

# The Data Center as a Grid Load Stabilizer

Hao Chen<sup>\*</sup>, Michael C. Caramanis<sup>\*\*</sup> and Ayse K. Coskun<sup>\*</sup>

<sup>\*</sup>Department of Electrical and Computer Engineering

<sup>\*\*</sup>Division of Systems Engineering

Boston University

{haoc, mcaraman, acoskun}@bu.edu



# Power Grid & Market

- Power supply = demand ? ( => blackouts )
- Renewable energy sources: intermittent



- Lack of reliable, large-scale, economical energy storage solutions
- Independent System Operator (ISO):
  - New power market features:
    - Demand side regulation service (RS)
  - Credits provided to the participant who modulates its power consumption dynamically so as to track the ***RS signal***

# Demand Side – Data Centers

- Electricity: **3%** of the overall consumption in the US<sup>[1]</sup>
- Power capping /management techniques
  - Enable **flexibility** in power consumption
- Workload flexibility

Data centers offer a unique opportunity for providing power regulation service (RS) reserves.

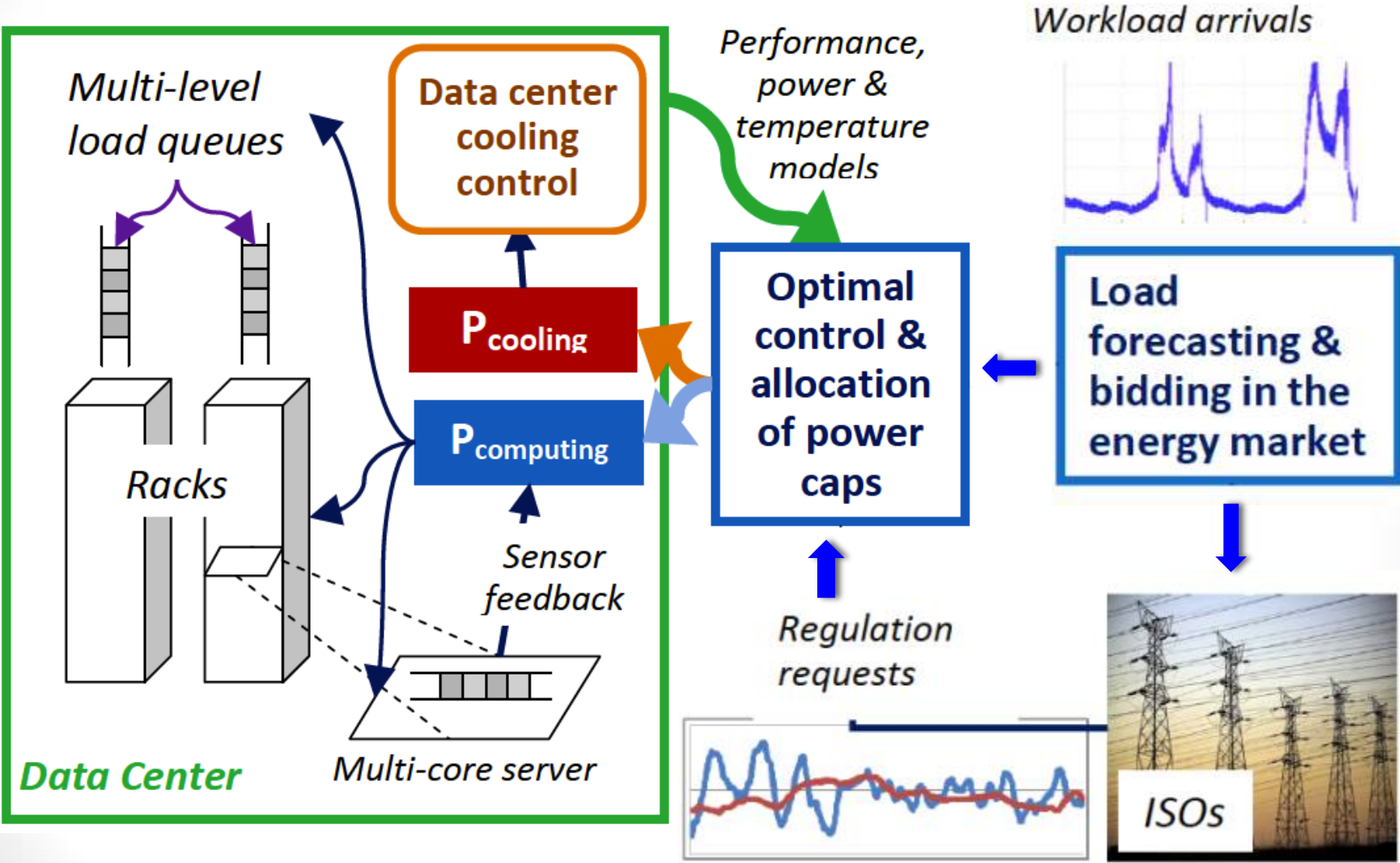


## Benefits of Participation

- Help solve unstable renewable energy problem
- Provide additional reserves to accommodate other less flexible uses of electricity
- Achieve significant monetary savings

[1]: J. Koomey. Growth in Data Center Electricity Use 2005 to 2010. Oakland, CA: Analytics Press. August, 1, 2010.

# Data Centers in Advanced Power Market



# Contributions

- A dynamic control policy for solving **server commitment** problem, leveraging:
  - Server-level power capping techniques
  - Information on server power states and overheads
  - Job scheduling & allocation decisions
- RS provision bidding value estimation
- **Data center level** (compared to previous work on a single server)
- Our solution is able to accurately track the ISO signal, and:
  - We achieve **50%+** monetary savings
  - The proposed policy does not cause major QoS degradation
  - Policy is agnostic of the specific type of workloads running
  - Significant improvement in both monetary savings and QoS compared to prior results based on a single server (*Chen et al. ICCAD 2013*)

# Outline

- **Background**
  - Data Center Power Management
  - Power Market and Data Center Participation
  - Regulation Service (RS)
- Data Center Model
- Dynamic Power Control Policy
- Regulation Reserves Bidding
- Results

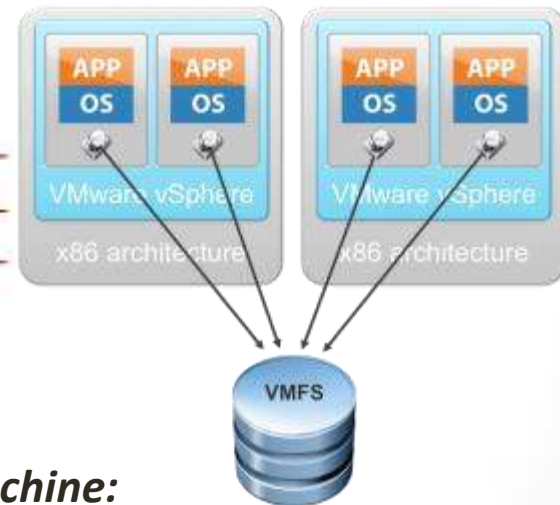
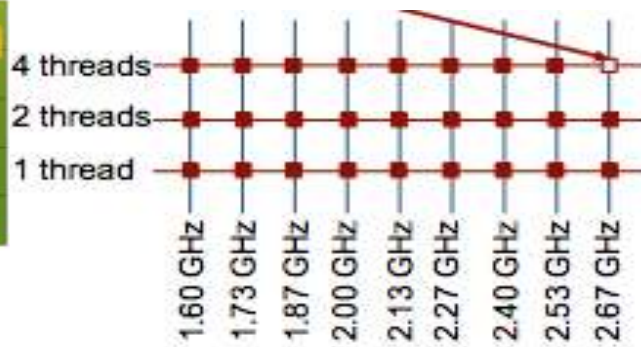
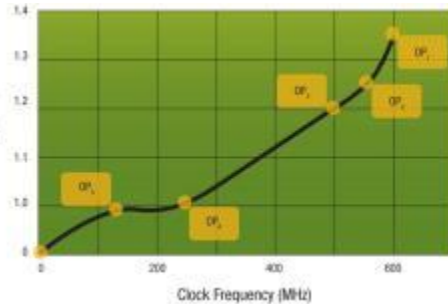


# Data Center Power Management



## Data Center Server Farms:

- Power and resource budgeting [Zhan DAC13][Gandhi SIGMETRICS09];
- Server Commitment: sleep and idle [Meisner Sigplan Not09][Ischi ISCA13][Gandhi IGCC12].



## Single Server Level:

- DVFS [Li HPCA06]
- Power Capping: DVFS + multi-thread allocation/migration [Cochran et al. Micro11][Rangan et al. ISCA09][Reda et al. Micro12]

## Virtual Machine:

- Power allocation [Nathuji et al. HPDC08]
- Resource consolidation policy [Hwang et al. ISLPED12]

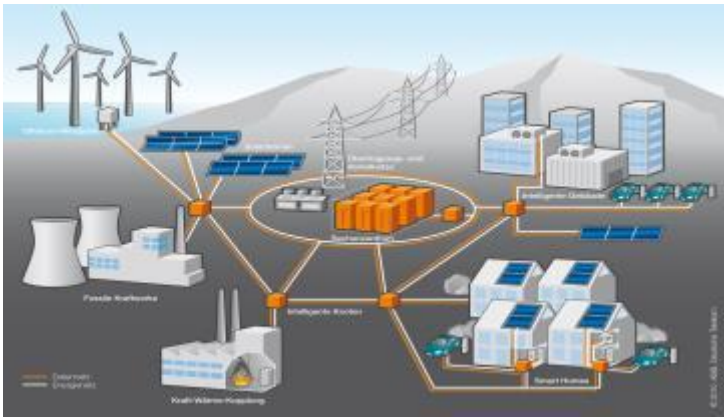
# Power Market and Data Center Participation

## Power Market:

- Dynamic pricing policy for RS bidding [Caramanis CDC12]
- Smart building RS provision [Paschalidis CDC-ECC11]

## Data Center Participation:

- Analytical profit model of data center participation [Ghamkhari SmartGridComm12]
- Analysis of different advanced power market for data centers to participate [Aikema IGCC12]
- Workload allocation among geographically distributed data centers [Wang ICDCS13][Wang SIGMETRICS13]



← Smart Grid

This work is **the first** to design policies for the **data center** for:

- *Power budgeting and management*
- *Server commitment*

to enable the data center to participate in the advanced power market programs.



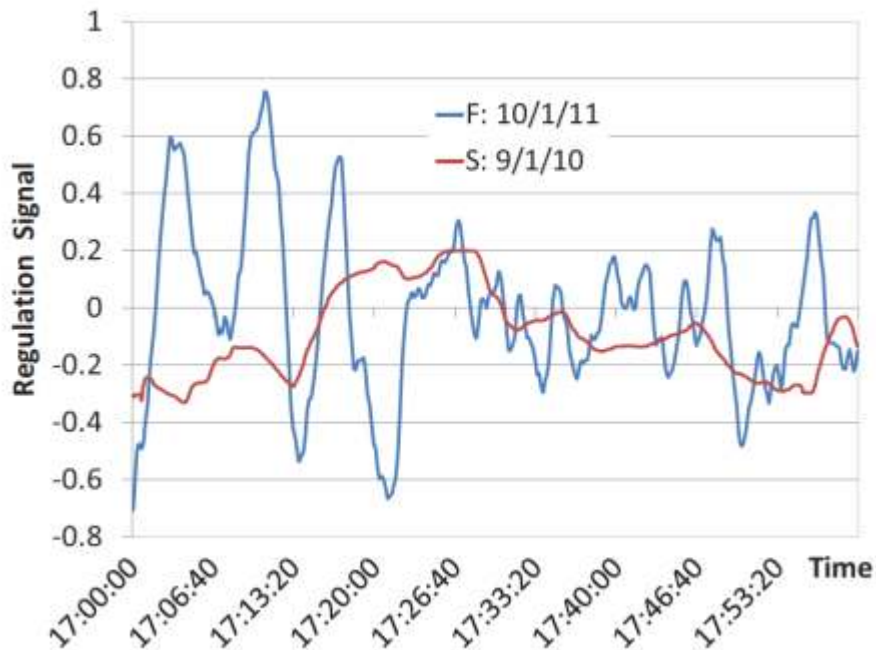
# Regulation Service (RS)

Bidding:  $(\bar{P}, R)$

Price Settling:  
Get contract

ISO: RS  
signal

Data Center  
Regulation



Typical PJM 150sec ramp rate (F) and 300sec ramp rate (S) regulation signal trajectories

$$P_{cap}(t) = \bar{P} + z(t)R$$

$$\text{Error: } \epsilon(t) = \frac{|P_{real}(t) - P_{cap}(t)|}{R}$$

$\epsilon(t)$  needs to be small:

$\epsilon(t) > \text{threshold} \Rightarrow \text{lose license}$

Costs:  $\Pi^E \bar{P} - \Pi^R R$  **Credit Earned**

- $\Pi^E$  and  $\Pi^R$  : market clearing prices
- Credits are reduced based on statistics of  $\epsilon(t)$

# Outline

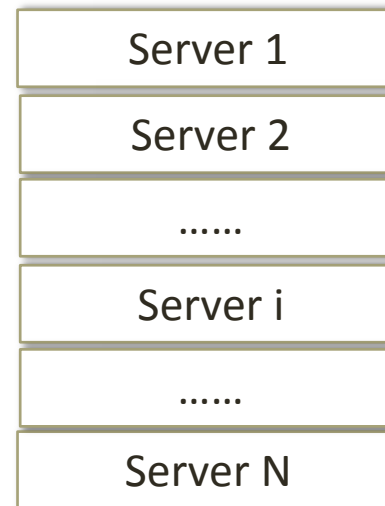
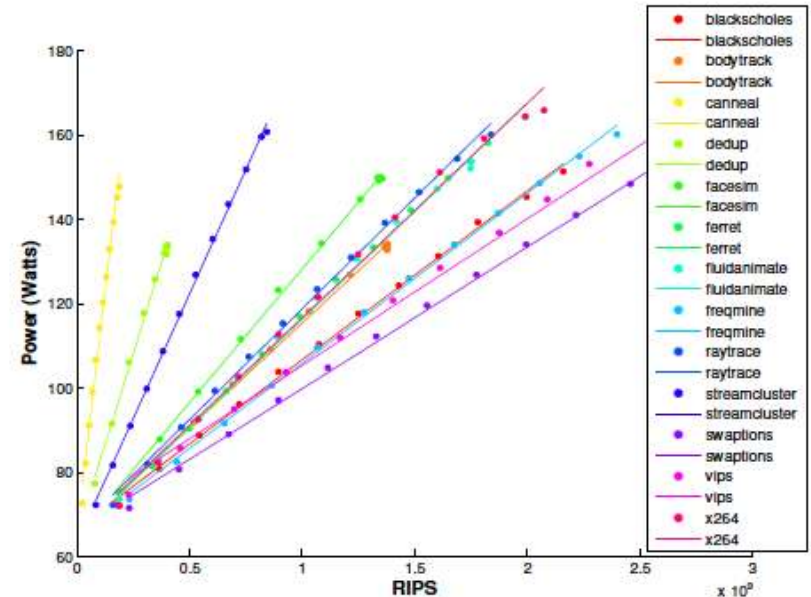
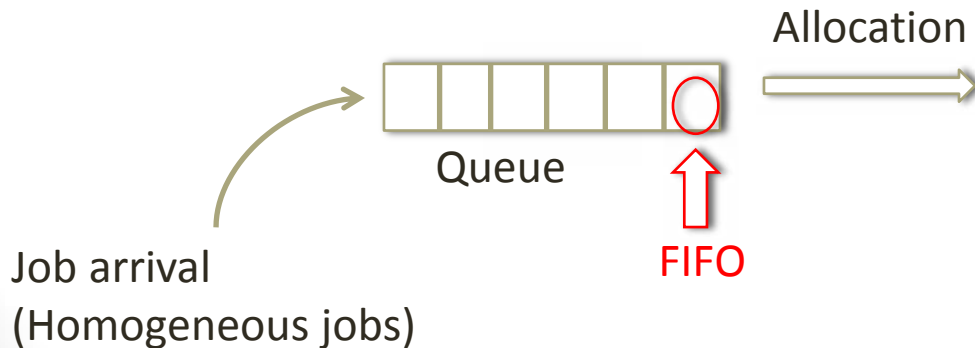
- Background
- **Data Center Model**
- Dynamic Power Control Policy
- Regulation Reserves Bidding
- Results

# Data Center Model

- Server States:

- Active:  $P_{\text{server}} = P_{\text{dyn}} + P_{\text{static}}$ 
  - $P_{\text{dyn}}$  can be modulated by DVFS or **CPU resource limits**
  - $P_{\text{dyn}} = k * \text{RIPS}$
- Idle:  $P_{\text{server}} = P_{\text{static}}$
- Sleep:  $P_{\text{server}} = P_{\text{sleep}}$ 
  - Constant low power, but resuming from sleep has time delay ( $t_{\text{res}}$ ) and energy cost ( $E_{\text{loss}}$ )

- Servicing Model:



Each server: 1 job at a time

# Outline

- Background
- Data Center Model
- **Dynamic Power Control Policy**
  - Goals and Optimization Problem
  - Designed Rules and Policies
- Regulation Reserves Bidding
- Results

# Dynamic Power Control Policy

- Goals:

- Reduce the tracking error  $(t) = \frac{|P_{real}(t) - P_{cap}(t)|}{R}$
- Improve the energy efficiency, including:
  - reduce the energy waste during the server state transition period
  - reduce the static energy waste
- Reduce the workload QoS performance degradation

- Optimization:

$$\min_{u(t)} J(x(t), u(t)) = \underbrace{1 |P_{real}(t) - P_{cap}(t)|}_{\text{Tracking Error}} + \underbrace{2 N_{tran}(t)}_{\text{Transition Energy Waste}} + \underbrace{3 N_{sleep}(t) + 4 N_{peak}(t)}_{\text{Static Energy Waste}}$$

- $x(t)$ : data center states at  $t$  (including server states and workload states);
- $u(t)$ : available control set at  $t$ ;
- $N_{tran}(t)$ : # of servers that are suspending to or resuming from the sleep state at  $t$ ;
- $N_{sleep}(t)$ : # of servers in sleep at  $t$ ;
- $N_{peak}(t)$ : # of servers running at their peak capacities at  $t$ .

# Dynamic Power Control Policy

## Additional Designed Rules:

- For a server that is running a job:  
*=> keep active at a power rate at least  $P_{min}$  until job finished, to guarantee QoS;*
- When no jobs are waiting in the queue:  
*=> no idle server is activated.*
- **Server State Transition Rules** [Gandhi IGCC12]:
  - A server that has been in idle  $> t_{out}$  (timeout threshold):  
*=> goes to sleep;*
  - When a new job arrives:  
*=> select the server with the smallest current  $t_{idle}(t)$  to activate;*
  - When we need to force servers to sleep:  
*=> select the servers with the largest current  $t_{idle}(t)$  to put to sleep.*

$t_{idle}(t)$ : the time that a server has been in the idle state at time  $t$ .



# Dynamic Power Control Policy

## Control Policy:

- **Case 1:**  $P_{\text{real}}(t-1) < P_{\text{cap}}(t)$

1. Active servers with  $P_{\text{server}} < P_{\text{peak}}$ :  $P_{\text{server}} \rightarrow P_{\text{peak}}$ ;
2. Existing waiting jobs and idle servers: **activate idle servers**  $\rightarrow P_{\text{peak}}$ ;
3. Sleeping servers: resume using *server state transition rules*.

Do the above three steps **in order** until  $P_{\text{real}}(t) = P_{\text{cap}}(t)$ .

- **Case 2:**  $P_{\text{real}}(t-1) > P_{\text{cap}}(t)$

1. Active servers with  $P_{\text{server}} < P_{\text{peak}}$ :  $P_{\text{server}} \rightarrow P_{\text{min}}$ ;
2. Active servers with  $P_{\text{server}} = P_{\text{peak}}$ :  $P_{\text{server}} \rightarrow P_{\text{min}}$ ;
3. Idle servers: suspend using *server state transition rules*.

Do the above three steps **in order** until  $P_{\text{real}}(t) = P_{\text{cap}}(t)$ .

# Outline

- Background
- Data Center Model
- Dynamic Power Control Policy
- **Regulation Reserves Bidding**
  - Estimate  $(\bar{P}, R)$
- Results

# Regulation Reserves Bidding

Average Power Consumption:

$$\bar{P} = \frac{\int_0^{1h} (\bar{P} + Rz(t)) dt}{1h} = \underbrace{\bar{N}_{active} * P_{active} + \bar{N}_{idle} * P_{idle} + \bar{N}_{sleep} * P_{sleep}}_{\text{Power of Servers in diff. states}} + \frac{E_{loss,1h}}{1h} \quad (1)$$

Avg. # of Servers in diff. states

Transition power waste

$$E_{loss,1h} = E_{loss} * \underbrace{N_{res}}_{\substack{\text{\# of state transitions in 1h} \\ \downarrow}} \quad (t_{res} * N_{tran}) \times (p_b * N_{dc}) \quad (2)$$

Energy waste of each transition

$$N_{dc} = \bar{N}_{active} + \bar{N}_{idle} + \bar{N}_{sleep} \quad (3)$$

$$\bar{N}_{active} = \frac{\int_0^{1h} N_{active}(t) dt}{1h} \quad \frac{E_{dyn}}{P_{dyn,max} * 1h} \quad \frac{* kI}{P_{dyn,max} * 1h} \quad (4)$$

Total dynamic energy for processing jobs

$\bar{N}_{idle} = \bar{N}_{sleep}$

Slack

**Regulation Reserve:**  $R \min\{ \underbrace{N_{dc} P_{peak}}_{\substack{\text{\# of servers in the data center} \\ \uparrow}}, \underbrace{\bar{P}, \bar{P}}_{\text{Min, Max power of servers}}, N_{dc} P_{sleep} \}$

# Outline

- Background
- Data Center Model
- Dynamic Power Control Policy
- Regulation Reserves Bidding
- **Results**
  - Methodology
  - Single Server vs. Data Center
  - Fast Sleep vs. Deep Sleep
  - Impact of Cluster Utilization
  - Impact of Different Workloads

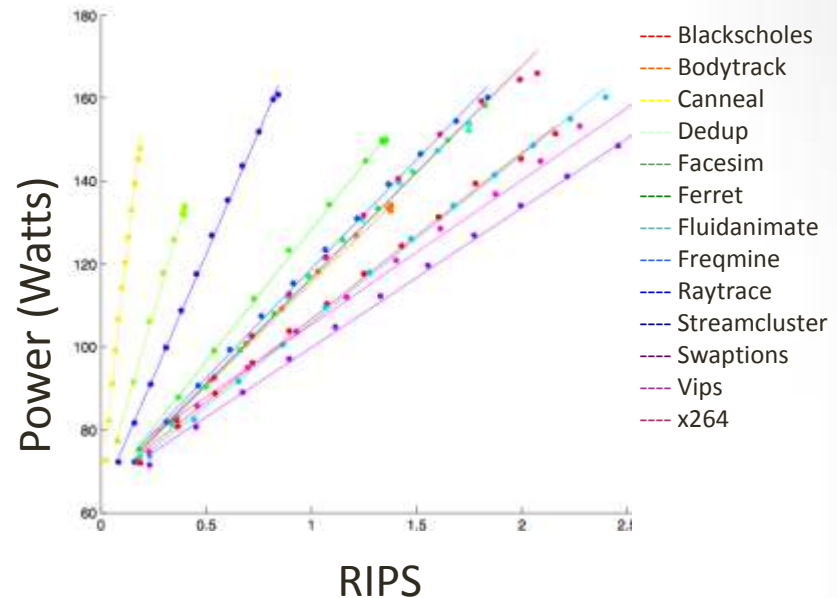
# Methodology

VMware vSphere 5.1

ESXi hypervisor

AMD Magny Cours  
(Opteron 6172)  
processor, 12 cores

Wattsup Power Meter



Linear Regression:  $P_{dyn,j}$

$$P_{server,j} = C_j * RIPS_j + P_{static}$$

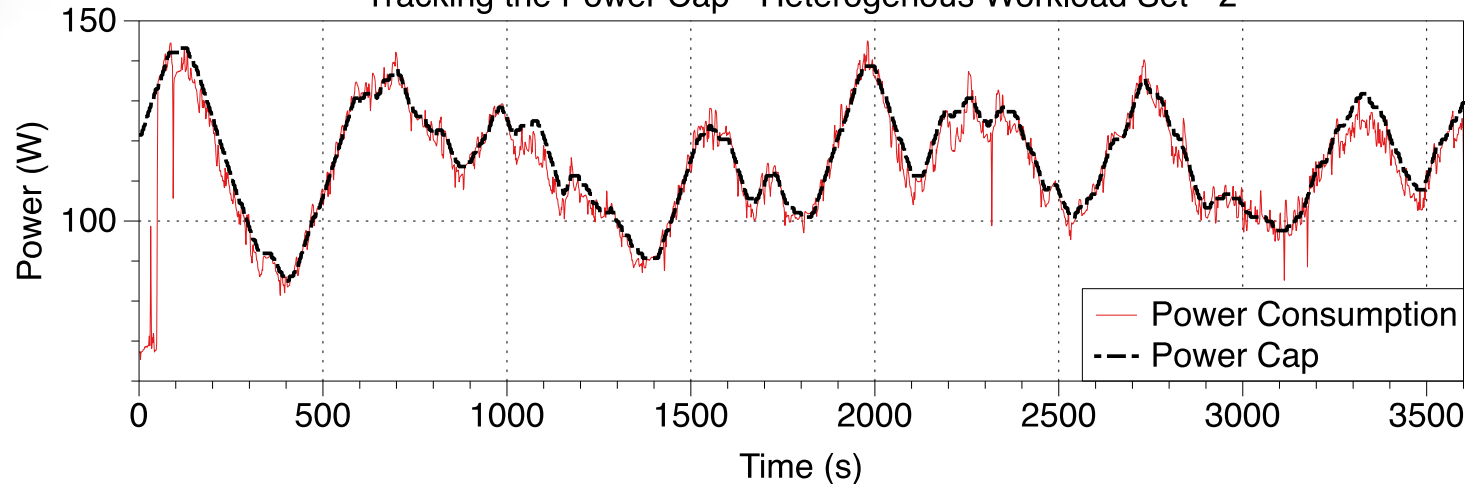
1-hour long HPC type workload (*run 10 times*)

- Applications from PARSEC 2.1 multi-threaded benchmark suite
- Job arrivals follow a Poisson process
- Generated by Monte Carlo method

Data Center: 100 Servers

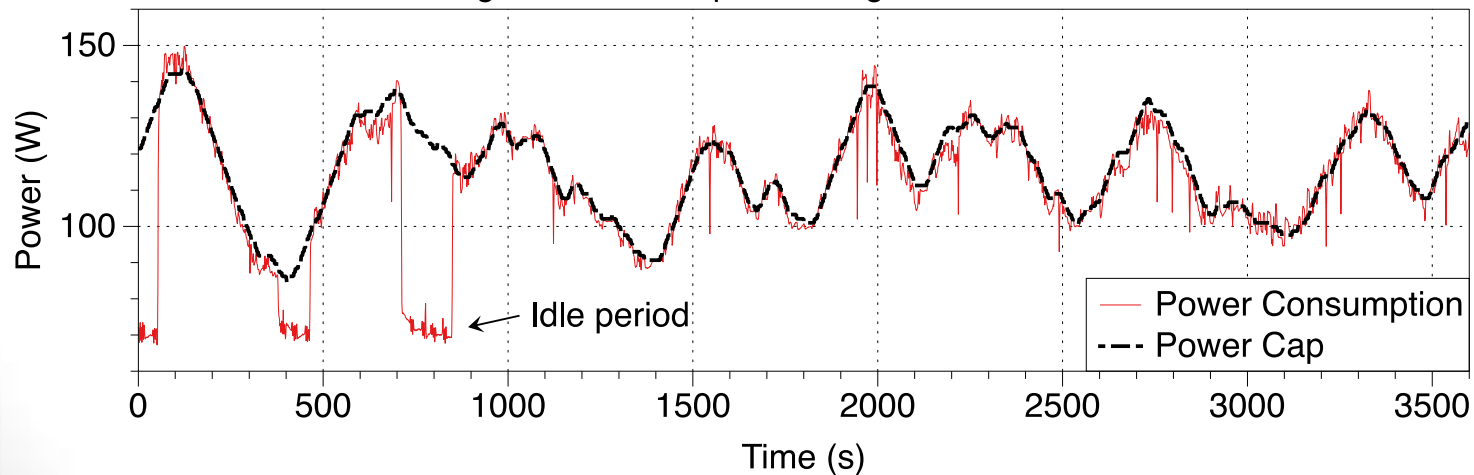
# Power Tracking – Single Server (ICCAD'13)

Tracking the Power Cap - Heterogenous Workload Set - 2



**Error  
7- 8%**

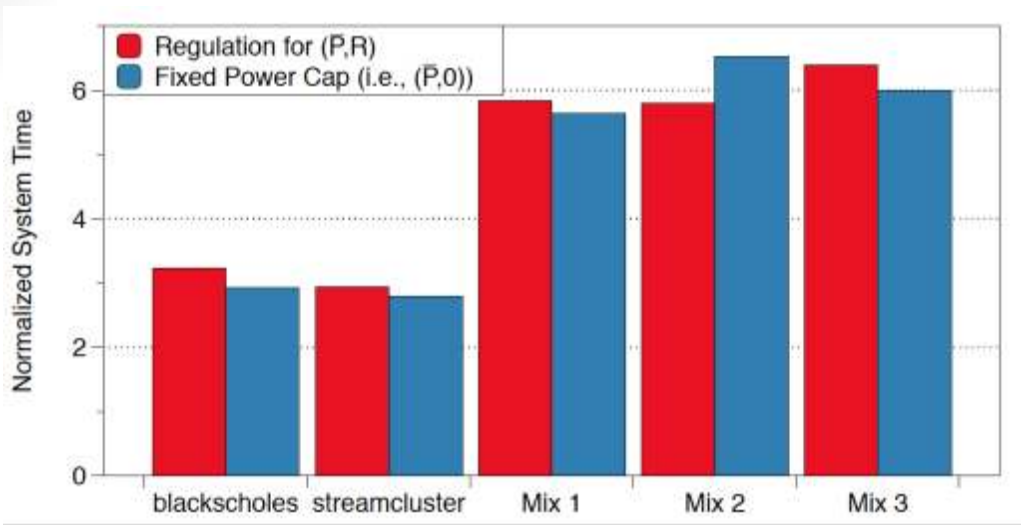
Tracking the Power Cap - Heterogenous Workload Set - 3



**Synthetic workload can fill in the idle periods.**



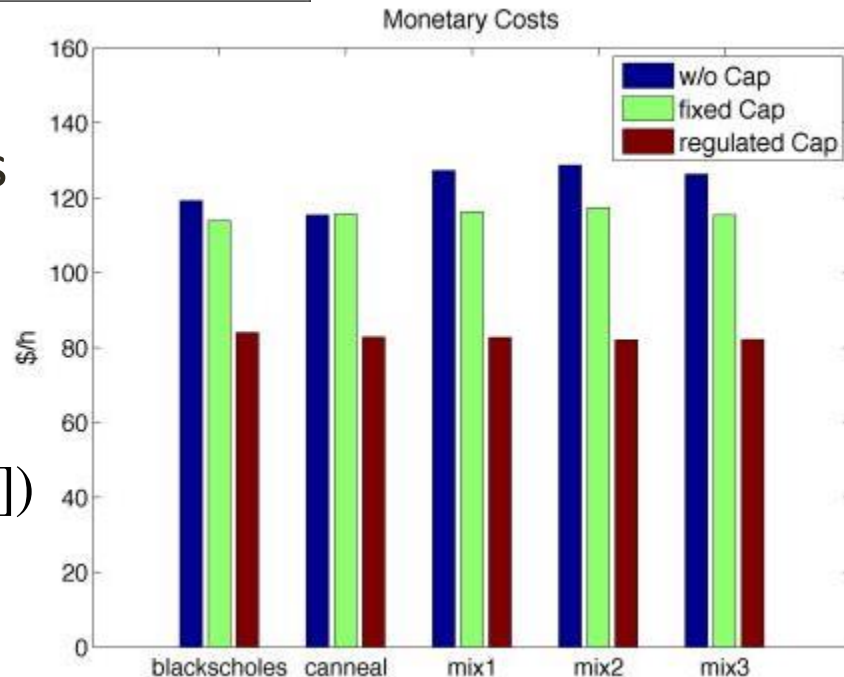
# QoS & Monetary Savings (ICCAD'13)



Similar  
Performance...

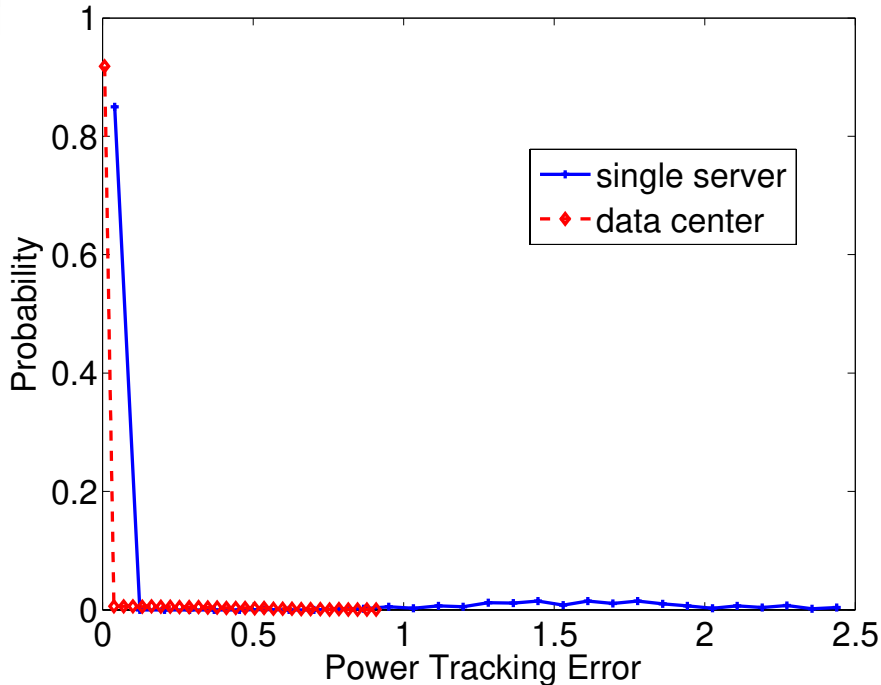
29% Monetary  
Savings!!!

- 10,000 identical servers
- w/o Cap:  $E P(t)$
- Fixed Cap:  $E \bar{P}$
- Regulation:  $E \bar{P} (R c [2 + (\bar{P})^2])$

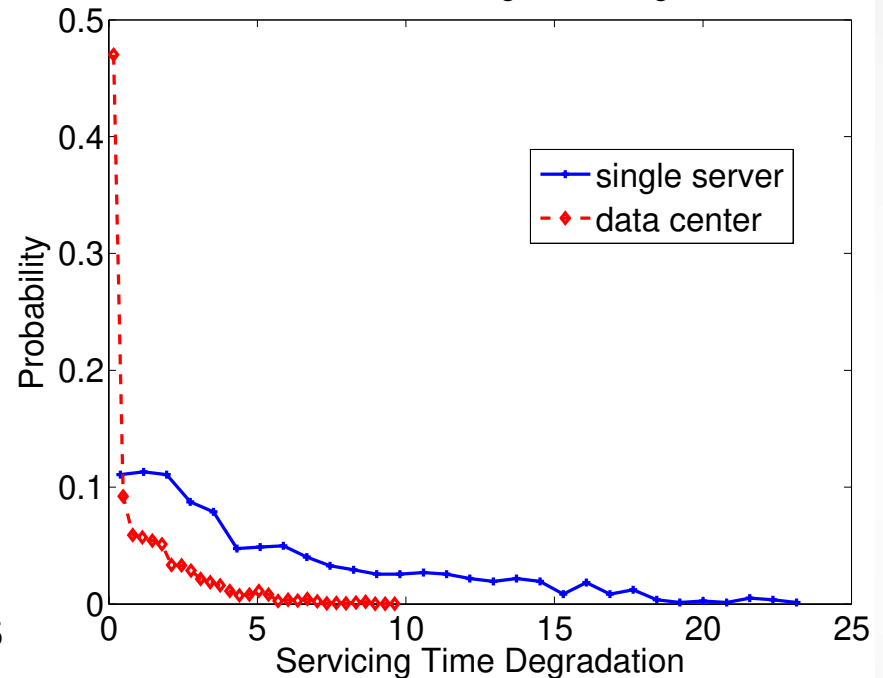


# Results - Single Server vs. Data Center

Distribution of Power Tracking Error



Distribution of Servicing Time Degradation

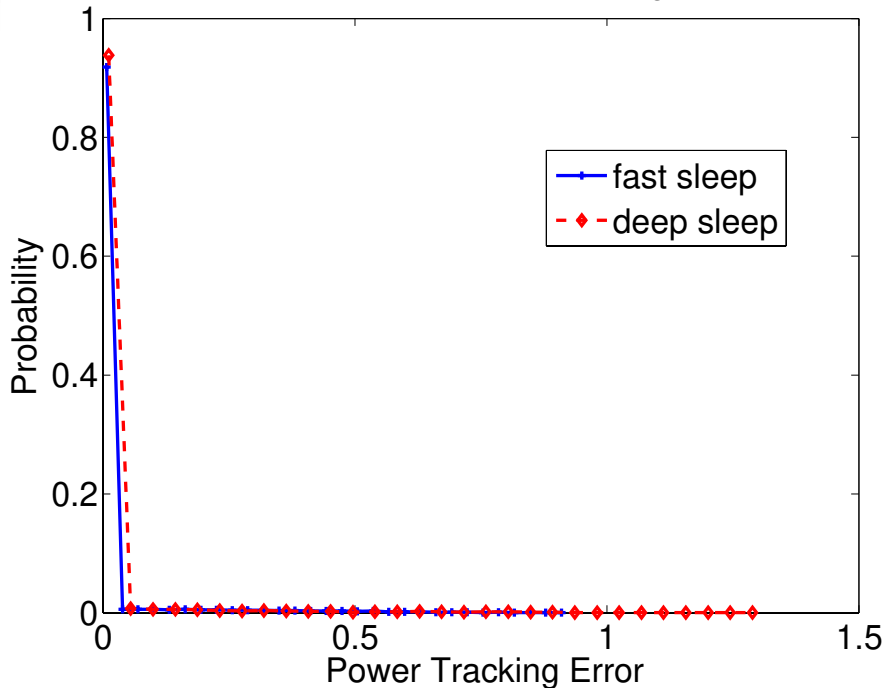


Regulation Reserves ( $R$ ) / Avg. Power Consumption ( $\bar{P}$ ):

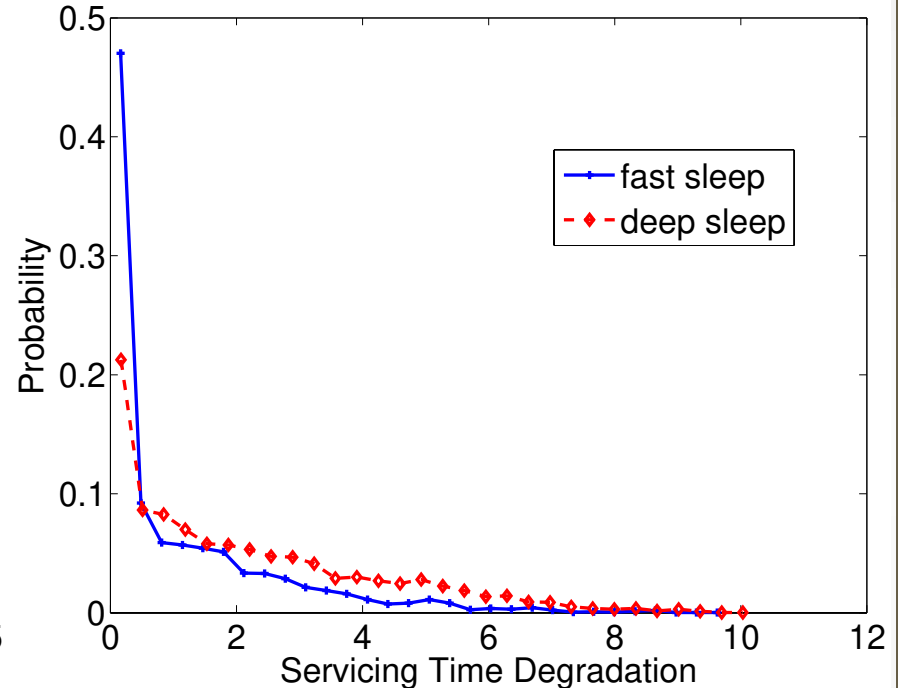
- Single Server: 29.7%
- Data Center: 56.8%

# Results - Fast Sleep vs. Deep Sleep

Distribution of Power Tracking Error



Distribution of Servicing Time Degradation

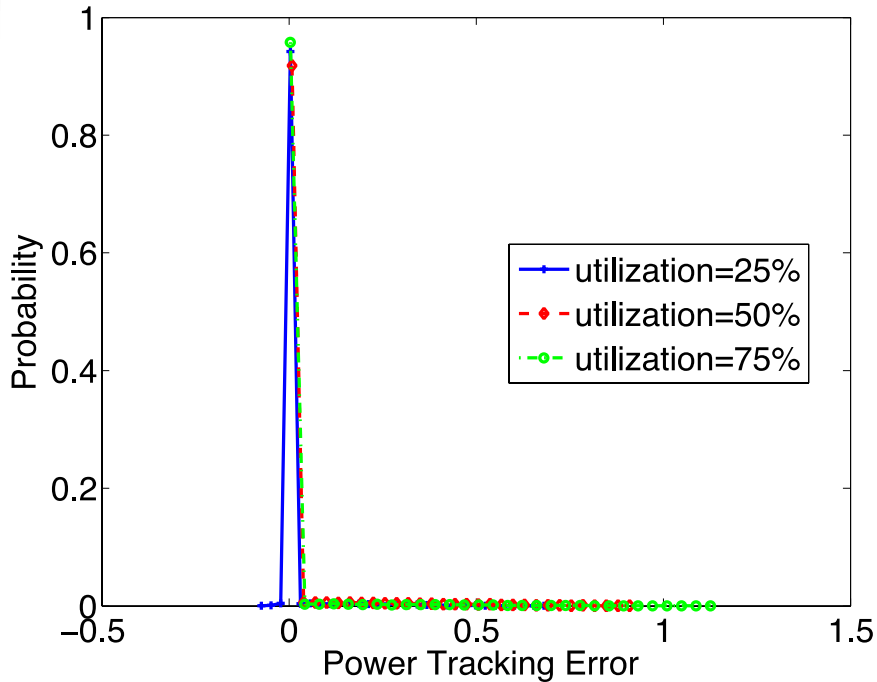


$R/\bar{P}$ :

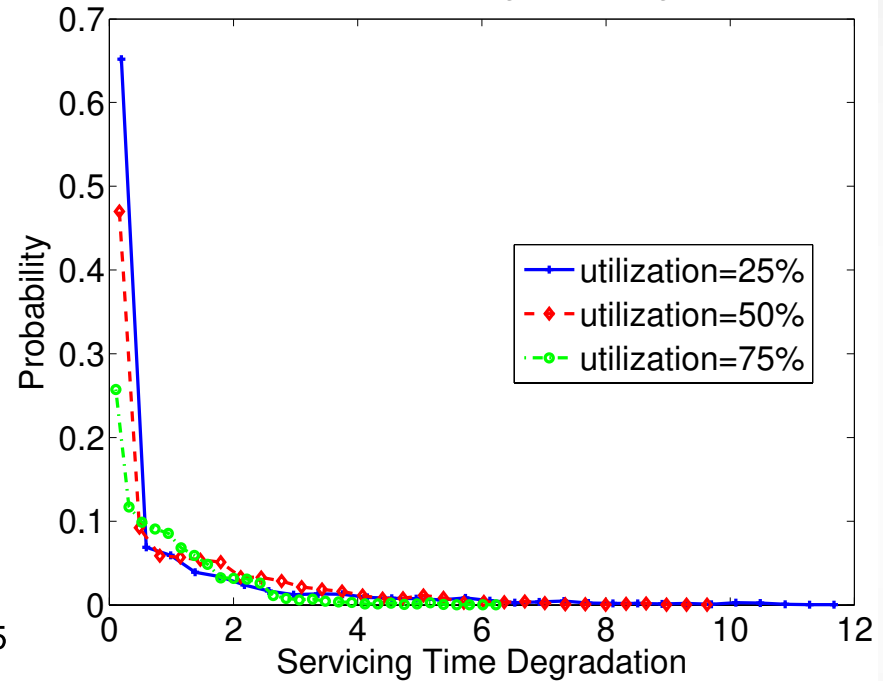
- Fast Sleep ( $t_{res}=10s$ ,  $P_{sleep}=10\%*P_{peak}$ ,  $P_{tran}=P_{peak}$ ): 56.8%
- Deep Sleep ( $t_{res}=200s$ ,  $P_{sleep}=5\%*P_{peak}$ ,  $P_{tran}=P_{peak}$ ): 36.9%

# Results - Impact of Cluster Utilization

Distribution of Power Tracking Error



Distribution of Servicing Time Degradation



$R/\bar{P}$ :

- 25% Utilization: 78.0%
- 50% Utilization: 56.8%
- 75% Utilization: 21.8%

# Results - Impact of Different Workloads

## CLUSTER LEVEL POWER REGULATION ON DIFFERENT WORKLOADS

	Blackscholes	Canneal	Streamcluster	Facesim	
$\bar{P}$	$9.75 * 10^3$	$9.71 * 10^3$	$9.84 * 10^3$	$9.84 * 10^3$	
$R$	$5.54 * 10^3$	$4.98 * 10^3$	$5.46 * 10^3$	$5.11 * 10^3$	
$\bar{D}$	1.13	1.13	0.21	0.22	} QoS Degradation
$\sigma_D$	1.54	0.69	0.26	0.27	
$\bar{\epsilon}$	0.03	0.03	0.03	0.03	} Tracking Error
$\sigma_\epsilon$	0.10	0.09	0.09	0.09	
$R/\bar{P}$	56.8%	51.3%	55.5%	52.0%	

<sup>a</sup> $\bar{D}$  and  $\sigma_D$  are mean and standard deviation of performance degradation;  $\bar{\epsilon}$  and  $\sigma_\epsilon$  are mean and standard deviation of tracking error.

# Conclusion & Future Work

- A **dynamic control policy** for the data center RS provision
- An **estimation method** to calculate the RS provision bidding value
- Data centers are promising candidates for RS provisioning:
  - Accurately track the RS signal;
  - Achieve 50%+ monetary savings;
  - With no major QoS degradation;
  - Regardless of types of workloads.
- Significant improvement of **data center** vs. prior single server results, taking **sleep states, utilization, etc.** into account

- Future work:

**Heterogeneous jobs & Power budgeting**

