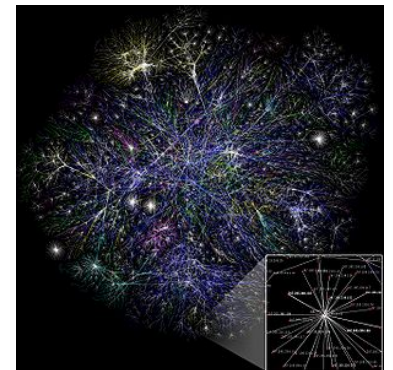


*Leana Golubchik, USC*

# Flexible Clouds

achieving the "right" cost-  
benefit balance through proper  
resource management



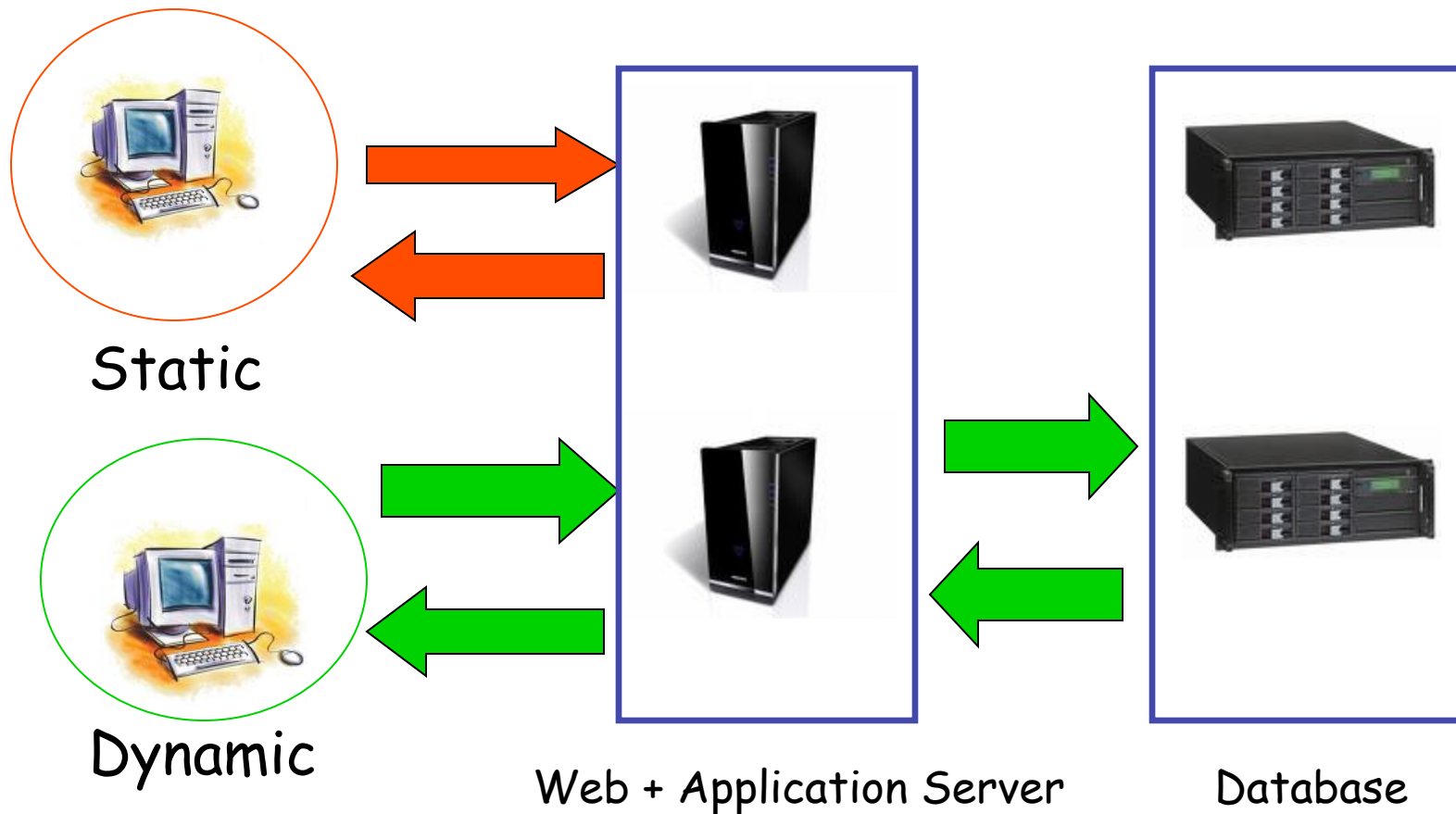
# Focus Today

- Benefits and possible sources of flexibility
  - Highly heterogeneous environment and workload
  - Maximizing flexibility through proper resource management
  - Facilitating an appropriate trade-off between benefits to the cloud users (performance, reliability, etc) and cost to the cloud providers (power, management, etc)
- Example problems
  - Workload characterization
  - Workload scheduling
  - QoS-aware power cost management

# Research Directions

- Management of energy costs
- Management of mixed workloads and heterogeneous resources
  - Data placement to facilitate flexible workload management
  - Ability to characterize and predict workloads
  - High server utilization (without performance degradation)
  - Greater opportunities for improved power usage
  - Resilience to failure and graceful degradation
- Ability to predict level of service delivered to users

# Workload Characterization



# Less Intrusive Approach

- Classify requests based on usage
- Without sophisticated monitoring/logging infrastructure



?



Request Categories  
Arrival rates/processes  
Per-category resource usage  
at each service tier

# Avoiding Invasive Instrumentation

- Use aggregate measurements (easier to obtain)
- Assume request categories correspond to high level transactions
- If can use an inference technique that doesn't assume knowledge of request categories, then can avoid invasive instrumentation

# Inference-based Approach

- Linear model of aggregate resource consumption

$$a_{i1}x_1 + a_{i2}x_2 + \dots + a_{im}x_m + c_i = r_i$$

$a_{i1}$ : # category 1 requests  
 $a_{im}$ : Resources consumed by one request  
 $c_i$ : "Noise"  
 $r_i$ : Aggregate resource usage

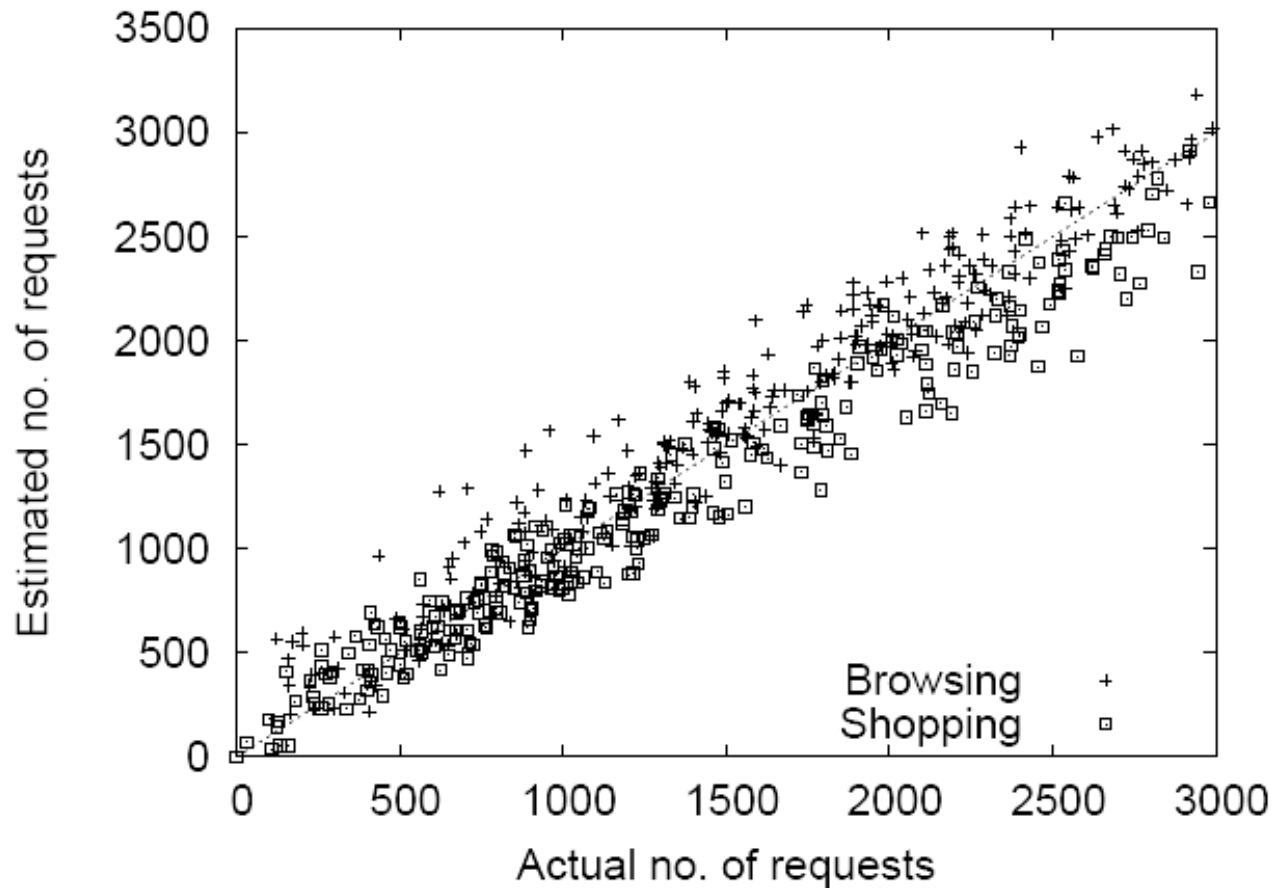
- System of equations for  $n$  resources,  $m$  categories and  $T$  measurements,  $\mathbf{AX} + \mathbf{C} = \mathbf{R}$

$$\begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} x_1 & \tilde{x}_1 \\ x_2 & \tilde{x}_2 \end{bmatrix} + \begin{bmatrix} c_1 \\ c_2 \end{bmatrix} = \begin{bmatrix} r_1 & \tilde{r}_1 \\ r_2 & \tilde{r}_2 \end{bmatrix}$$

$(n \times m)$        $(m \times T)$        $(n \times 1)$        $(n \times T)$

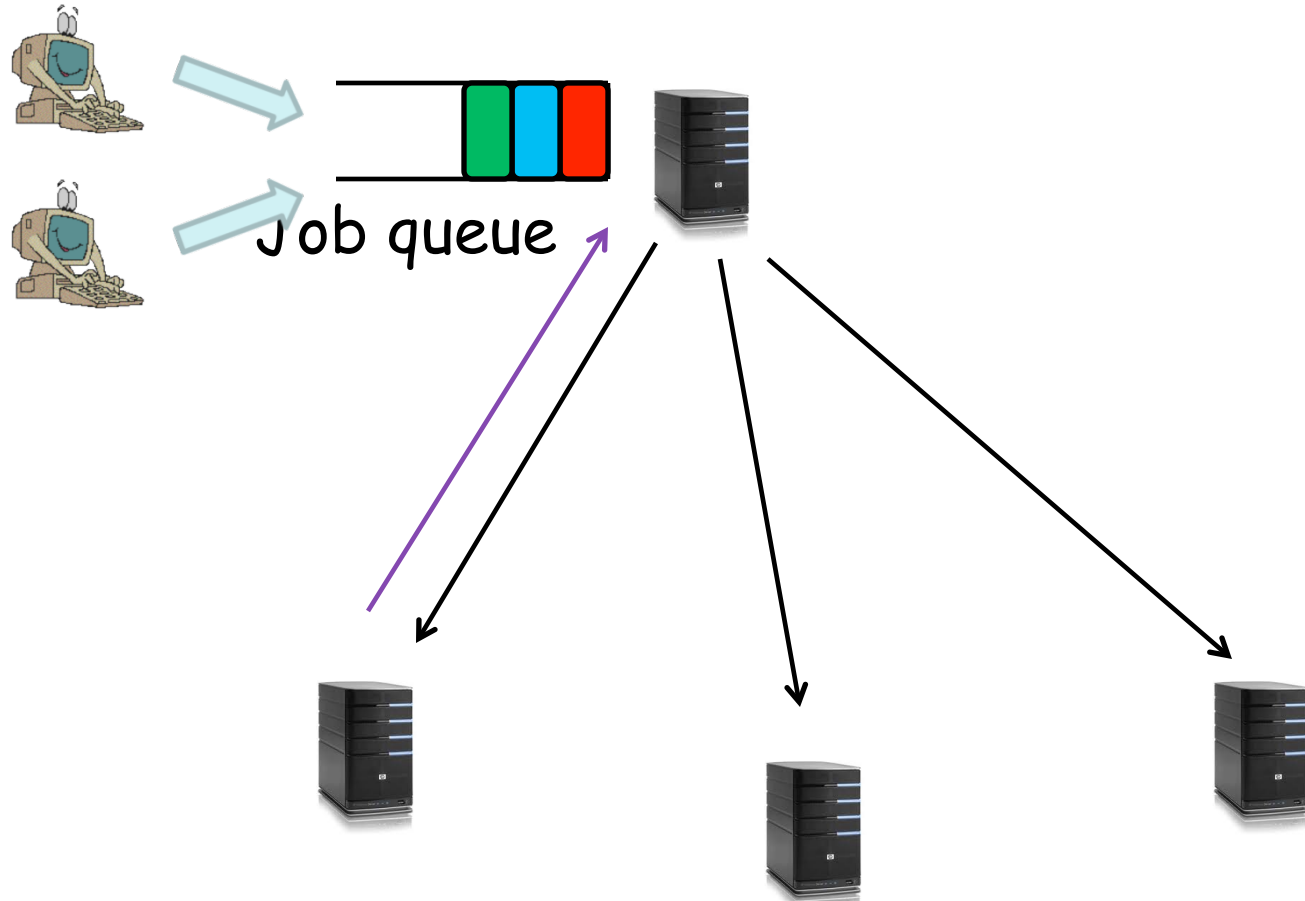
“blind source separation”/Independent Component Analysis


# Request Classification





# Workload Scheduling



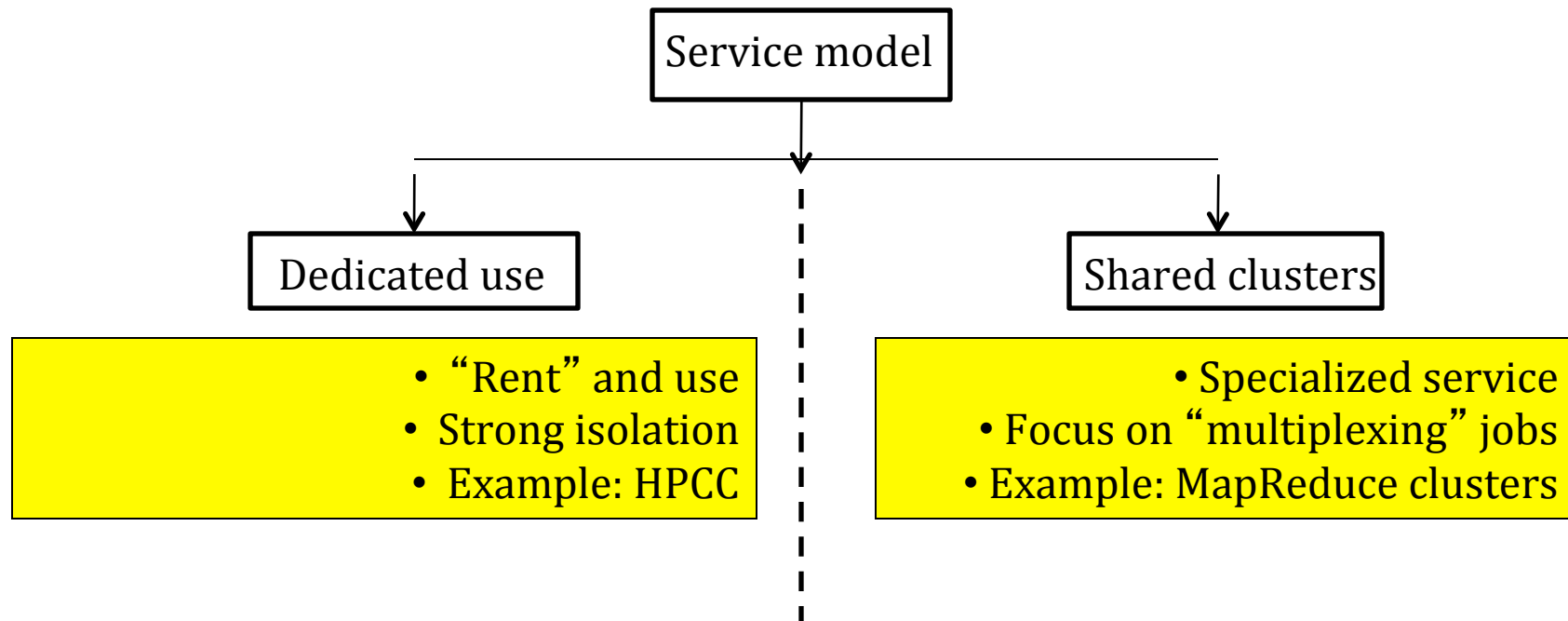
 Node-hours

# Resource Management

- “Match” the available resources to request demands
  - Resources: CPU, memory, storage, bandwidth
- Two “classic” scenarios for resource management:
  - Capacity planning
  - Resource allocation/arbitration
- Performance goals
  - Utilization, quality of service, operating costs
  - Competing goals: Users vs. Administrators
    - Response time vs. Utilization

# Limitations

- No predictability in job service times
- Limited service differentiation

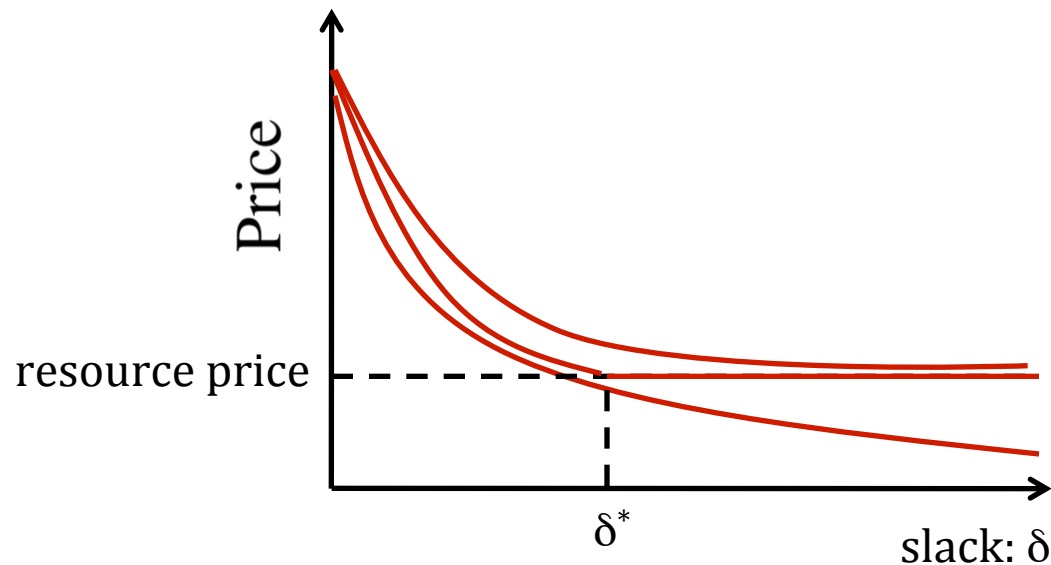


# Predictability and Service differentiation

- Associate a deadline with each job.
- *Predictability*: ensure that deadlines are not violated.
- *Service differentiation*: tighter/earlier deadlines for delay sensitive jobs.
- How do we determine the deadline for a job?
  - Let users select deadlines, with restrictions (earliest finish time is feasible) - FCFS?
  - Incentivize "sufficient slack" in the system (facilitate service differentiation)

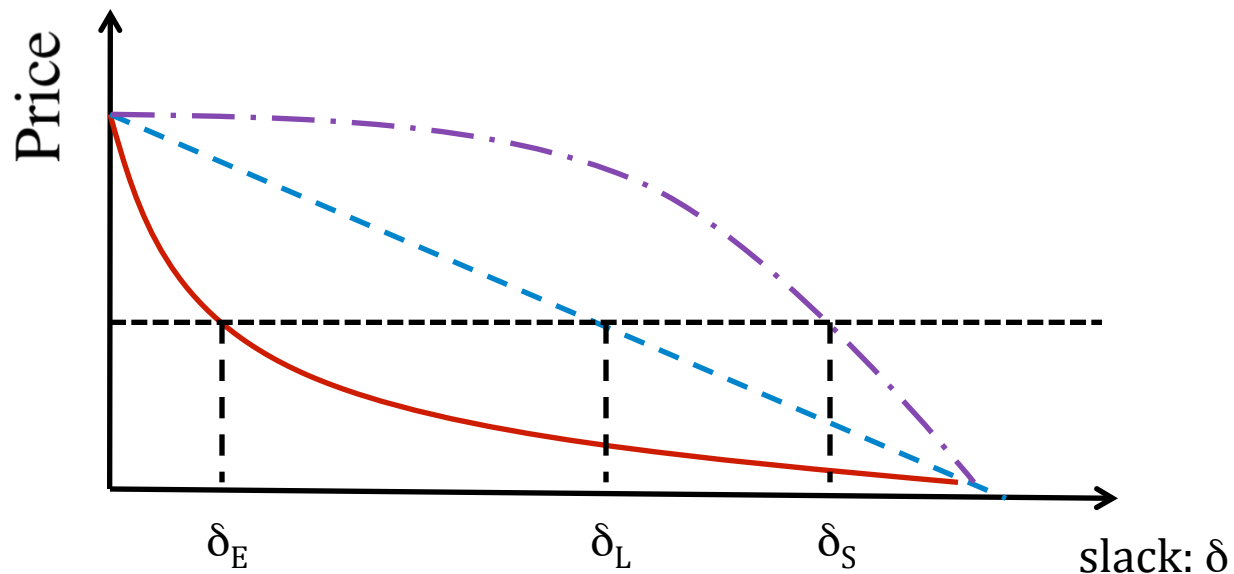
# Incentives for offering slack

- Charge a price for executing a job
  - Price of a job depends on its deadline
  - Administrator offers a price-deadline curve
  - Goal is not to maximize revenue but incentivize users to offer slack



# Choosing the price-deadline curve?

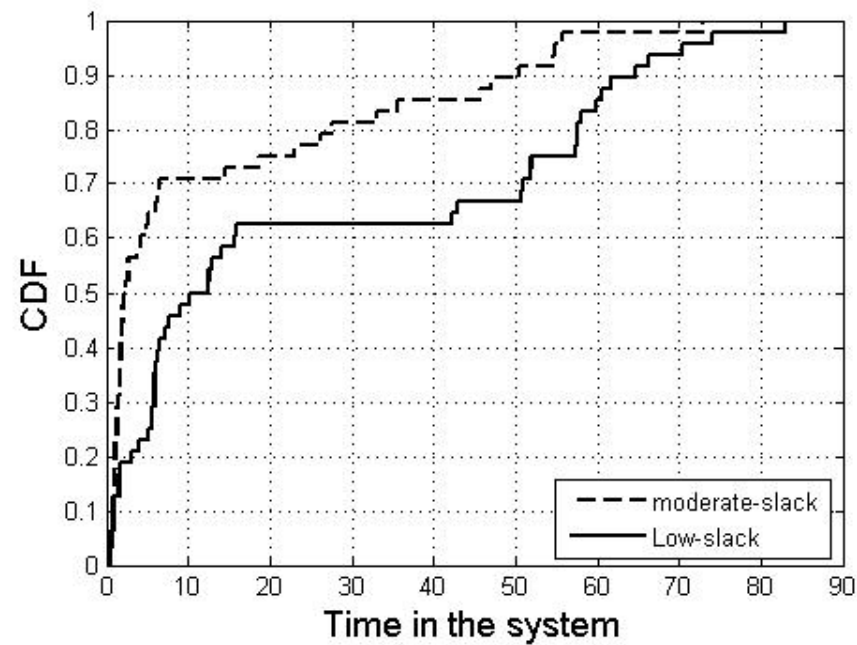
- Design choices:
  - “Decay” shape and rate: How does the price decrease with slack?
  - “Premium” price: How much to charge for not offering any slack?



# Deadline violations with low slack

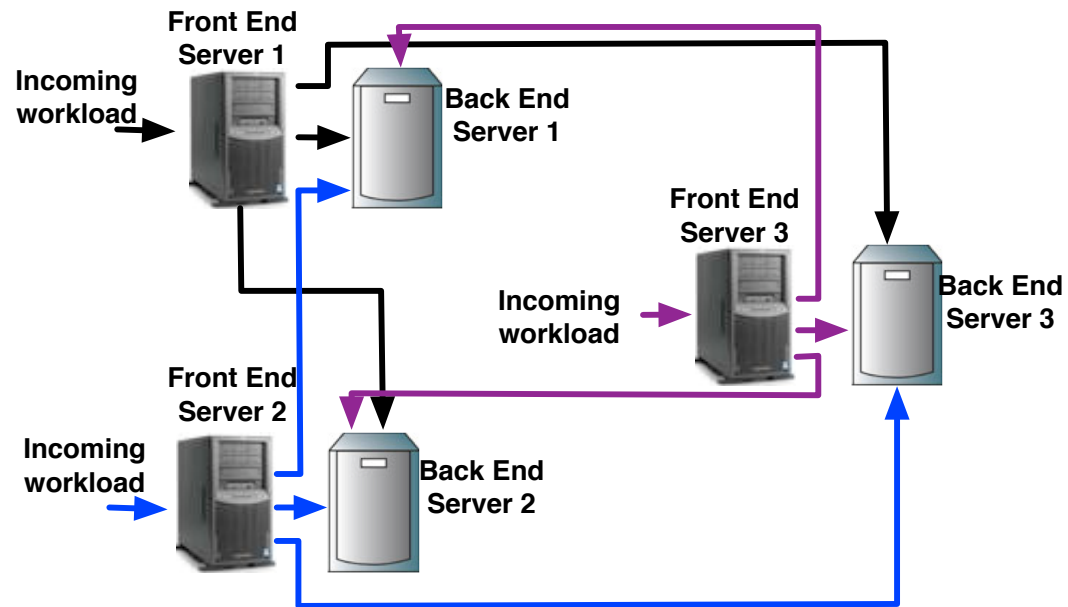
	Shared cluster		Private cluster	
# jobs	% tardy jobs	Tardiness 95th perc.	% tardy jobs	Tardiness 95th perc.
125	0%	0	8%	5.5 mins
100	2%	20 s	0%	0
75	2.7%	24 s	6.7%	2.5 mins
50	0%	0	3%	4.3 mins
25	0%	0	0%	0

# Service Differentiation



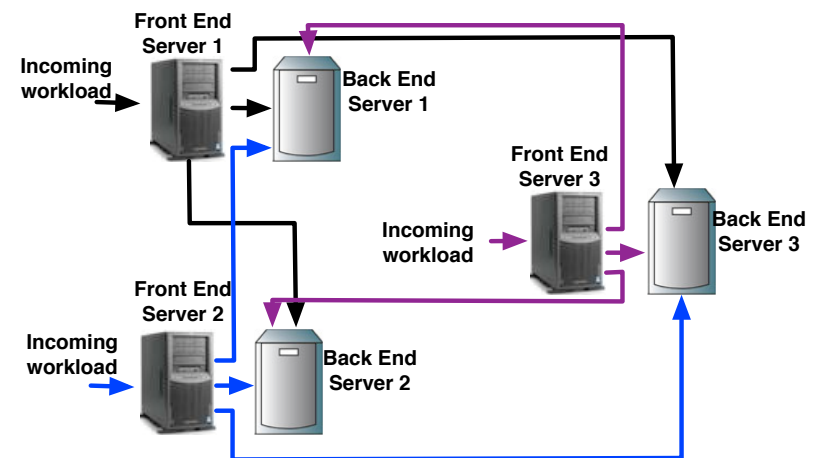


# QoS-aware Power Cost Management

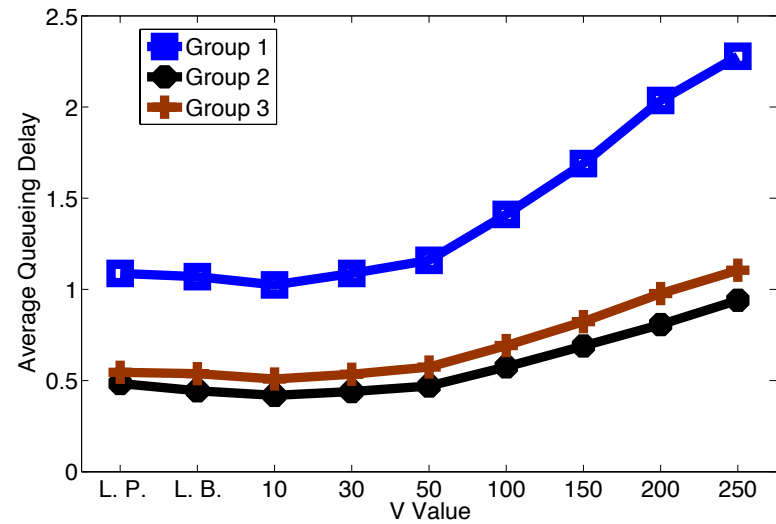
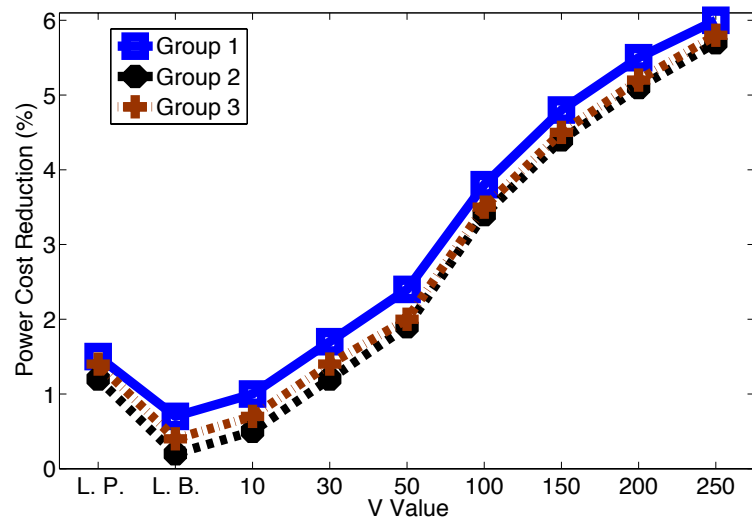


# Reducing Power Cost

- Workload scheduling in geographically distributed centers
- Maintaining QoS
  - Reducing power cost and usage



# Power cost vs Delay



# Summary

- Flexibility in resource management techniques (scheduling, data placement, etc) will lead to a “good” tradeoff between cost (e.g., power) and benefits (e.g., performance) of clouds

# Collaborators

- Yuan Yao, USC
- Abhishek Sharma, USC
- Longbo Huang, USC
- Ramesh Govindan, USC
- Fei Sha, USC
- Mike Neely, USC
- Ranjita Bhagwan, Microsoft
- Monojit Choudhury, Microsoft
- Geoffrey Voelker, UCSD
- Harsha Madhyastha , UCR
- Srikanth Kandula, Microsoft