

# Energy-efficient Dynamic Capacity Provisioning in Server Farms

**VARUN GUPTA**  
Carnegie Mellon University

Partly based on joint work with:

Anshul Gandhi  
(CMU)

Mor Harchol-Balter  
(CMU)

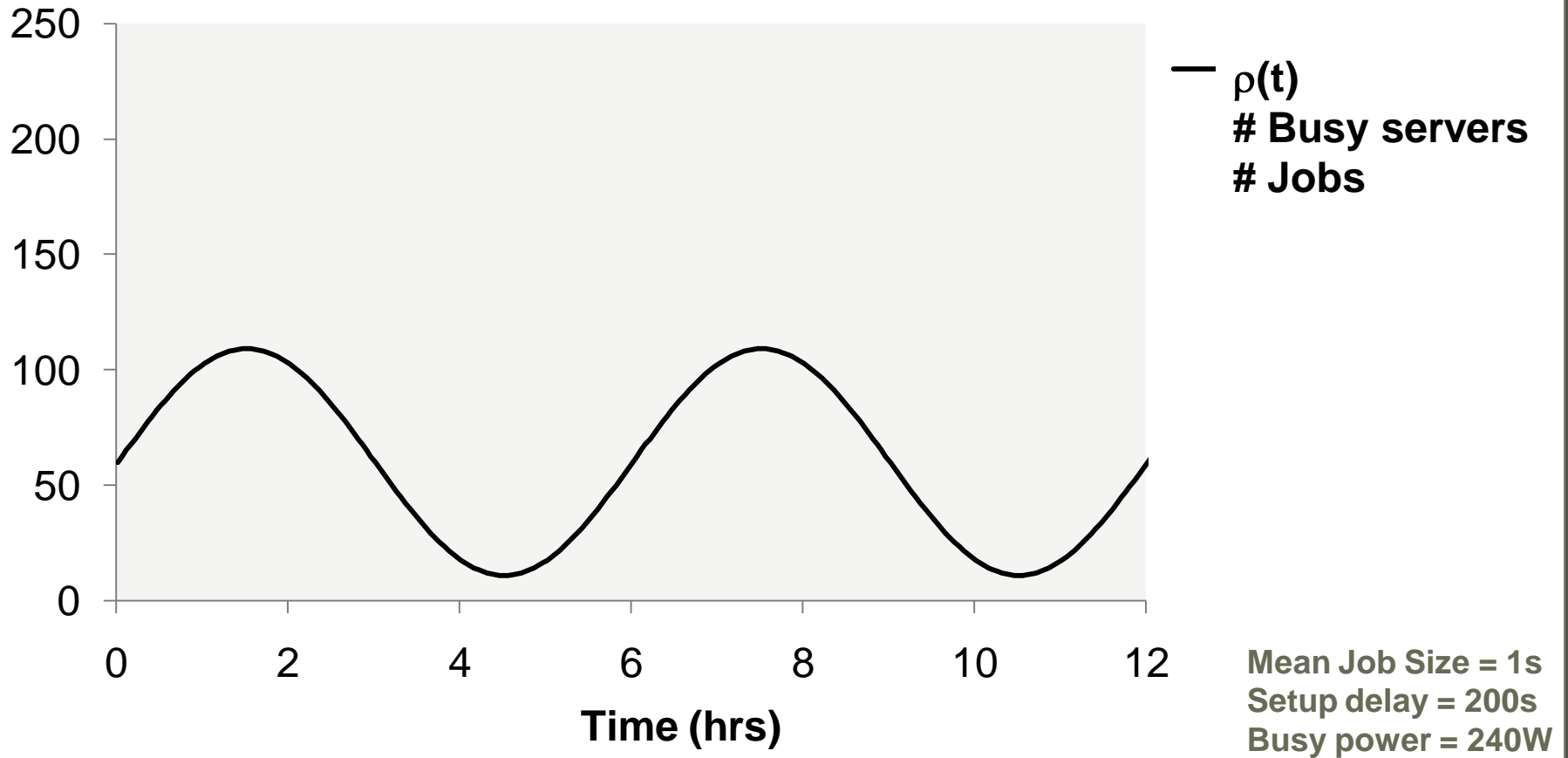
Mike Kozuch  
(Intel Research)

# The “provisioning for peak” problem

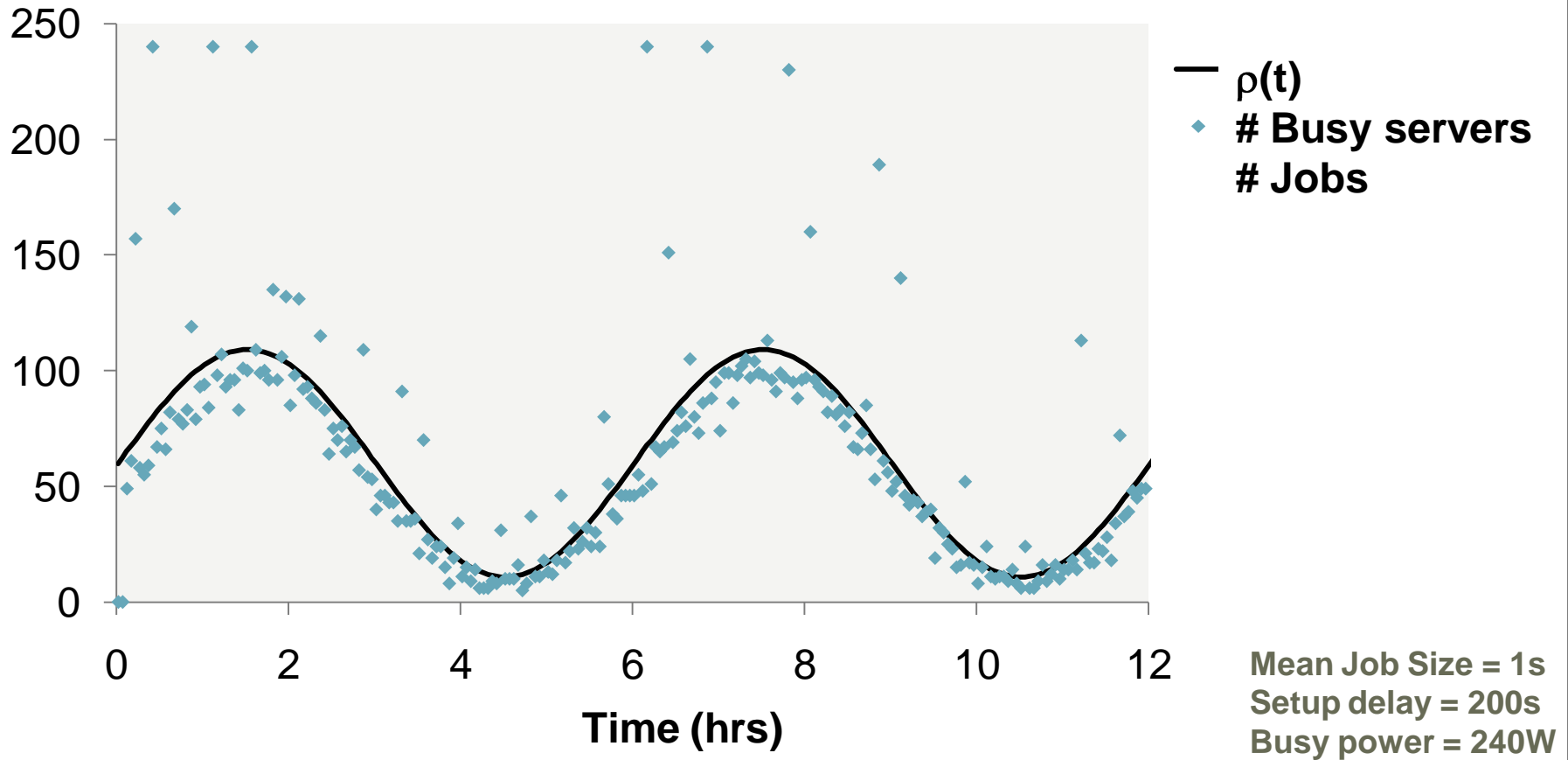


**PROBLEM: Want to turn servers OFF/ON to match  $\rho(t)$ ...**  
**... and also minimize setup penalties!**

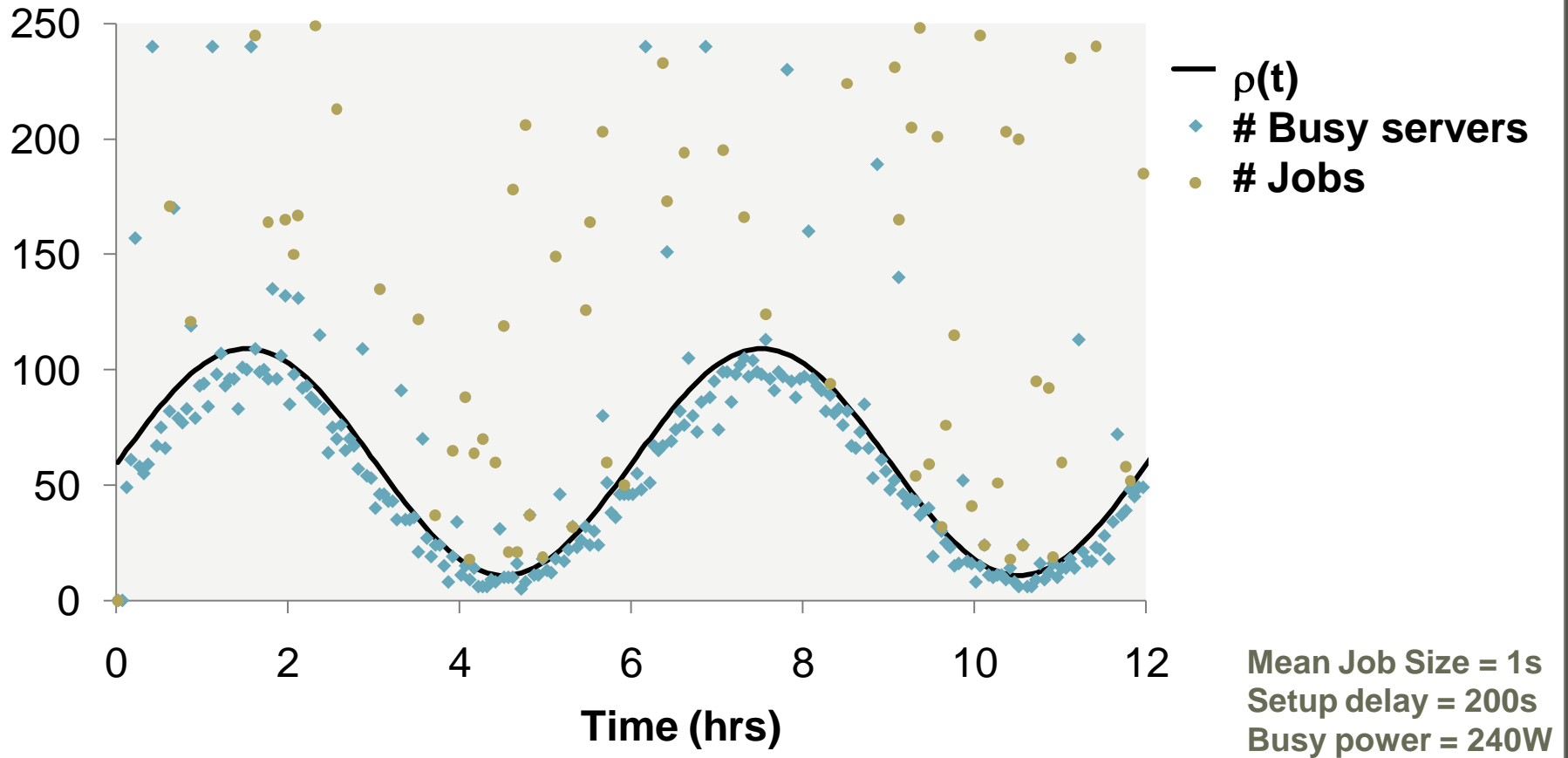
# First Attempt: ON/OFF policy



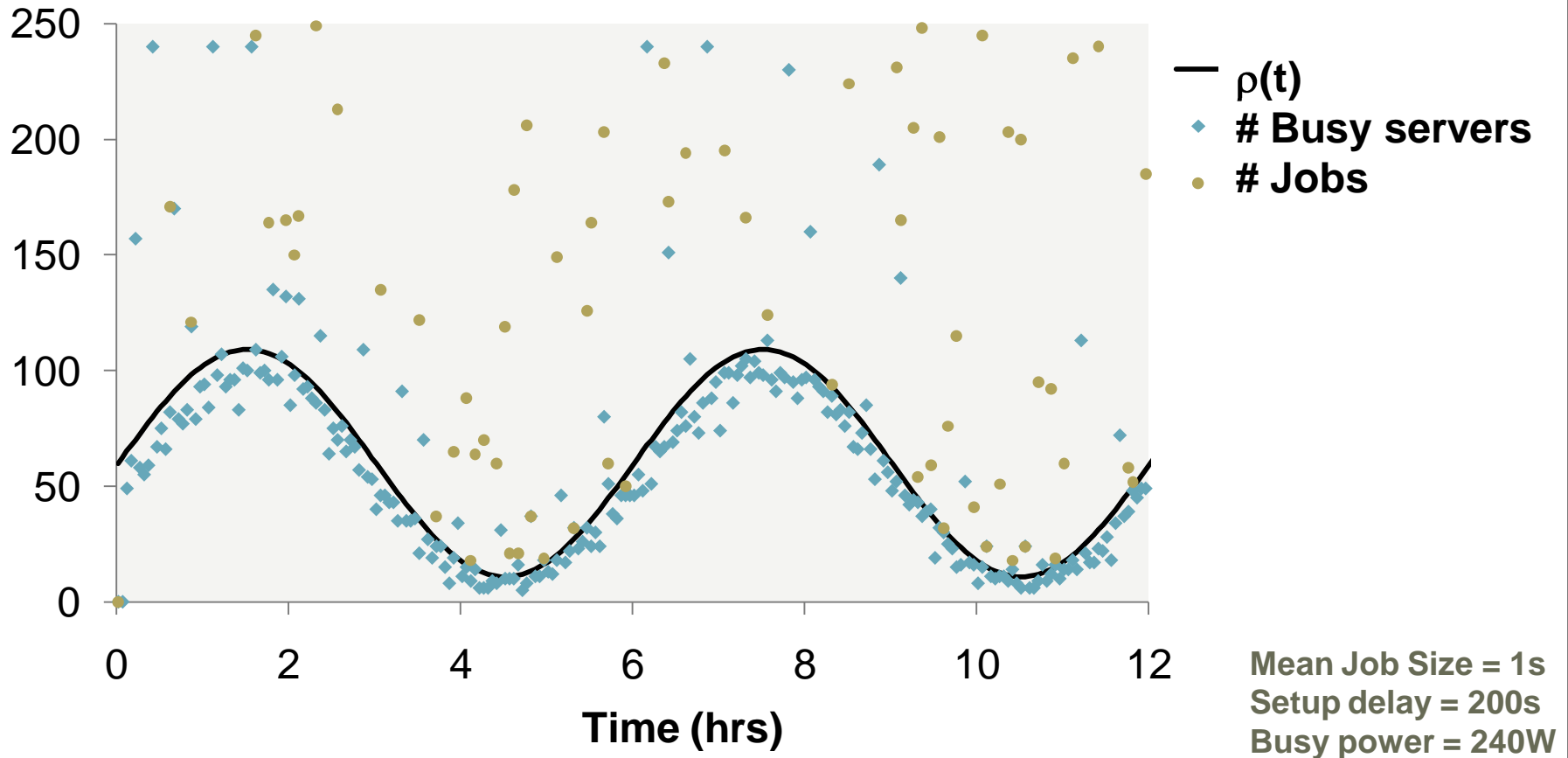
# First Attempt: ON/OFF policy



# First Attempt: ON/OFF policy

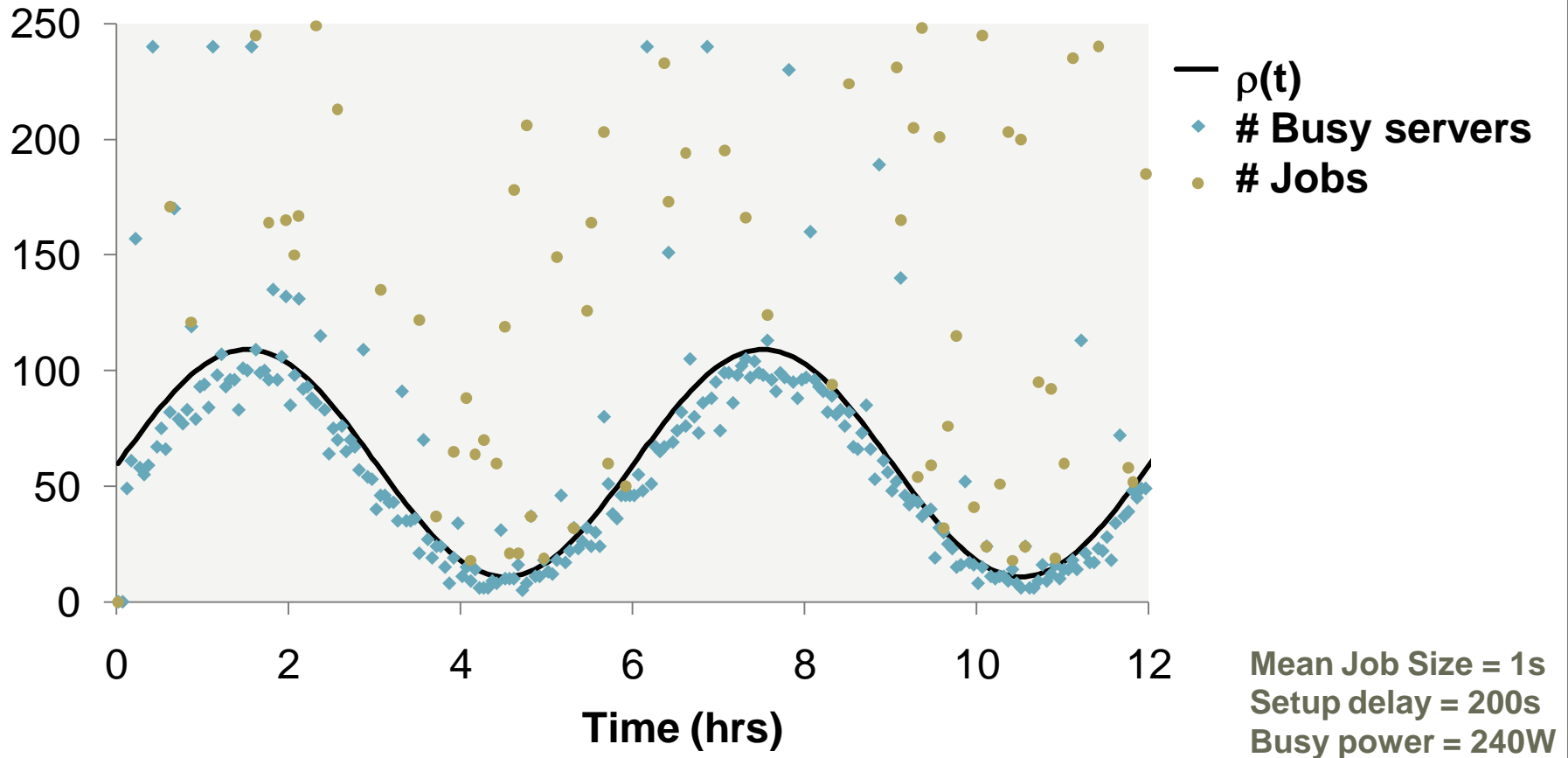


# First Attempt: ON/OFF policy



**WITH SETUP:**     $E[\text{Power}] = 51.9 \text{ kW}$      $E[\text{Response time}] = 13\text{s}$   
**NO SETUP:**     $E[\text{Power}] = 14.4 \text{ kW}$      $E[\text{Response time}] = 1\text{s}$

# First Attempt: ON/OFF policy



**Can we do better than ON/OFF?**

**Add inertia while turning servers OFF**

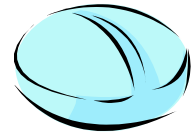
# Our Prescription: DELAYEDOFF



Turn OFF a server after idle for  $t_{wait}$

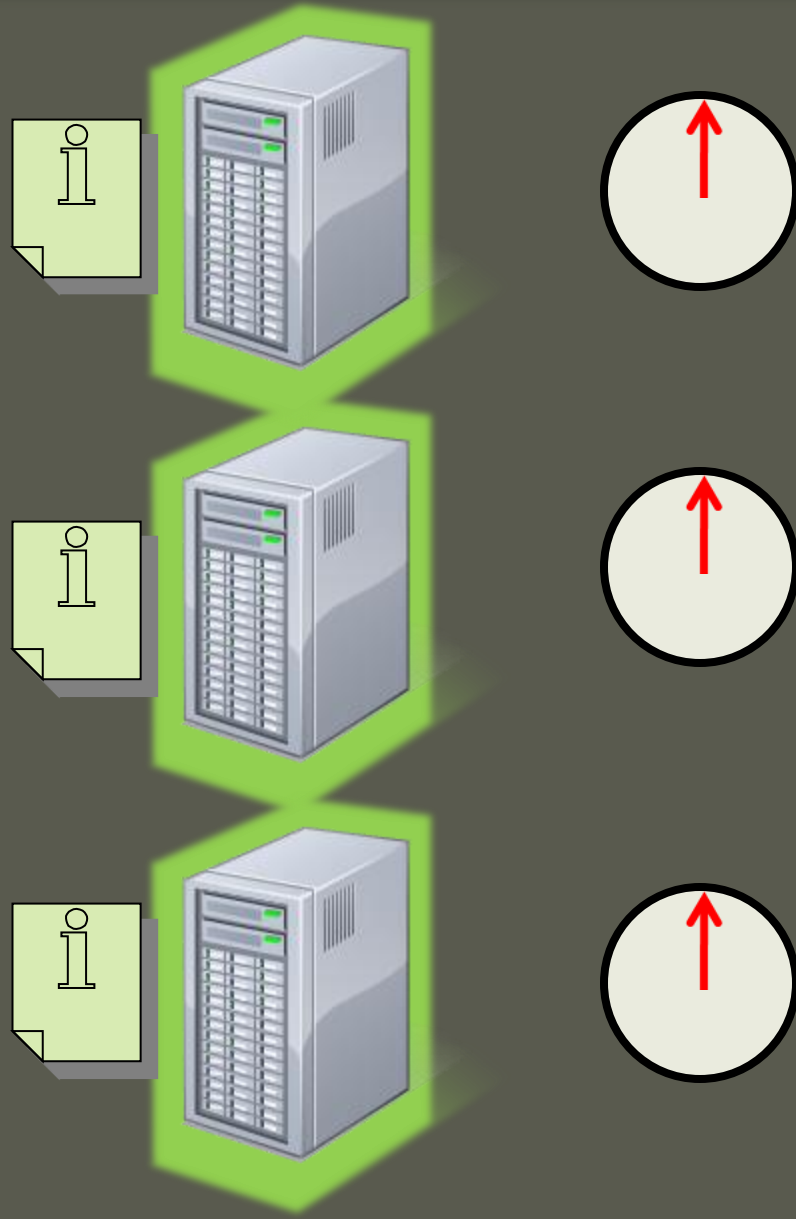
$t_{wait}$  independent of  $\rho(t)$ !

new arrival routed to the  
*most recently busy* (MRB) server

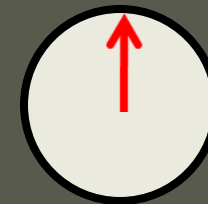
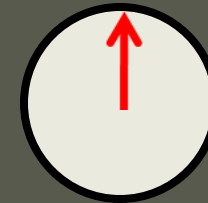
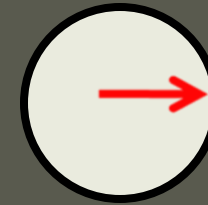
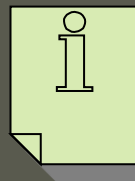
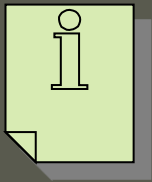


arriving job turns a server ON  
if all servers busy

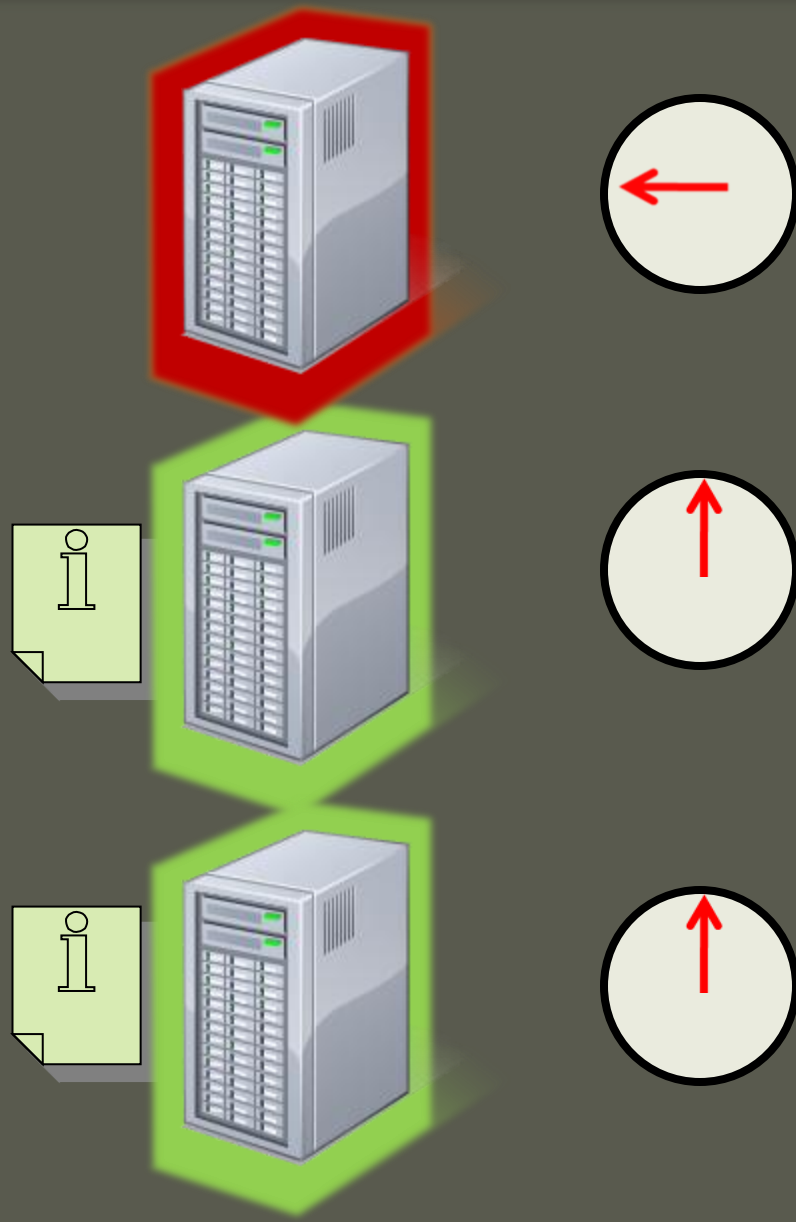




$t_{wait}$  timers



$t_{wait}$  timers

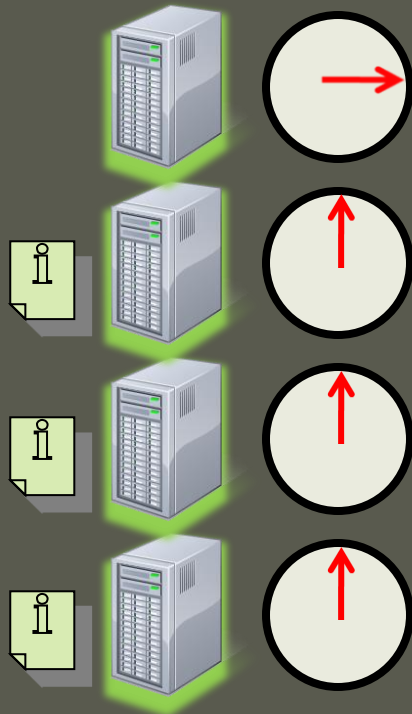


$t_{wait}$  timers

# Our Prescription: DELAYEDOFF

**THEOREM:** For Poisson arrivals, as the load  $\rho \rightarrow \infty$ , the number of ON servers is concentrated around  $\rho + \sqrt{\rho \log \rho}$

--- Proof Intuition ---



Step 1: An equivalent system view

$k$  jobs run on the “first”  $k$  servers

Step 2: Analysis of idle periods of  $M/G/\infty$

Time from a  $k \rightarrow (k-1)$  to the next  
 $(k-1) \rightarrow k$  transition

## Intuition for idle periods



$$\rho + c\sqrt{\rho \log \rho}$$

1. # jobs  $\approx$  Normal with mean and variance  $\rho$

$$2. \Pr[\text{jobs} \geq \rho + c\sqrt{\rho \log \rho}] \propto \frac{1}{\sqrt{\rho}} e^{-\frac{c^2 \log \rho}{2}}$$

$$\propto \frac{1}{\rho^{\frac{c^2+1}{2}}}$$

3. Events happen at rate  $\rho$

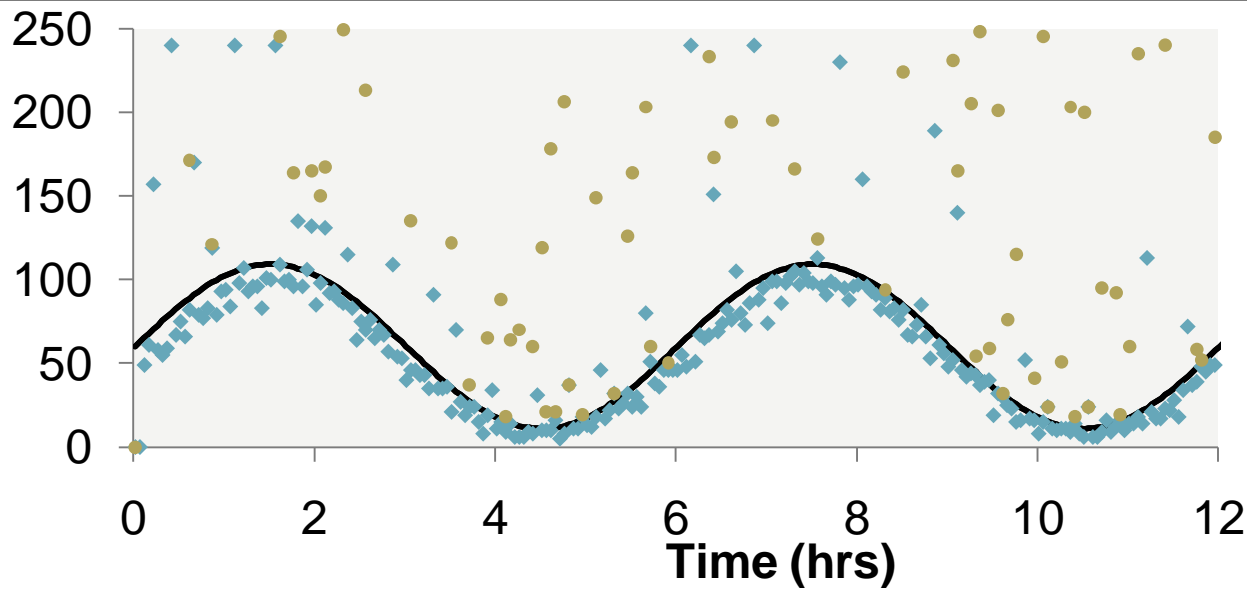
$$4. \text{Mean idle period of } \rho + c\sqrt{\rho \log \rho} \text{ server} \propto \frac{\rho^{\frac{c^2+1}{2}}}{\rho}$$

$$\propto \rho^{\frac{c^2-1}{2}}$$



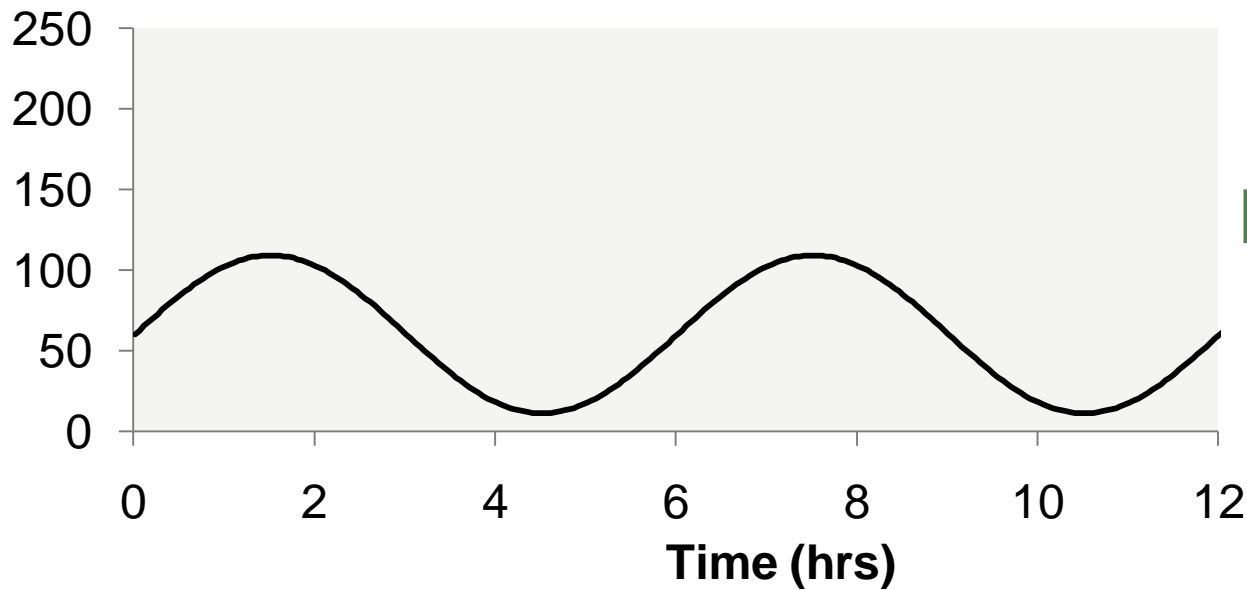
**MRB**  $\Rightarrow \rho + (1 \pm \epsilon)\sqrt{\rho \log \rho}$  servers for any constant  $t_{\text{wait}}$  !

In practice, we choose  $t_{\text{wait}}$  to amortize setup cost



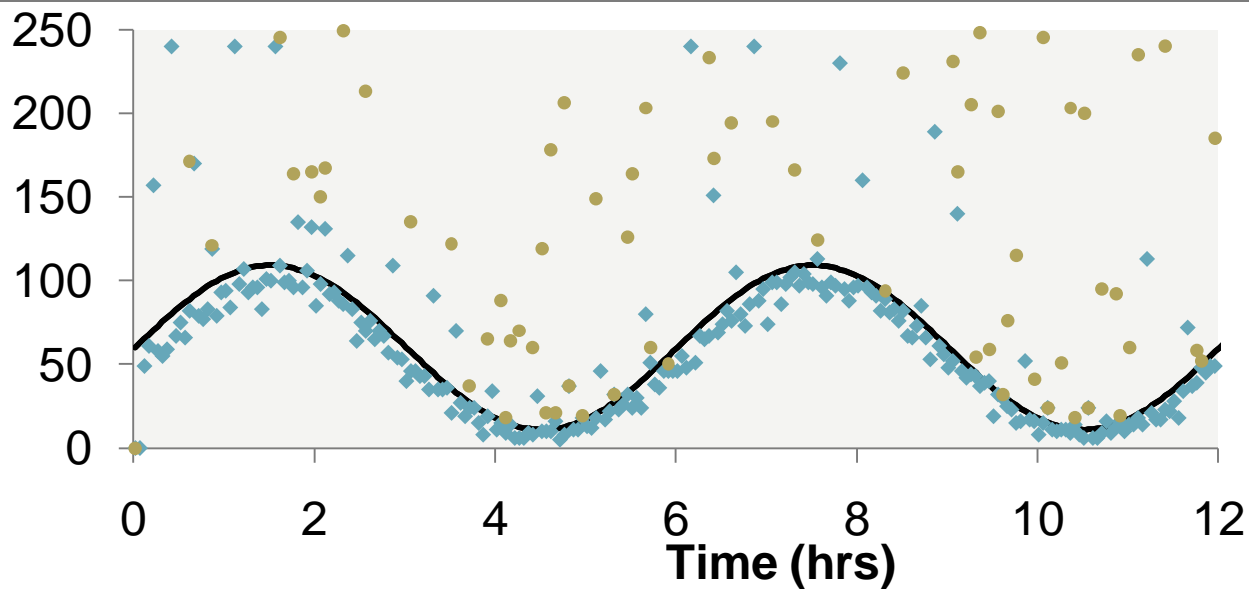
**ON/OFF**

- $\rho(t)$
- ◆ # Busy+idle servers
- # Jobs



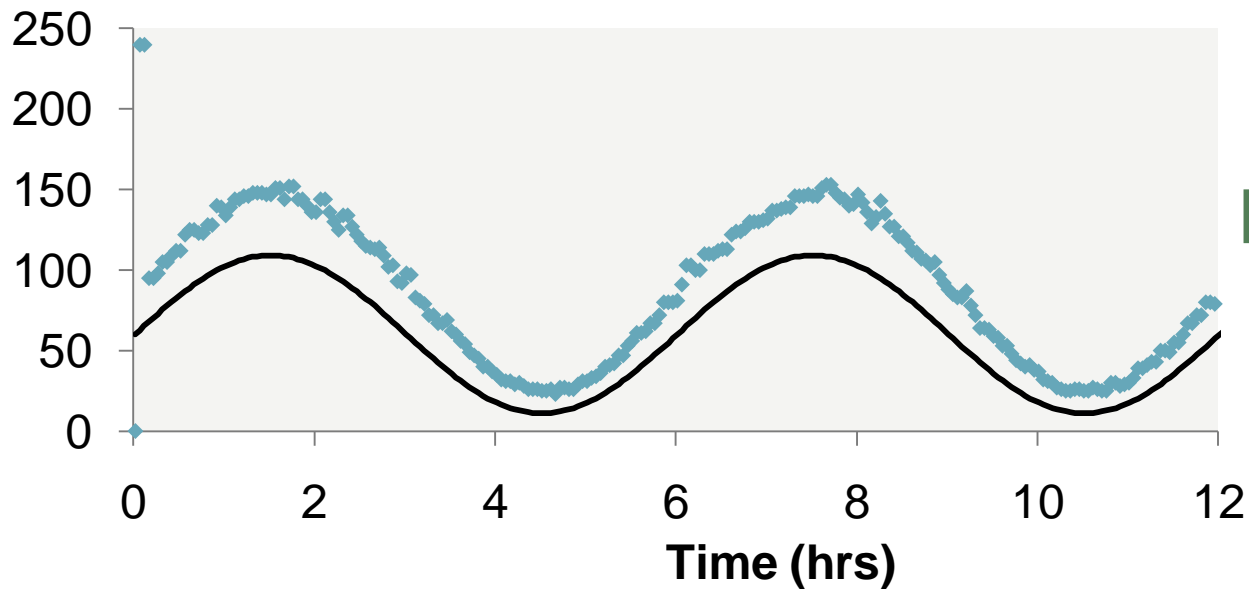
**DELAYEDOFF**

Mean Job Size = 1s  
Setup delay = 200s  
Busy power = 240W



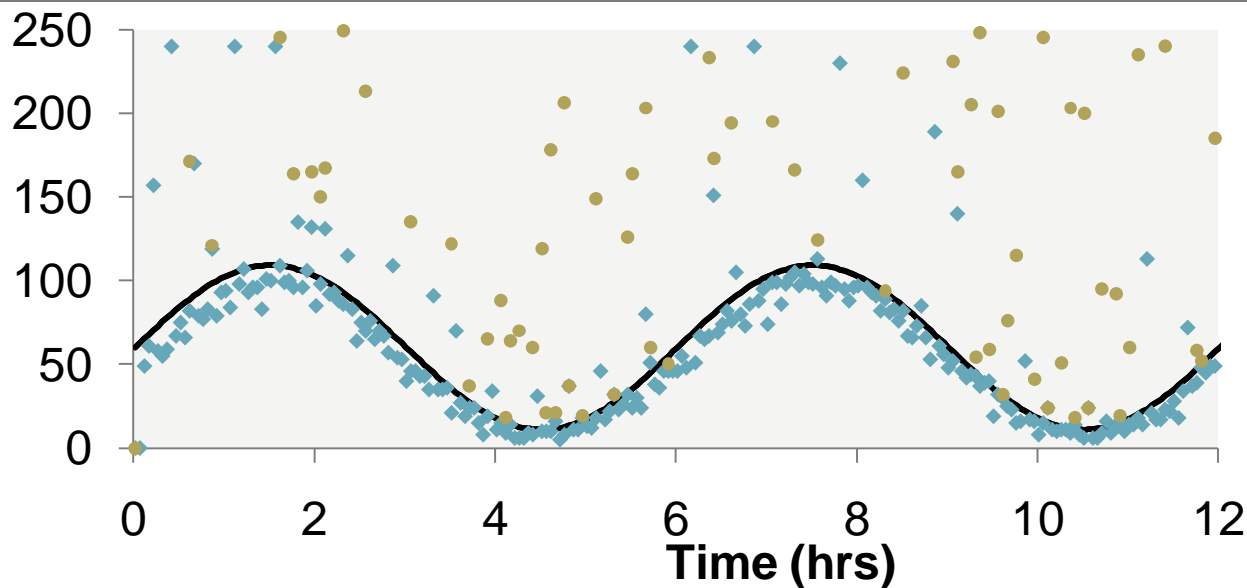
**ON/OFF**

- $\rho(t)$
- ◆ # Busy+idle servers
- # Jobs



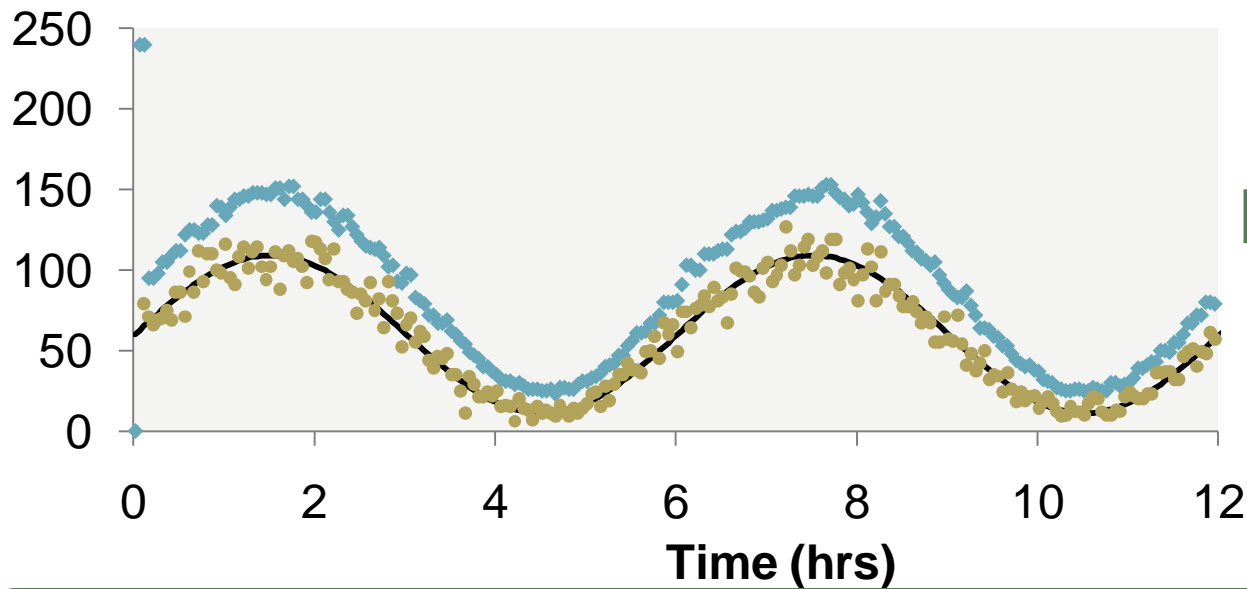
**DELAYEDOFF**

Mean Job Size = 1s  
Setup delay = 200s  
Busy power = 240W



**ON/OFF**

—  $\rho(t)$   
 ◆ # Busy+idle servers  
 ● # Jobs

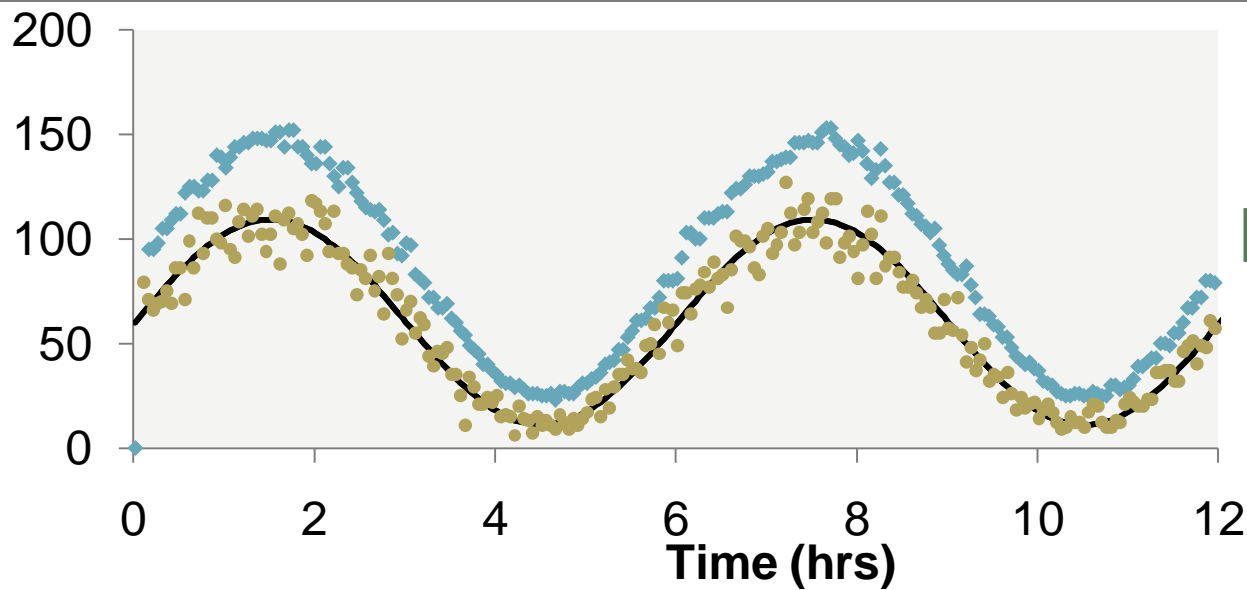


**DELAYEDOFF**

Mean Job Size = 1s  
 Setup delay = 200s  
 Busy power = 240W

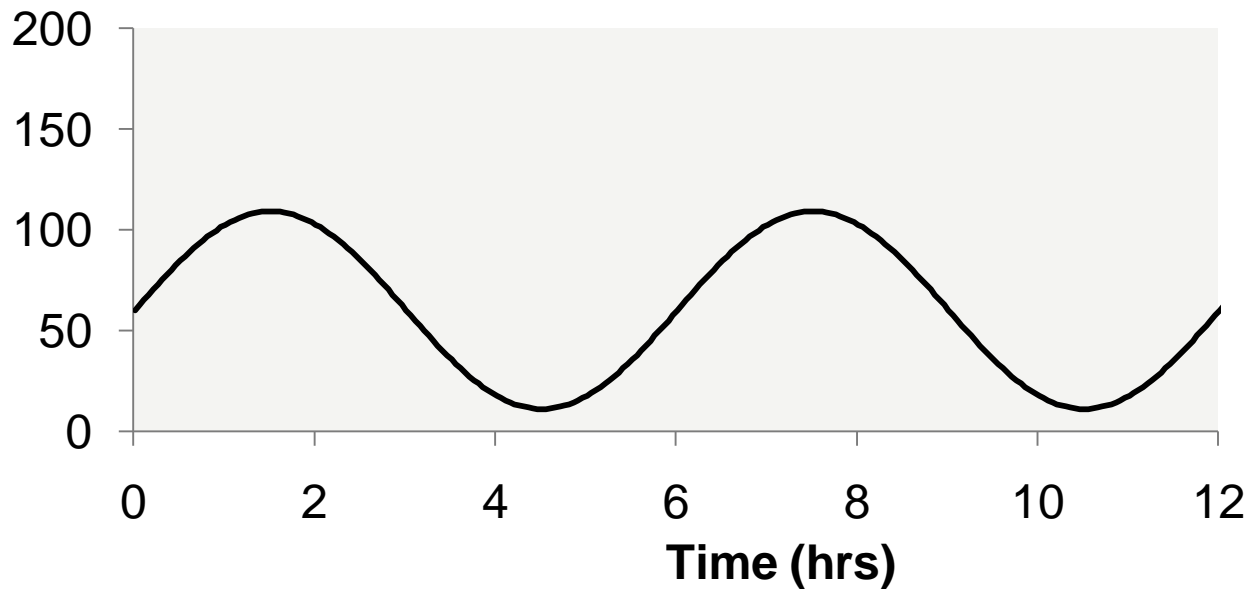
**DELAYEDOFF:**  $E[\text{Power}] = 18.9 \text{ kW}$      $E[\text{Response time}] = 1.002\text{s}$   
**NO SETUP:**  $E[\text{Power}] = 14.4 \text{ kW}$      $E[\text{Response time}] = 1\text{s}$





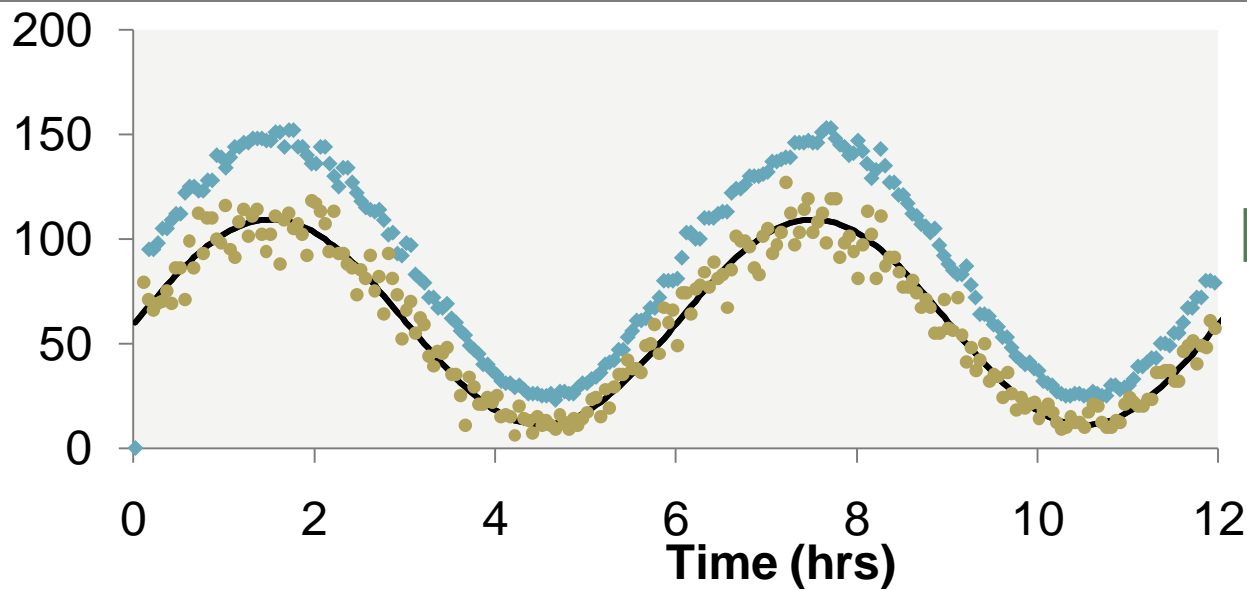
**DELAYEDOFF**

- $\rho(t)$
- ◆ # Busy+idle servers
- # Jobs



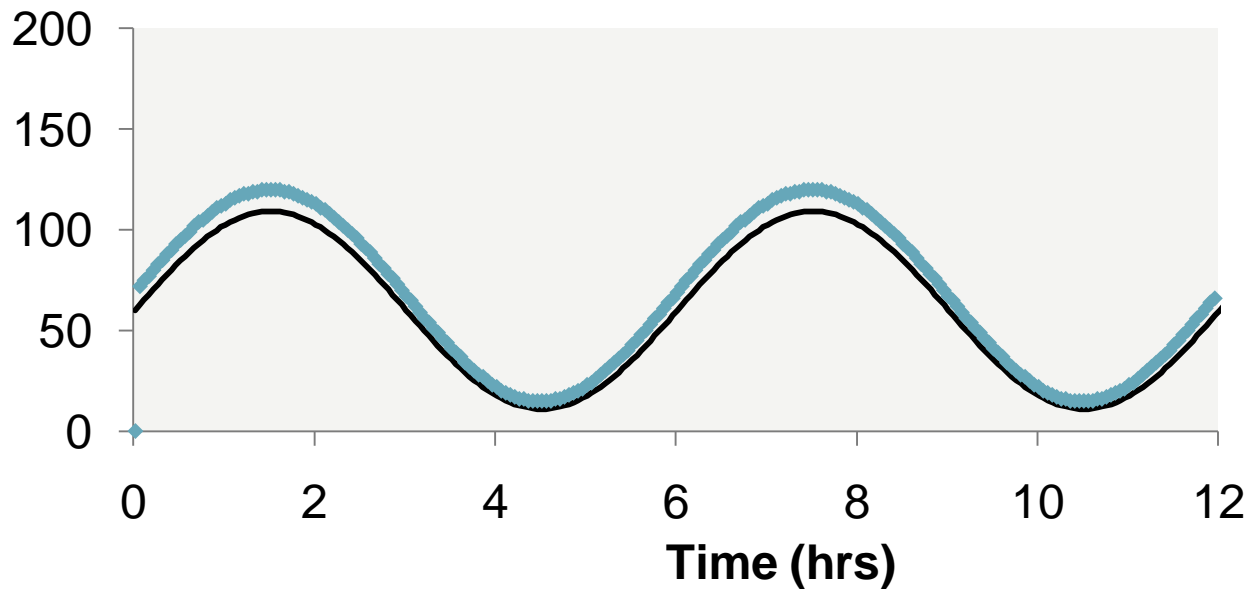
**SQ. ROOT W/  
LOOKAHEAD**  
("optimal offline")

Mean Job Size = 1s  
Setup delay = 200s  
Busy power = 240W



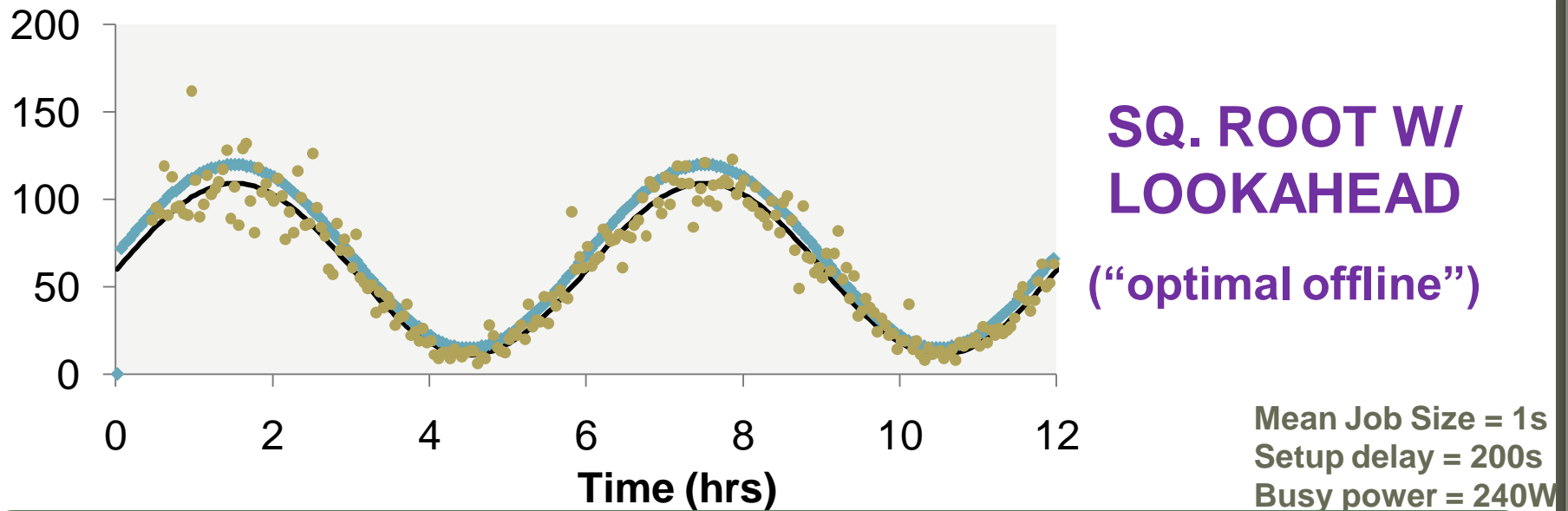
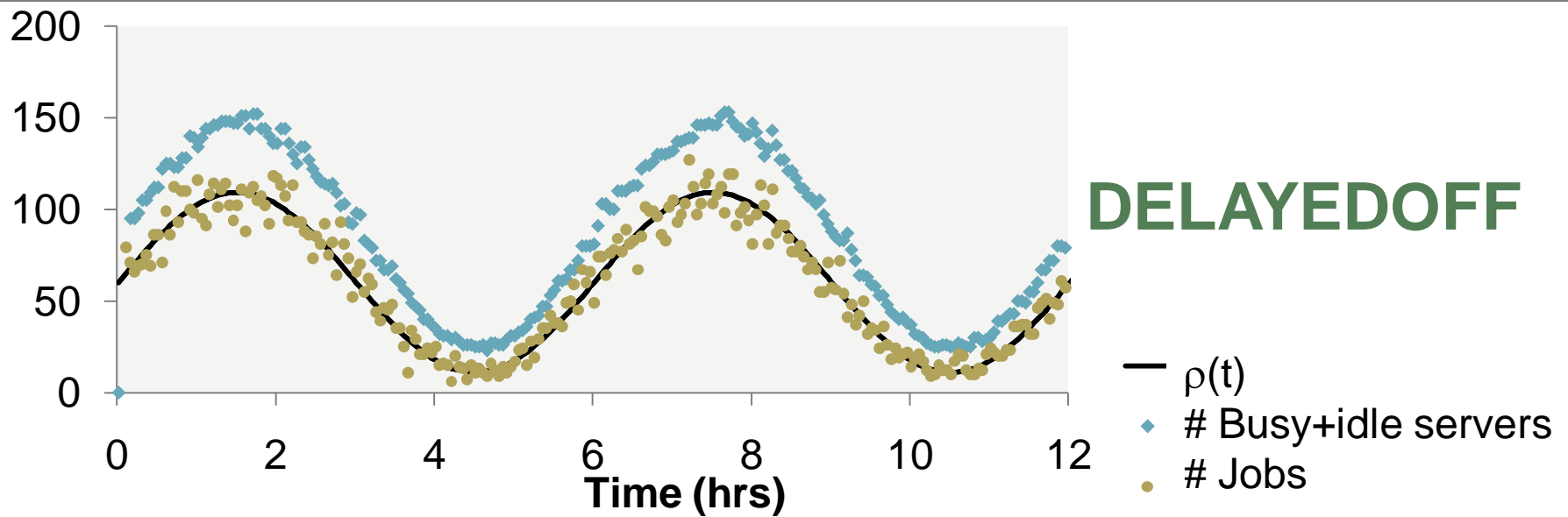
**DELAYEDOFF**

- $\rho(t)$
- ◆ # Busy+idle servers
- # Jobs



**SQ. ROOT W/  
LOOKAHEAD**  
("optimal offline")

Mean Job Size = 1s  
Setup delay = 200s  
Busy power = 240W



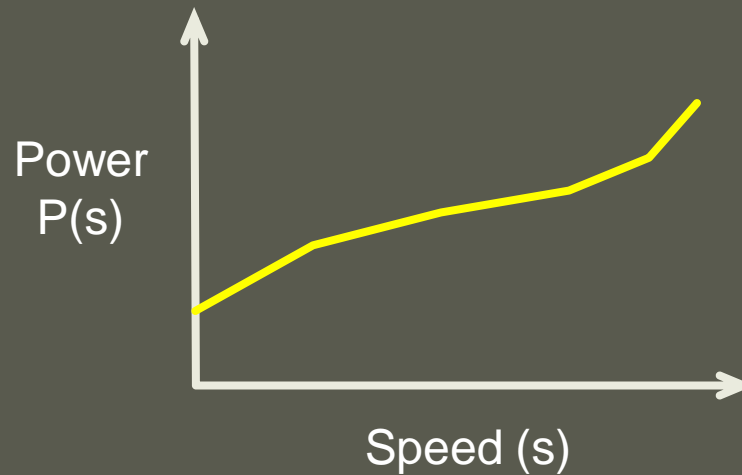
DELAYEDOFF:  $E[\text{Power}] = 18.9 \text{ kW}$      $E[\text{Response time}] = 1.002\text{s}$   
 ✓ w/LOOKAHEAD:  $E[\text{Power}] = 15.8 \text{ kW}$      $E[\text{Response time}] = 1.036\text{s}$

## DELAYEDOFF mods

---

1. Speed scaling algorithms
2. A simple proxy for MRB routing
3. Heterogeneous servers
4. Managing Virtual Machines in the Cloud
5. Wear-leveling/Performance tradeoffs

# Effect of DELAYEDOFF on speed-scaling



**Q:** Optimal speed to balance energy-performance?

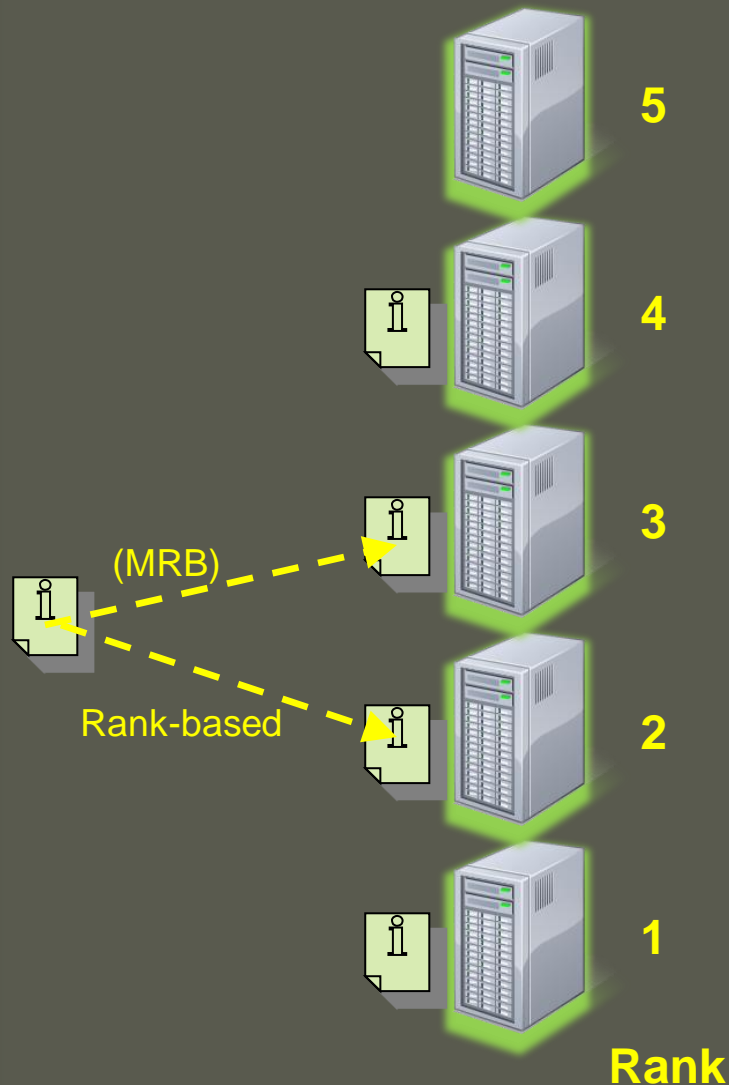
**A:** your favorite metric =  $\mathbf{F}(\mathbf{E}[\text{energy/job}], \mathbf{E}[\text{response time}])$

$$\mathbf{E}[\text{energy/job}] \xrightarrow{MRB} \frac{P(s)}{s}$$

$$\mathbf{E}[\text{response time}] \xrightarrow{MRB} \frac{1}{s}$$

$$s^* = \operatorname{argmin}_s \mathbf{F}\left(\frac{P(s)}{s}, \frac{1}{s}\right)$$

# A simple proxy for MRB routing



MRB requires a lot of state updates



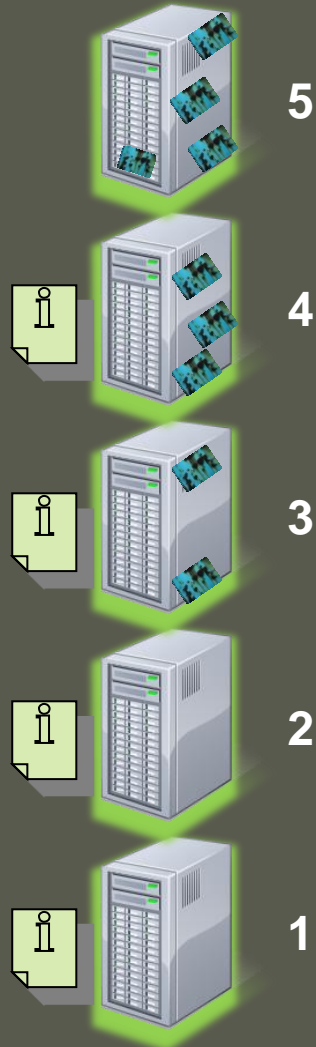
Proxy policy

- Assign static ranks to servers
- Route a new arrival to *highest ranked idle* server

Almost the same performance as MRB  
+ easier to implement than MRB  
+ easy to extend

# DELAYEDOFF for heterogeneous servers

less efficient

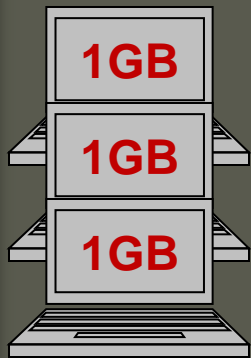


Assign ranks to servers based on efficiency (1 = most efficient)

Route a new arrival to *highest ranked idle* server

more efficient

# Managing Virtual Machines in the Cloud



**2GB RAM**



**2GB RAM**



**2GB RAM**



# Managing Virtual Machines in the Cloud



Split physical servers into “virtual” servers



Assign static ranks to virtual servers

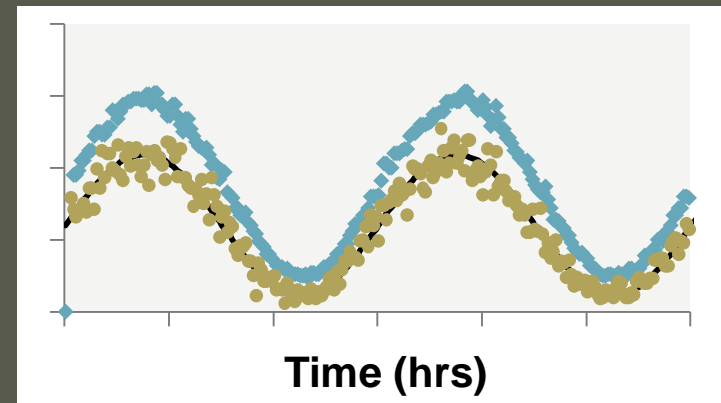
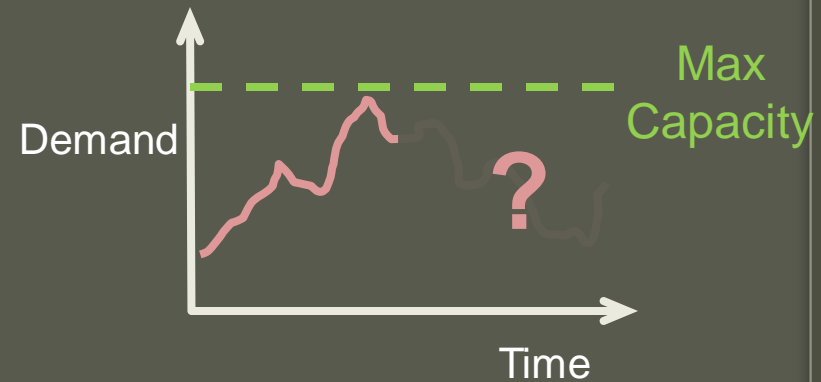


Route a new VM request to *highest ranked idle* virtual server

Turn the physical server OFF after each virtual server has idled for  $t_{wait}$

# Conclusions

- Problem: unpredictable demand and non-trivial setup costs
- DELAYEDOFF : A new traffic-oblivious capacity scaling scheme
- Extensions to real-world scenarios



# The Importance of Being MRB

**THEOREM [MRB]:** As the load  $\rho \rightarrow \infty$ , the number of ON servers is concentrated around  $\rho + \sqrt{\rho \log \rho}$

**THEOREM [Round-Robin]:** As the load  $\rho \rightarrow \infty$ , for constant job sizes, the number of ON servers is  $\rho \left(1 + \frac{t_{wait}}{\text{job size}}\right)$

