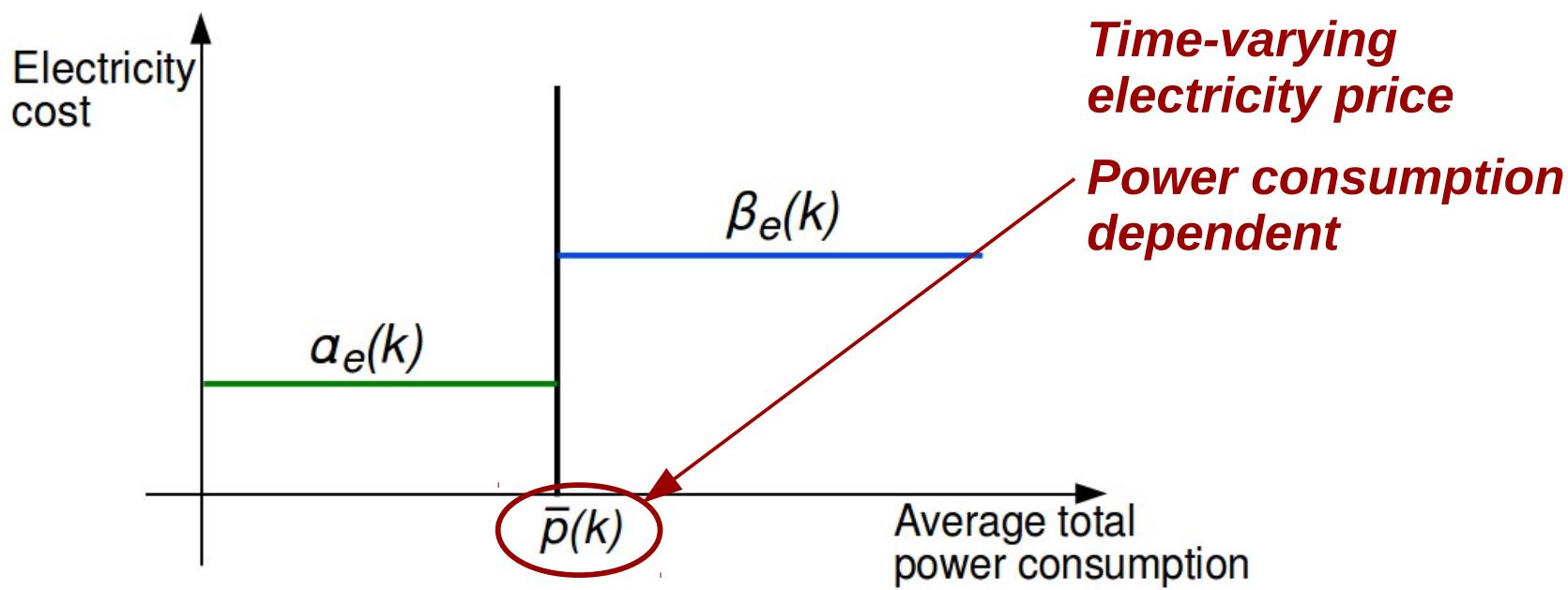
Interaction with the smart-grid

- Goal of the data-center-level controller
 - Minimize run-time operating cost
- Run-time operating cost
 - Cost of powering the data center
 - Defined by the service level agreement (SLA) with the power-grid, SLA_G



Interaction with the smart-grid

- Goal of the data-center-level controller
 - Minimize run-time operating cost
- Run-time operating cost
 - Cost of powering the data center
 - Defined by the service level agreement (SLA) with the power-grid, SLA_G



Interaction with the smart-grid

- **Goal of the data-center-level controller**
 - Minimize run-time operating cost
- **Run-time operating cost**
 - Cost of powering the data center
 - Defined by the service level agreement (SLA) with the power-grid, SLA_G
 - Revenue induced by executing the jobs with a certain QoS
 - Defined by the SLA with the users, SLA_U
 - Cost of migrating jobs among the zones
- **Controllers have to find the best trade-off between executing jobs with high QoS and reducing the powering cost**



Simulation set-up

- Neglects the baseline controller
 - Baseline controller does not consider a cost function in the synthesis of its control strategy
- Two cases are considered
 - 1) A coordinated and an uncoordinated controller
 - Controllers are forced to execute jobs with the highest QoS, e.g., neglect SLA_U
 - The uncoordinated controller only considers the reduced electricity price ($\alpha_e(k)$)

Simulation set-up

- Neglects the baseline controller
 - Baseline controller does not consider a cost function in the synthesis of its control strategy
- Two cases are considered
 - 1) A coordinated and an uncoordinated controller
 - Controllers are forced to execute jobs with the highest QoS, e.g., neglect SLA_U
 - The uncoordinated controller only considers the reduced electricity price ($\alpha_e(k)$)
 - 2) Two coordinated controllers
 - The first coordinated controller is forced to execute jobs at the highest QoS
 - The parameters in the cost function of the second coordinated controller make it more favorable to drop jobs

Simulation set-up

- Neglects the baseline controller

- Baseline controller does not consider a cost function in the synthesis of its control strategy

- Two cases are considered

- 1) A coordinated and an uncoordinated controller
 - Controllers are forced to execute jobs with the highest QoS, e.g., neglect SLA_U
 - The uncoordinated controller only considers the reduced electricity price ($\alpha_e(k)$)

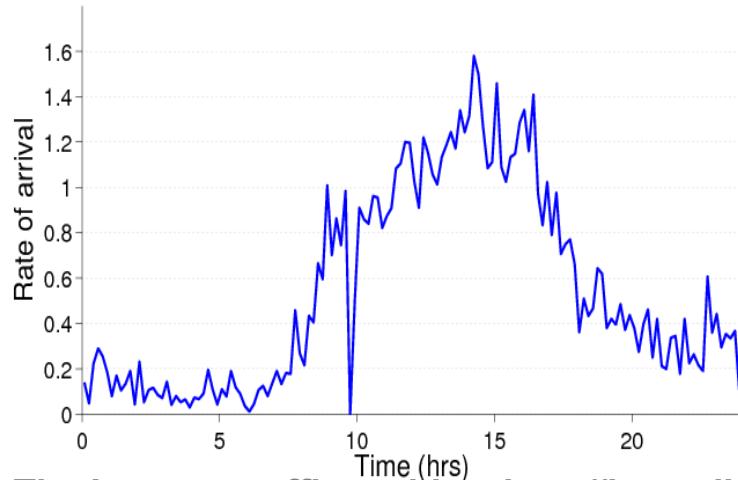
- 2) Two coordinated controllers

- The first coordinated controller is forced to execute jobs at the highest QoS
 - The parameters in the cost function of the second coordinated controller make it more favorable to drop jobs

Job arrival rate and electricity costs

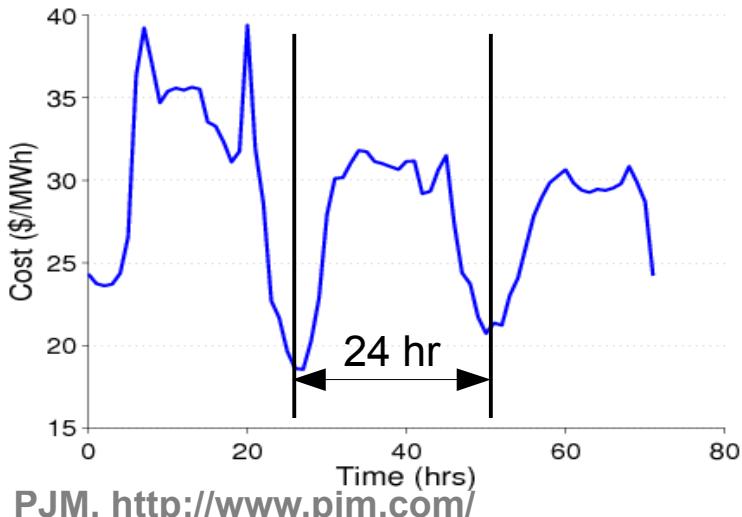
■ Real data

- Job arrival rate



The Internet traffic archive, <http://ita.ee.lbl.gov/>

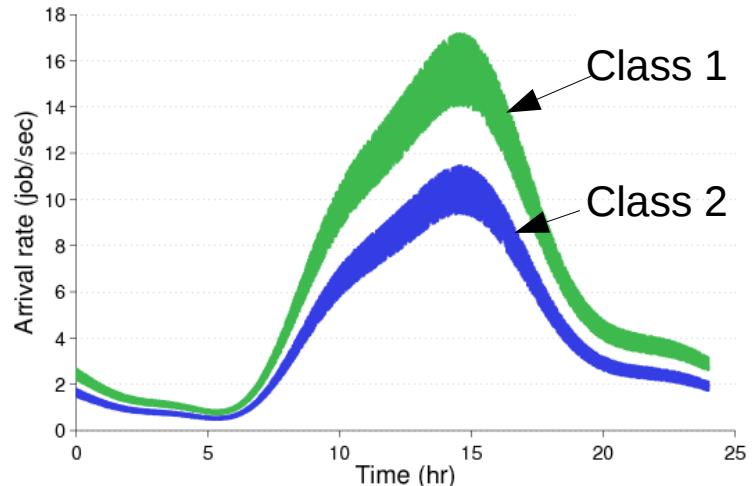
- Day-ahead electricity cost



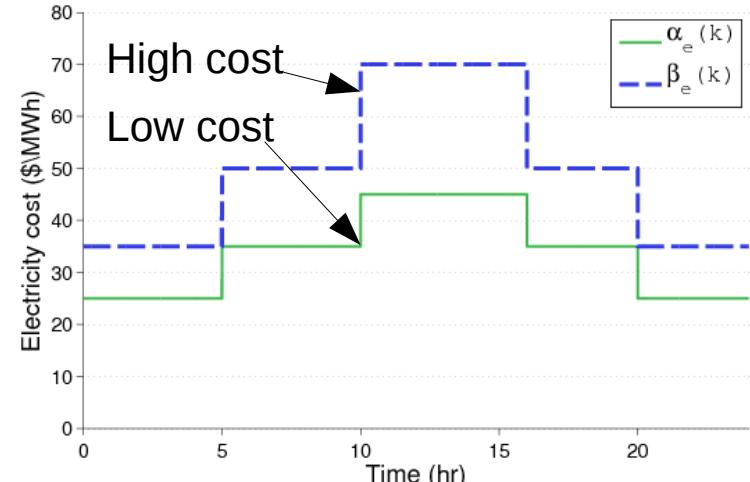
PJM, <http://www.pjm.com/>

■ Parameters of the simulation

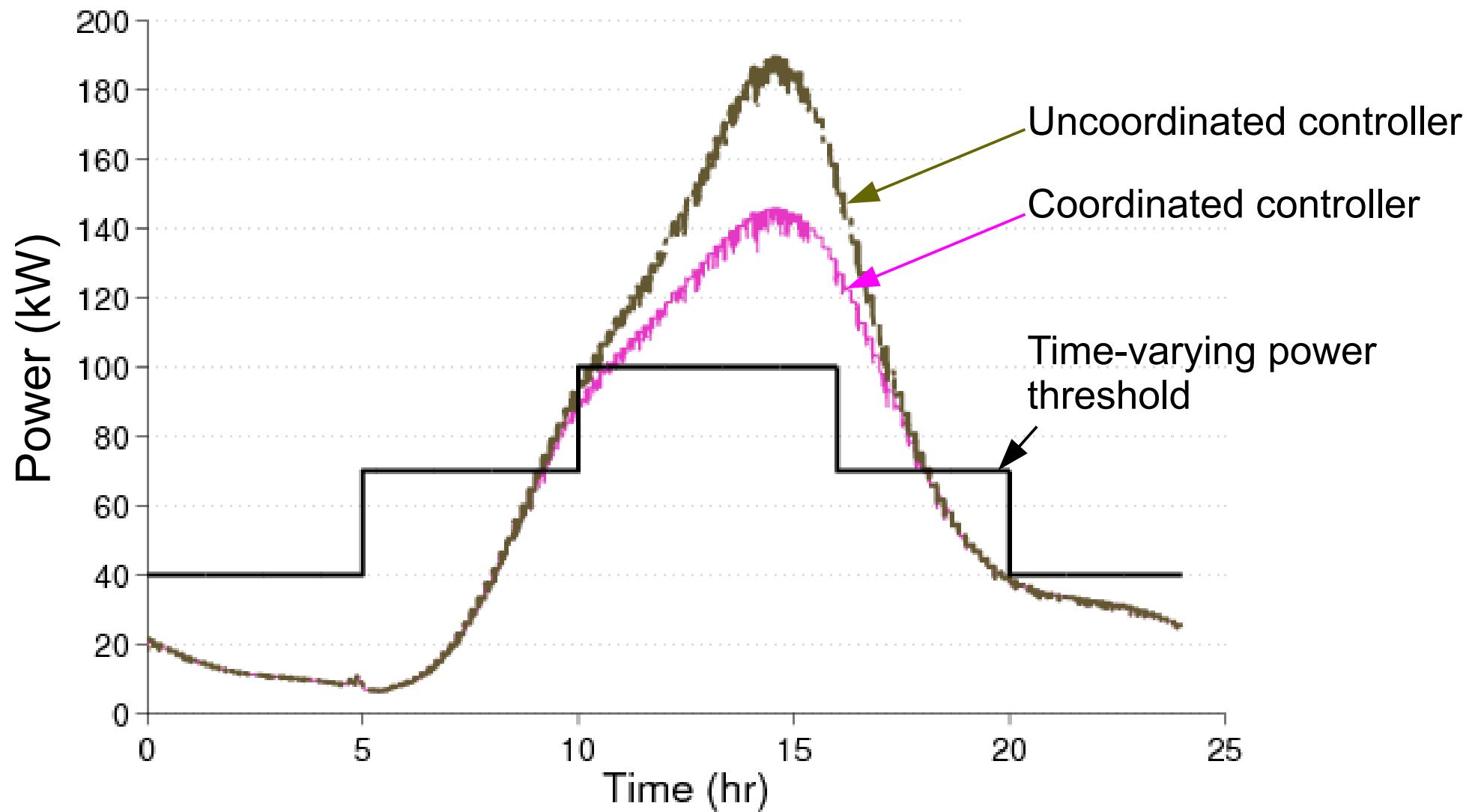
- Job arrival rate, 2 classes of jobs



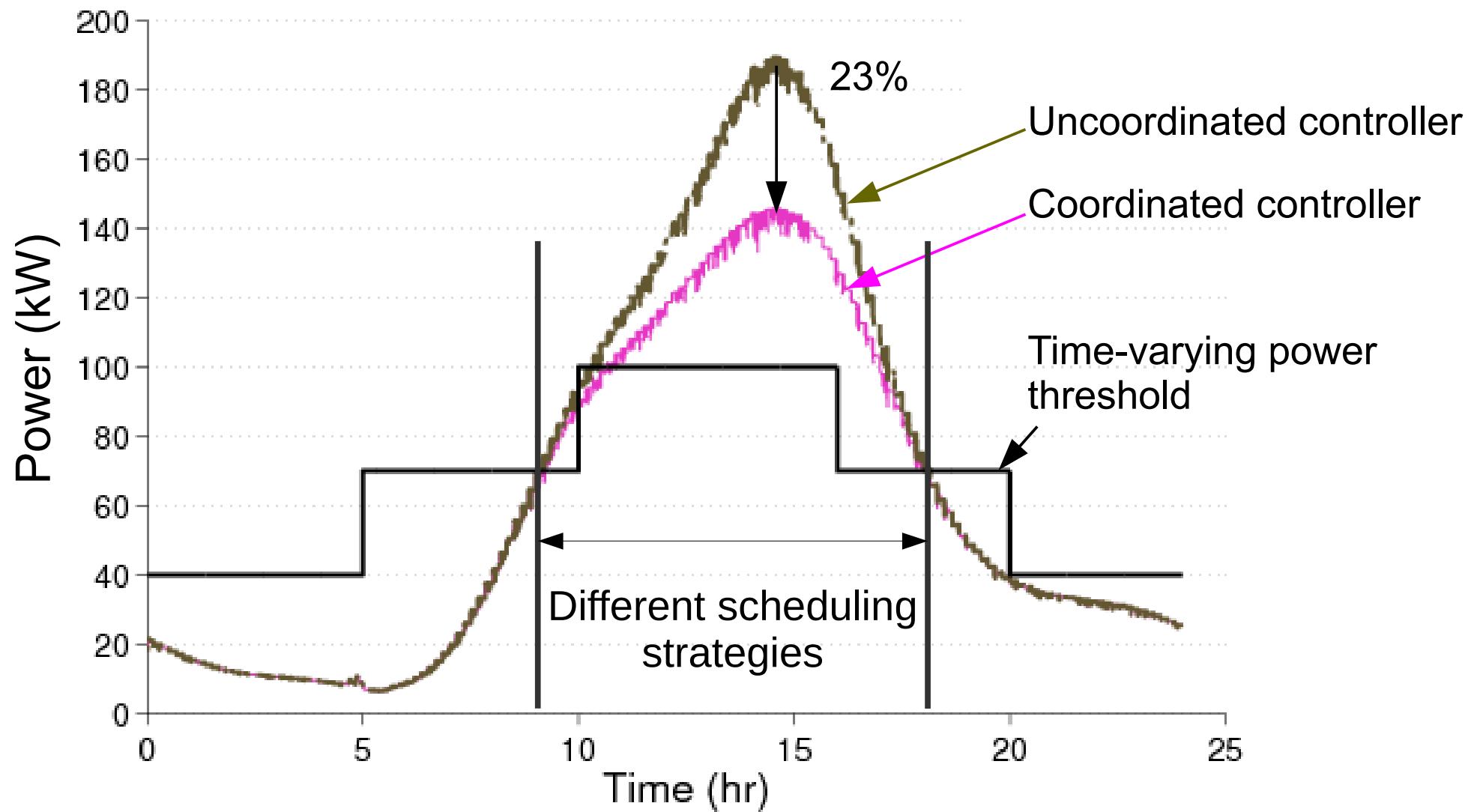
- Electricity costs



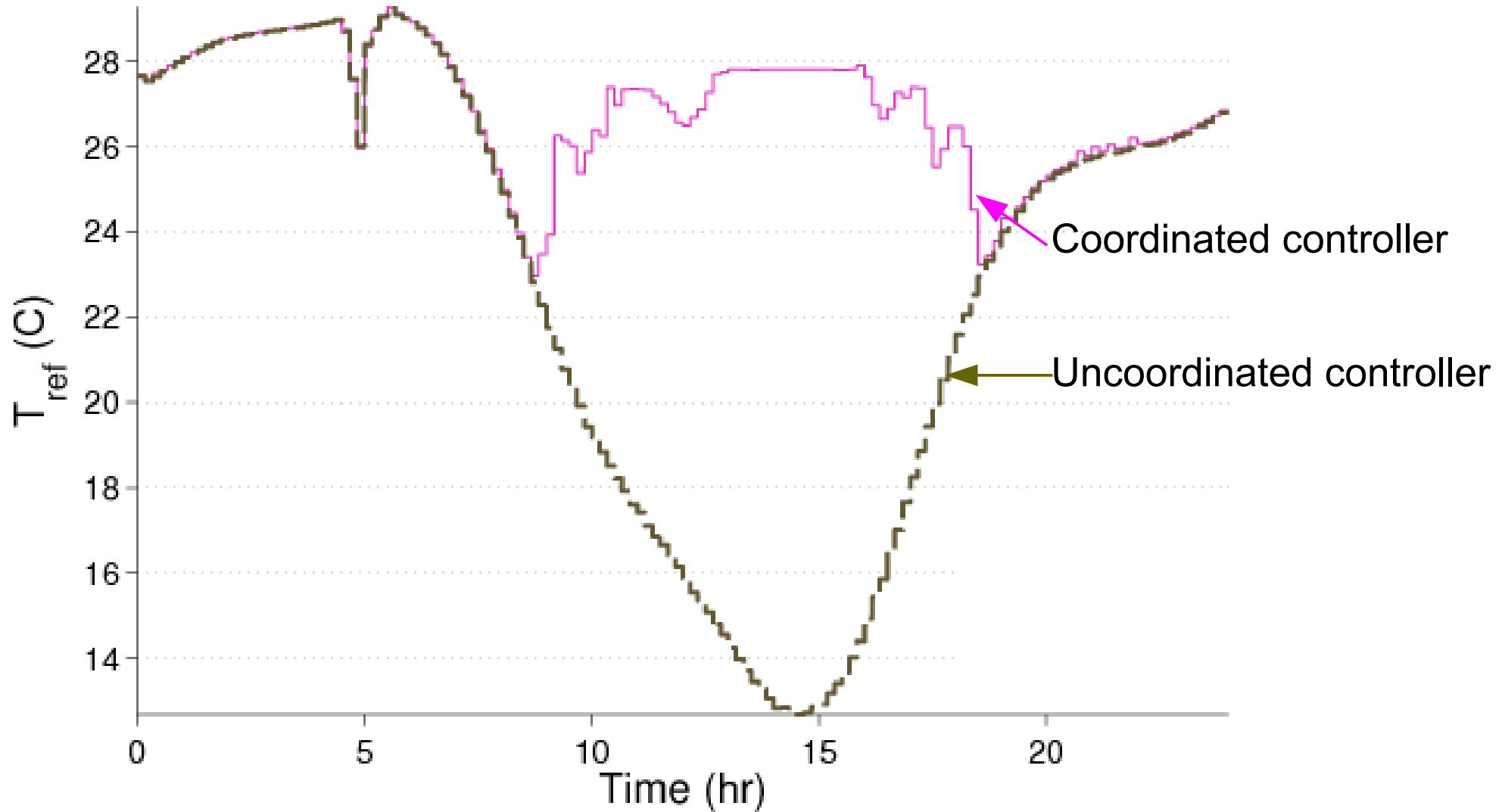
Total data center power consumption



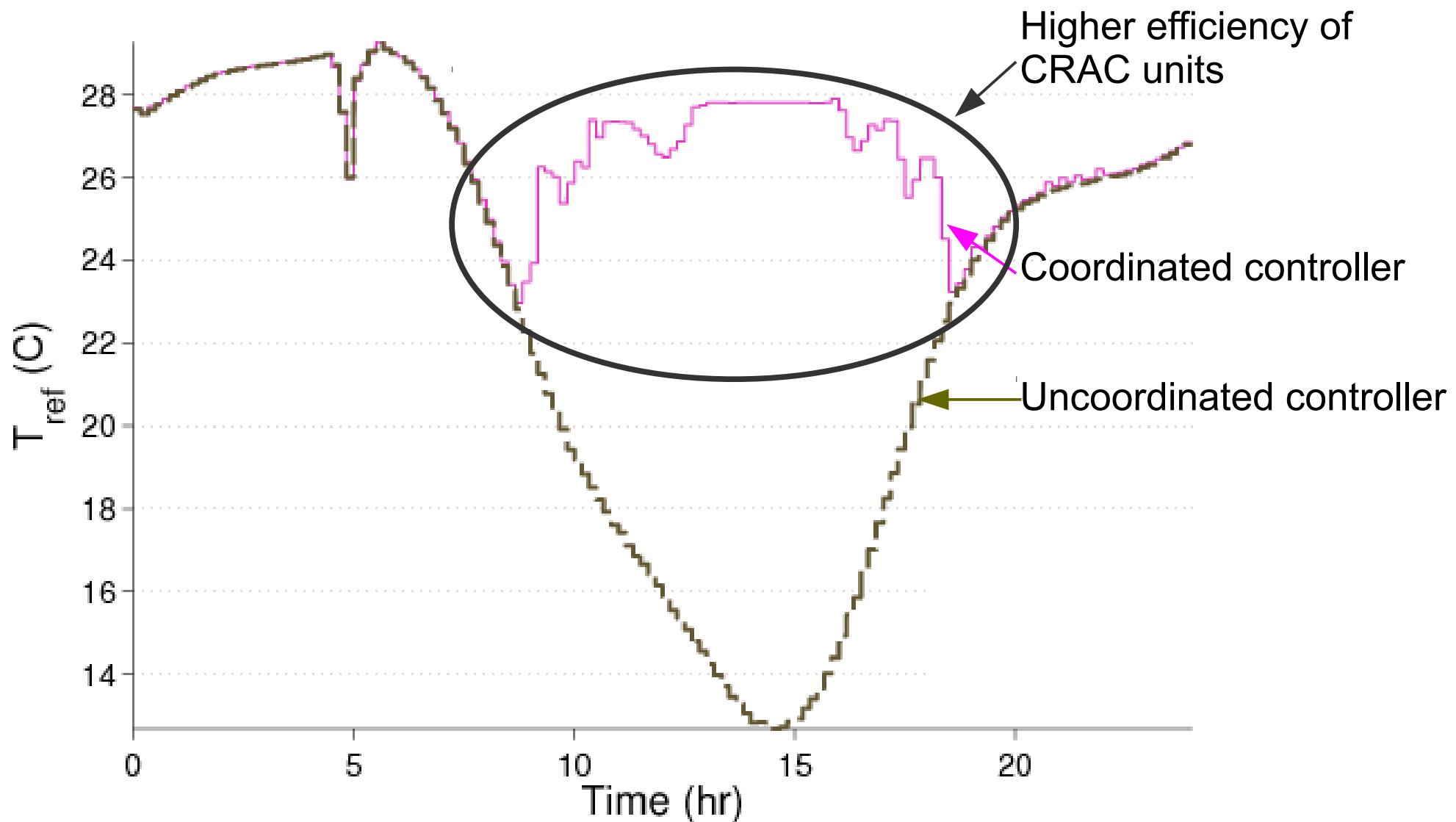
Total data center power consumption



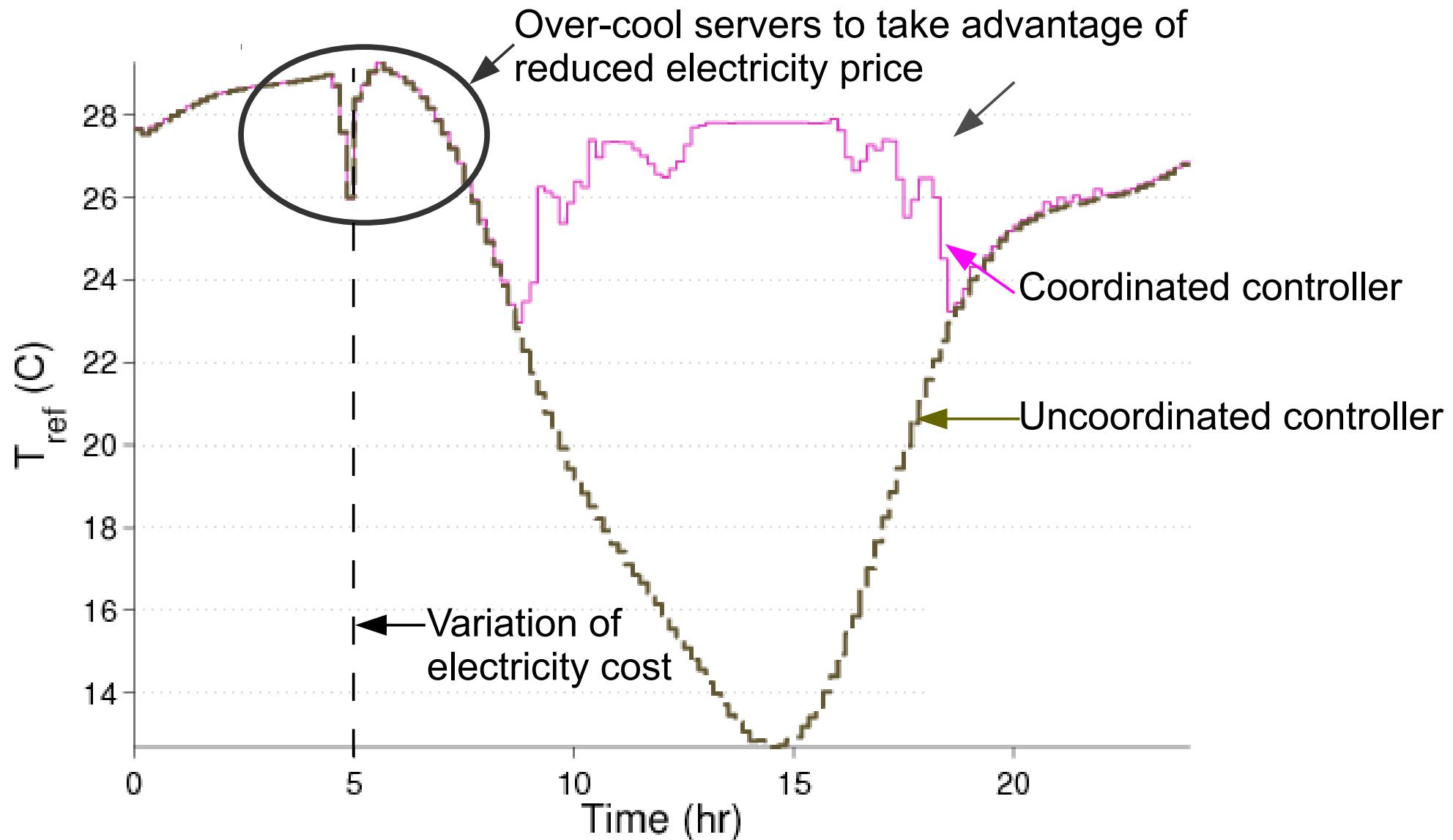
Average reference temperatures



Average reference temperatures



Average reference temperatures



Simulation set-up

- Neglects the baseline controller

- Baseline controller does not consider a cost function in the synthesis of its control strategy

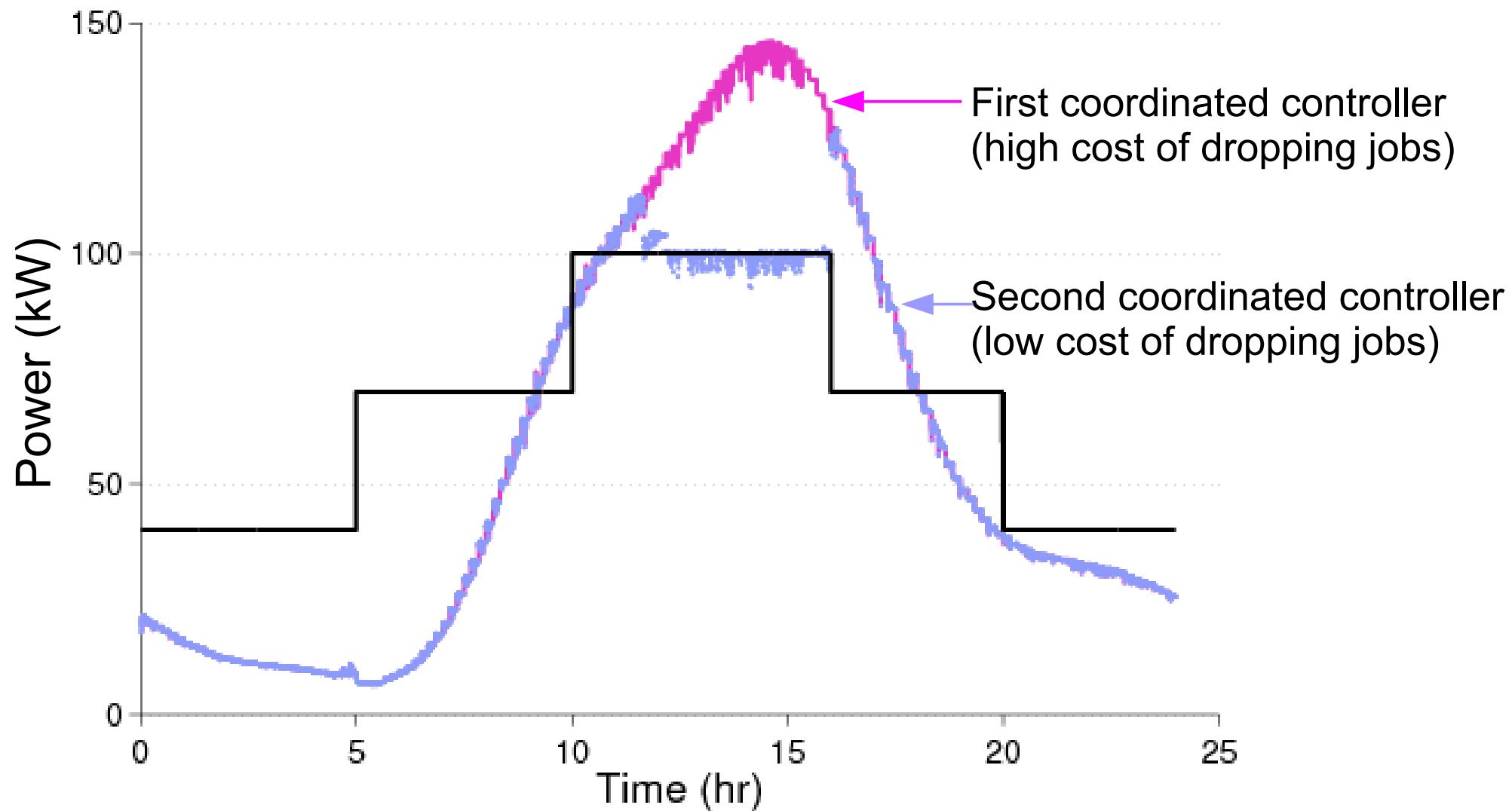
- Two cases are considered

- 1) A coordinated and an uncoordinated controller
 - Controllers are forced to execute jobs with the highest QoS, e.g., neglect SLA_U
 - The uncoordinated controller only considers the reduced electricity price ($\alpha_e(k)$)

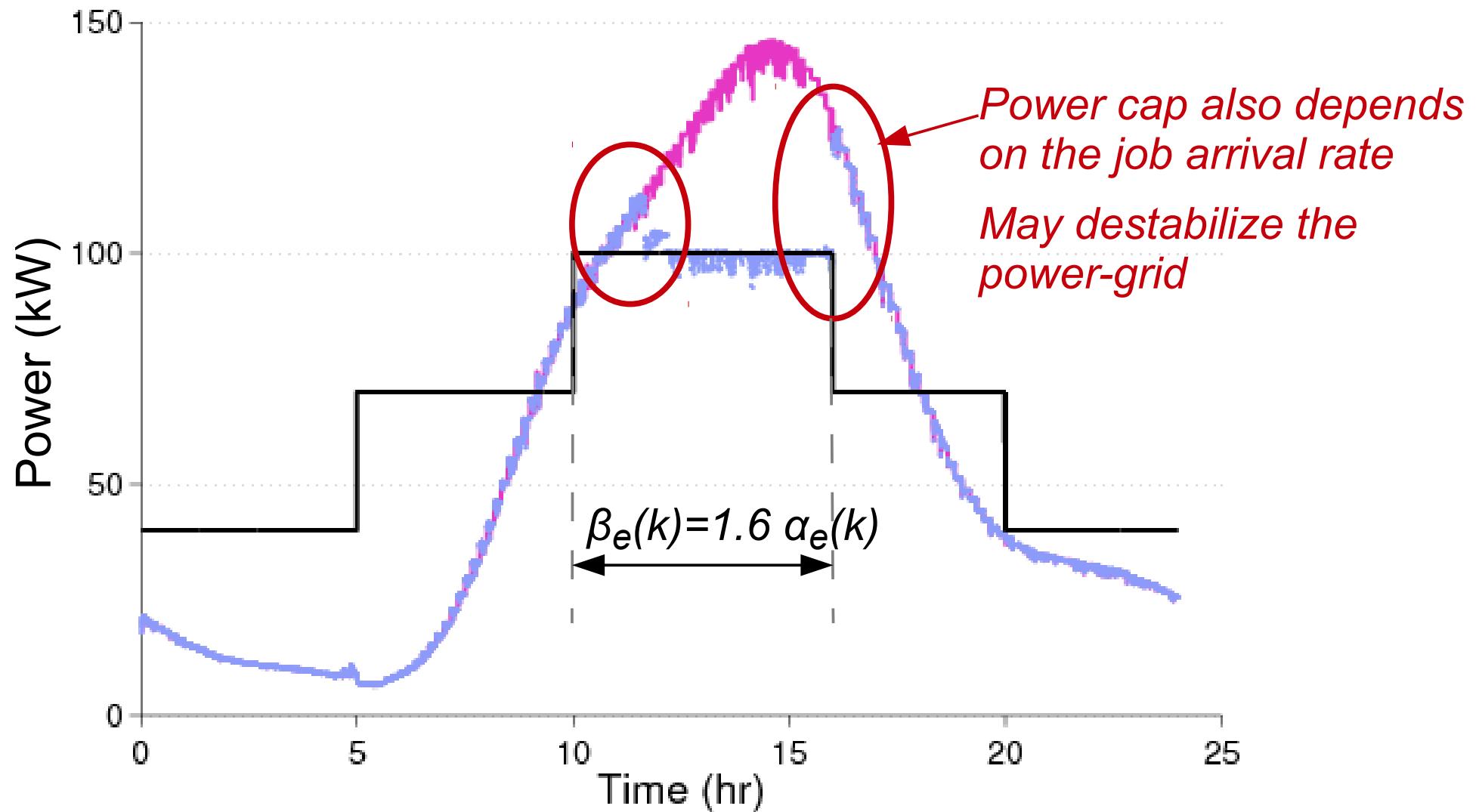
- 2) Two coordinated controllers

- The first coordinated controller is forced to execute jobs at the highest QoS
- The parameters used in the cost function of the second coordinated controller make it more favorable to drop jobs

Total data center power consumption



Total data center power consumption



Outline

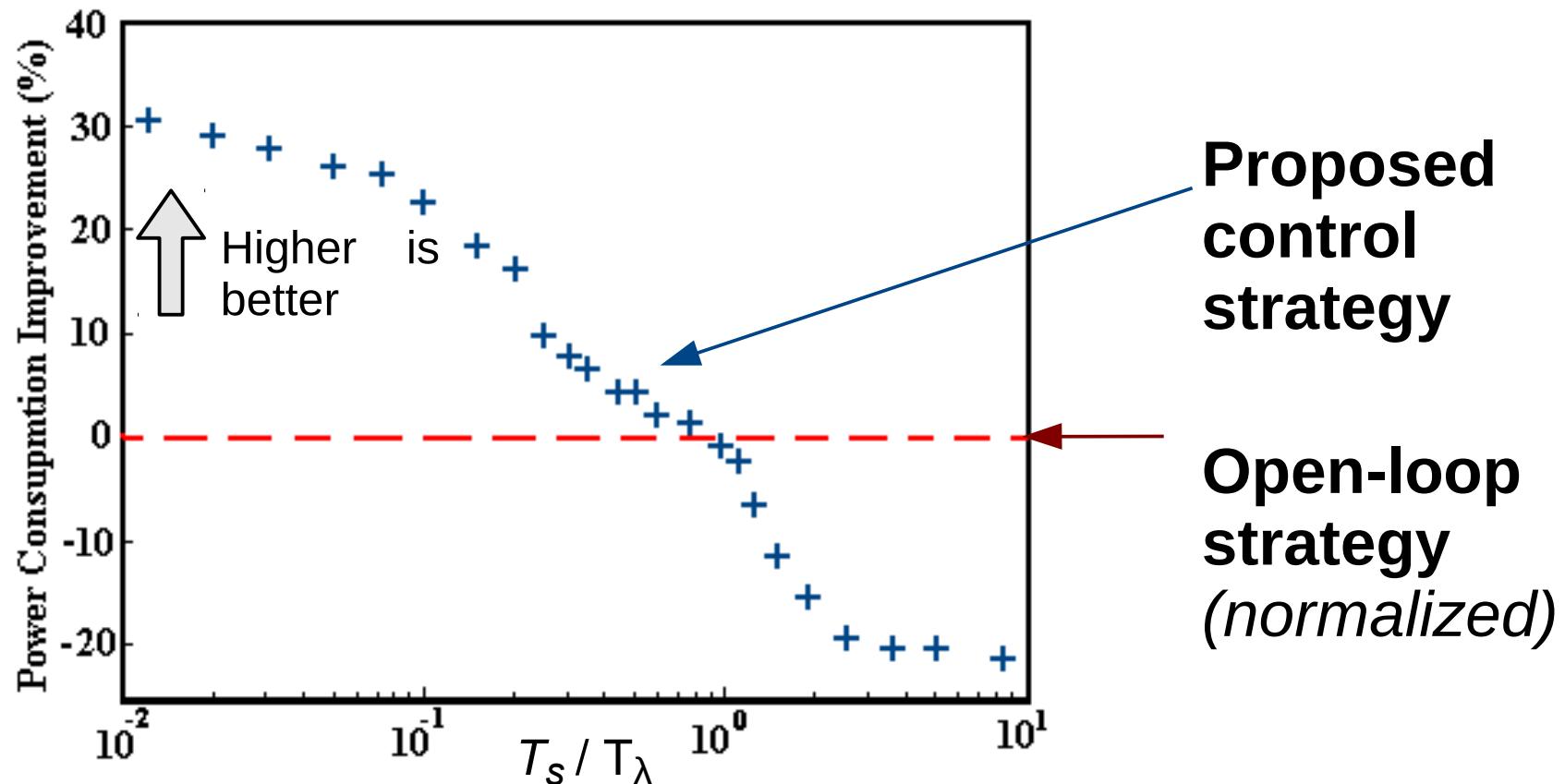
- **Introduction**
- **Proposed control strategy**
- **Performance analysis for constant job arrival rate**
- **Interaction with the smart-grid**
- **Zone-level control**
- **Conclusion and future work**

Zone-level controller

- Operates on the minutes time-scale
- Decides how many servers in the zone should be turned on
 - Servers take time to turn on
 - When turning on, servers consume power and do not execute jobs
- Control actions based on
 - Desired jobs execution rate (predictive control)
 - Current use of resources in the zone (reactive control)
- Considers
 - Time to turn servers on: T_s
 - Variability of workload arrival rate: $1/T_\lambda$
- Controller details in the dissertation

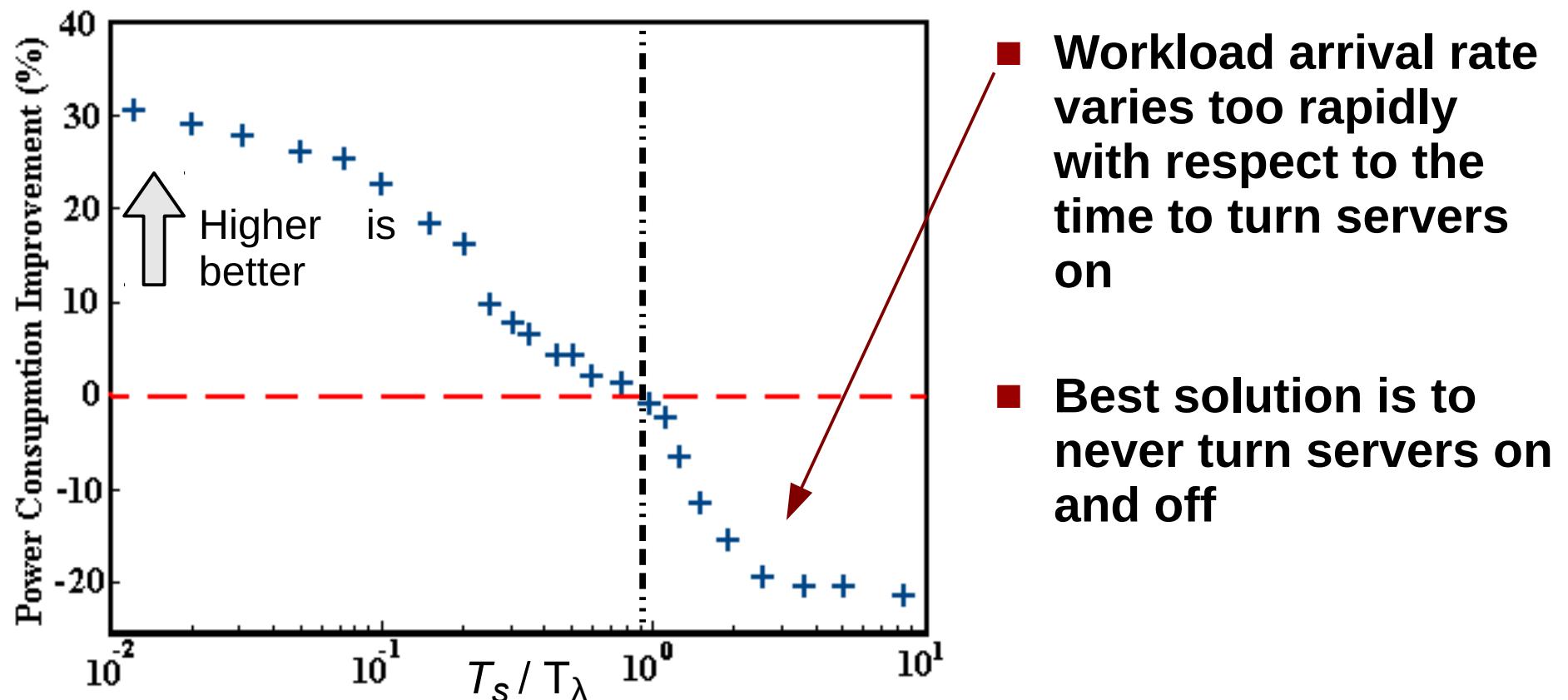
Simulation results

- Compare the proposed control approach against optimal open-loop strategy
 - The number of ON servers is *fixed* and based on the long-term average job arrival rate



Simulation results

- Compare the proposed control approach against optimal open-loop strategy
 - The number of ON servers is *fixed* and based on the long-term average job arrival rate



Outline

- **Introduction**
- **Proposed control strategy**
- **Performance analysis for constant job arrival rate**
- **Interaction with the smart-grid**
- **Zone-level control**
- **Conclusion and future work**

Conclusion

- Data centers have operating regimes for which the baseline and the uncoordinated controller are as optimal as the coordinated controller
- There are cases where the uncoordinated controller is *always* as optimal as the coordinated controller
- Proposed the cyber-physical index
 - First attempt to characterize the thermal and computational characteristics of the data center with a single index
- Data centers can take advantage of SLAs with the power grid
 - Depending on the SLA, the data center may destabilize the power-grid
- If the job arrival rate changes too quickly, then never turn servers off is the best strategy

Future work

- **Model validation**
 - Multiple data center cases have to be considered in order to understand which model is useful in which condition
- **Data collection and sharing with the research community**
 - The lack of public available data is one of the major issue related to data center modeling and control
- **Interaction with the power-grid**
 - How should the SLA with the power-grid be stipulated, in order to avoid sudden changes in the power demand?
- **Management of data storage**
 - Including management of data storage in the optimal control problem may improve the overall performance
- **Further development of the data center simulator**
 - A data center simulator used and recognized by the whole research community would allow the comparison of different control algorithms under the same test cases

Publications

- *L. Parolini*, B. Sinopoli, B. H. Krogh. Models and Control Strategies for Data Centers in the Smart Grid. In Control and Optimization Theory for Electric Smart Grids, Springer
- *L. Parolini*, B. Sinopoli, B. H. Krogh, Z. Wang. A Cyber-Physical-System Approach to Data Center Modeling and Control for Energy Efficiency. Proceedings of the IEEE, Special Issue on Cyber-Physical Systems
- *L. Parolini*, B. Sinopoli, B. H. Krogh. Model predictive control of data centers in the smart grid scenario. 18th World Congress of the International Federation of Automatic Control. August 2011
- *L. Parolini*, E. Garone, B. Sinopoli, B. H. Krogh. A Hierarchical Approach to Energy Management in Data Centers. 49th IEEE Conference on Decision and Control. December 2010
- M. Aghajani, *L. Parolini*, B. Sinopoli. Dynamic Power Allocation in Server Farms: a Real Time Optimization Approach. 49th IEEE Conference on Decision and Control. December 2010
- *L. Parolini*, N. Tolia, B. Sinopoli, B. H. Krogh. A Cyber-Physical Systems Approach to Energy Management in Data Centers. First International Conference on Cyber-Physical Systems. April 2010
- *L. Parolini*, B. Sinopoli, B. H. Krogh. A Unified Thermal-Computational Approach to Data Center Energy Management. Fourth International Workshop on Feedback Control Implementation and Design in Computing Systems and Networks. April 2009
- *L. Parolini*, B. Sinopoli, B. H. Krogh. Reducing Data Center Energy Consumption via Coordinated Cooling and Load Management. HotPower. December 2008

Thank you!

Additional slides

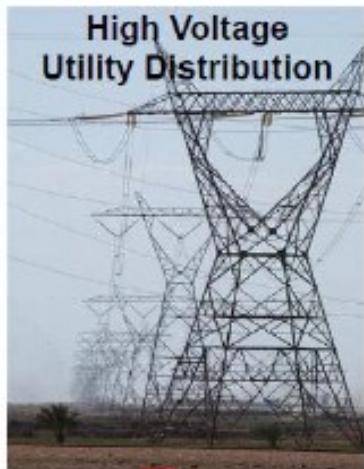
Introduction to Data centers

What is a data center?

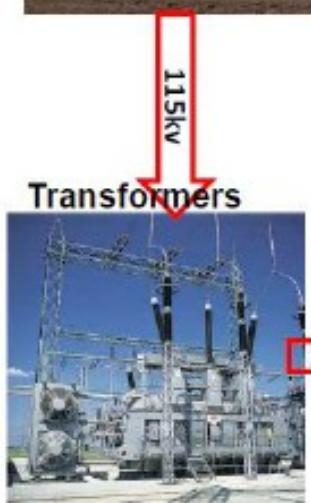
- A data center is a facility used to house computer systems and associated components



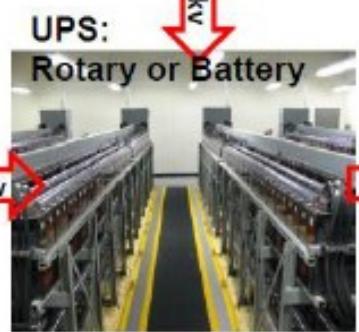
Power distribution



8% distribution loss
 $.997^3 \cdot .94 \cdot .99 = 92.2\%$



115kV



13.2kV

Transformers

13.2kV

Transformers

480V

~1% loss in switch gear & conductors

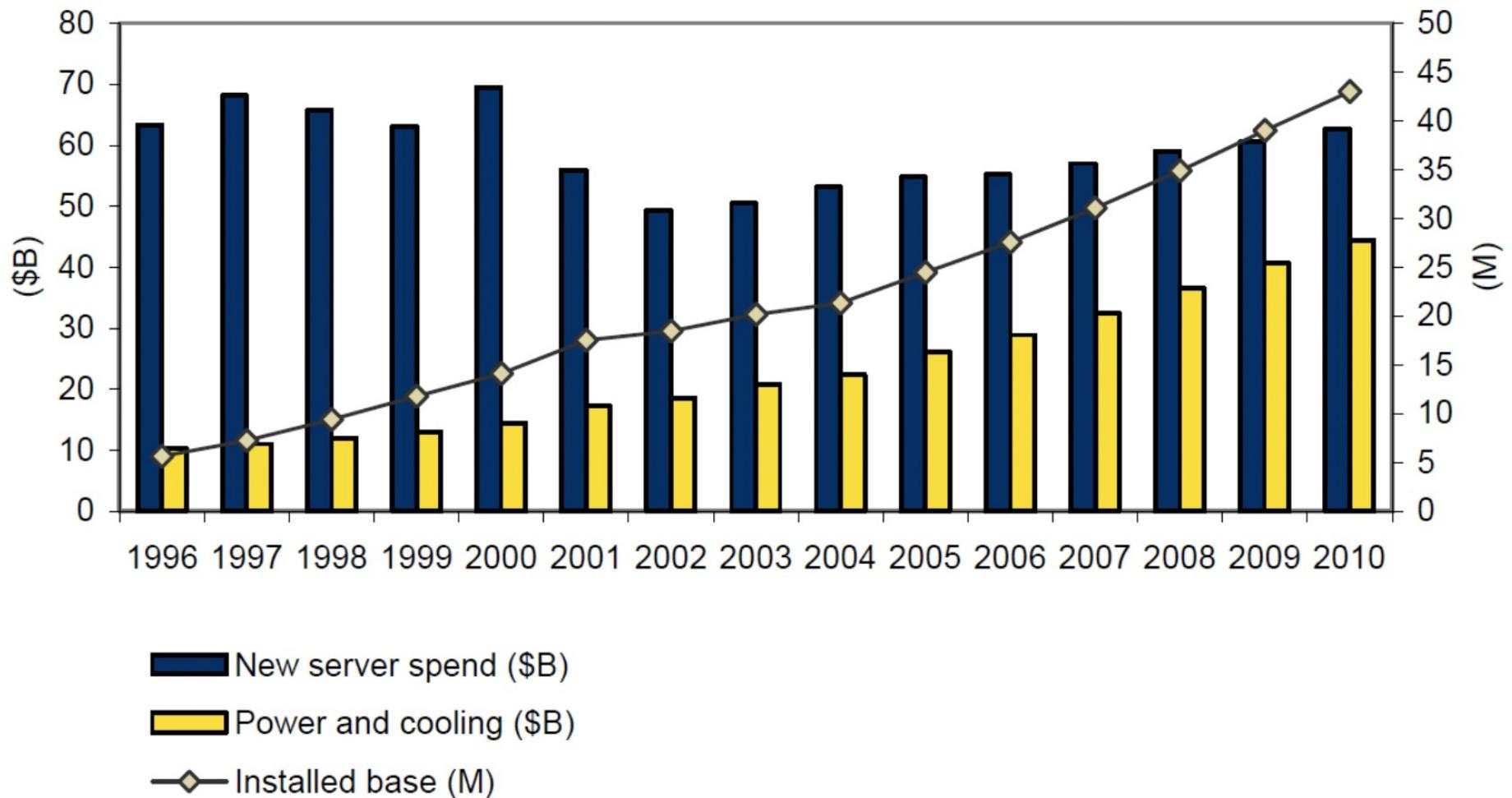
0.3% loss
99.7% efficient

6% loss
94% efficient, ~97% available

0.3% loss
99.7% efficient

0.3% loss
99.7% efficient

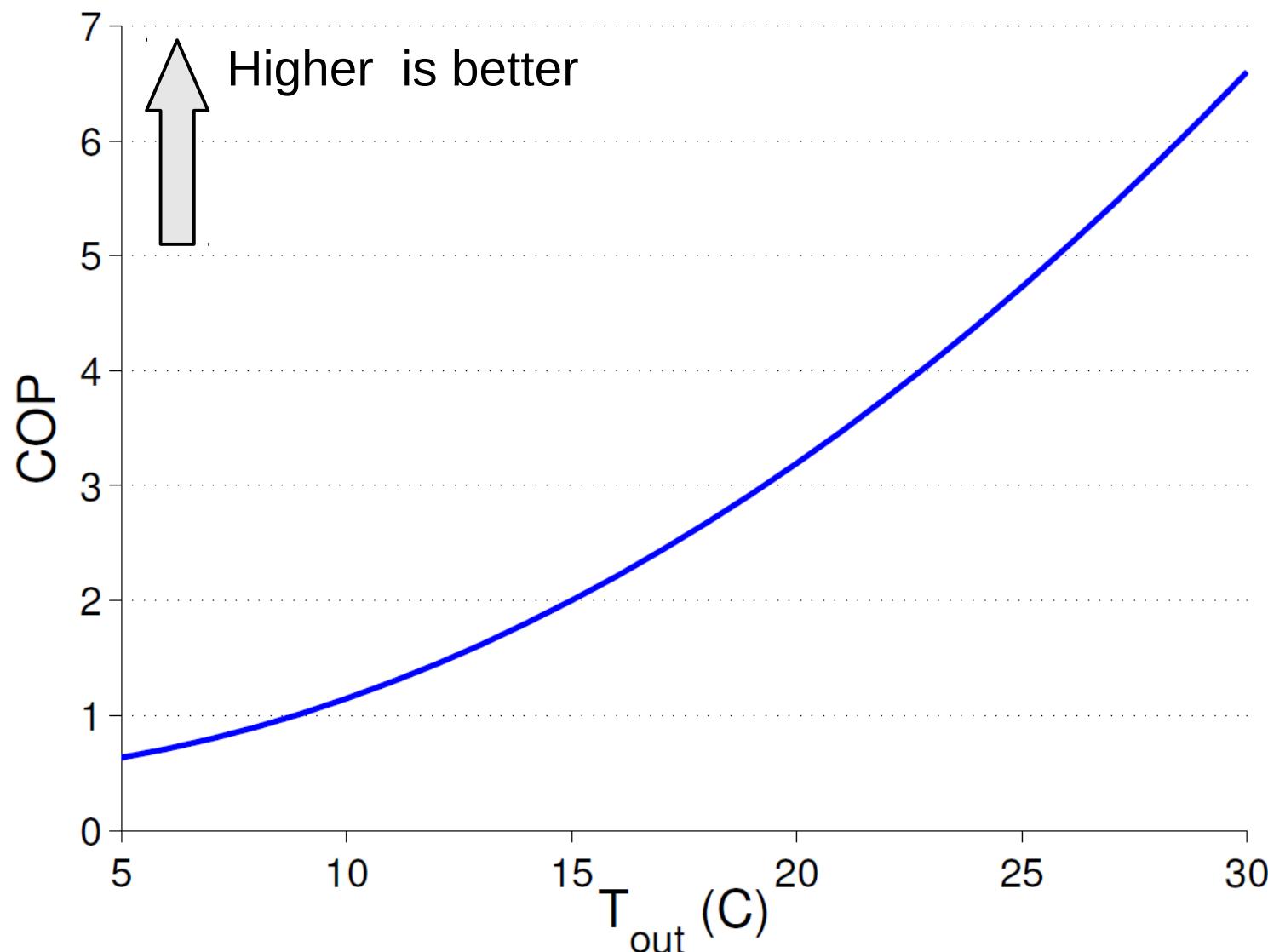
Spending trend



Data centers and cyber-physical systems

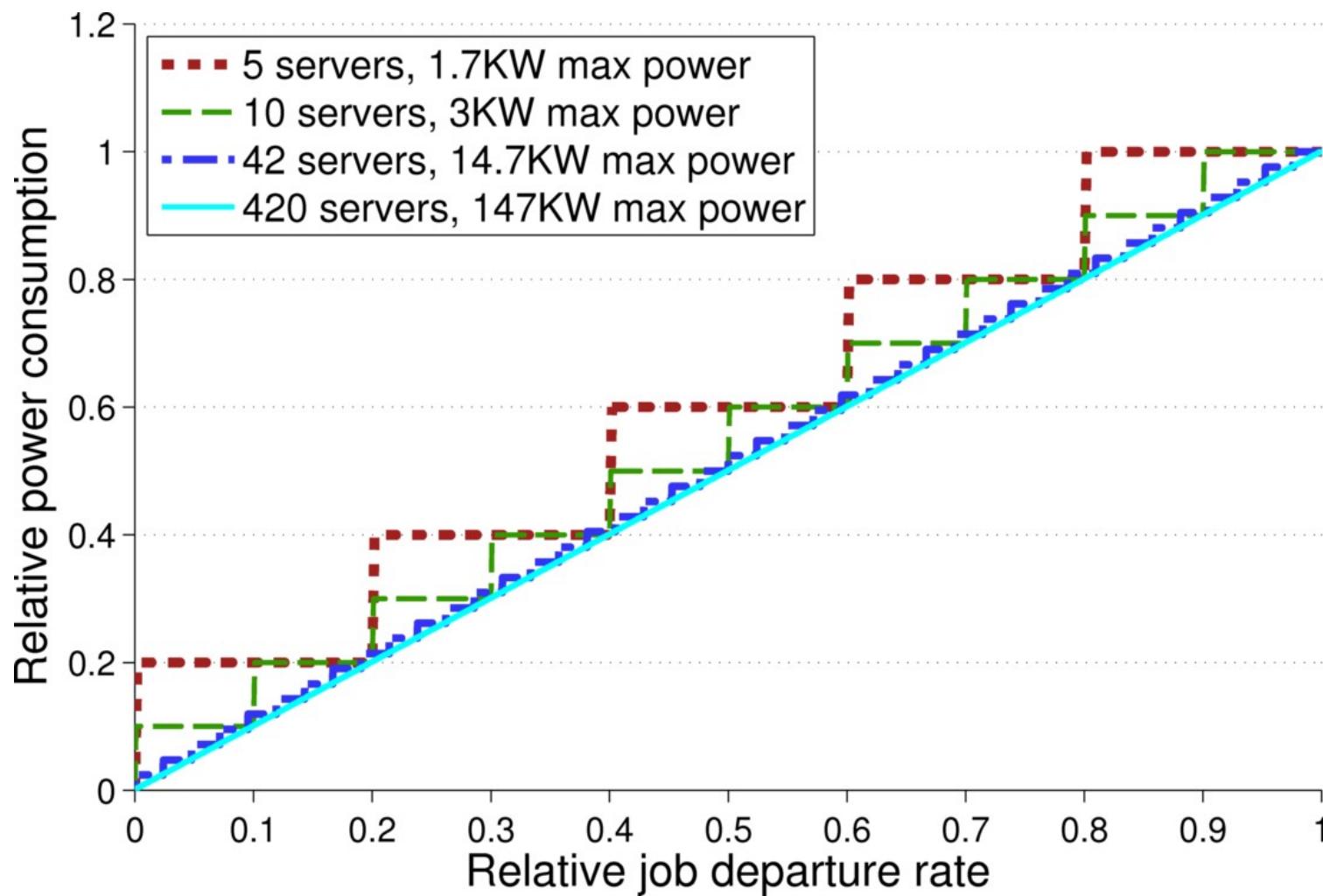
- We have designed CPS since the advent of the computer
 - A digitally controlled system is a CPS
- Is it any different now? If so what is different?
 - No separation of concerns
 - Coupling
 - Sensitivity
- Data Centers are a prime example of a CPS

Efficiency of CRAC units



Ensemble power consumption

- Single server max power consumption: 350 W



Examples of SLAs

- **Amazon Elastic Compute Cloud (EC2)**
 - “[...] If the Annual Uptime Percentage for a customer drops below 99.95% for the Service Year, that customer is eligible to receive a Service Credit equal to 10% of their bill [...]”
- **Google 500ms additional delay**
 - +20% reduction in search
- **Amazon 100ms additional delay**
 - +1% loss in sales

Examples of SLAs

- **Amazon Elastic Compute Cloud (EC2)**
 - “[...] If the Annual Uptime Percentage for a customer drops below 99.95% for the Service Year, that customer is eligible to receive a Service Credit equal to 10% of their bill [...]”
- **Google 500ms additional delay**
 - +20% reduction in search
- **Amazon 100ms additional delay**
 - +1% loss in sales

Issues in controlling a data center

- Large-scale nonlinear system
- Time scale of the controlled processes
- Cyber-physical coupling
- Lack of a widely accepted model

Issues in controlling a data center

■ Large-scale nonlinear system

- Example of nonlinearities
 - Response time with respect to the amount of available hardware resources
 - Server idle power consumption ~60 % of peak power
 - Cooling efficiency increases with temperature
 - Power conversion loss decreases with power consumption
 - Fan power consumption increases cubically with the speed of the fan
- Example of the system scale
 - Control the voltage and frequency of every CPU (e.g., 4 CPU per server, 50k servers → 200k CPUs)
 - On/Off state of every server
 - Speed of server and CRAC fan (100k fans)
 - Performance indicators, (e.g., 2 per server → ~100k)

Issues in controlling a data center

- **Time scale of the controlled processes**
 - Server response time to HTTP request ~100ms
 - Server turn-on time ~5min
 - Data center thermal dynamic in the hour time scale

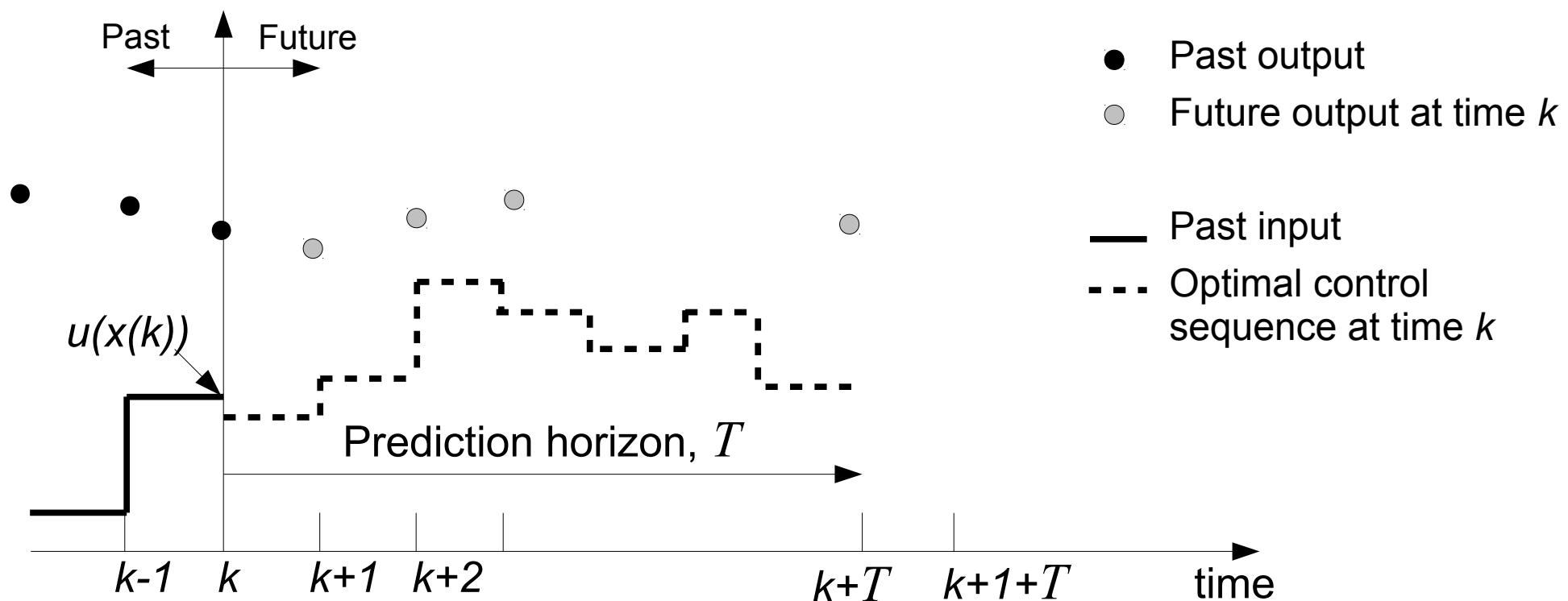
Issues in controlling a data center

- **Time scale of the controlled processes**
 - Server response time to HTTP request ~100ms
 - Server turn-on time ~5min
 - Data center thermal dynamic in the hour time scale
- **Cyber-physical coupling**
 - The way workload is executed and where it is executed affects the cooling power consumption
 - The capacities of the cooling and of the power distribution subsystem set limits to the computational subsystem

Brief introduction to Model Predictive Control (MPC)

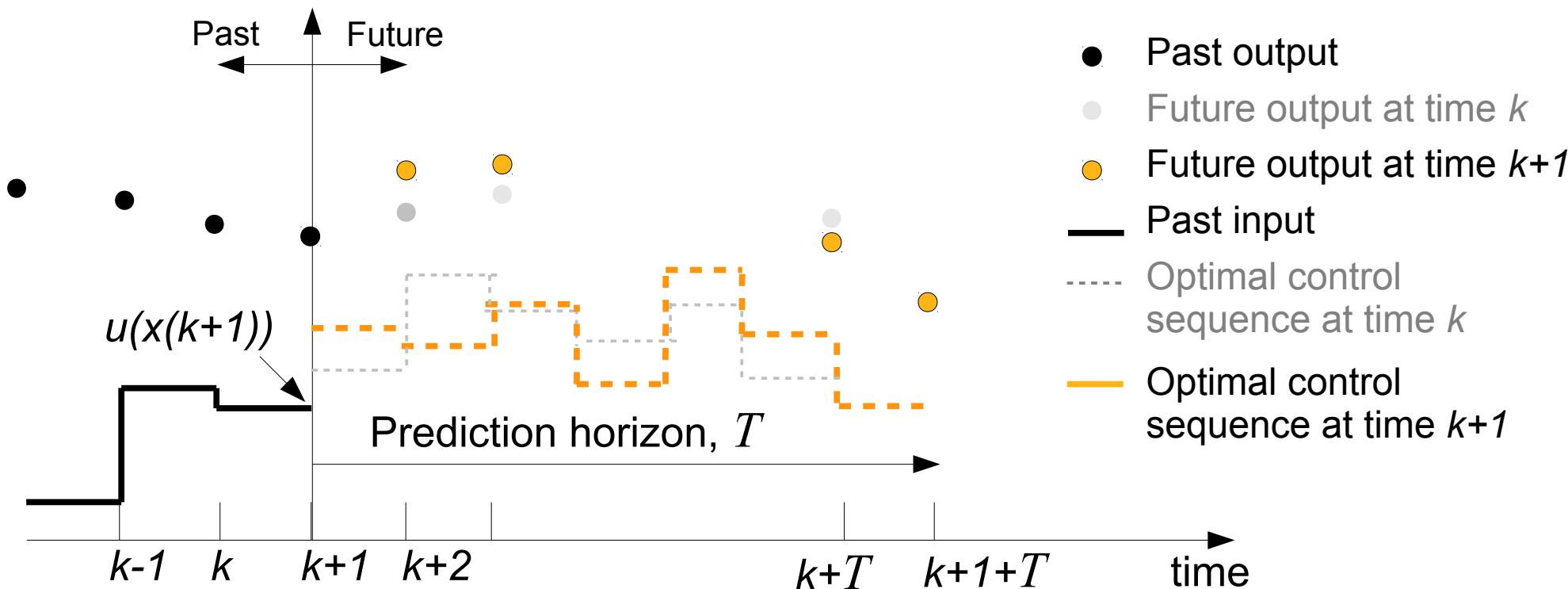
Model predictive control (MPC)

- The controller uses an *internal model of the plant* to predict future states and outputs of the plant
- *The optimal control sequence may change at every step*



Model predictive control (MPC)

- The controller uses an *internal model of the plant* to predict future states and outputs of the plant
- *The optimal control sequence may change at every step*



Thermal and computational network

Thermal network

■ Linear model

$$\hat{T}_{\text{out}}(\nu + 1|k) = A_{T,D} \hat{T}_{\text{out}}(\nu|k) + B_{T,D} [\hat{\mathbf{p}}_N(\nu|k)^T \hat{\mathbf{T}}_{\text{ref}}(\nu|k)^T]^T$$

Output temperature of zones and CRAC units
Power consumption of zones *Reference temperature of CRAC units*

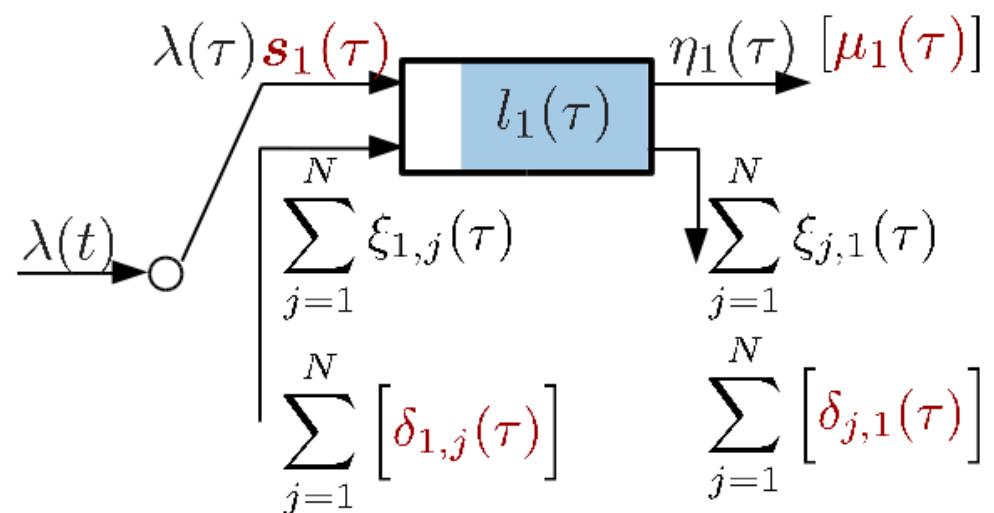
$$\hat{T}_{\text{in}}(\nu|k) = \Psi \hat{T}_{\text{out}}(\nu|k) \leq \overline{T}_{\text{in}}$$

■ CRAC power consumption depends on the coefficient of performance (COP)

$$p_i(t) = \frac{\dot{Q}_i(t)}{COP_i(T_{\text{out},i}(t))}$$

Heat removed rate (W)

Computational network



$$\xi_{j,1}(\tau) = \begin{cases} \delta_{j,1}(\tau) & \text{if } l_1(\tau) > 0 \\ \frac{\delta_{j,1}(\tau)}{N} (a_1(\tau) - \nu_1(\tau)) & \text{or } a_1(\tau) > \nu_1(\tau) \\ \sum_{j=1}^N \delta_{j,1} & \text{otherwise} \end{cases}$$

$$i(\tau) = a_1(\tau) - d_1(\tau)$$

$$d_1(\tau) = \eta_1(\tau) + \sum_j^N \xi_{j,i}(\tau)$$

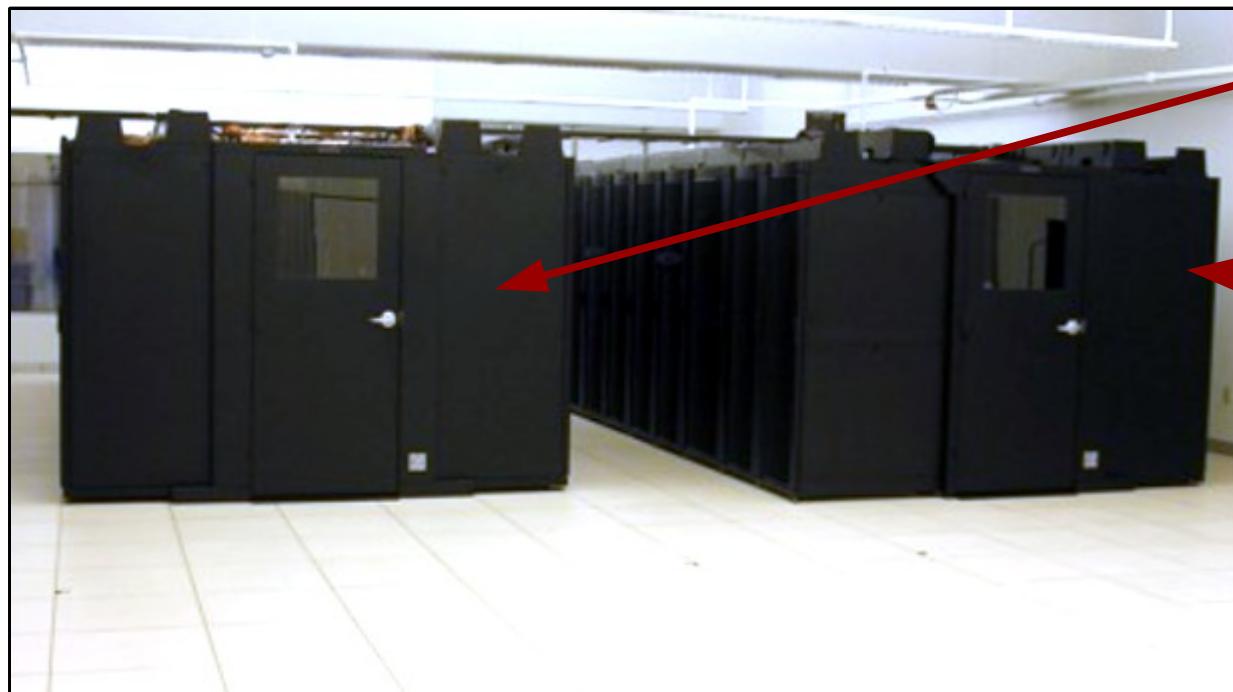
$$a_1(\tau) = \lambda(\tau)s_1(\tau) + \sum_{j=1}^N \xi_{i,j}(\tau)$$

$$\eta_1(\tau) = \begin{cases} \mu_1(\tau) & \text{if } l_1(\tau) > 0 \\ a_1(\tau) & \text{or } a_1(\tau) > \mu_1(\tau) \\ a_1(\tau) & \text{otherwise} \end{cases}$$

Data collection in the
Data Center Observatory (DCO)
Carnegie Mellon University
Pittsburgh, Pa

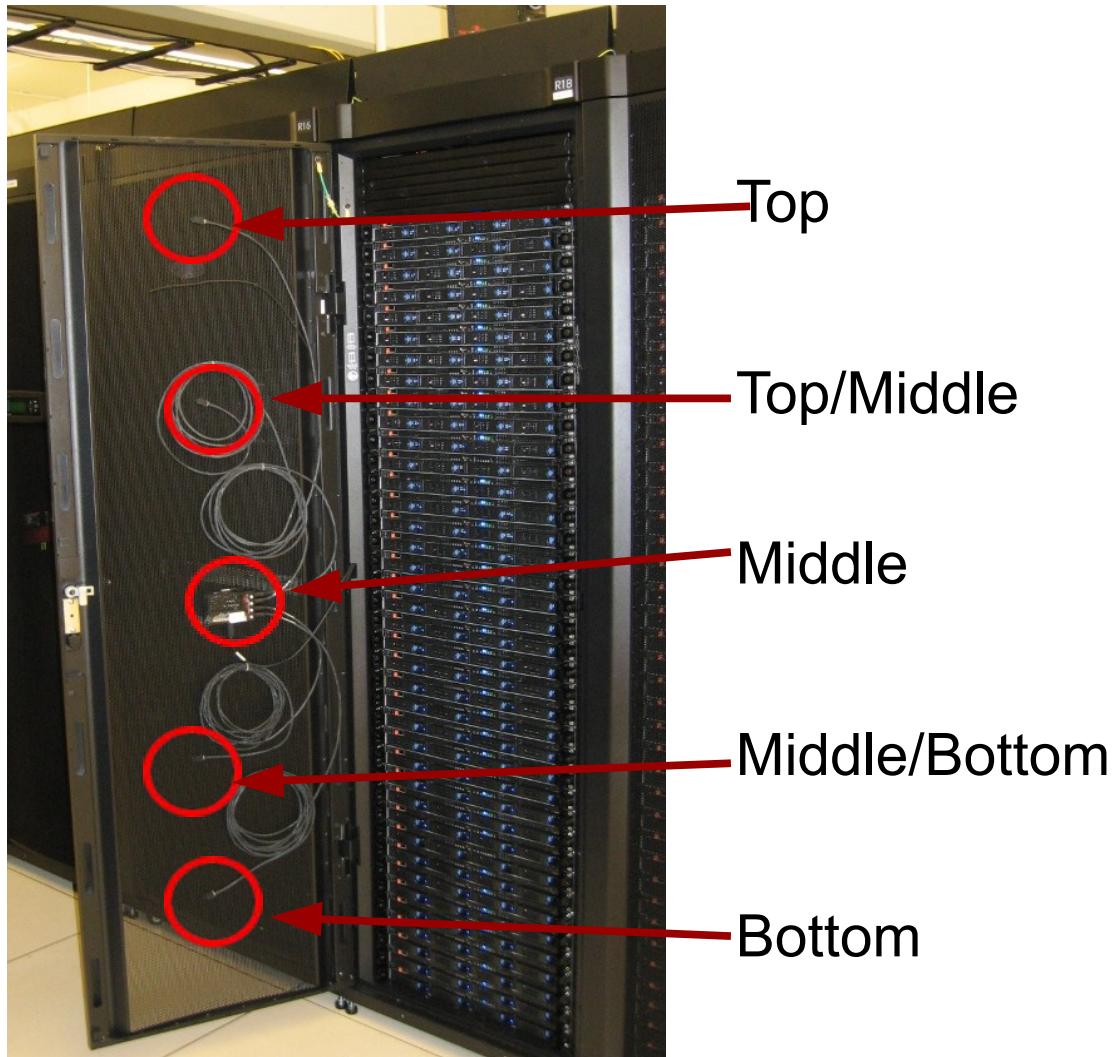
Data collection in the DCO

- **Data center observatory (DCO)**
 - Testbed for research on data-intensive workload

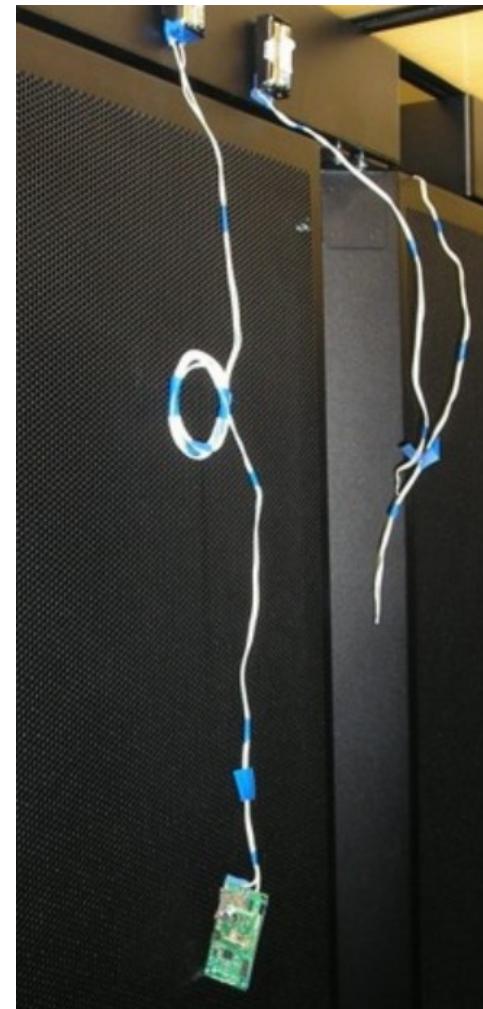


Temperature and humidity sensors

■ Wired sensors

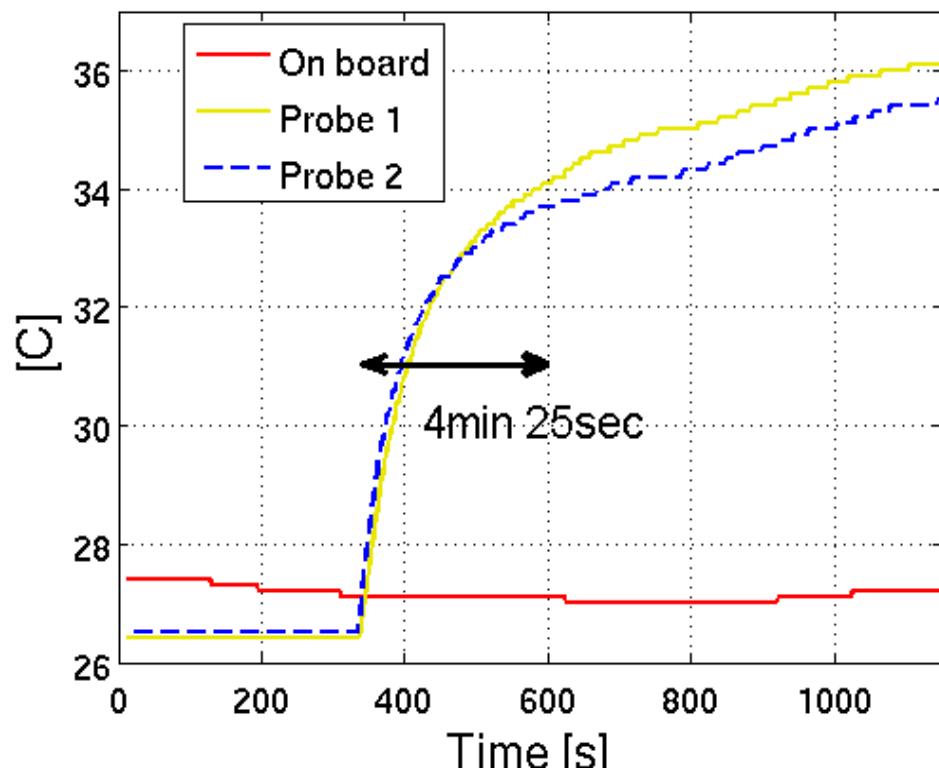


■ Wireless sensors



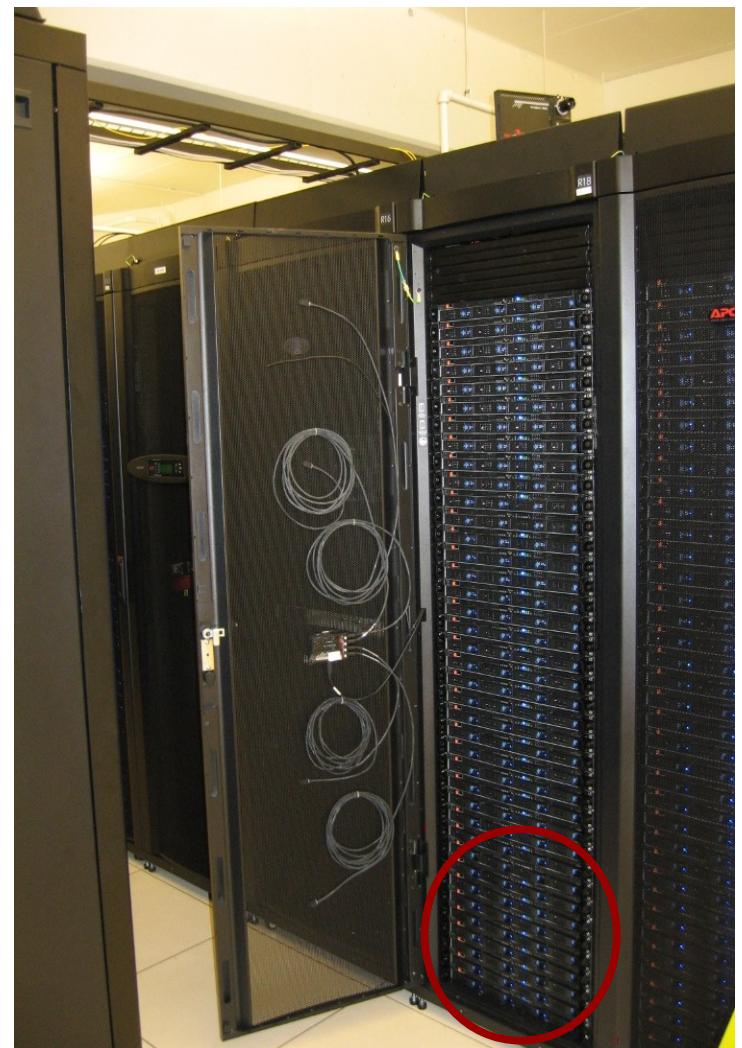
Wired sensors

- APC NetBotz sensor
- Time constant ~53 sec
 - Wireless sensors have time constant ~13sec



Thermal test on 6 servers

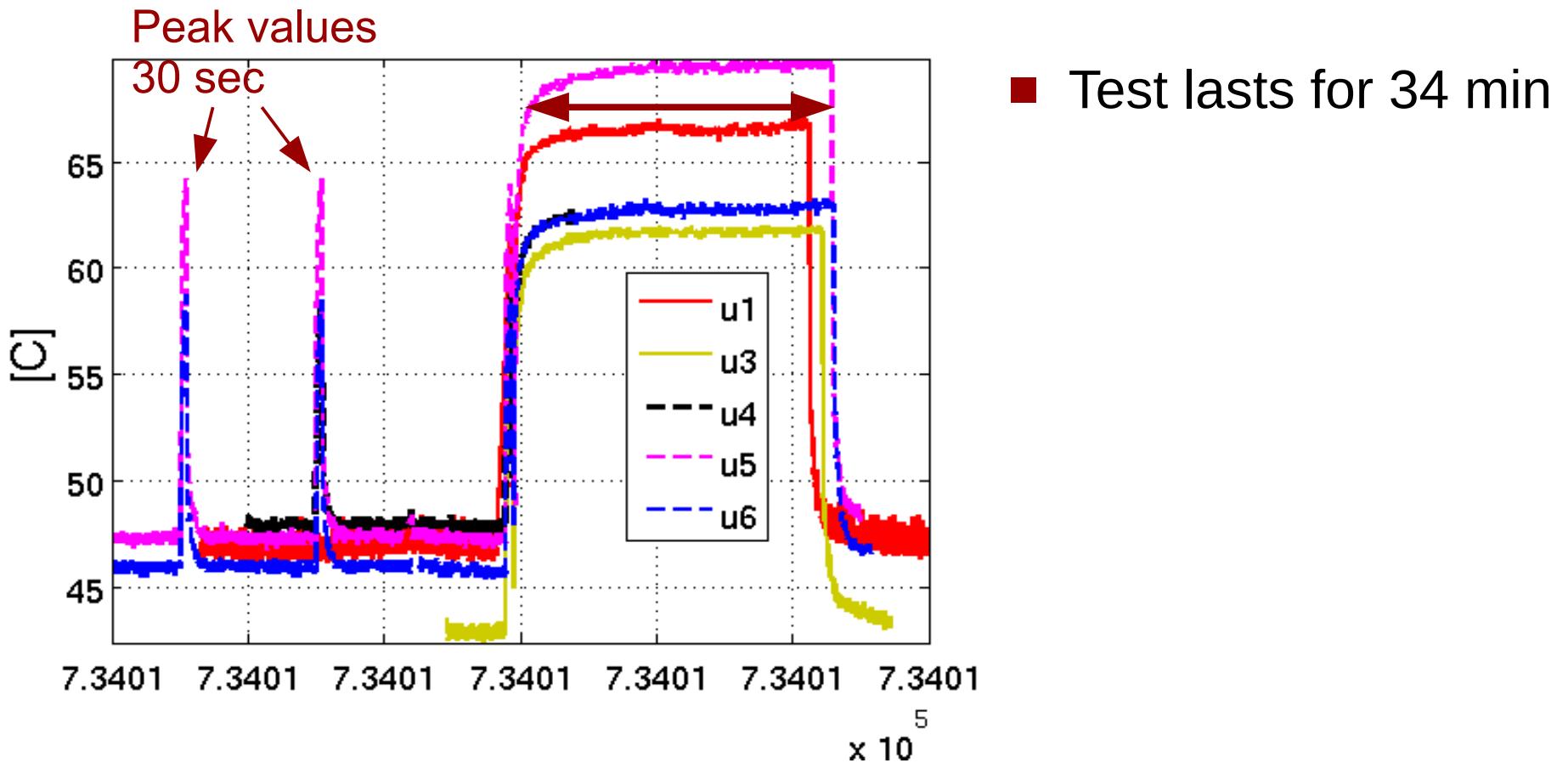
- **Time varying workload**
 - Measure server power consumption
 - CPU temperature (on-board sensors)
 - Inlet and outlet temperatures



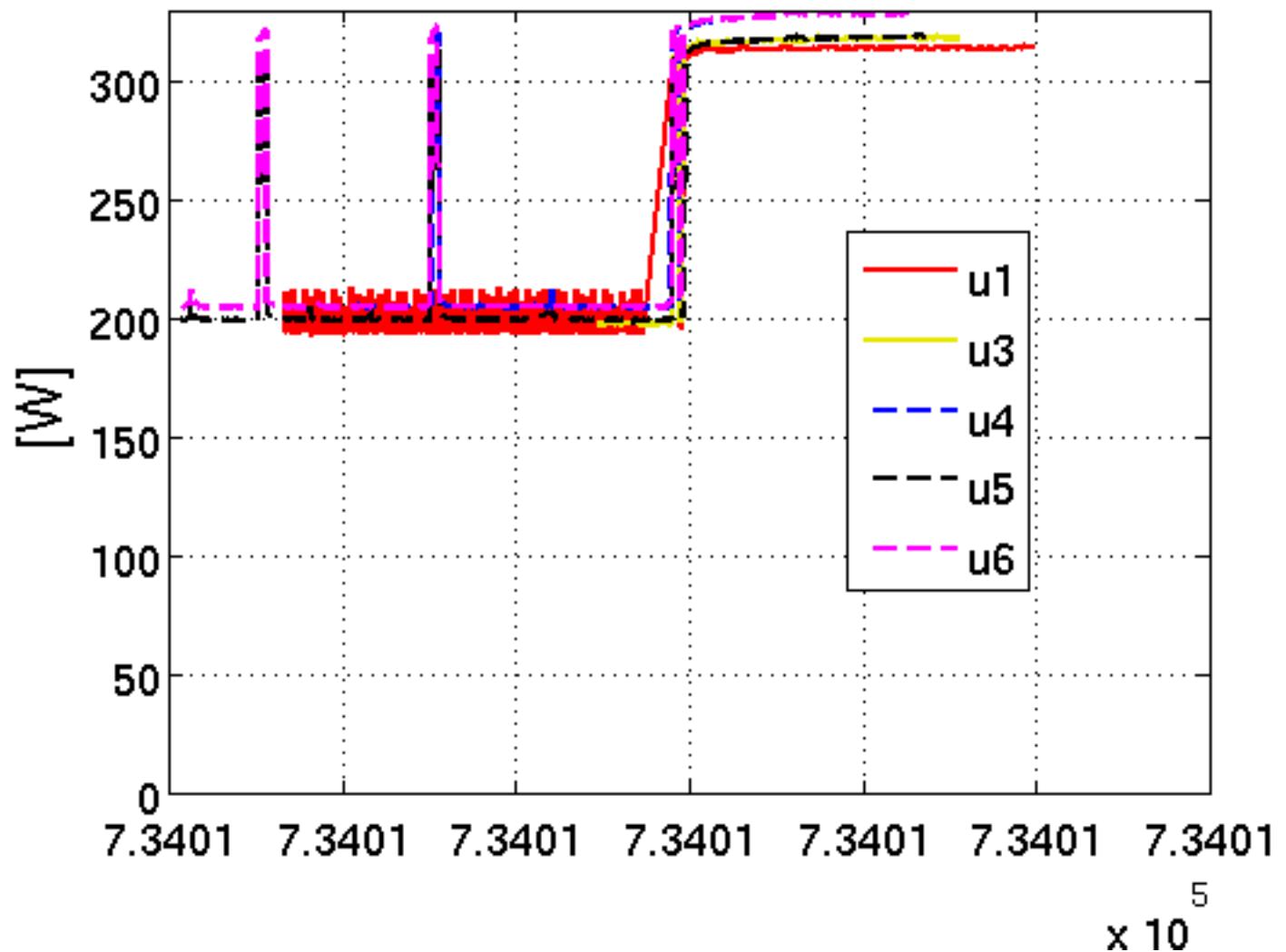
CPU temperature

- **8 Cores per CPU**

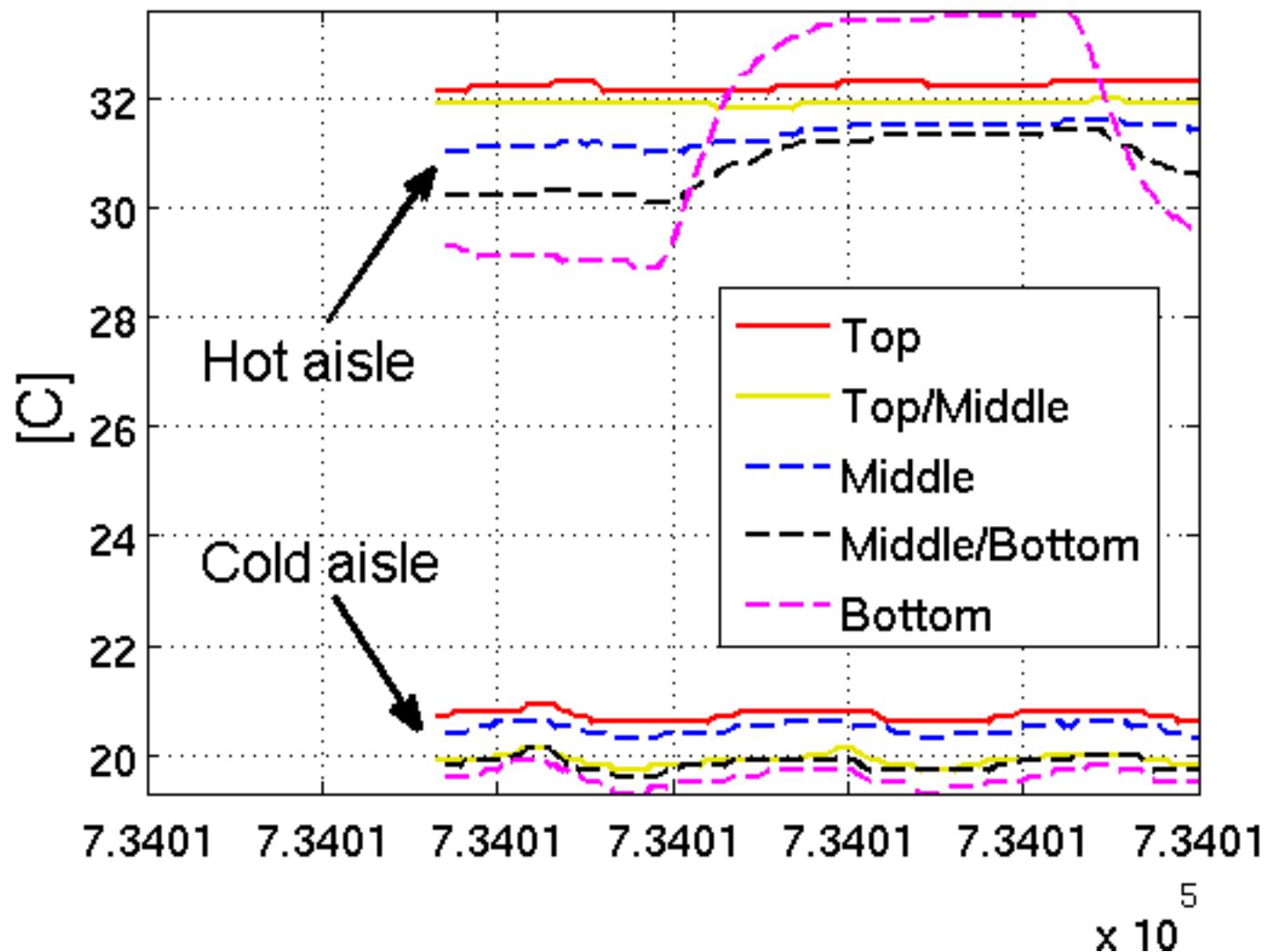
- We consider the average value of the 8 cores



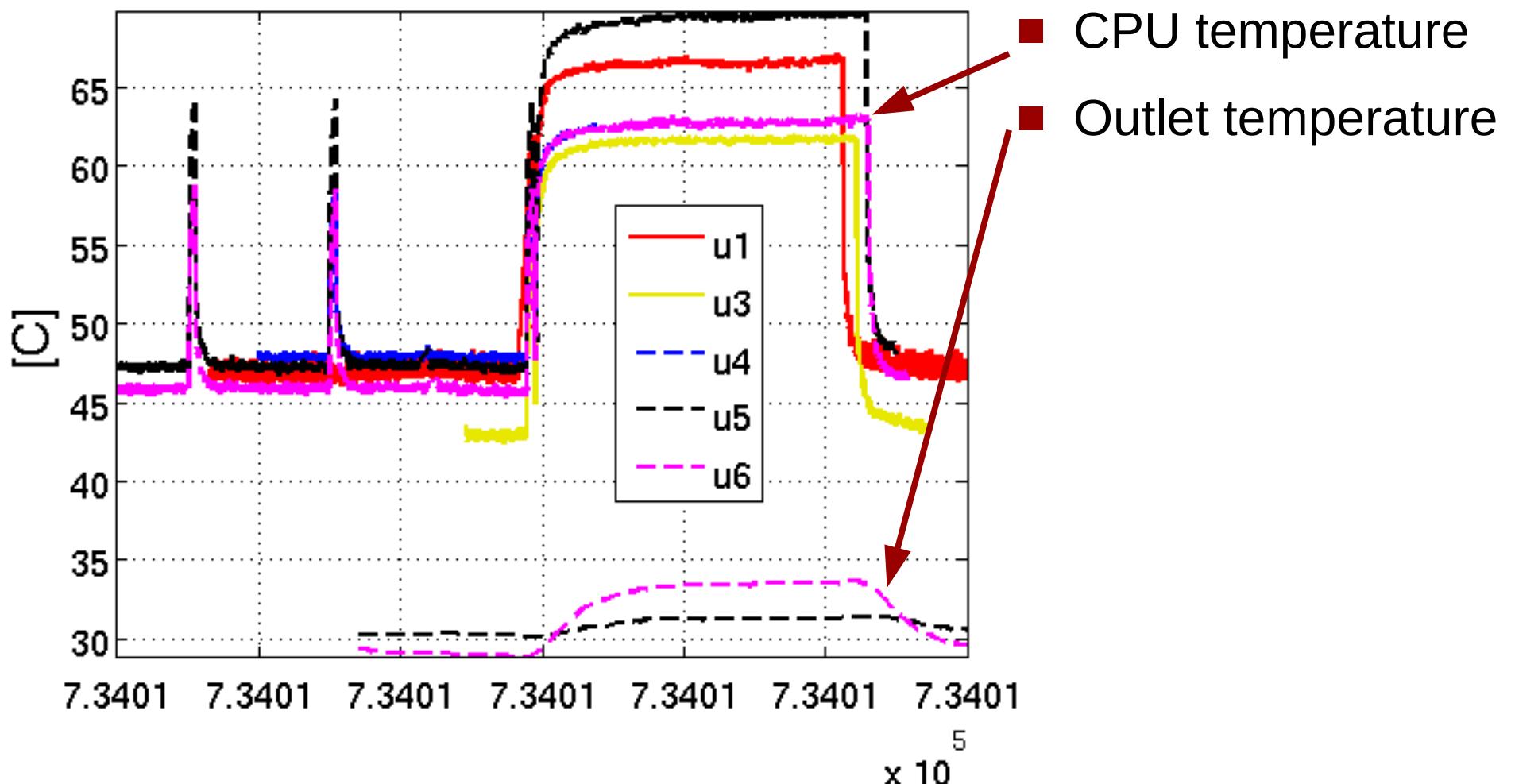
Server power consumption



Inlet and outlet temperatures



Outlet and CPU temperatures



Estimation of server power consumption

Joint work with
Anshul Gandhi and Kevin Woo

Servers

■ Dave's cn007

- 8 core Intel Xeon X5355 2.67 GHz

■ PH server

- 2 core Intel Xeon 3.00 GHz

Tools

■ Software

- SAR
 - CPU, network, disk
- Perfmon
 - Cache misses
- Powertop
 - P-states and C-states

■ Hardware

- Watts-Up Pro
 - Power consumption

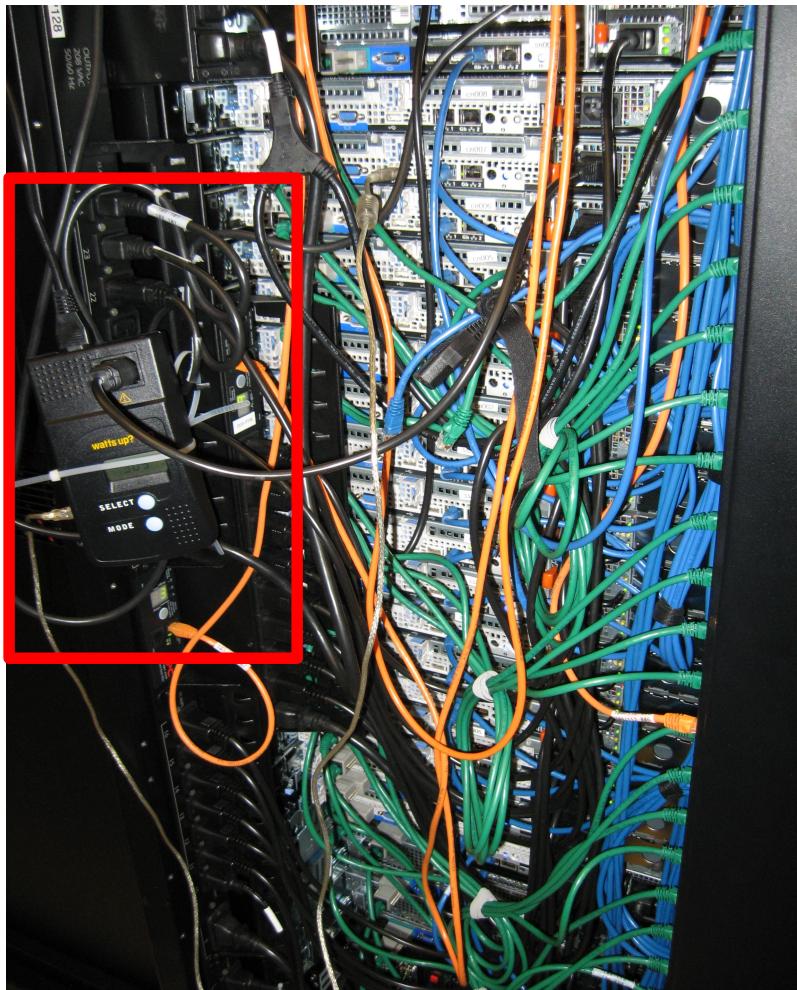
Perfmon

- Tool developed by HP
 - Exposes hardware performance counters
- Composed of 3 components
 - Kernel Patch
 - Kernel Module
 - User Interface
- Requires kernel recompilation
 - Unstable
- Not tried on cn007, tested on PH

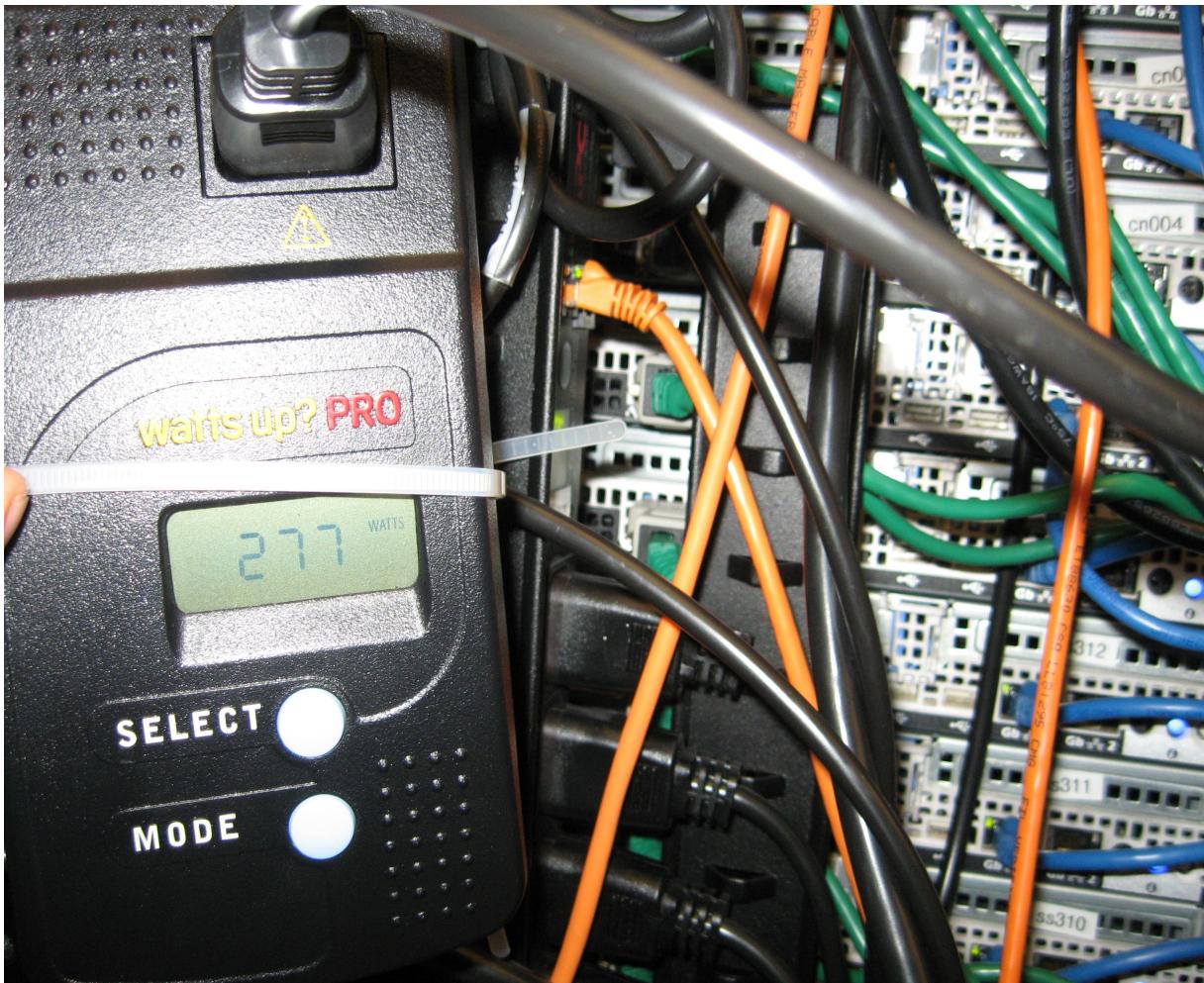
cn007



Watts-Up Pro Installation

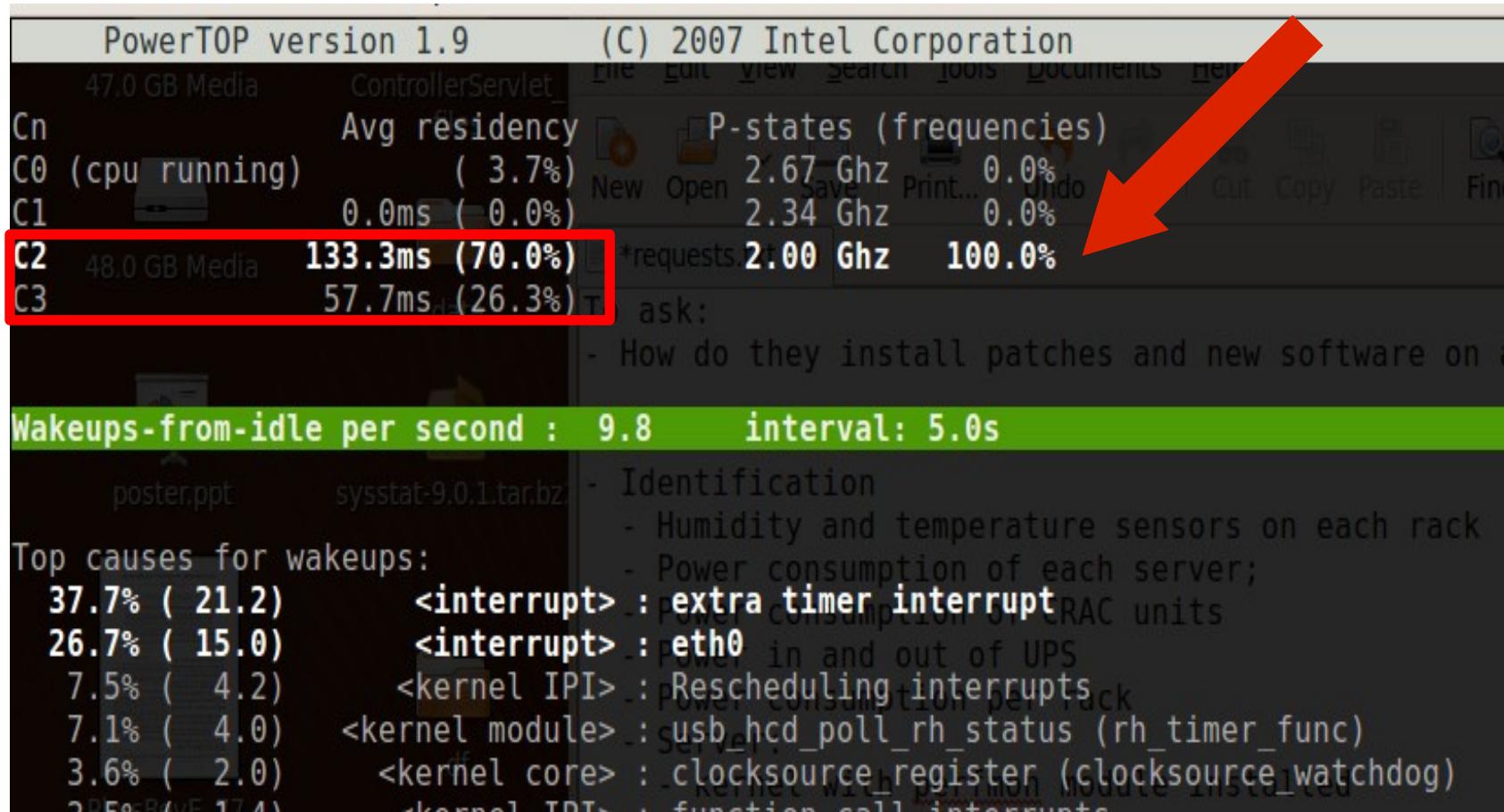


Watts-Up Pro Monitoring



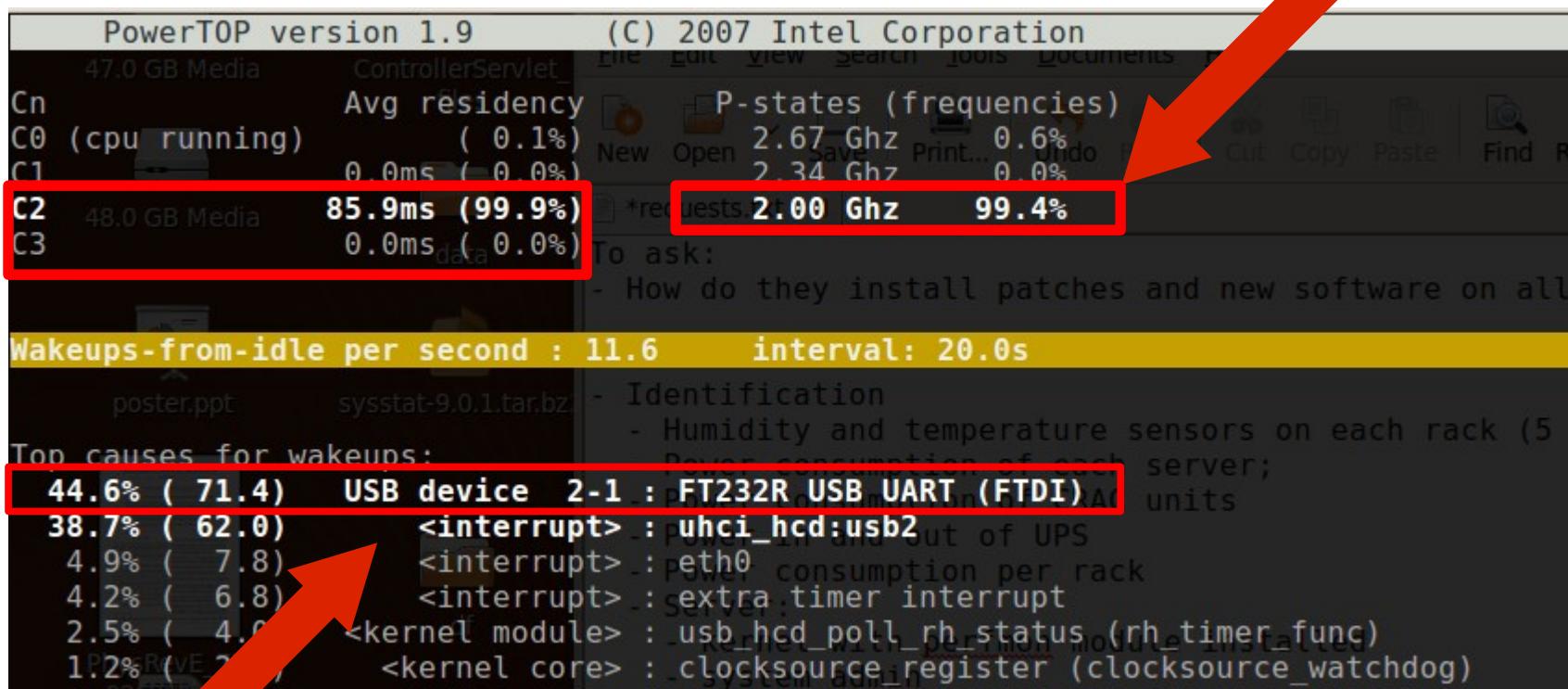
Script Overhead

Idle System



Script running

Low script overhead



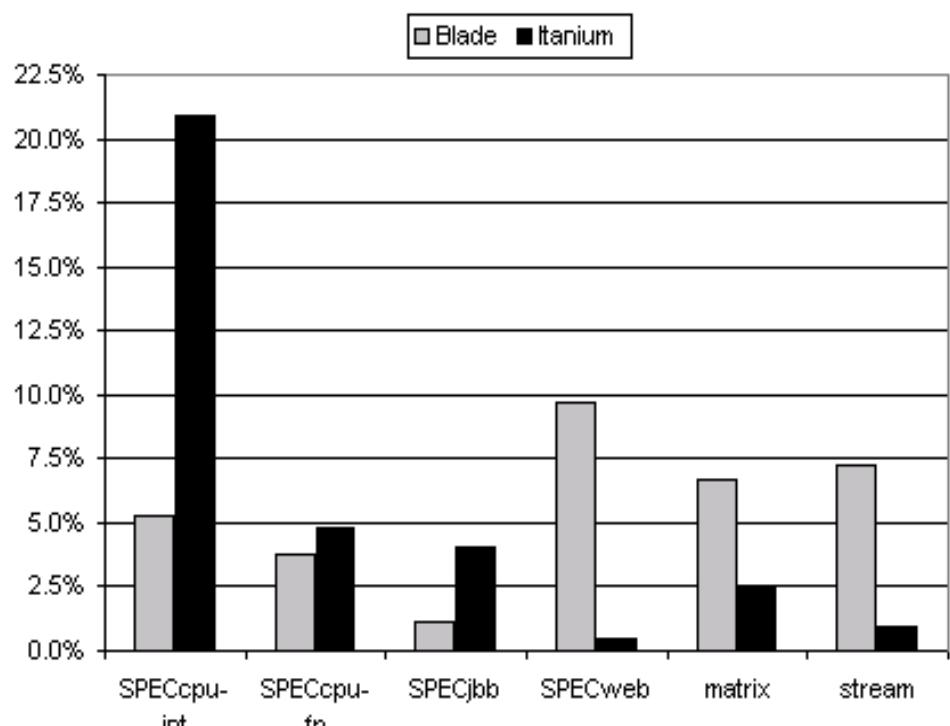
Watts-Up Communication

No Polling

```
File Edit View Terminal Help
PowerTOP version 1.9          (C) 2007 Intel Corporation

Cn          Avg residency      P-states (frequencies)
C0 (cpu running)    ( 4.6%)    2.67 Ghz    0.0%
C1           0.0ms ( 0.0%)    2.34 Ghz    0.0%
C2           177.6ms (73.8%)   2.00 Ghz   100.0%
C3           94.2ms (21.7%)
Wakeups-from-idle per second : 6.5      interval: 5.0s
Top causes for wakeups:
 39.2% ( 14.2) <interrupt> · extra timer interrupt
```

Previous work



90th percentile estimation error.

[Economou06]

- Server power consumption is a (affine) linear combination of:

- CPU utilization
- Off-chip memory access
- HD I/O rate
- Network I/O rate

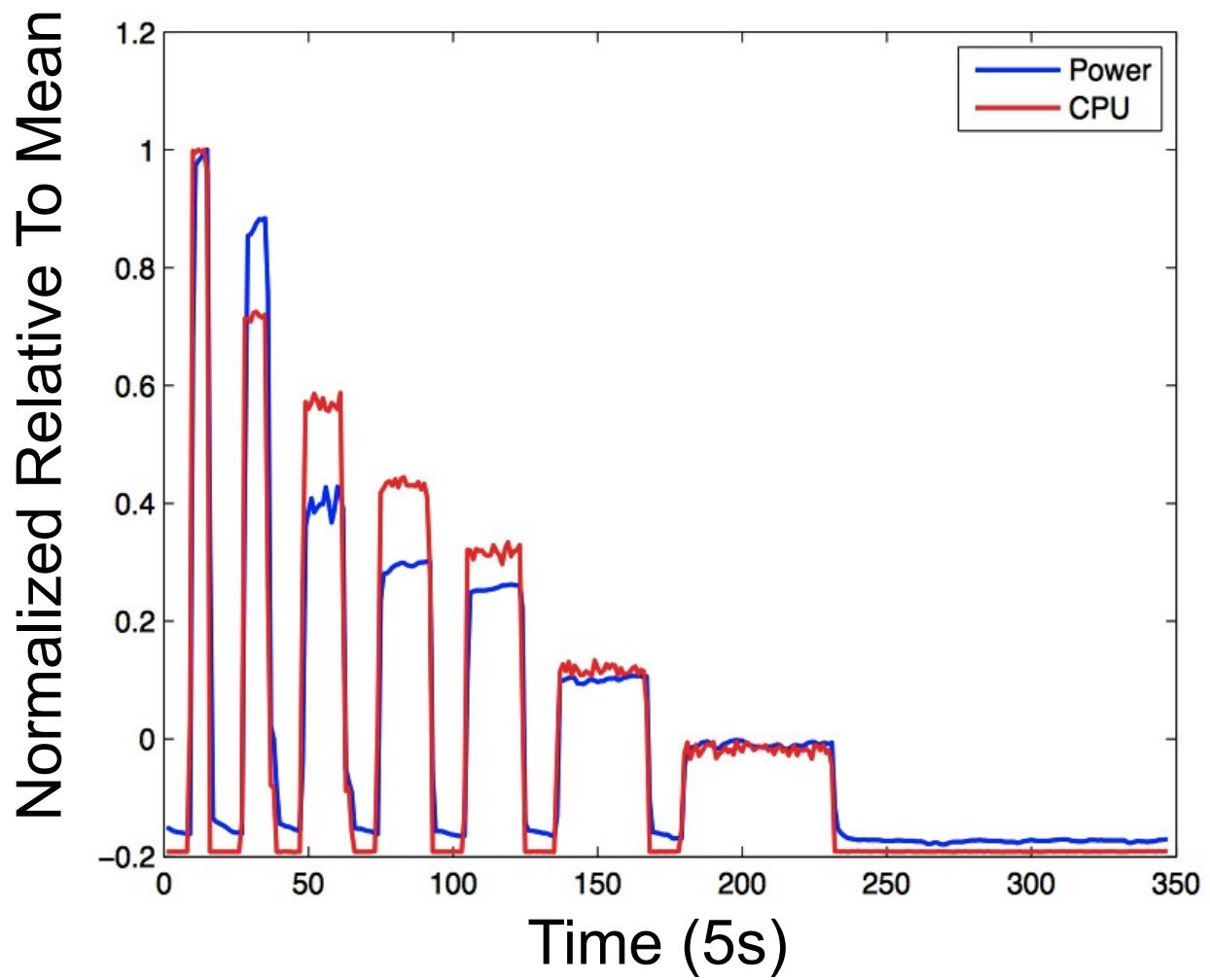
Model

■ Hypothesis

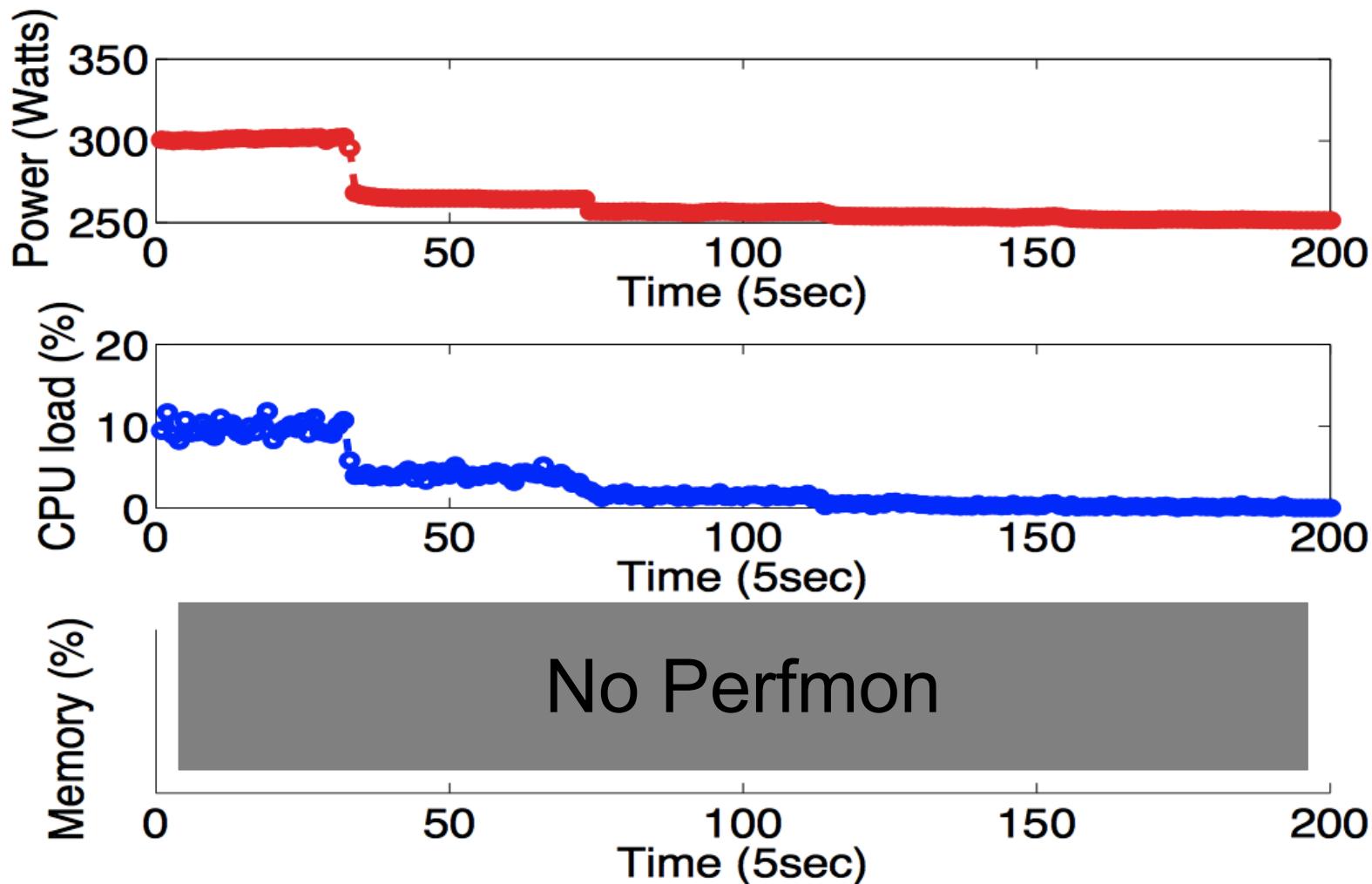
- Server power consumption only depends on current data values
- Server power consumption is a (affine) linear combination of ?

$$Ax = \begin{bmatrix} ? \end{bmatrix}^b \begin{bmatrix} x \end{bmatrix} = \begin{bmatrix} \text{power} \end{bmatrix}$$

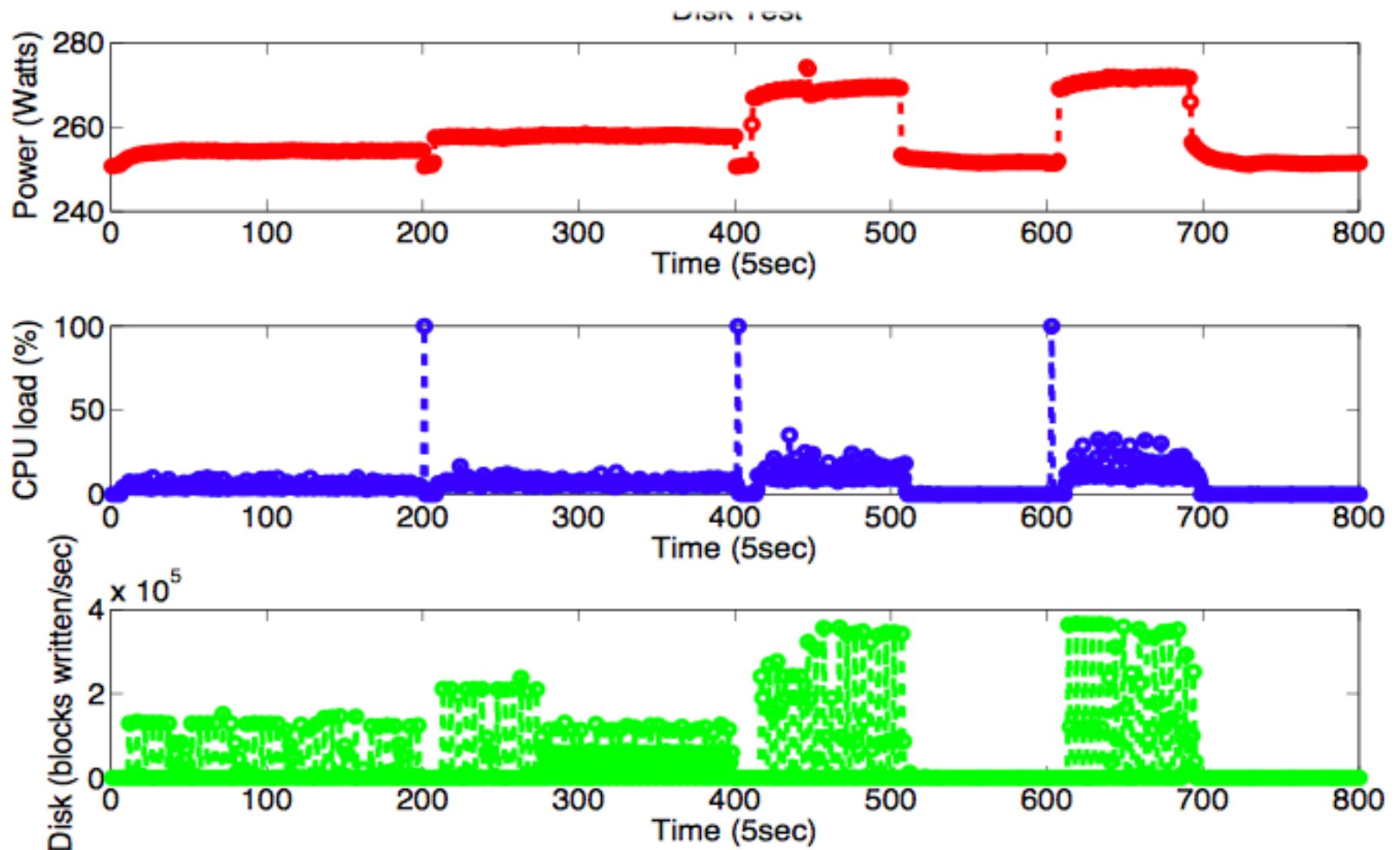
CPU Test



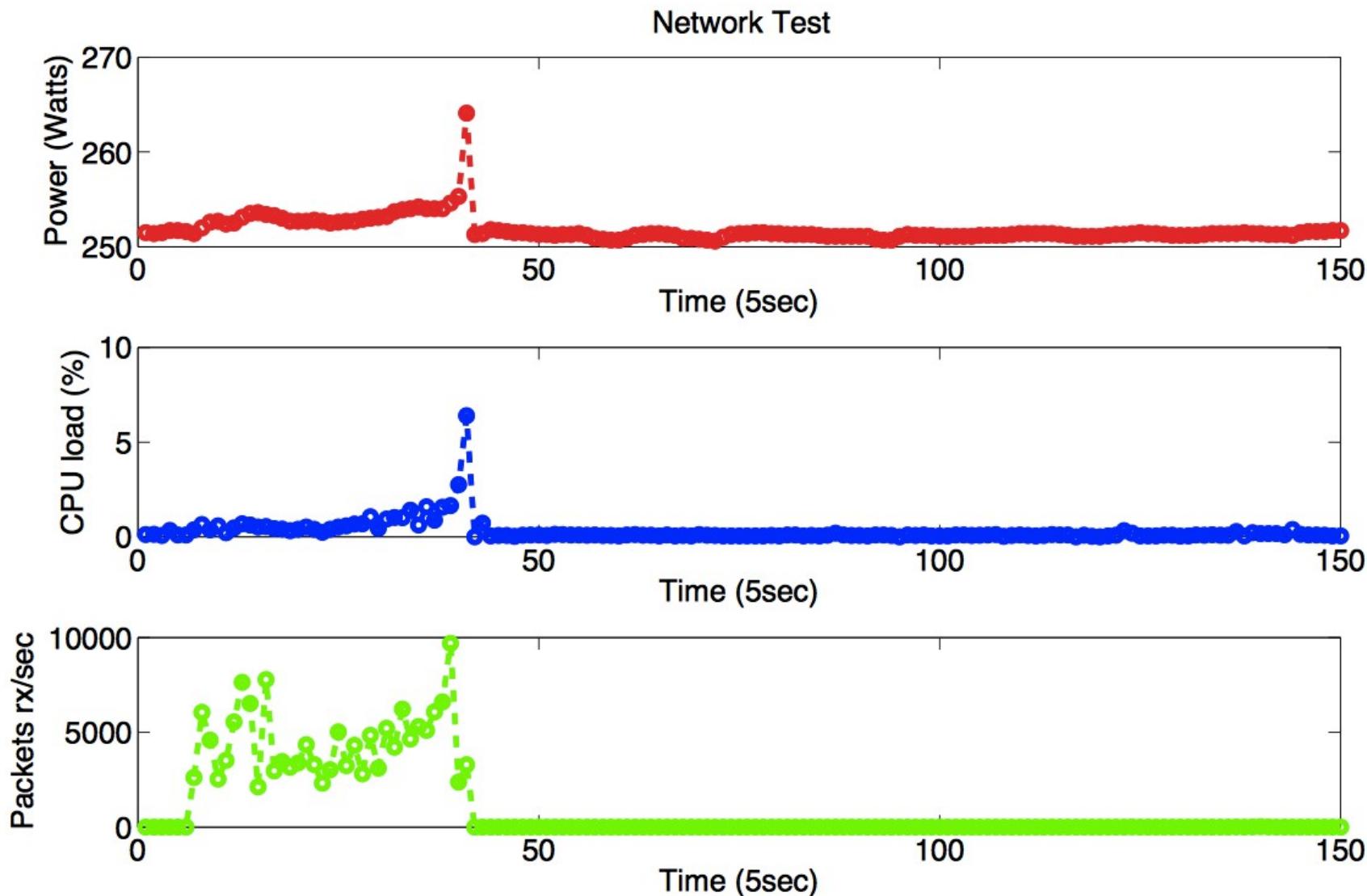
Memory Test



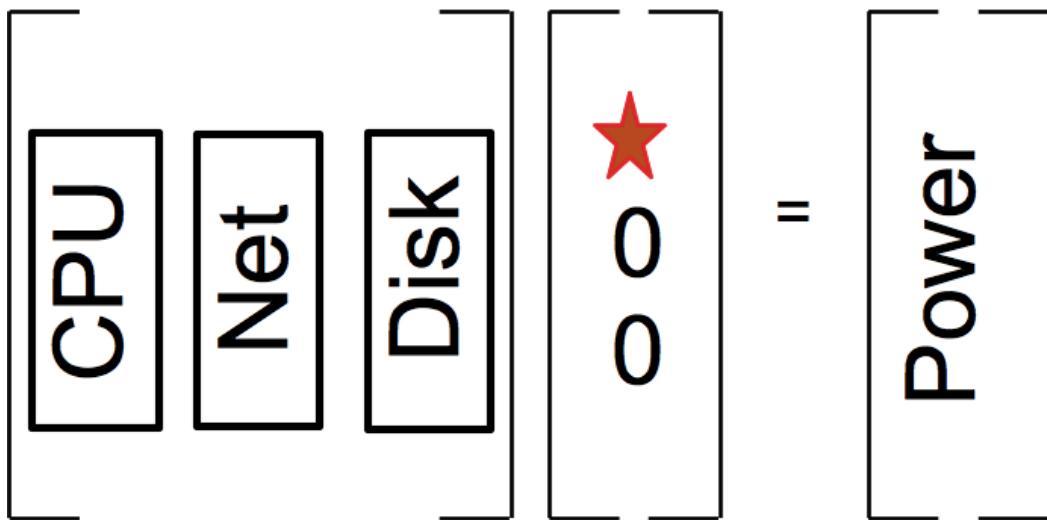
Disk Test



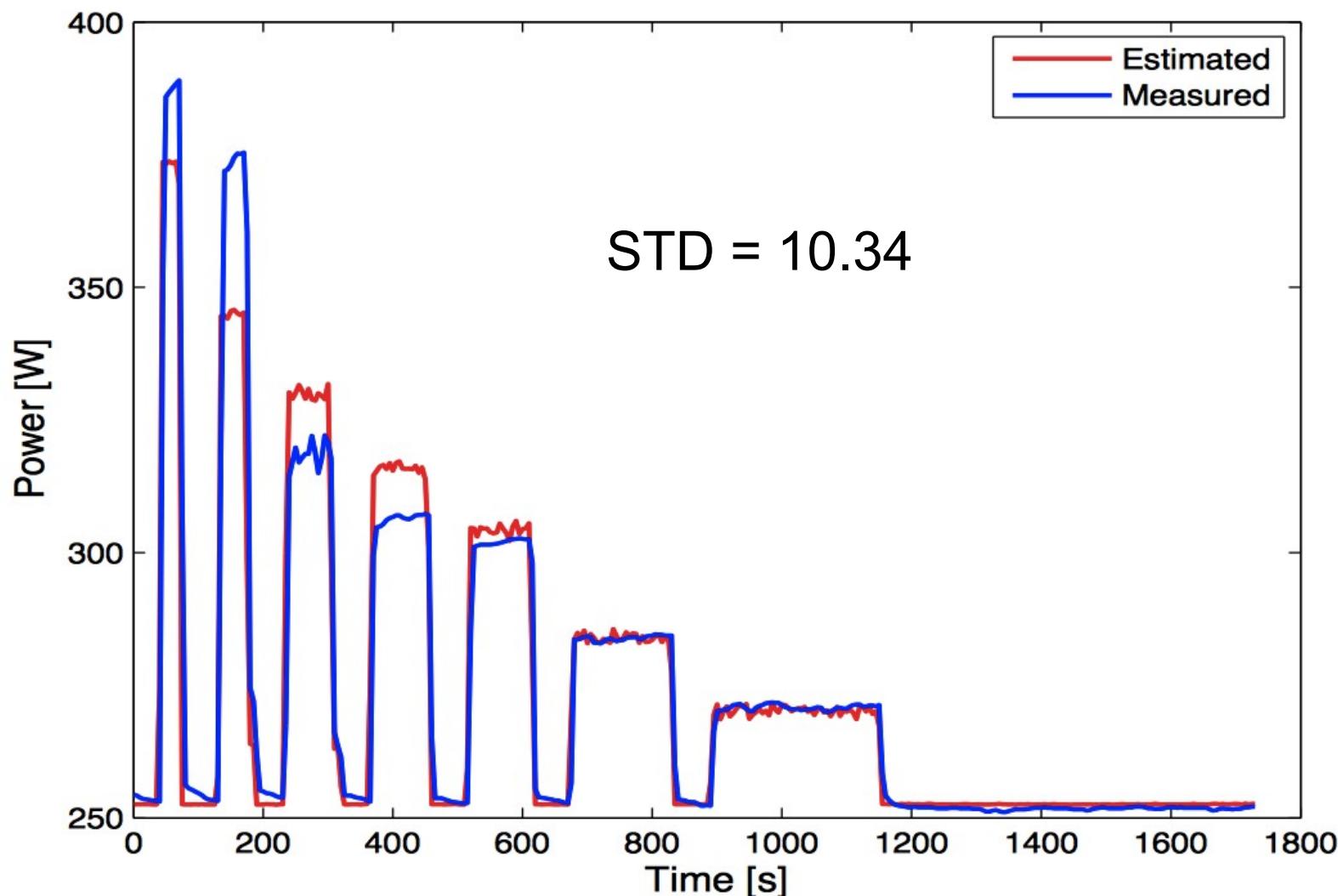
Network Test



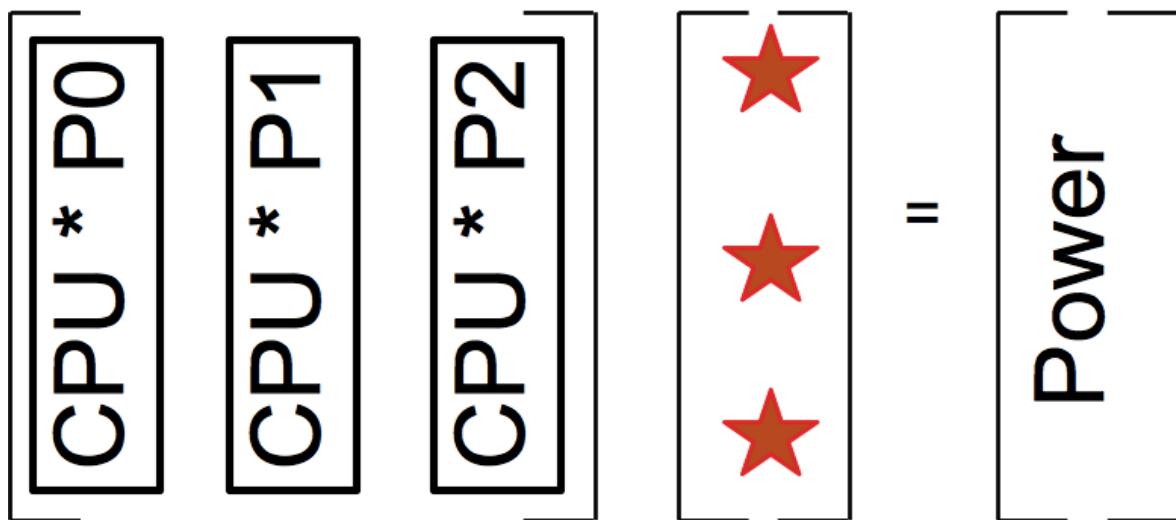
Model Analysis



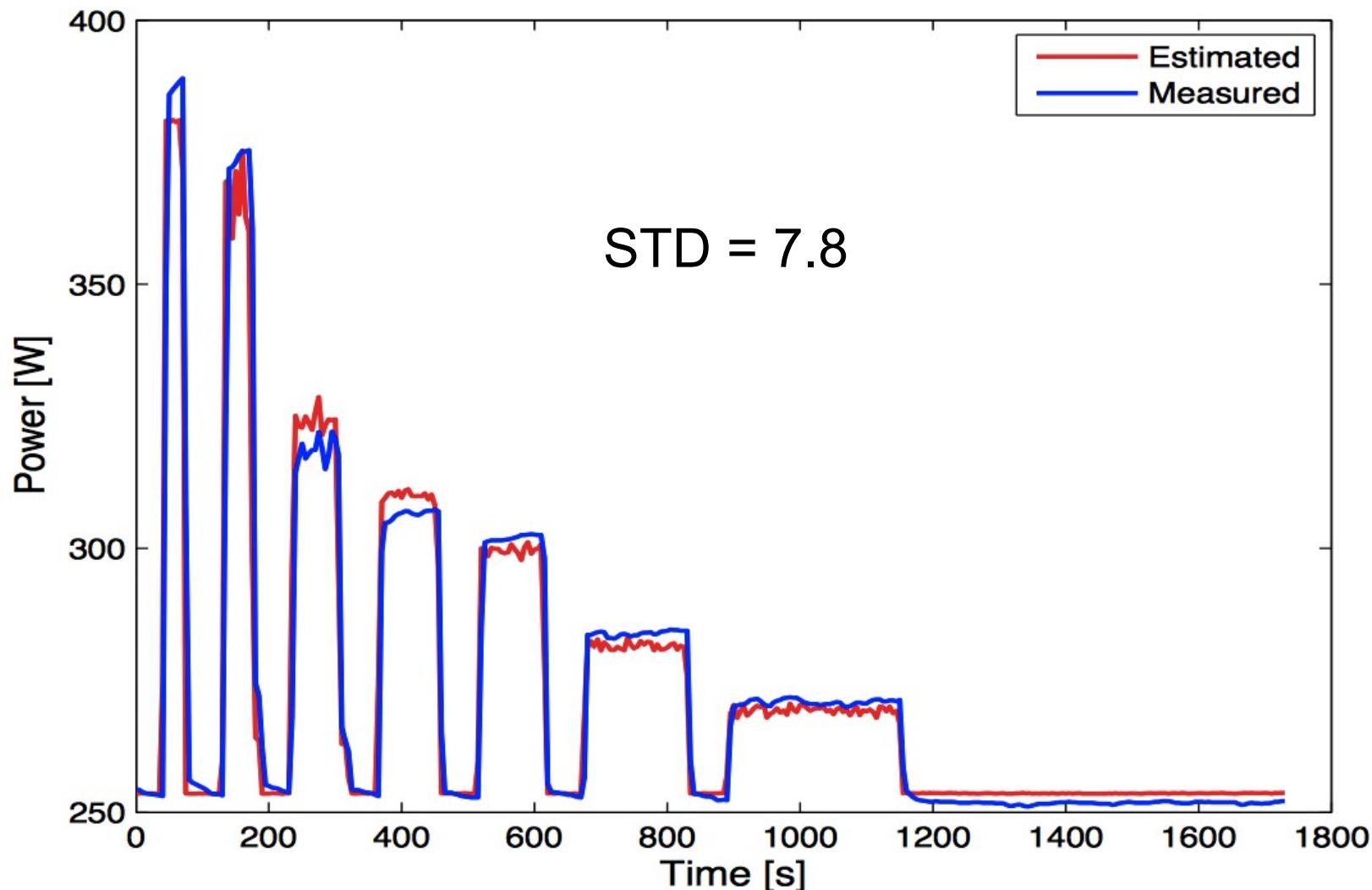
Estimation - CPU data only



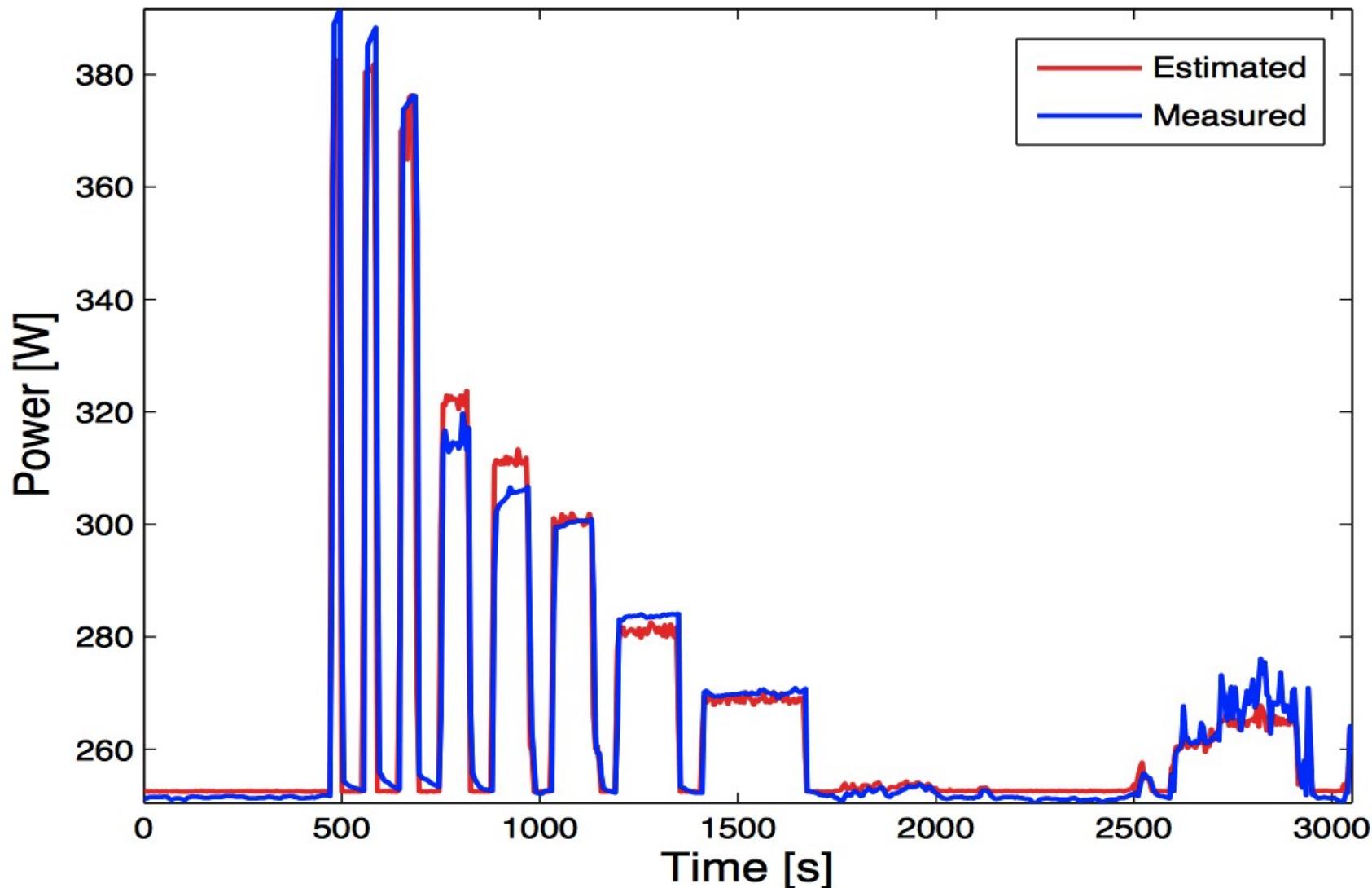
New Model



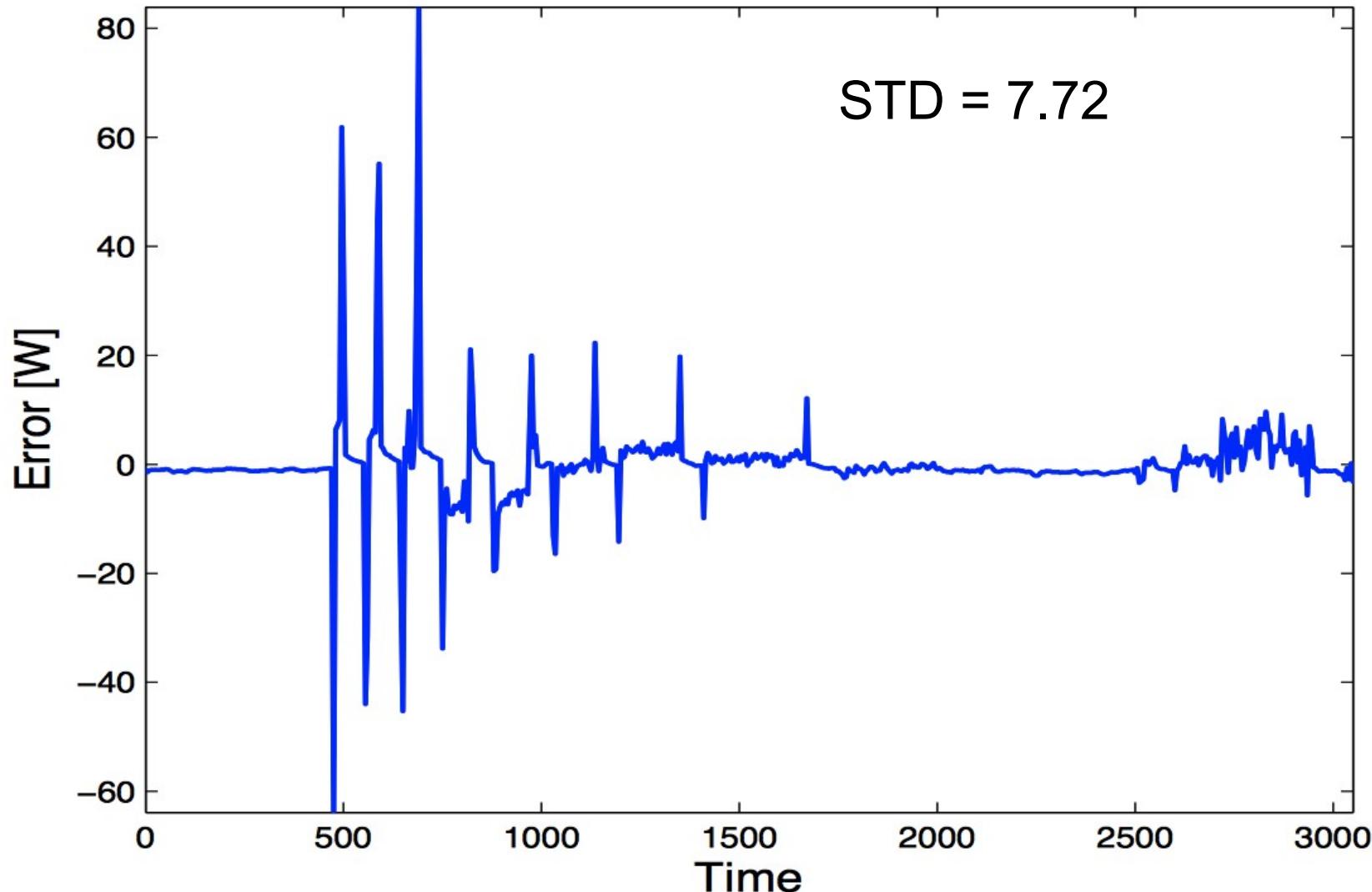
Estimation – P-States



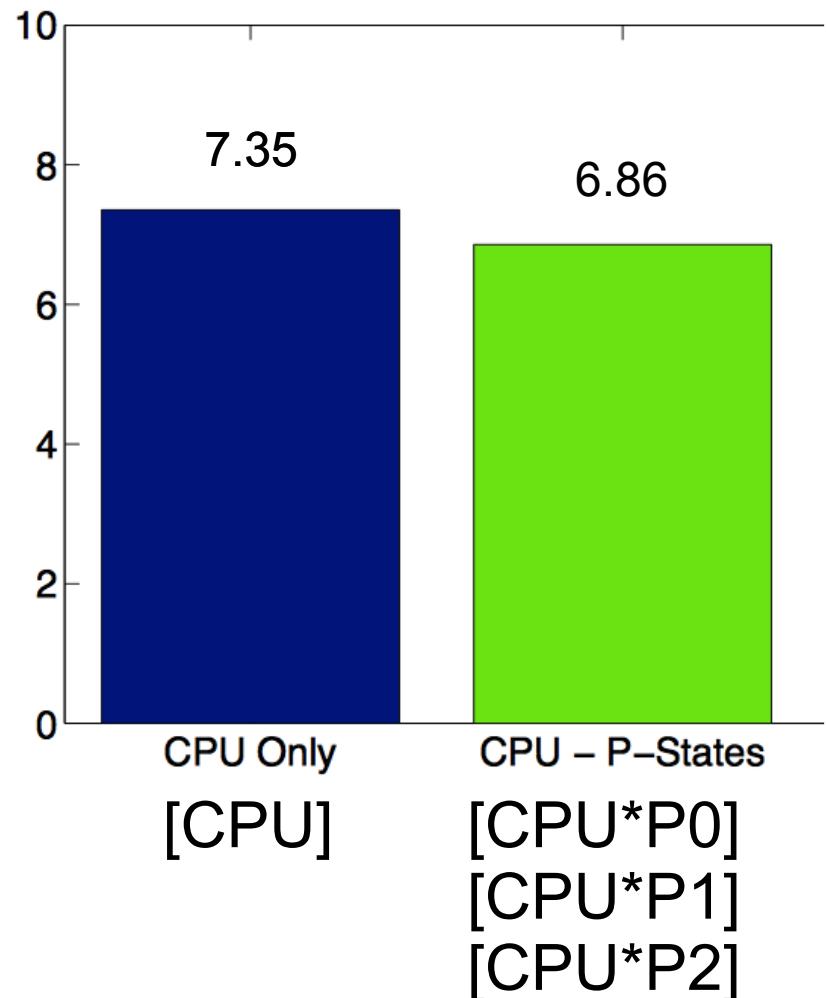
Model Validation



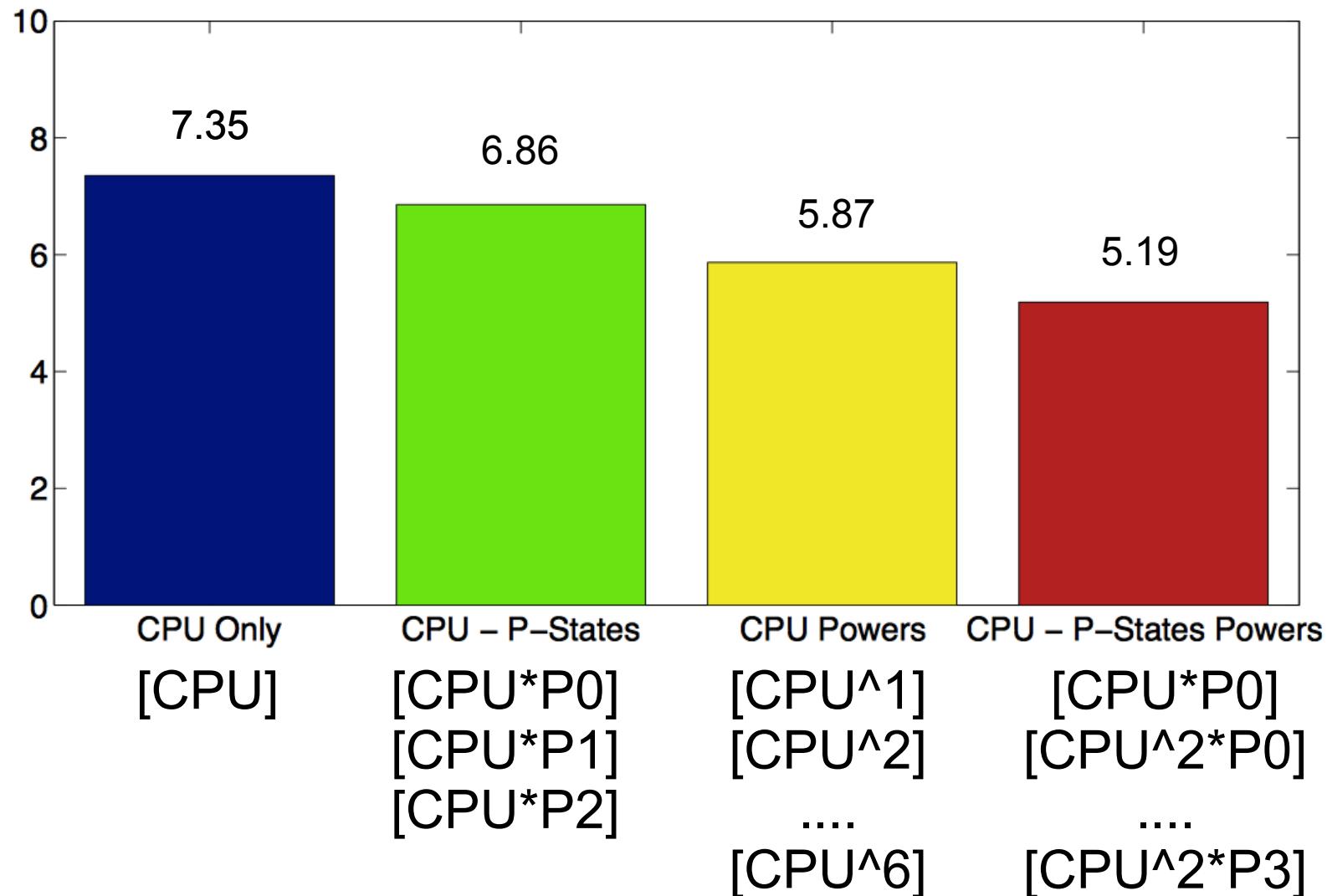
Error Analysis



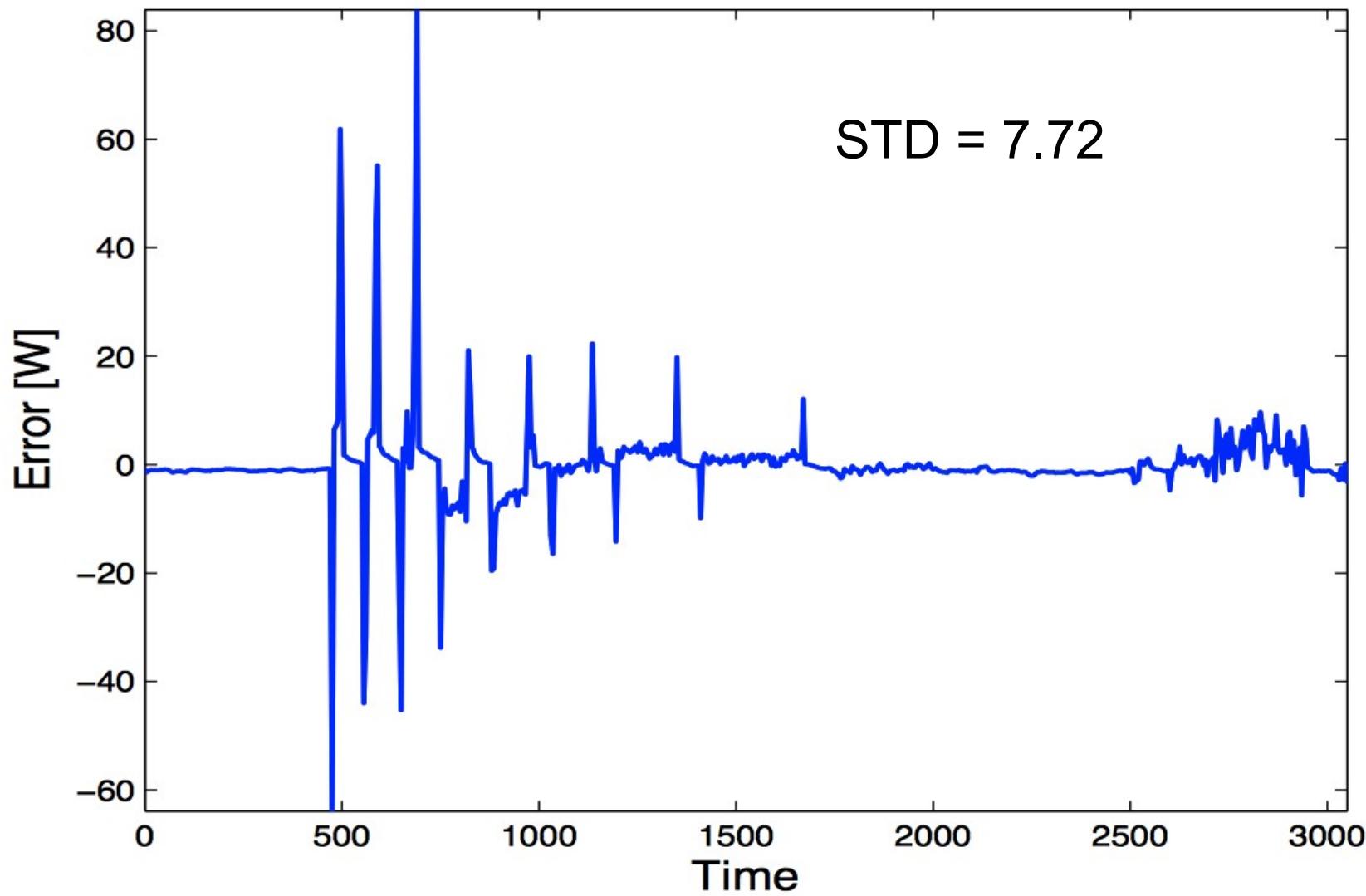
90th Percentile



90th Percentile

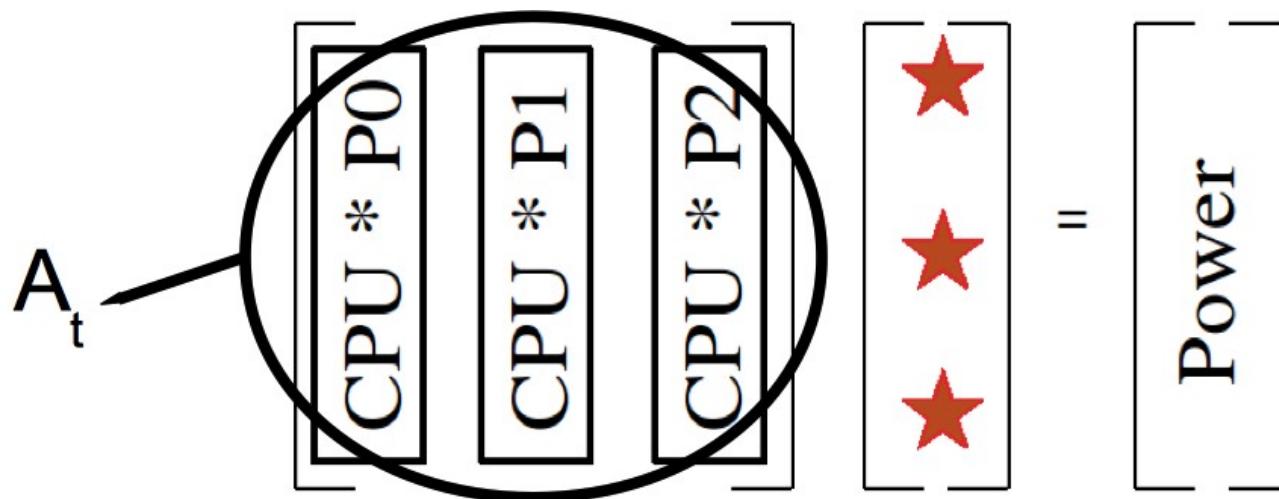


Error Analysis



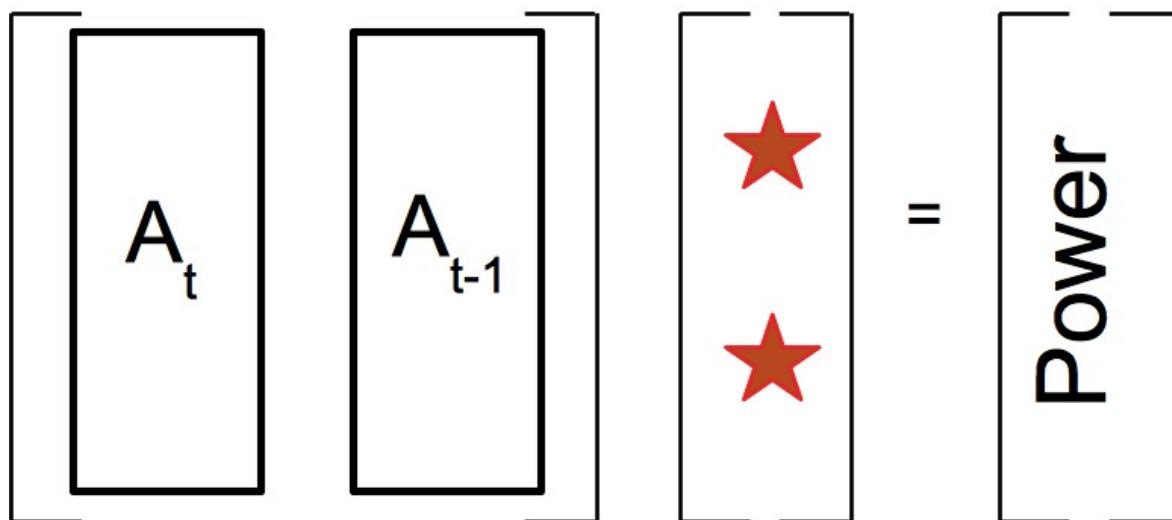
Dynamic Model

- Main error during CPU load transitions
 - Possibly due to power dynamics caused by capacitors in power circuitry etc.
- Solution: use history

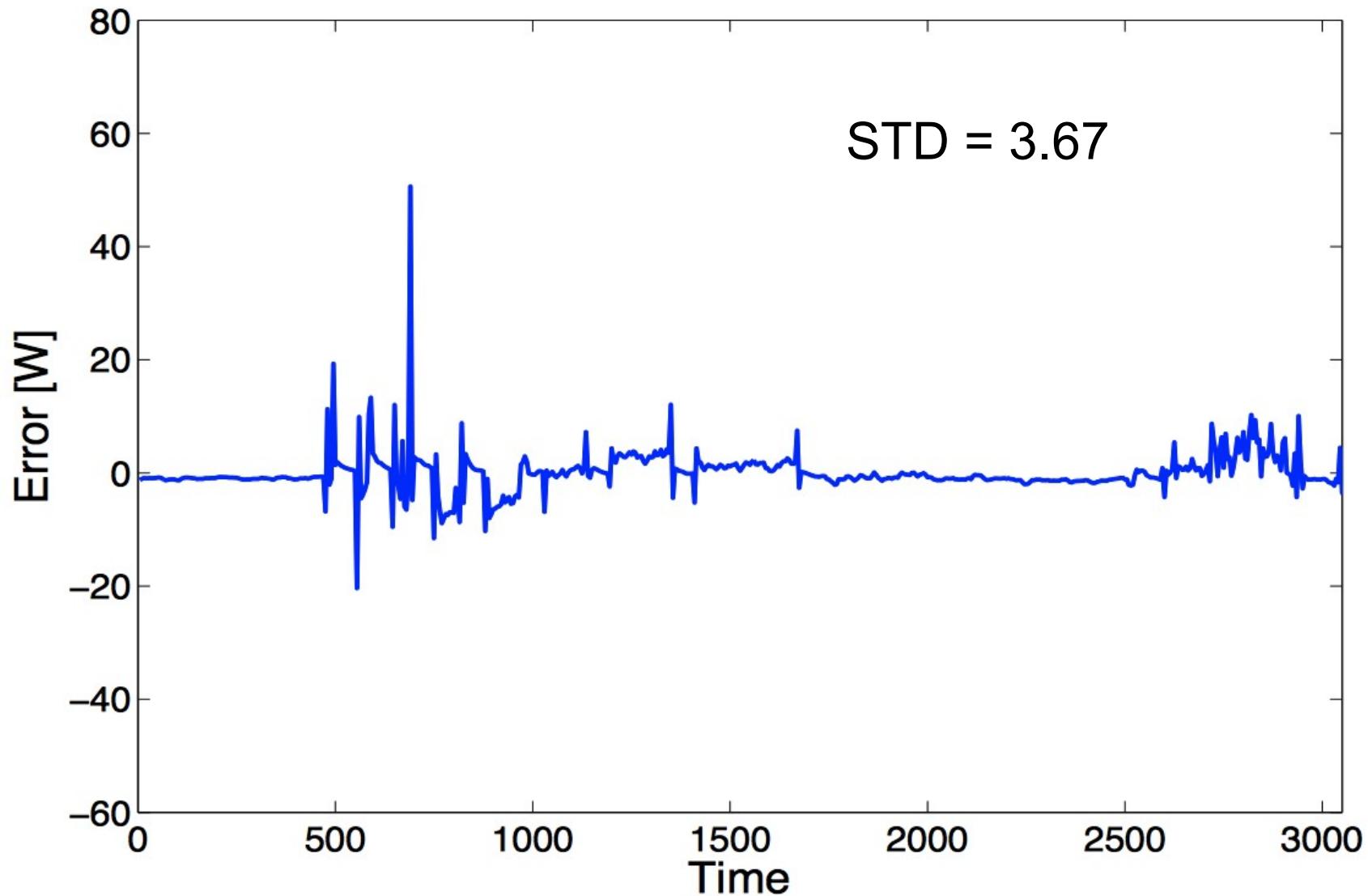


Dynamic Model

- Main error during CPU load transitions
 - Possibly due to power dynamics caused by capacitors in power circuitry etc.
- Solution: use history

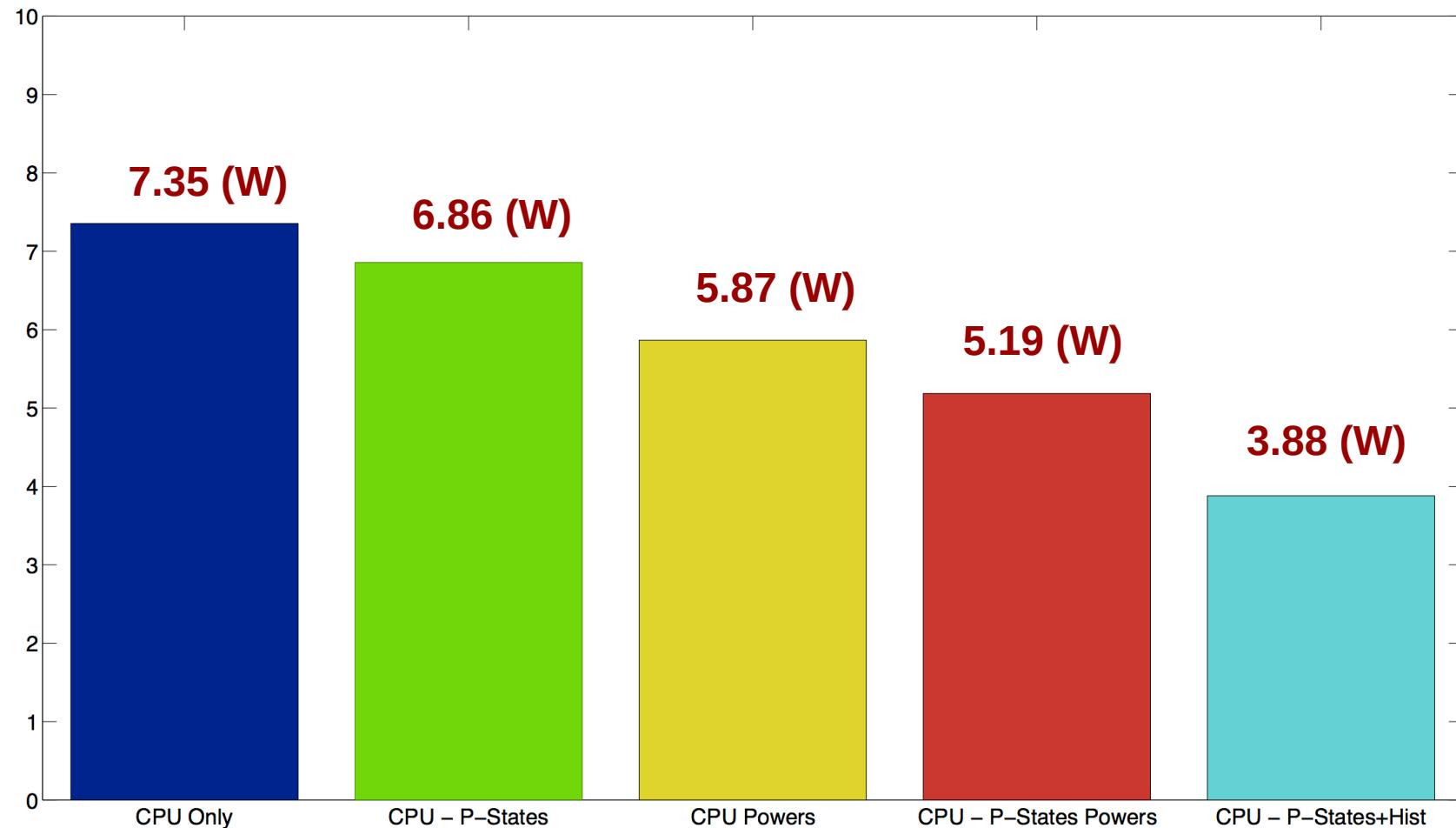


Error Analysis – Dynamic Model



Estimation of server power consumption

- Bars shows the 90th percentile of the estimation error

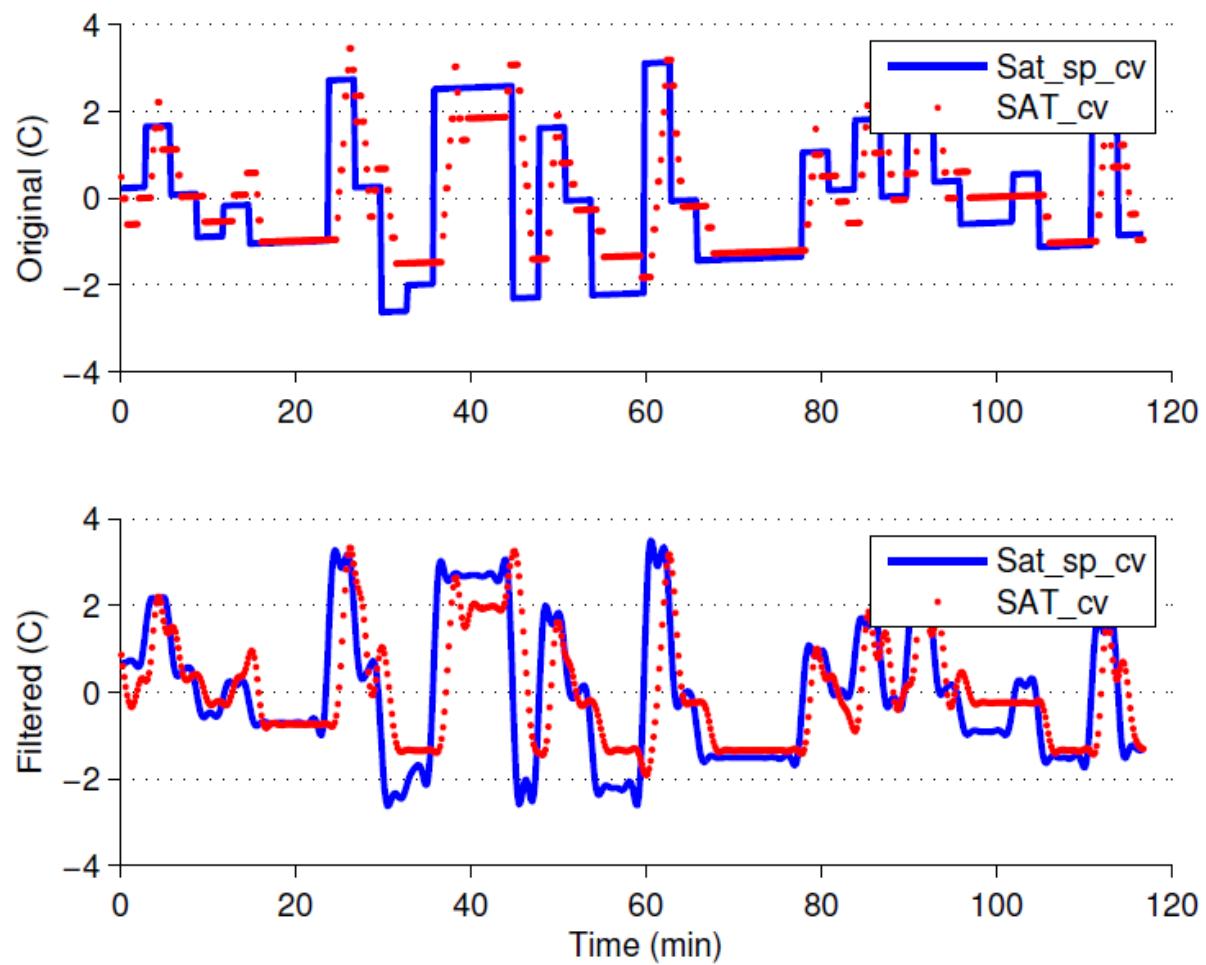


Conclusions

- For cn007, CPU Load + P-State are the informative data
 - Network and Disk are insignificant
 - Memory?
- Main error during CPU load transition
 - Possible Causes:
 - Synchronization between tools
 - Effects of power circuitry
 - Solution: use history (dynamic model)
- Non-linear models seem to fit data better

I/O models of CRAC units
Data collected by researchers at
HP Labs, Palo Alto, CA

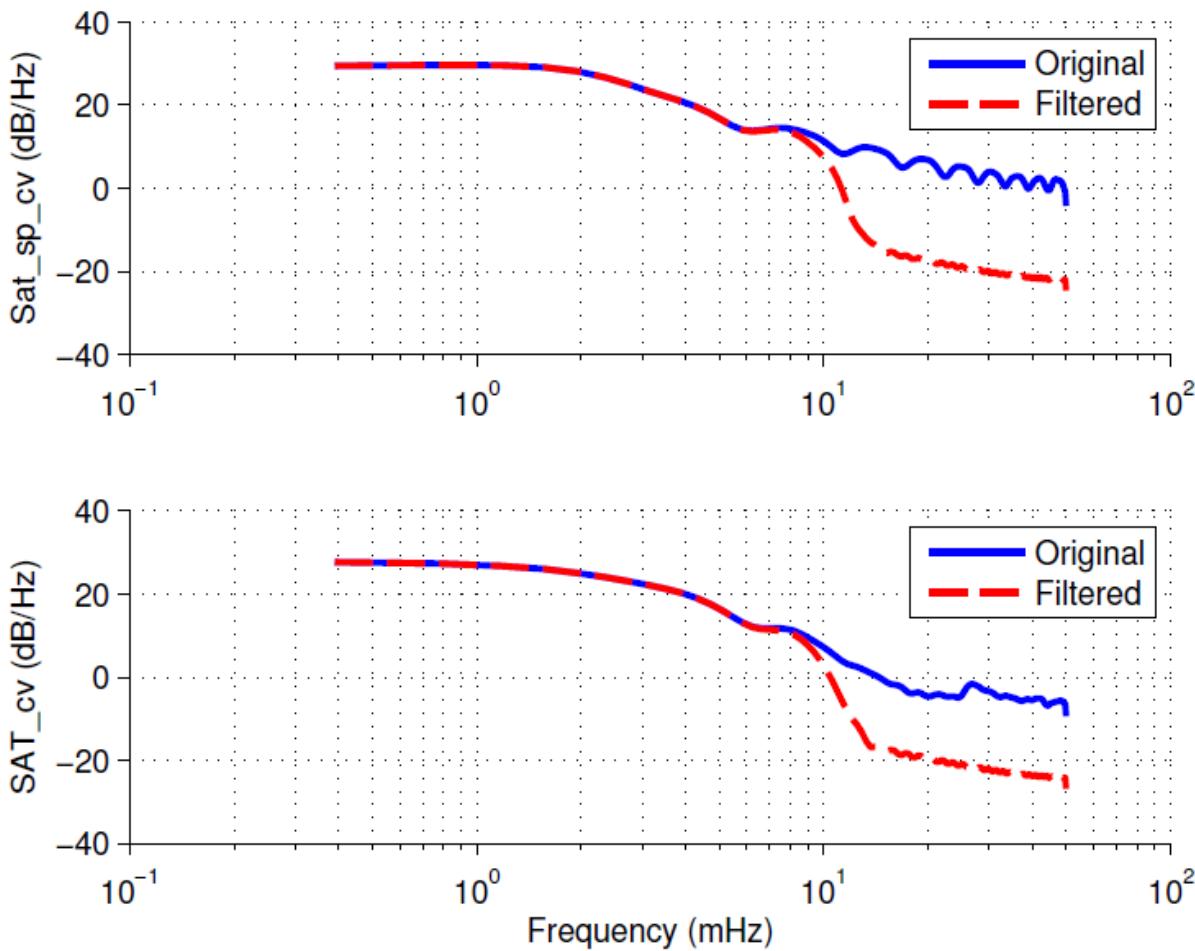
Input and output data over time



- **Original data**
 - Set point (SAT_cv)
 - Output temperature (SAT)

- **Filtered data**
 - Low-pass filter at 10(mHZ)

Power spectral density of I/O data



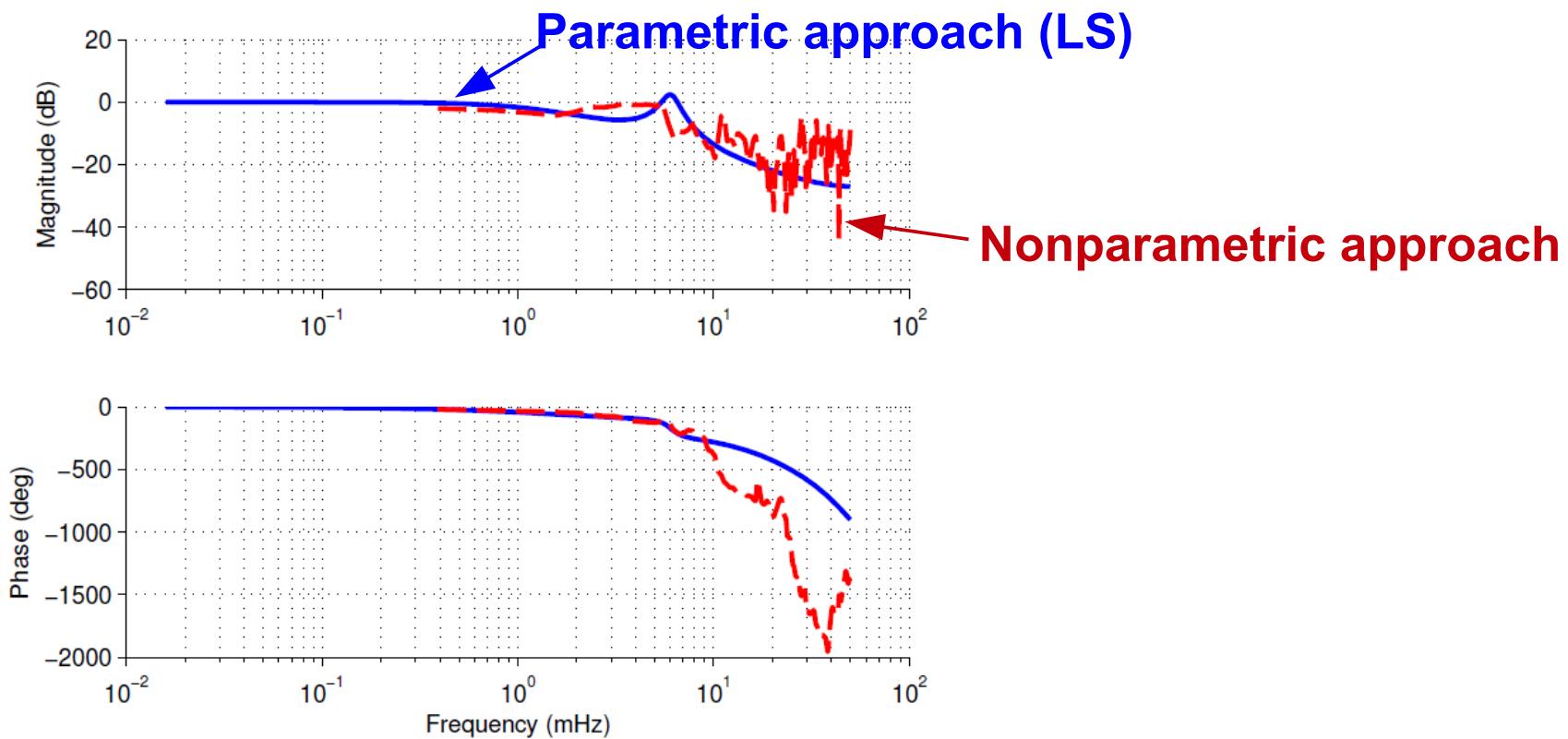
- Reference temperature
 - Set point

- Output temperature
 - Supply air temperature (SAT)

Estimated transfer function

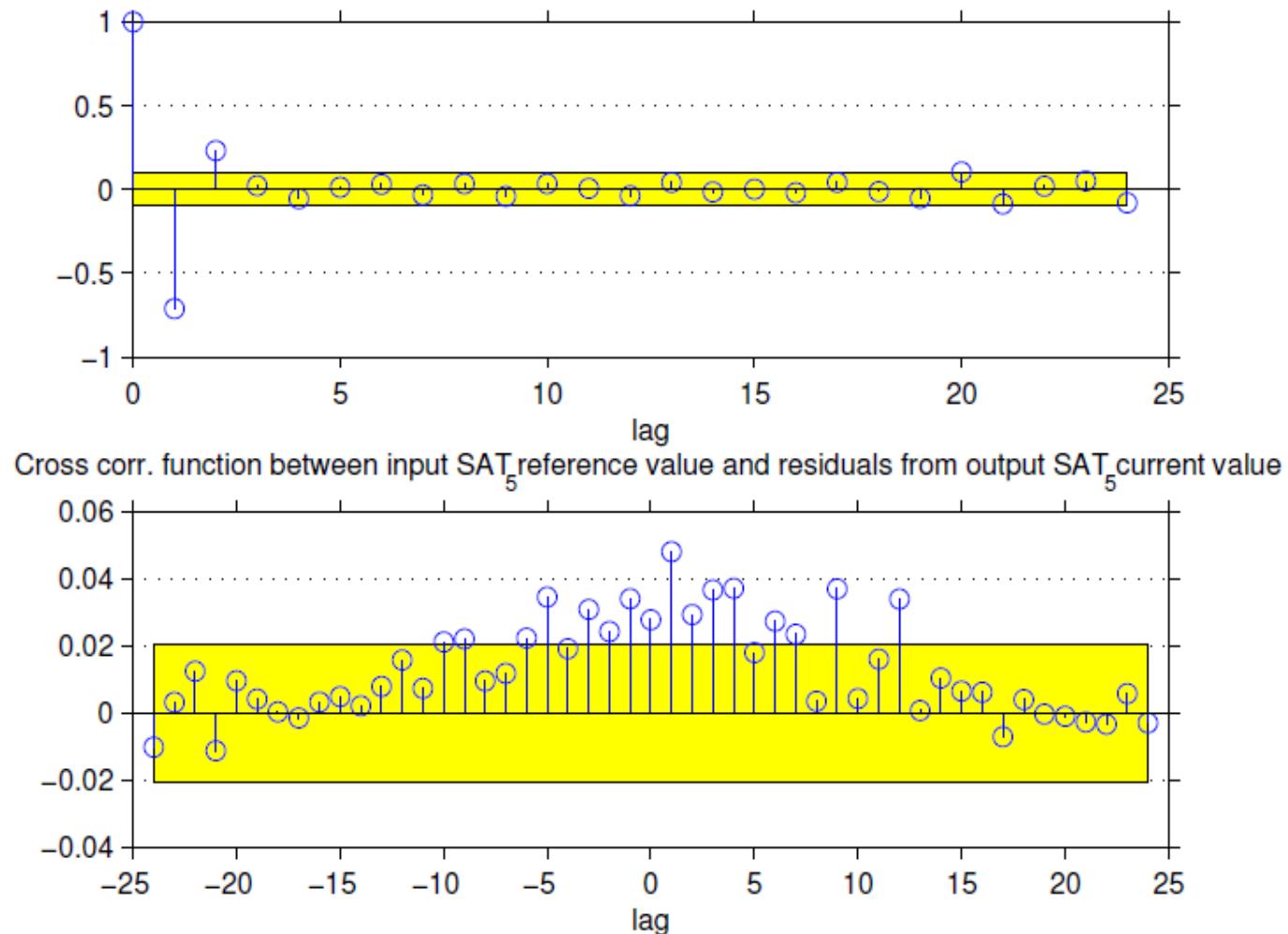
■ ARX model

$$ARX433_5(s) = e^{-40.40s} \frac{0.0088743(s^2 + 2 \cdot 0.6028 \cdot 0.0375s + 0.0375^2)}{(s + 0.008609)(s^2 + 2 \cdot 0.1014 \cdot 0.0381s^2 + 0.0381^2)}$$



Residual analysis

- Fit on the validation data-set: 84 %



Notes on the cyber-physical index (CPI)

Sensitivity index

■ Given a data center

- How much energy can be saved by a coordinated controller, with respect to an uncoordinated controller?

■ Zone input temperature at the equilibrium

$$\mathbf{T}_{\text{in},\mathcal{N}} = \mathcal{L}\boldsymbol{\eta} + \Psi_{[\mathcal{N},\mathcal{C}]} \mathbf{T}_{\text{ref}}$$

- Relative sensitivity index of the i^{th} zone

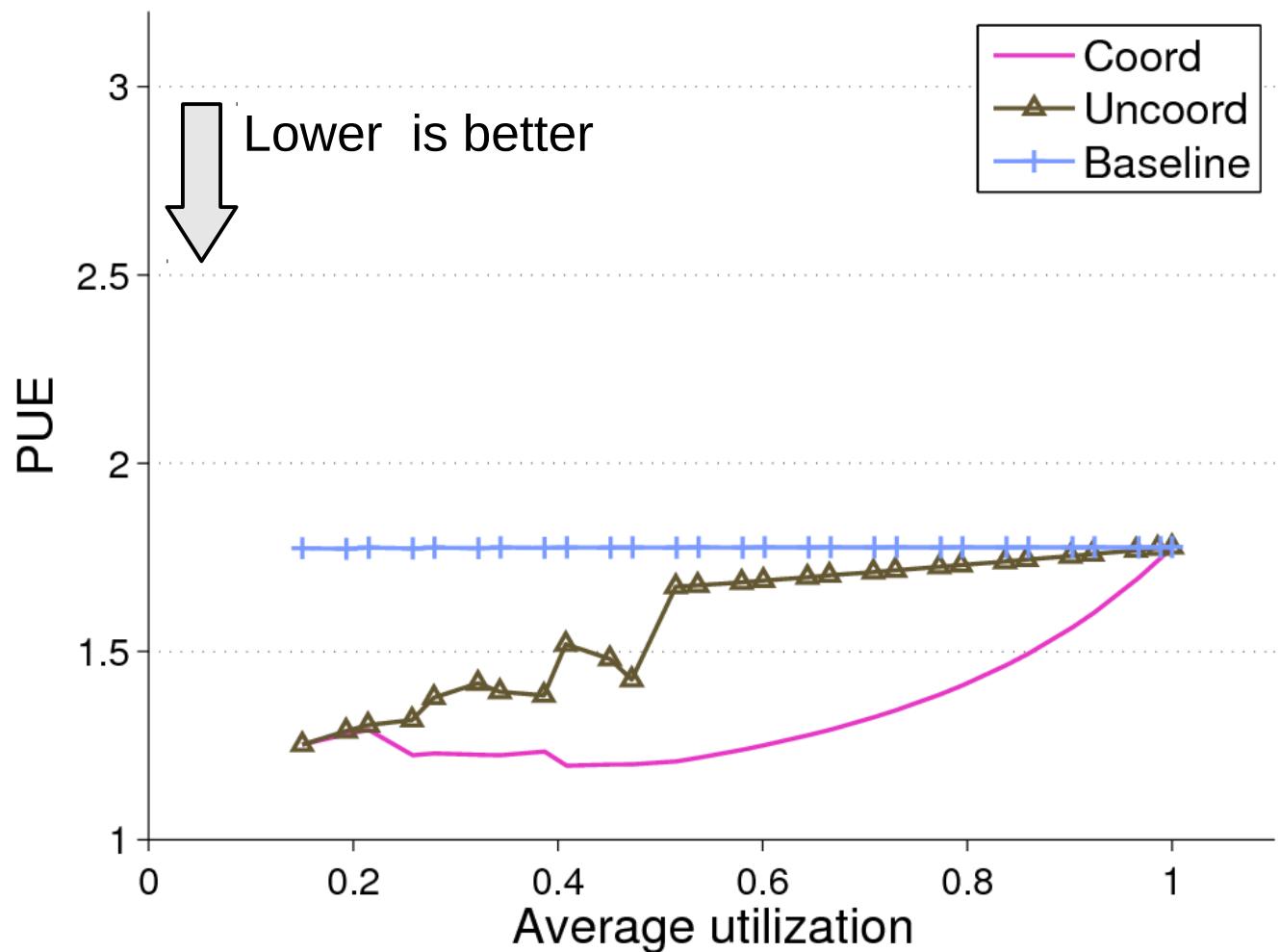
$$\mathcal{S}_i = \left\| \frac{\partial T_{\text{in},i}}{\partial \mathbf{T}_{\text{ref}}} \right\|_2 / \left\| \frac{\partial T_{\text{in},i}}{\partial \mathbf{z}} \right\|_2 \quad \mathbf{z} = [\mathbf{T}_{\text{ref}}^T \boldsymbol{\eta}^T]^T$$

■ Data center *cyber-physical index*

$$\textcolor{red}{CPI} = k \text{ std} ([S_1 \dots S_N]) \quad \text{where } k = \sqrt{\sum_{i=1}^N S_i^2}$$

Power usage effectiveness

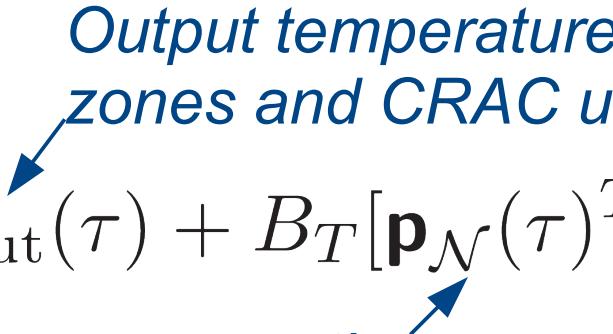
- $\text{PUE} = \frac{\text{Total data center power consumption}}{\text{Server power consumption}}$



Thermal network

■ Linear model

$$\dot{T}_{\text{out}}(\tau) = A_T T_{\text{out}}(\tau) + B_T [\mathbf{p}_N(\tau)^T \quad \mathbf{T}_{\text{ref}}(\tau)^T]^T$$

Output temperature of zones and CRAC units

Power consumption of zones *Reference temperature of CRAC units*

$$T_{\text{in}}(\tau) = \Psi T_{\text{out}}(\tau) \leq \overline{T}_{\text{in}}$$

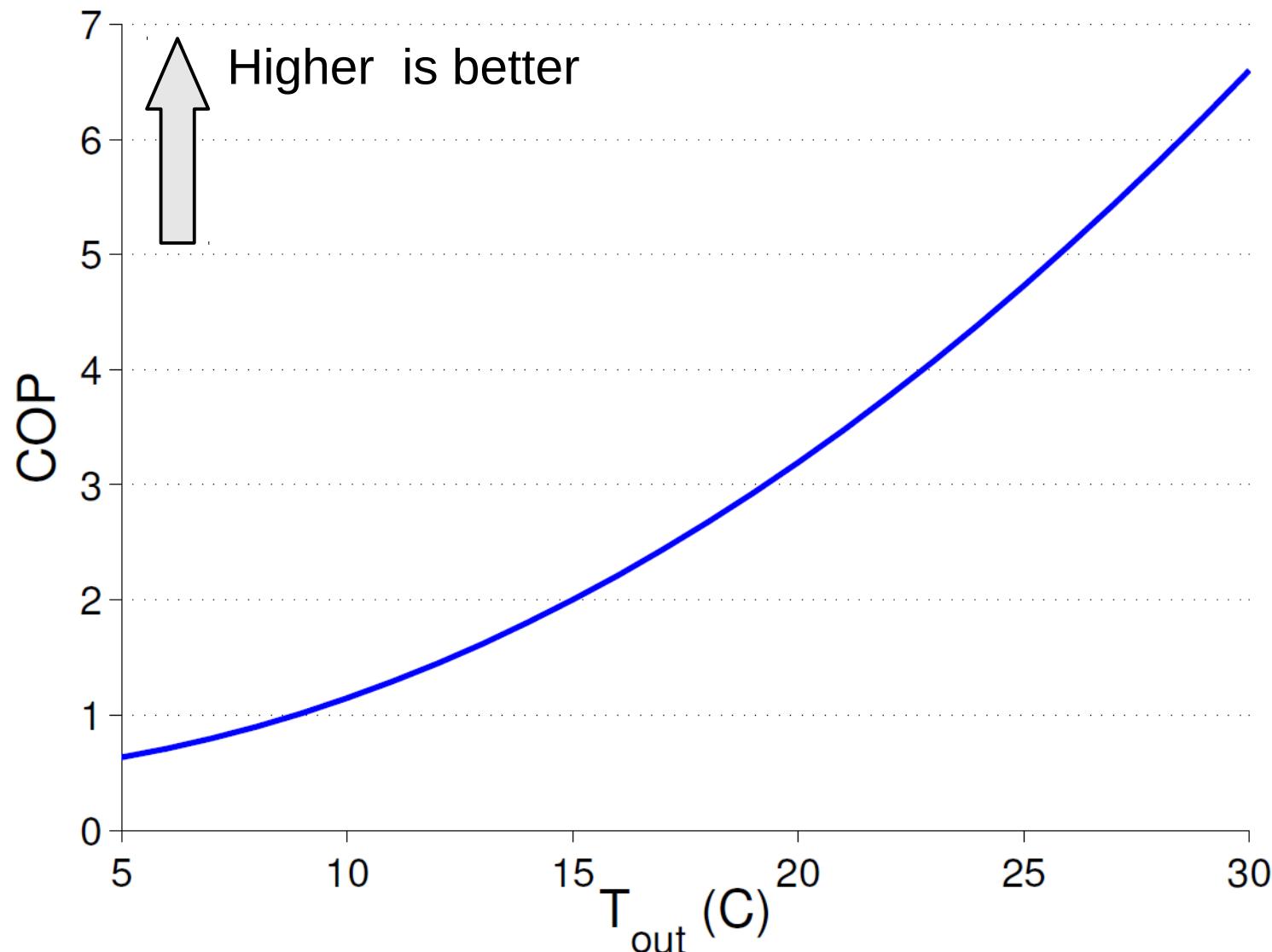

Inlet temperature constraint

■ CRAC power consumption depends on the coefficient of performance (COP)

$$p_i(t) = \frac{\dot{Q}_i(t)}{COP_i(T_{\text{out},i}(t))}$$


Heat removed rate (W)

Efficiency of CRAC nodes



J. Moore et al. "Making scheduling "cool": temperature-aware workload placement in data centers." 2005

Sensitivity index

- Given a data center
 - How much energy can be saved by a coordinated controller, with respect to an uncoordinated controller?
- Relative sensitivity of the i^{th} zone at the equilibrium
 - $\mathcal{S}_i = \left\| \frac{\partial T_{\text{in},i}}{\partial \mathbf{T}_{\text{ref}}} \Big|_{(\mathbf{l}, \mathbf{T}_{\text{out}})} \right\|_2 / \left\| \frac{\partial T_{\text{in},i}}{\partial \mathbf{z}} \Big|_{(\mathbf{l}, \mathbf{T}_{\text{out}})} \right\|_2 \quad \mathbf{z} = [\mathbf{T}_{\text{ref}}^T \boldsymbol{\eta}^T]^T$

Sensitivity index

■ Given a data center

- How much energy can be saved by a coordinated controller, with respect to an uncoordinated controller?

■ Relative sensitivity of the i^{th} zone at the equilibrium

$$\blacksquare \quad S_i = \left\| \frac{\partial T_{\text{in},i}}{\partial \mathbf{T}_{\text{ref}}} \Big|_{(\mathbf{l}, \mathbf{T}_{\text{out}})} \right\|_2 / \left\| \frac{\partial T_{\text{in},i}}{\partial \mathbf{z}} \Big|_{(\mathbf{l}, \mathbf{T}_{\text{out}})} \right\|_2 \quad z = [\mathbf{T}_{\text{ref}}^T \boldsymbol{\eta}^T]^T$$

$$\begin{bmatrix} \frac{\partial T_{\text{in},i}}{\partial \eta_1} \Big|_{(\mathbf{l}, \mathbf{T}_{\text{out}})} \\ \vdots \\ \frac{\partial T_{\text{in},i}}{\partial \eta_1} \Big|_{(\mathbf{l}, \mathbf{T}_{\text{out}})} \end{bmatrix} = \frac{\partial T_{\text{in},i}}{\partial \boldsymbol{\eta}} \quad \begin{array}{l} \nearrow \\ \searrow \end{array} \quad \frac{\partial T_{\text{in},i}}{\partial \mathbf{T}_{\text{ref}}} = \begin{bmatrix} \frac{\partial T_{\text{in},i}}{\partial T_{\text{ref},1}} \Big|_{(\mathbf{l}, \mathbf{T}_{\text{out}})} \\ \vdots \\ \frac{\partial T_{\text{in},i}}{\partial T_{\text{ref},C}} \Big|_{(\mathbf{l}, \mathbf{T}_{\text{out}})} \end{bmatrix}$$

Sensitivity index

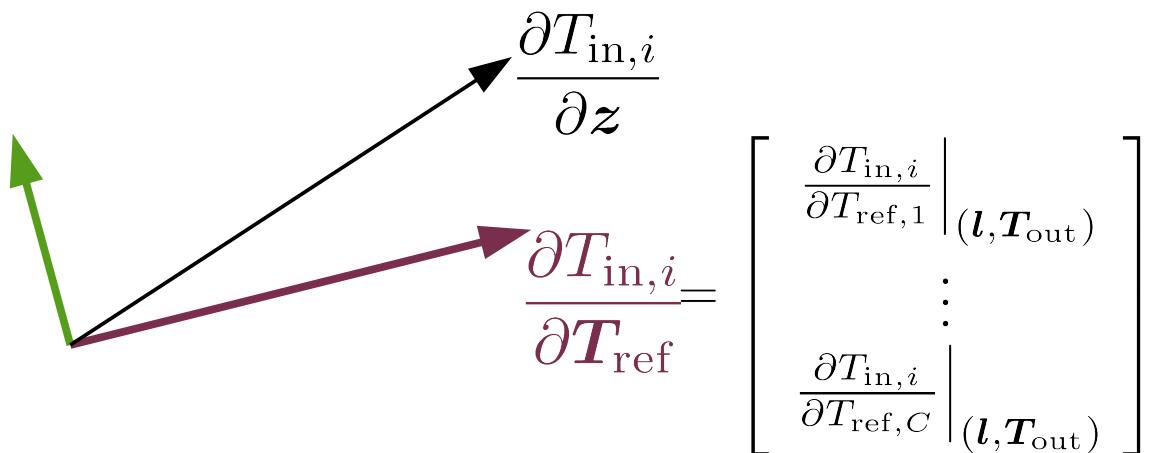
■ Given a data center

- How much energy can be saved by a coordinated controller, with respect to an uncoordinated controller?

■ Relative sensitivity of the i^{th} zone at the equilibrium

- $\mathcal{S}_i = \left\| \frac{\partial T_{\text{in},i}}{\partial \mathbf{T}_{\text{ref}}} \Big|_{(\mathbf{l}, \mathbf{T}_{\text{out}})} \right\|_2 / \left\| \frac{\partial T_{\text{in},i}}{\partial \mathbf{z}} \Big|_{(\mathbf{l}, \mathbf{T}_{\text{out}})} \right\|_2 \quad \mathbf{z} = [\mathbf{T}_{\text{ref}}^T \boldsymbol{\eta}^T]^T$

$$\begin{bmatrix} \frac{\partial T_{\text{in},i}}{\partial \eta_1} \Big|_{(\mathbf{l}, \mathbf{T}_{\text{out}})} \\ \vdots \\ \frac{\partial T_{\text{in},i}}{\partial \eta_1} \Big|_{(\mathbf{l}, \mathbf{T}_{\text{out}})} \end{bmatrix} = \frac{\partial T_{\text{in},i}}{\partial \boldsymbol{\eta}}$$



The diagram illustrates the decomposition of the relative sensitivity. A green arrow points from the vector of partial derivatives with respect to $\boldsymbol{\eta}$ to a red arrow representing the partial derivative with respect to \mathbf{T}_{ref} . This red arrow is shown as the hypotenuse of a right triangle, where the vertical leg is labeled $\frac{\partial T_{\text{in},i}}{\partial \mathbf{z}}$ and the horizontal leg is labeled $\frac{\partial T_{\text{in},i}}{\partial \mathbf{T}_{\text{ref}}}$. To the right, a matrix is shown where each row is a partial derivative with respect to a specific element of \mathbf{T}_{ref} , specifically $\frac{\partial T_{\text{in},i}}{\partial T_{\text{ref},1}}, \dots, \frac{\partial T_{\text{in},i}}{\partial T_{\text{ref},C}}$.

$$\frac{\partial T_{\text{in},i}}{\partial \mathbf{z}} = \sqrt{\left(\frac{\partial T_{\text{in},i}}{\partial \mathbf{T}_{\text{ref},1}} \Big|_{(\mathbf{l}, \mathbf{T}_{\text{out}})} \right)^2 + \dots + \left(\frac{\partial T_{\text{in},i}}{\partial \mathbf{T}_{\text{ref},C}} \Big|_{(\mathbf{l}, \mathbf{T}_{\text{out}})} \right)^2}$$

Sensitivity index

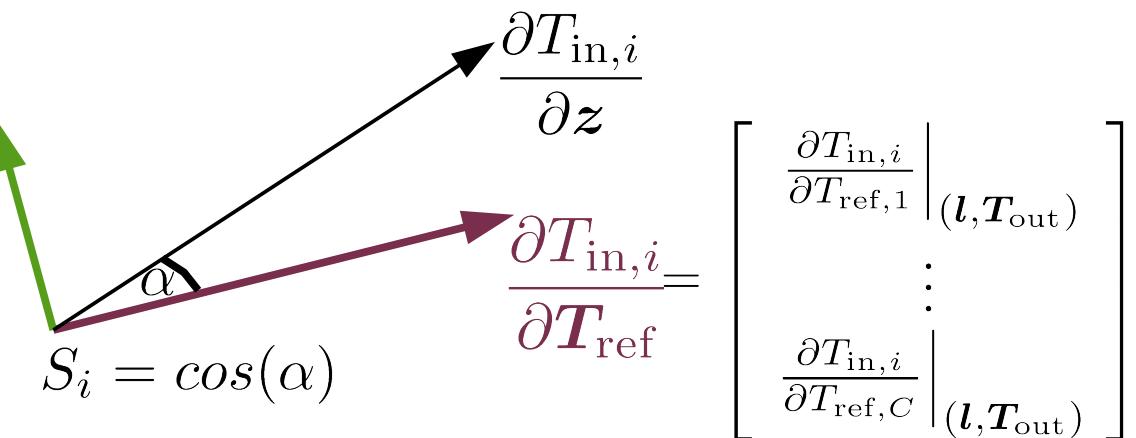
■ Given a data center

- How much energy can be saved by a coordinated controller, with respect to an uncoordinated controller?

■ Relative sensitivity of the i^{th} zone at the equilibrium

$$S_i = \left\| \frac{\partial T_{\text{in},i}}{\partial \mathbf{T}_{\text{ref}}} \Big|_{(\mathbf{l}, \mathbf{T}_{\text{out}})} \right\|_2 / \left\| \frac{\partial T_{\text{in},i}}{\partial \mathbf{z}} \Big|_{(\mathbf{l}, \mathbf{T}_{\text{out}})} \right\|_2 \quad z = [\mathbf{T}_{\text{ref}}^T \boldsymbol{\eta}^T]^T$$

$$\begin{bmatrix} \frac{\partial T_{\text{in},i}}{\partial \eta_1} \Big|_{(\mathbf{l}, \mathbf{T}_{\text{out}})} \\ \vdots \\ \frac{\partial T_{\text{in},i}}{\partial \eta_1} \Big|_{(\mathbf{l}, \mathbf{T}_{\text{out}})} \end{bmatrix} = \frac{\partial T_{\text{in},i}}{\partial \boldsymbol{\eta}}$$



Sensitivity index for the i^{th} zone

- Given a data center

- How much energy can be saved by a coordinated controller, with respect to an uncoordinated controller?

- Cyber-physical index (CPI)

- Relative sensitivity of the i^{th} zone at the equilibrium

$$S_i = \left\| \frac{\partial T_{\text{in},i}}{\partial \mathbf{T}_{\text{ref}}} \Big|_{(\mathbf{l}, \mathbf{T}_{\text{out}})} \right\|_2 / \left\| \frac{\partial T_{\text{in},i}}{\partial \mathbf{z}} \Big|_{(\mathbf{l}, \mathbf{T}_{\text{out}})} \right\|_2 \quad \mathbf{z} = [\mathbf{T}_{\text{ref}}^T \boldsymbol{\eta}^T]^T$$

- Easy to compute when using the proposed model

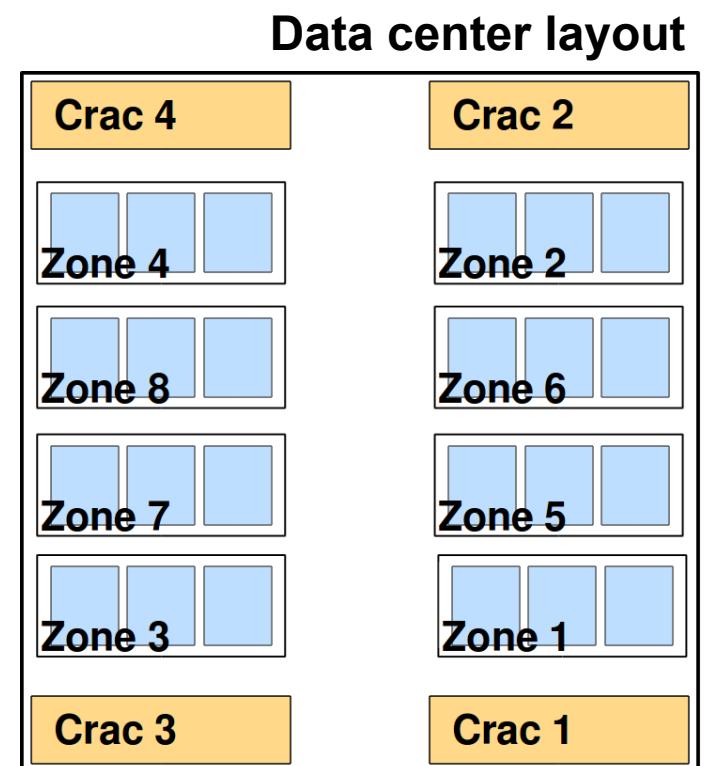
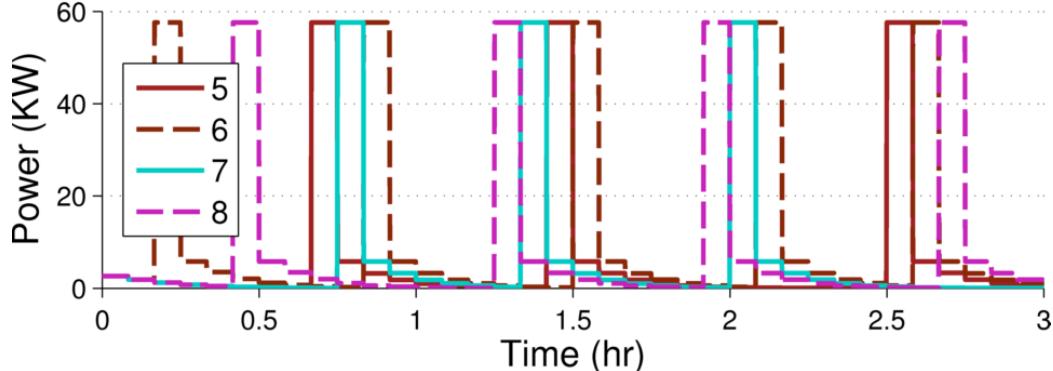
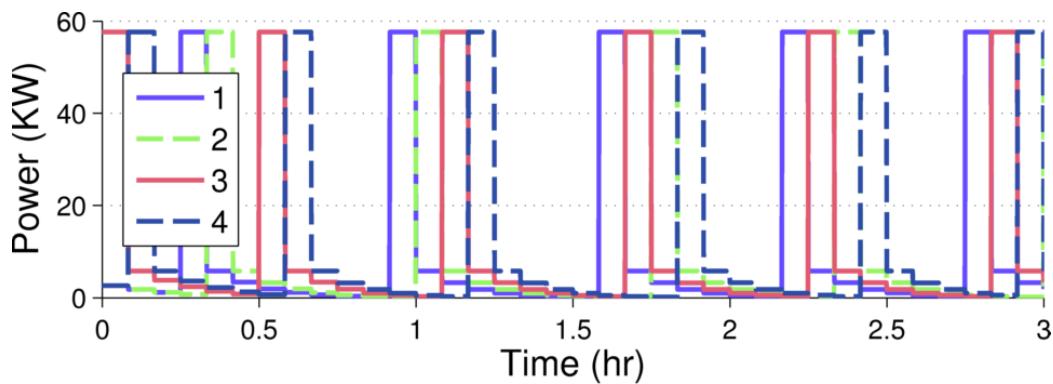
- Data center *cyber-physical index*

$$\text{CPI} = k \text{ std} ([S_1 \quad \dots \quad S_N])$$

Stability of the coordinated controller

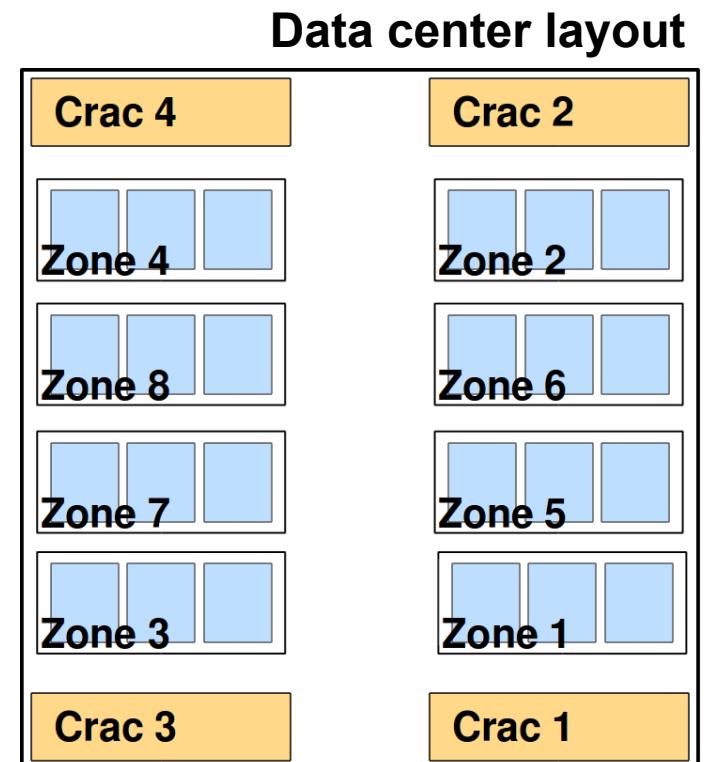
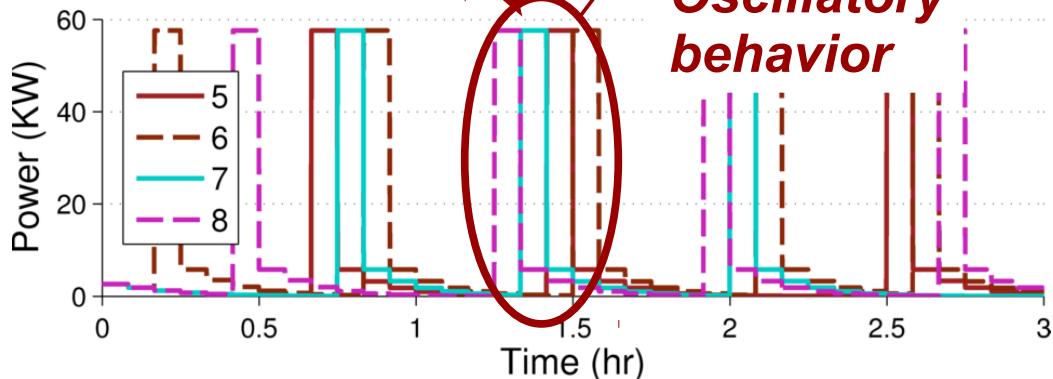
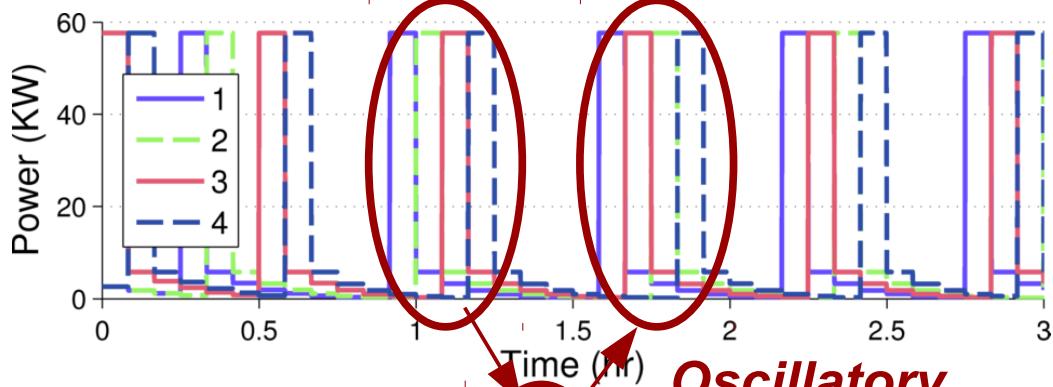
Economic MPC and data center

- Coordinated control approach
 - Perfect predictive model
 - Constant workload arrival rate
- Zone power consumption



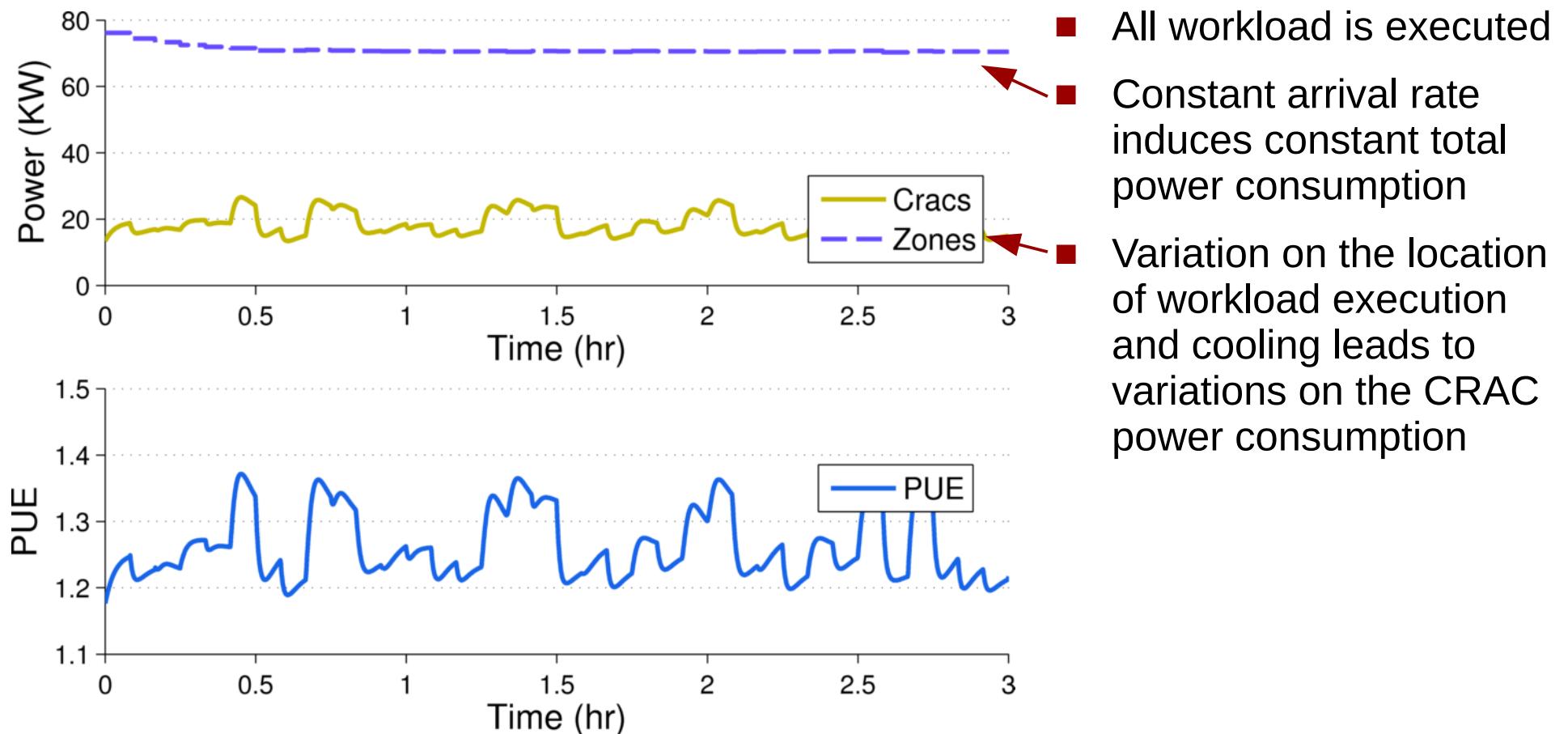
Economic MPC and data center

- Coordinated control approach
 - Perfect predictive model
 - Constant workload arrival rate
- Zone power consumption



Economic MPC and data center

■ Total server and CRAC power consumption



Effects of model mismatch

Control-oriented model

- Considers each rack and each CRAC as a dynamic subsystem

$$\begin{aligned} \mathbf{x}_i(k+1) &= A_i \mathbf{x}_i(k) + B_i(k) \mathbf{u}_i(k) + \sum_{\substack{j=1 \\ j \neq i}}^{R+C} B_{i,j} \mathbf{z}_j(k) \\ \mathbf{z}_i(k) &= G_i \mathbf{x}_i(k) \end{aligned}$$

	$\mathbf{x}_i(k)$	$\mathbf{u}_i(k)$	$\mathbf{z}_i(k)$
Racks	Outlet temperature of every server in the rack (<i>vector</i>)	Relative amount of resources assigned to different jobs (<i>vector</i>)	Average rack outlet temperature
CRACs	Supplied air temperature (<i>scalar</i>)	Desired supplied air temperature (<i>scalar</i>)	Supplied air temperature

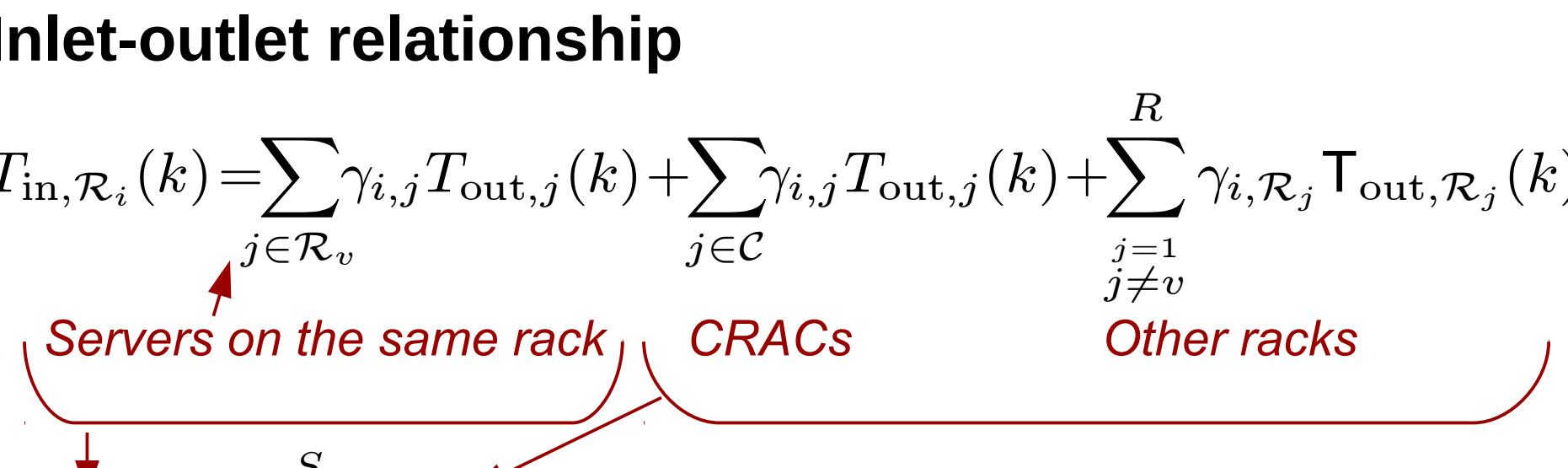
Thermal constraints

- Bound the inlet temperature of each rack

$$T_{\text{in},\mathcal{R}_i}(k) \leq \bar{T}_{\text{in},\mathcal{R}_i}, \quad i = 1, \dots, R \leftarrow \text{Number of racks}$$

- Inlet-outlet relationship

$$T_{\text{in},\mathcal{R}_i}(k) = \sum_{j \in \mathcal{R}_v} \gamma_{i,j} T_{\text{out},j}(k) + \sum_{j \in \mathcal{C}} \gamma_{i,j} T_{\text{out},j}(k) + \sum_{\substack{j=1 \\ j \neq v}}^R \gamma_{i,\mathcal{R}_j} T_{\text{out},\mathcal{R}_j}(k)$$


 $F_i x_i(k) + \sum_{\substack{j=1 \\ j \neq i}}^S F_{i,z_j} z_j(k) \leq Z_i \quad S: \text{number of subsystems } (S=R+C)$

One step ahead coordinated control

- The coordinated control problem can be written as

$$\min_{\mathcal{U}_1} \sum_{i=1}^S \mathbf{c}_{i,u}^T(k) \hat{\mathbf{u}}_i(k|k)$$

s.t.

$$\begin{aligned} \underline{\mathbf{u}}_i &\leq \hat{\mathbf{u}}_i(k|k) \leq \bar{\mathbf{u}}_i, \\ H_i \hat{\mathbf{u}}_i(k|k) &\leq 1 \end{aligned}$$

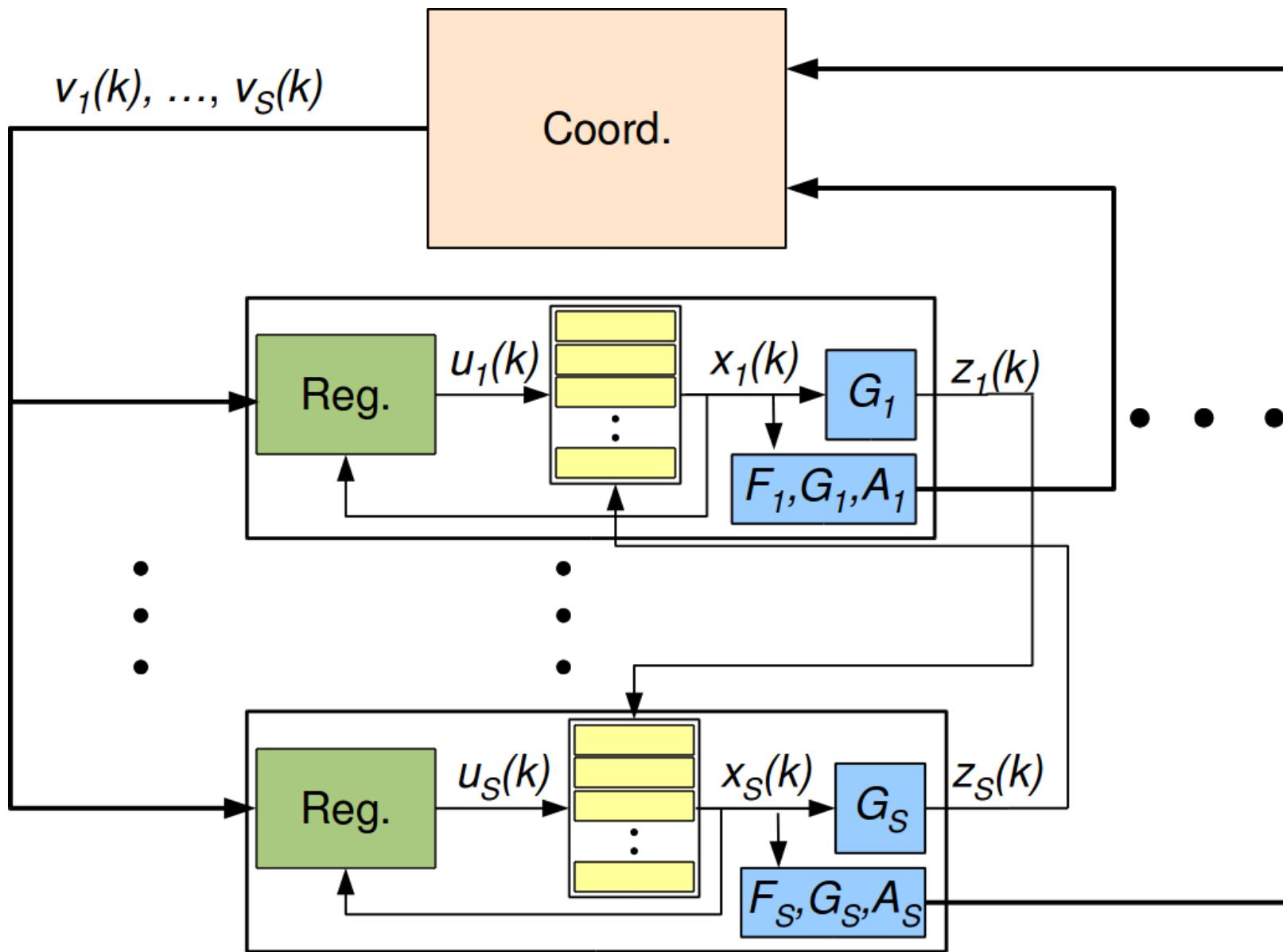
$$F_i B_i(k) \hat{\mathbf{u}}_i(k|k) + \underbrace{\sum_{\substack{j=1 \\ j \neq i}}^S F_{i,z_j} \mathbf{G}_j B_j(k) \hat{\mathbf{u}}_j(k|k)}_{\text{Coupling term}} + \mathbf{k}_i(k) \leq \bar{\mathbf{z}}_i,$$

$$\mathcal{U}_1 = \left\{ \hat{\mathbf{u}}_1(k|k), \dots, \hat{\mathbf{u}}_S(k|k) \right\},$$

Problem analysis

- **Data centers are large-scale systems**
 - It may be impossible to collect data from all of the sensors and compute the control actions with a single controller
- **Every subsystem affects other subsystems only through its output**
 - *For racks, the output has a much lower dimension than the state*
- *Is it possible to exploit the particular coupling, in order to derive an effective hierarchical control strategy?*

Proposed hierarchy strategy

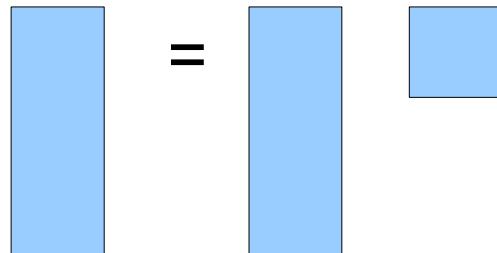


Coordinator optimization problem

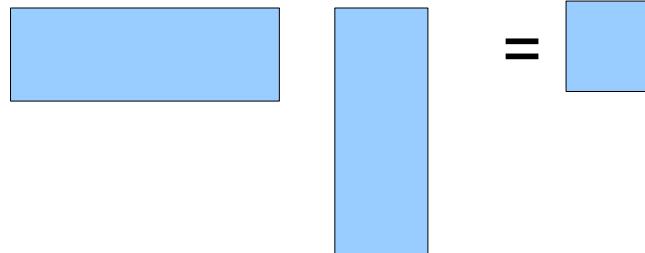
- Stated in terms of the virtual control inputs $\{v_i(k)\}$
- A collection of matrices $\{M_i(k)\}$ approximate the actions of local controllers

- The coordinator considers the following relationship

$$u_i(k) = M_i(k)v_i(k)$$



$$G_i B_i(k) M_i(k) = I$$



- Multiple collections of $\{M_i(k)\}$ matrices can be defined, but not all of them yield an overall optimal control strategy

Stage one – coordinator problem

- Global optimization problem in the $\{v_i(k)\}$ variables

$$\min_{\mathcal{V}_1} \sum_{i=1}^S \mathbf{c}_{i,u}^T(k) M_i(k) \hat{\mathbf{v}}_i(k|k)$$

s.t.

$$\underline{\mathbf{u}}_i \leq M_i(k) \hat{\mathbf{v}}_i(k|k) \leq \overline{\mathbf{u}}_i,$$

$$H_i M_i(k) \hat{\mathbf{v}}_i(k|k) \leq 1$$

$$F_i B_i(k) M_i(k) \hat{\mathbf{v}}_i(k|k) + F_{i,z} \hat{\mathbf{v}}(k|k) + \mathbf{k}_i(k) \leq \overline{\mathbf{z}}_i.$$

$$\mathcal{V}_1 = \{\hat{\mathbf{v}}_1(k|k), \dots, \hat{\mathbf{v}}_S(k|k)\}$$

Stage two – regulator problem

- Every regulator solves a local optimization problem

$$\min_{\hat{\mathbf{u}}_i(k|k)} \mathbf{c}_{i,u}^T(k) \hat{\mathbf{u}}_i(k|k)$$

s.t.

$$\underline{\mathbf{u}}_i \leq \hat{\mathbf{u}}_i(k|k) \leq \overline{\mathbf{u}}_i$$

$$H_i \hat{\mathbf{u}}_i(k|k) \leq 1$$

$$F_i B_i(k) \hat{\mathbf{u}}_i(k|k) + \sum_{\substack{l=1 \\ l \neq i}}^S F_{i,z_l} \hat{\mathbf{v}}_l(k|k) + \mathbf{k}_i(k) \leq \overline{\mathbf{z}}_i$$

$$\hat{\mathbf{v}}_i^*(k|k) = G_i B_i(k) \hat{\mathbf{u}}_i(k|k)$$



Ensures coherence among optimization problems

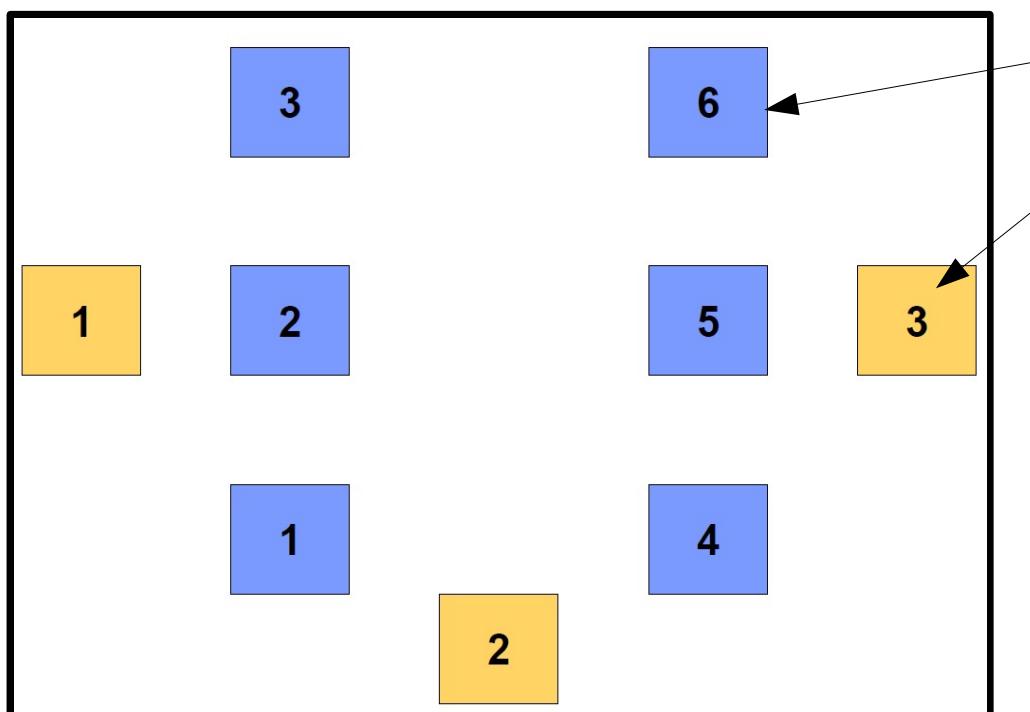
Properties of the control strategy

- If the coordinator problem has an admissible solution, then all of the local regulator problems have an admissible solution
- The condition $\hat{v}_i^*(k|k) = G_i B_i(k) \hat{u}_i(k|k)$
 - can be weakened to
$$F_{j,z_i} \hat{v}_i^*(k|k) \geq F_{j,z_i} M_i(k) \hat{u}_i(k|k) \quad i, j = 1, \dots, S, \quad i \neq j$$
- There always exists a collection of matrices $\{M_i(k)\}$ such that the performance obtained by the hierarchy control strategy equals the performance of the centralized coordinated controller

Choice of the $\{M_i(k)\}$ matrices

- Matrices $\{M_i(k)\}$ are used by the coordinator to approximate the actions of local controllers
 - Choice based on the a priori information about local subsystems
- A partial characterization of the $\{M_i(k)\}$ matrices is given in the paper
- How relevant is the effect of the $\{M_i(k)\}$ matrices on the overall system performance?

Simulation

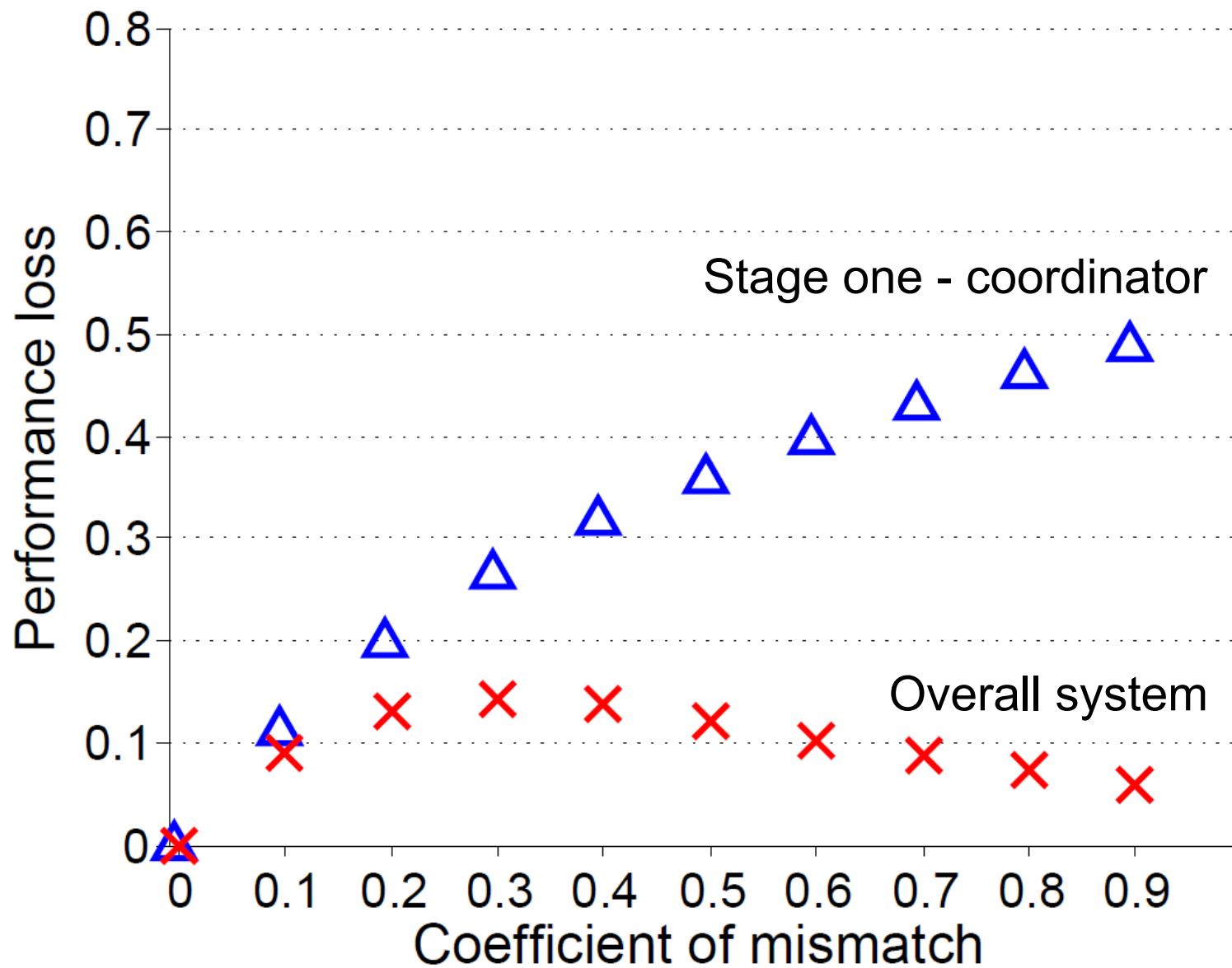


- **6 Racks**
- **3 CRACs**
- **6 job classes**
- **42 Server per rack**
 - Total of 1515 control variables

Simulation study

- **Objective:** analyze the effects of a bad choice of the $\{M_i(k)\}$ matrices
- **Analysis based on**
 - the performance loss estimated at the first stage of the hierarchical strategy
 - the performance loss obtained at the second stage of the hierarchical strategy (overall system performance)
- **Coefficient of mismatch**
 - When the coefficient is 0, the $\{M_i(k)\}$ matrices are optimally chosen
 - Higher values of the coefficient imply worse choices of the $\{M_i(k)\}$ matrices
- **Analysis obtained via a Monte Carlo approach**
 - Each point is obtained averaging over 500 simulations

Performance loss



Conclusion and future work

- The particular coupling among the subsystems and the large dimension of the control inputs makes hierarchical control an interesting approach
 - In the particular case we discussed the hierachal approach is effective
- We are currently working in order to extend this approach for a multi-step MPC and for nonlinear cost functions
- Techniques to take advantage of the particular form of the input constraints are also being developed