



| IBM Research

Coordinating multiple managers to achieve specified power-performance tradeoffs

IBM Watson: [Jeff Kephart](#), [Hoi Chan](#), [Raja Das](#),
[David Levine](#), [Gerry Tesauro](#)

IBM Austin: [Charles Lefurgy](#), [Freeman Rawson](#)

June 13, 2007

© 2007 IBM Corporation

Team Members



Charles Lefurgy



Freeman Rawson



David Levine



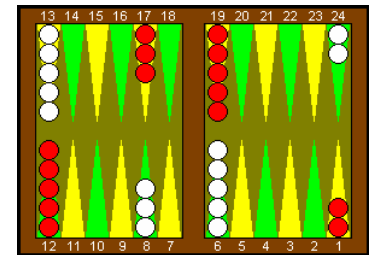
Jeff Kephart



Hoi Chan



Raja Das



Gerry Tesauro

Austin Lab

Watson Lab

Background

- The costs and constraints associated with electrical power are becoming an increasing concern for data center operation
 - 50% of data centers will have insufficient power and cooling capabilities by 2008 (Gartner)
 - Power will be the second-highest operating cost (after labor) in 70% of data centers
 - According to Berkeley RAD lab, 30% reduction in power would save \$15B and 100M metric tons of CO₂ emissions per year in US (1.7% of total emissions)

- Power efficiency is becoming the subject of energy and environmental regulations by governments around the world
 - EPA report to Congress due in mid-2007; Energy Star standards for systems, data centers
 - Other agencies such as DoE and other governments involved as well

- Non-governmental organizations & industry groups actively studying the problem
 - Green Grid: consortium of IT companies developing best practices for reducing power consumption in data centers
 - AMD, Dell, HP, IBM, Sun Microsystems, Microsoft, VMWare, ...
 - SPEC working on a power/performance benchmark

- Industry is aggressively developing and marketing power-conserving hardware and software solutions

Agenda

- **Background**
- **Power-Performance Research**
 - Algorithms
 - Results
- **Commercialization**

Universal algorithms for managing power and performance

- **Universally Optimal Power Management Algorithm**

Algorithm MaxPowerSavings ()

For each object O in DataCenter

TurnOff (O)

- **Universal Performance Management Algorithm**

Algorithm MaxPerformance ()

For each object O in DataCenter

TurnOnCompletely (O)

PerformanceManage (O)

Now all we need to do
is combine these
algorithms
somehow...?!?

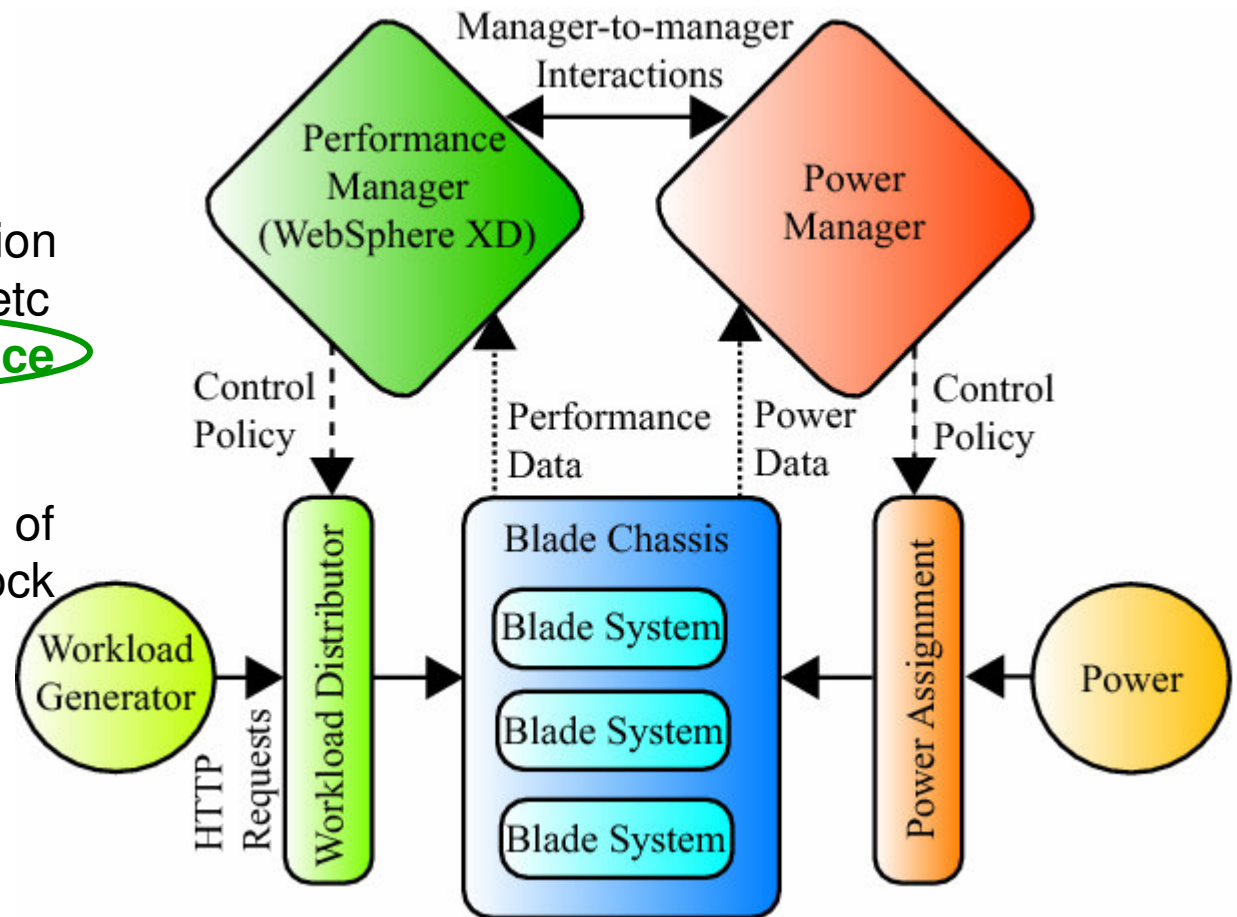
Power and Performance Management

■ Performance agent

- IBM WebSphere middleware adjusts load balancing, CPU application placement parameters, etc to **maximize performance**

■ Power agent

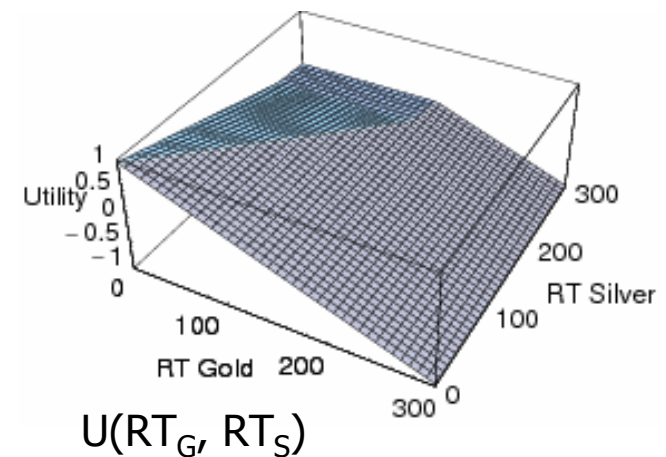
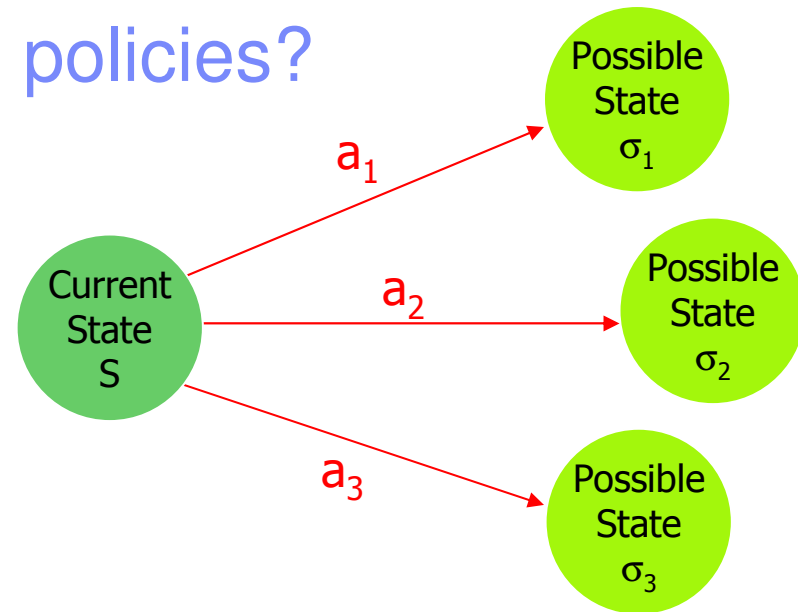
- Autonomic Management of Energy throttles CPU clock to slow down processor and **save power**



An inherent conflict!

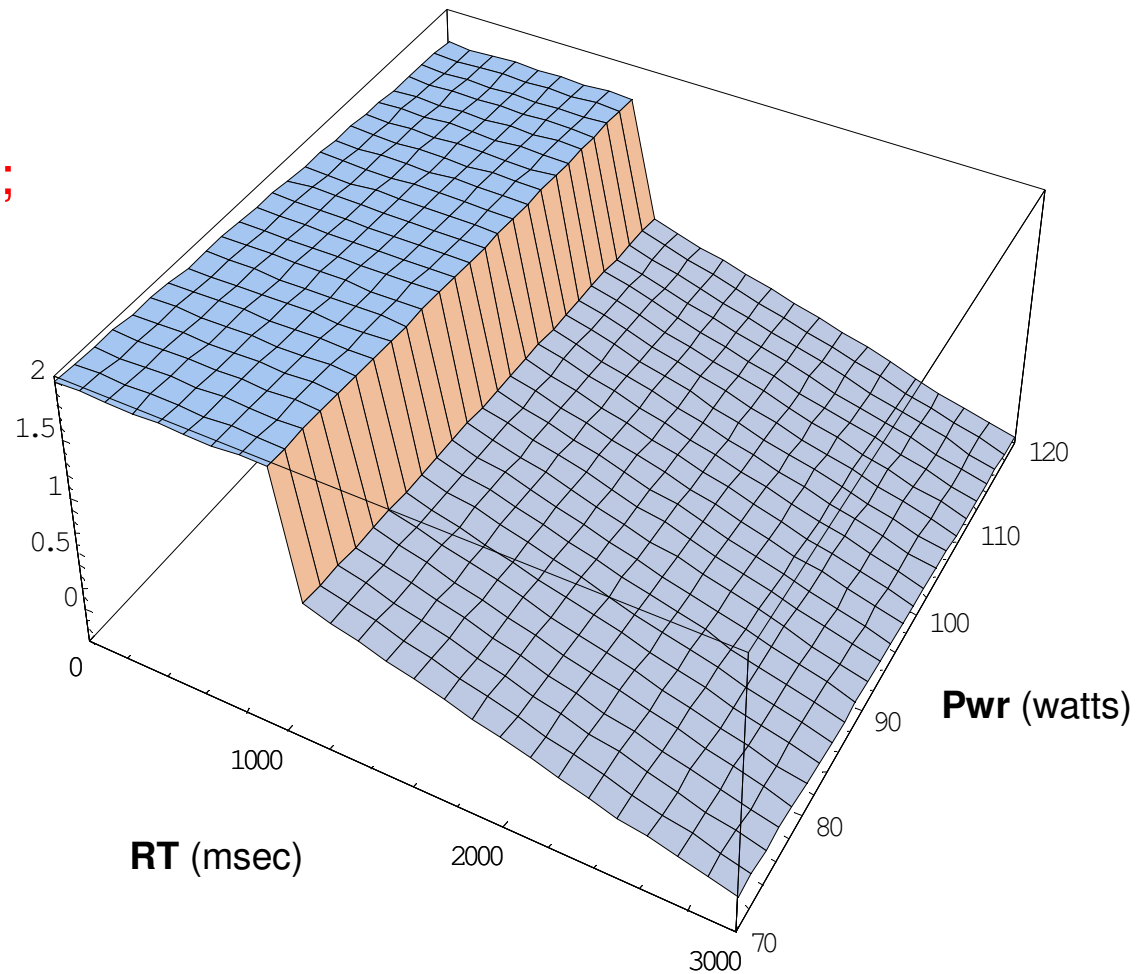
How to represent high-level policies?

- **Utility functions** map any possible state of a system to a scalar value
- They can be derived from
 - a service level agreement
 - preference elicitation techniques
 - simple templates, e.g. specify response time thresholds and “importance” levels
- They are a generally useful representation for high-level objectives, e.g.
 - Minimize power while meeting SLA
 - Maximize performance while meeting power constraint
 - Range in between

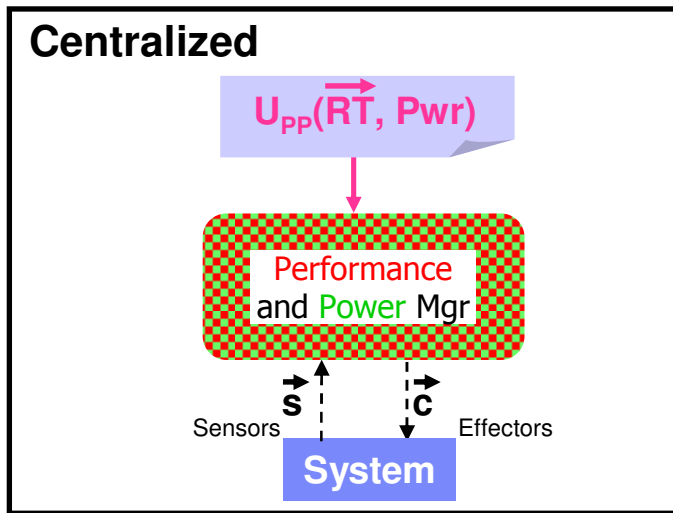


Power-performance utility functions

- $U(\text{perf}, \text{pwr}) = U(\text{perf}) - \varepsilon \text{Pwr}$;
 $\text{Pwr} < \text{Pwr}_{\text{Max}}$
- $U(\text{perf}, \text{pwr}) = U(\text{perf})/\text{Pwr}$;
 $\text{Pwr} < \text{Pwr}_{\text{Max}}$



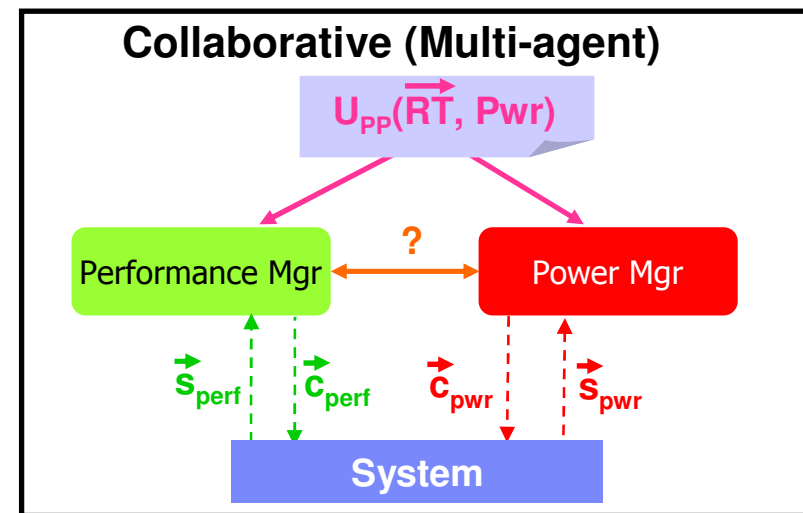
Multiagent approach to Power-Performance optimization



Set effectors c to maximize U_{PP}

Conceptually easiest

Not very practical!



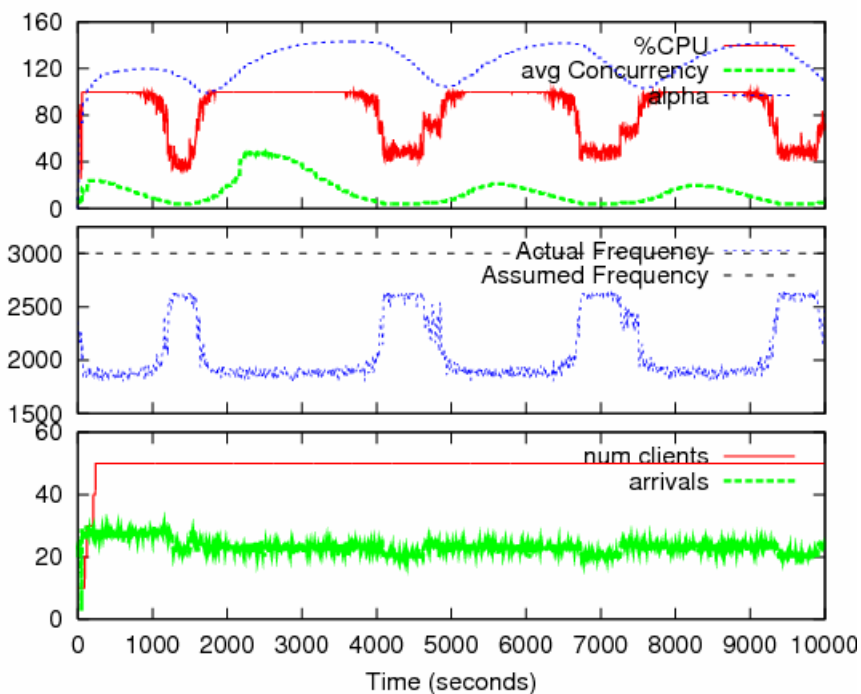
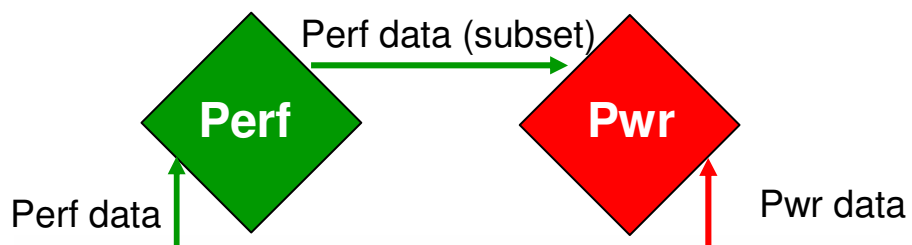
What info should be exchanged, and how?

- Do we need negotiation?
- Do we need mediation?
- What are the right power control knobs?

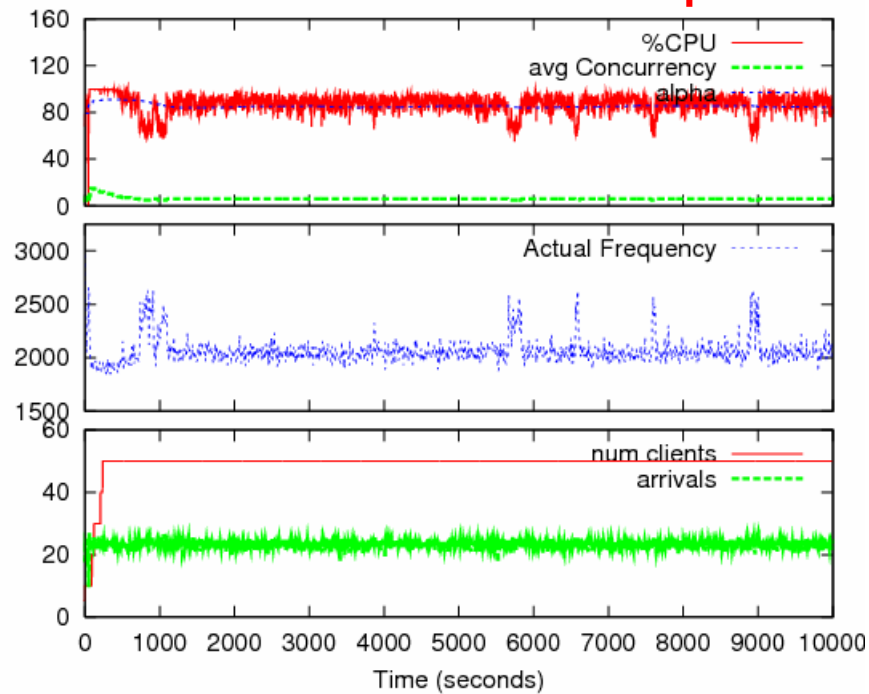
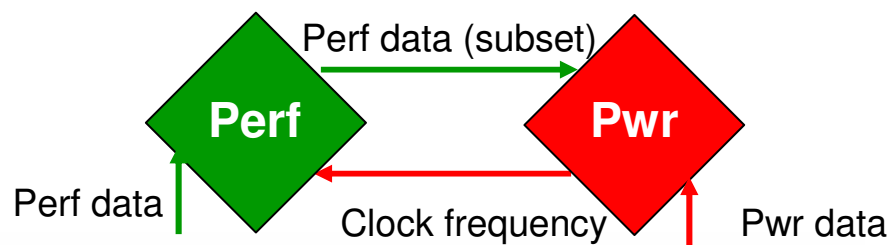
Strategy: Try the simplest method first

- No negotiation or mediation
- Power control knobs = power cap settings
- Minimize changes to WXD
- Add complexity only if/when really needed

What needs to be communicated?

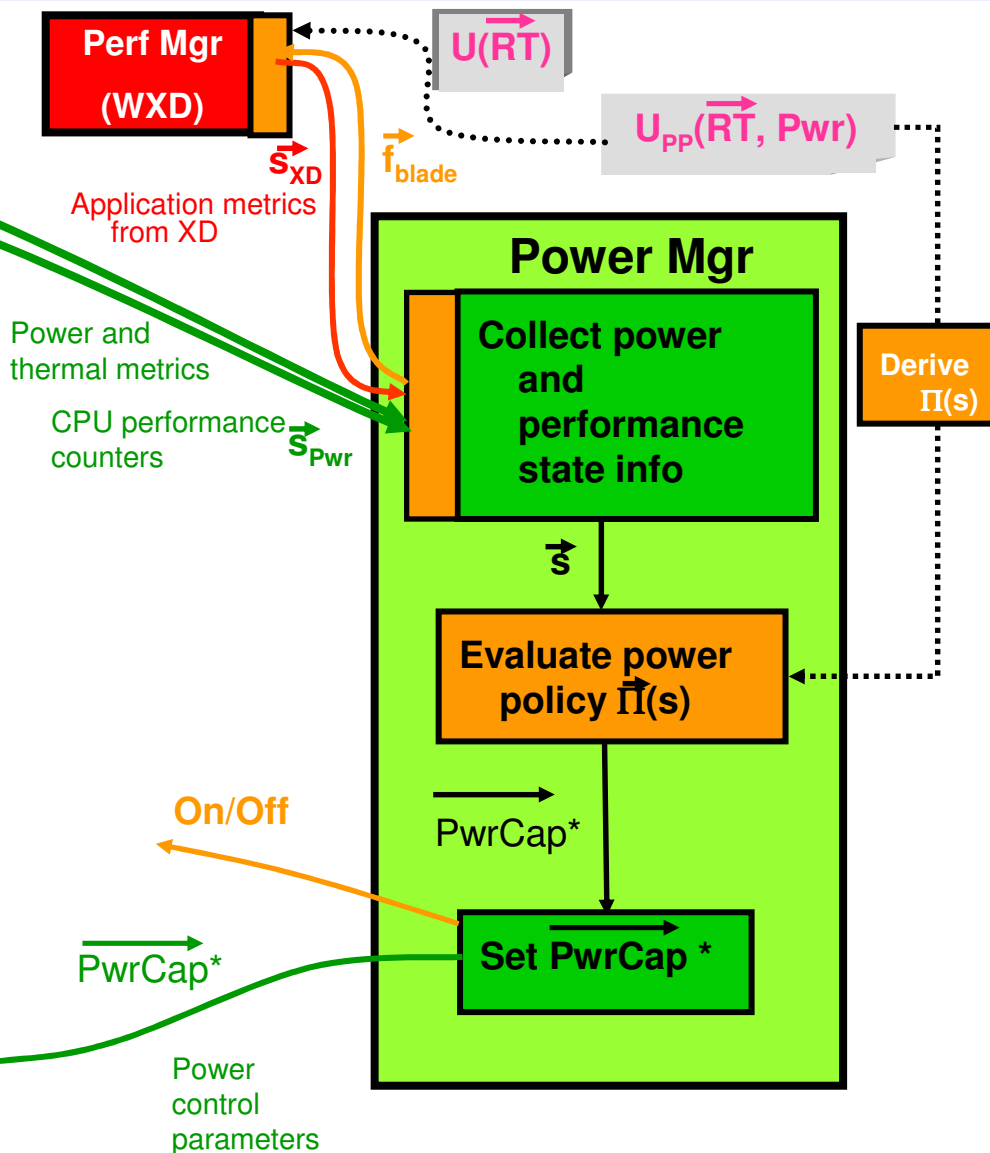
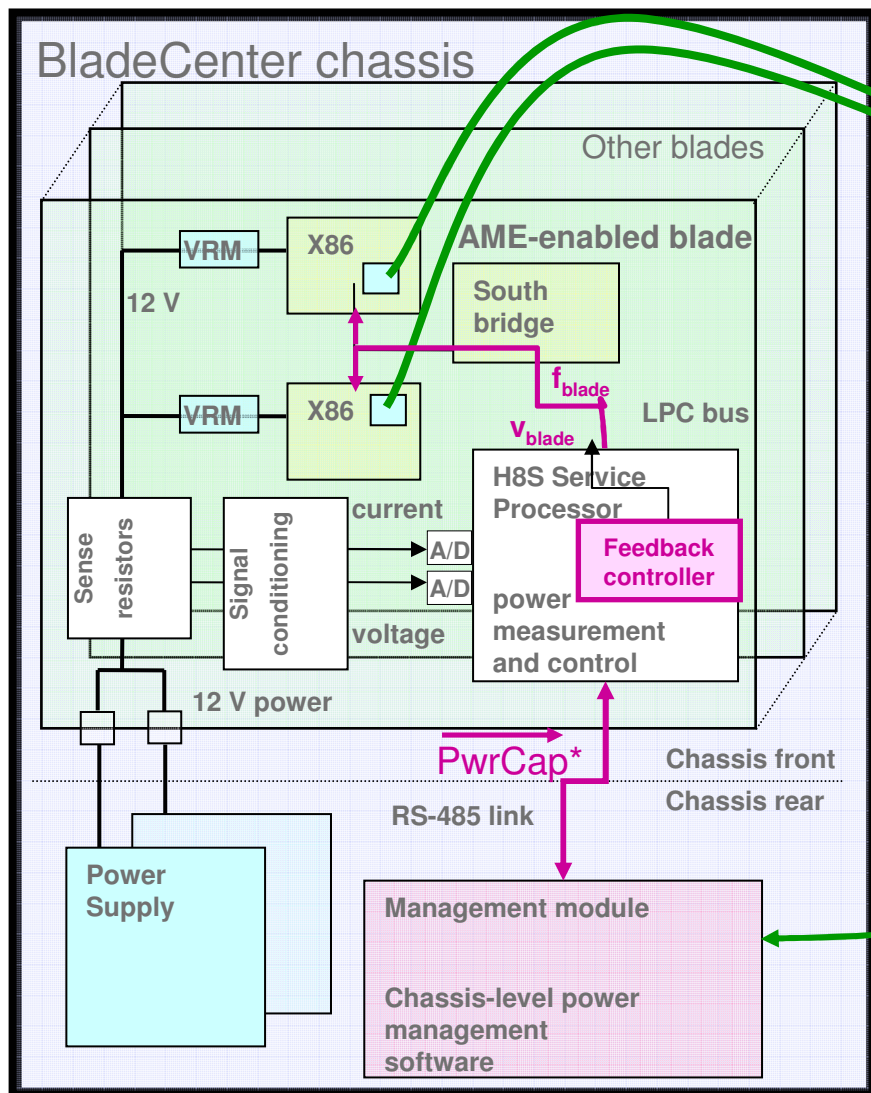


No frequency feedback



With frequency feedback

Experimental setup



BladeCenter product with added software agents

Deriving a Powercap Policy

- Measure offline how performance and power use depend upon system state s and power cap P_{cap} for each blade
 - ~1 week of data collection

- Use interpolation or regression to generate models

- RT(s , P_{cap})
- Pwr(s , P_{cap})

Models

- Substitute models into utility function

$$U_{pp}(RT(s, P_{cap}), Pwr(s, P_{cap})) = U'(s, P_{cap})$$

- For given state s , determine P_{cap} that maximizes U'

$$P_{cap}^*(s) = \operatorname{argmax}_{P_{cap}} U'(s, P_{cap})$$

Optimization

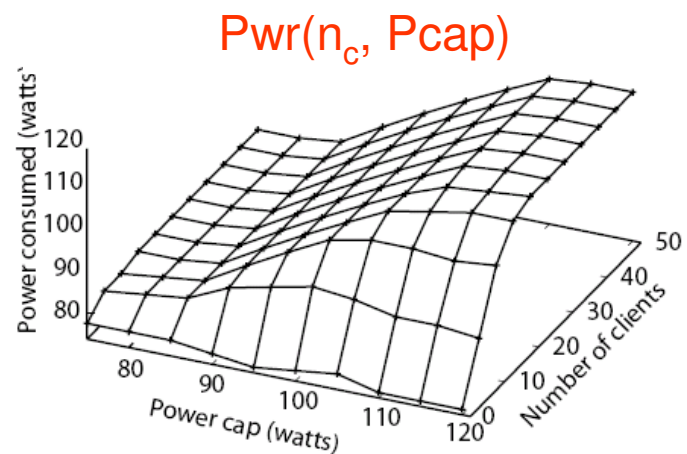
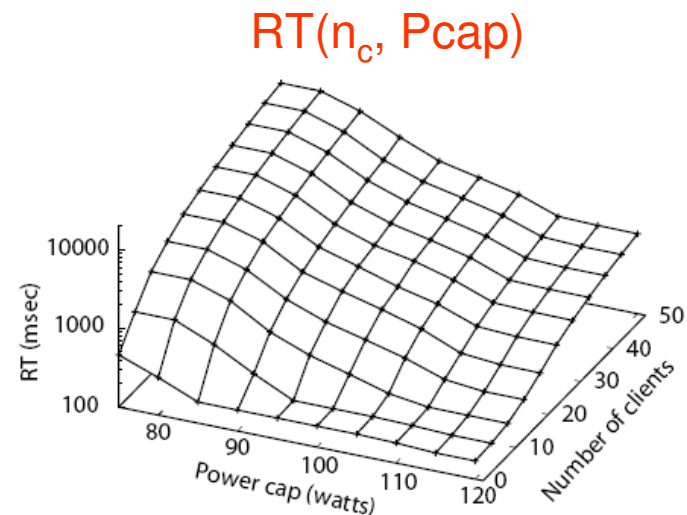
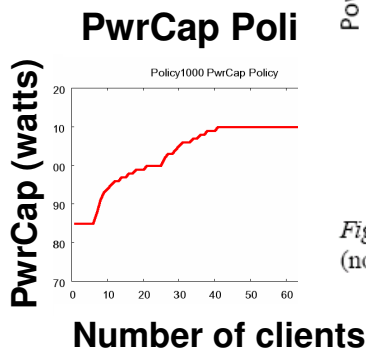
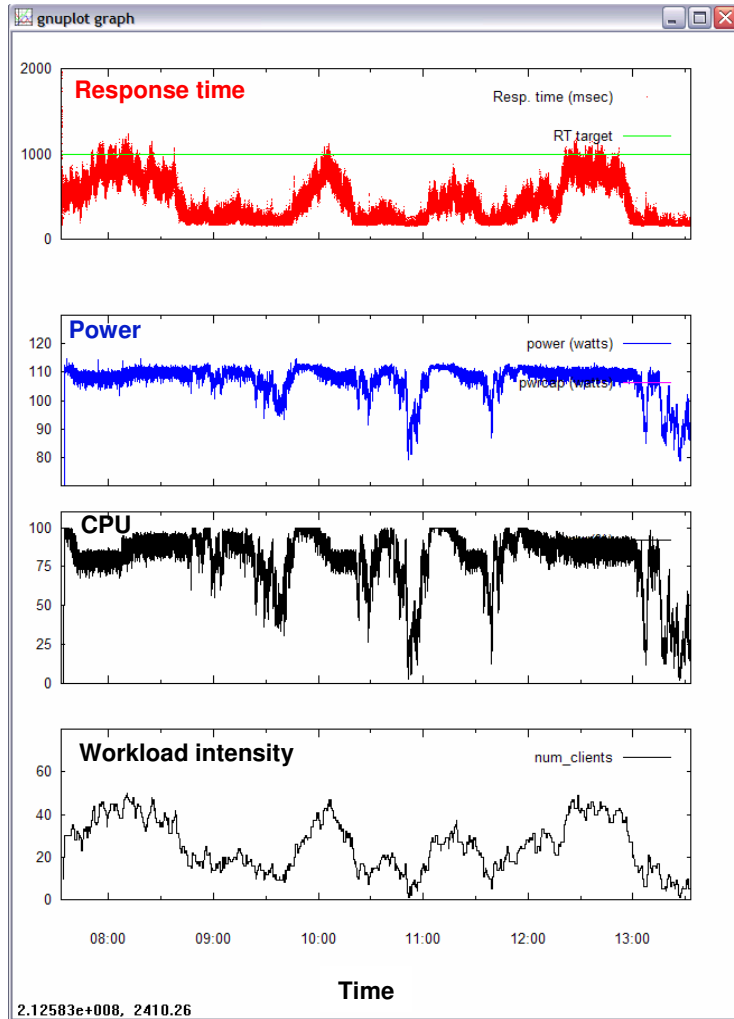


Figure 4: a) Experimentally-determined models of a) $RT(p_{\kappa}, n_c)$ (note log scale) and b) $Pwr(p_{\kappa}, n_c)$ for GB.

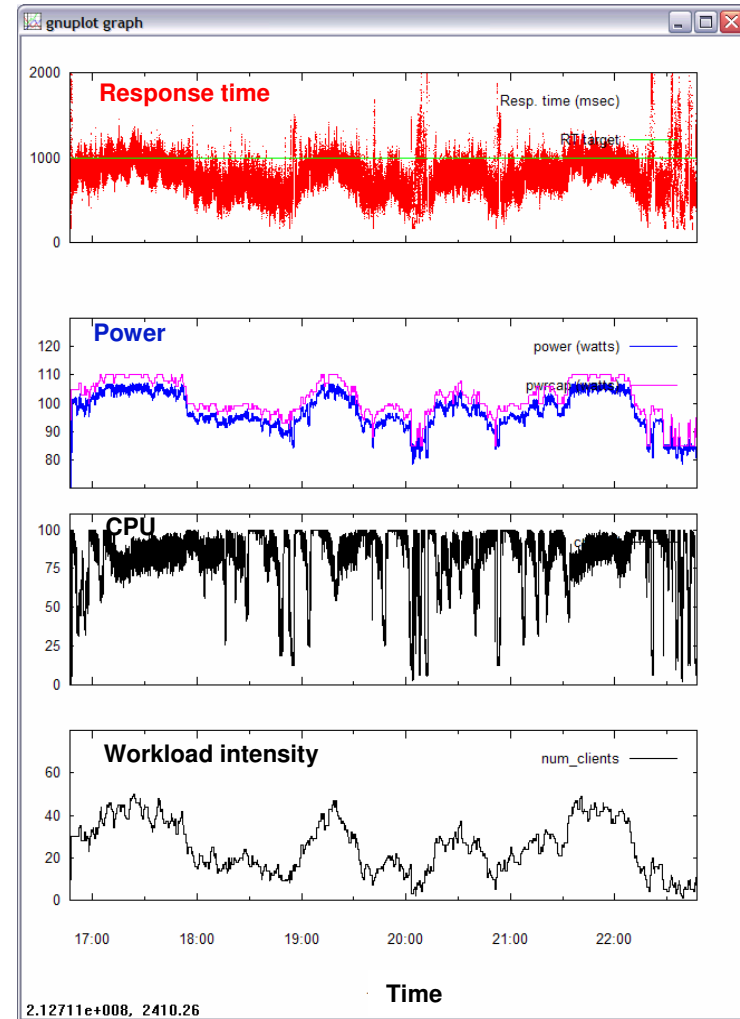
Experiment: Hand-crafted Power Policy

No power management



Avg power = 107.9 watts

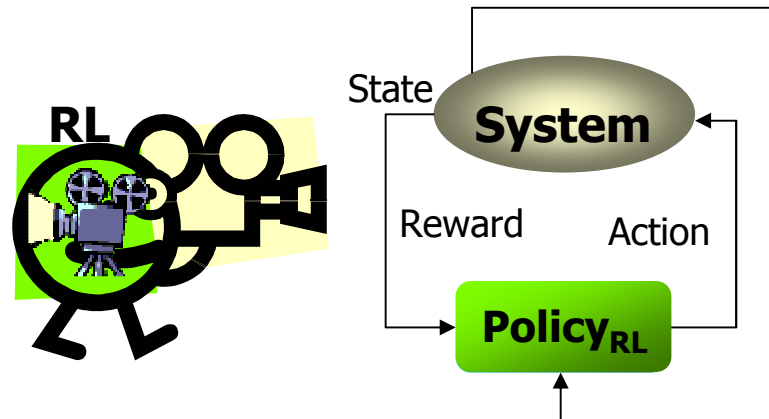
Power management, using Policy HC01



Avg power = 96.6 watts (savings: 11.3 watts = 10.5%)

Hybrid Reinforcement Learning

- Reinforcement learning methods learn to make good decisions by
 - observing $\langle \text{state}, \text{action}, \text{reward} \rangle$ tuples
 - learning long-range value functions $V(\text{state}, \text{action})$
 - Abiding by optimal policy: when in state s , take action a that maximizes $V(s, a)$
- Typically, they learn by updating $V(\text{state}, \text{action})$ starting from random assumptions
- This can take a long time, and performance can be very poor during the learning phase
- We invented a new RL technique, Hybrid RL, that starts from an existing policy, and improves upon it
- Very general method that automatically improves **any** existing systems management policy
 - No knowledge engineering needed



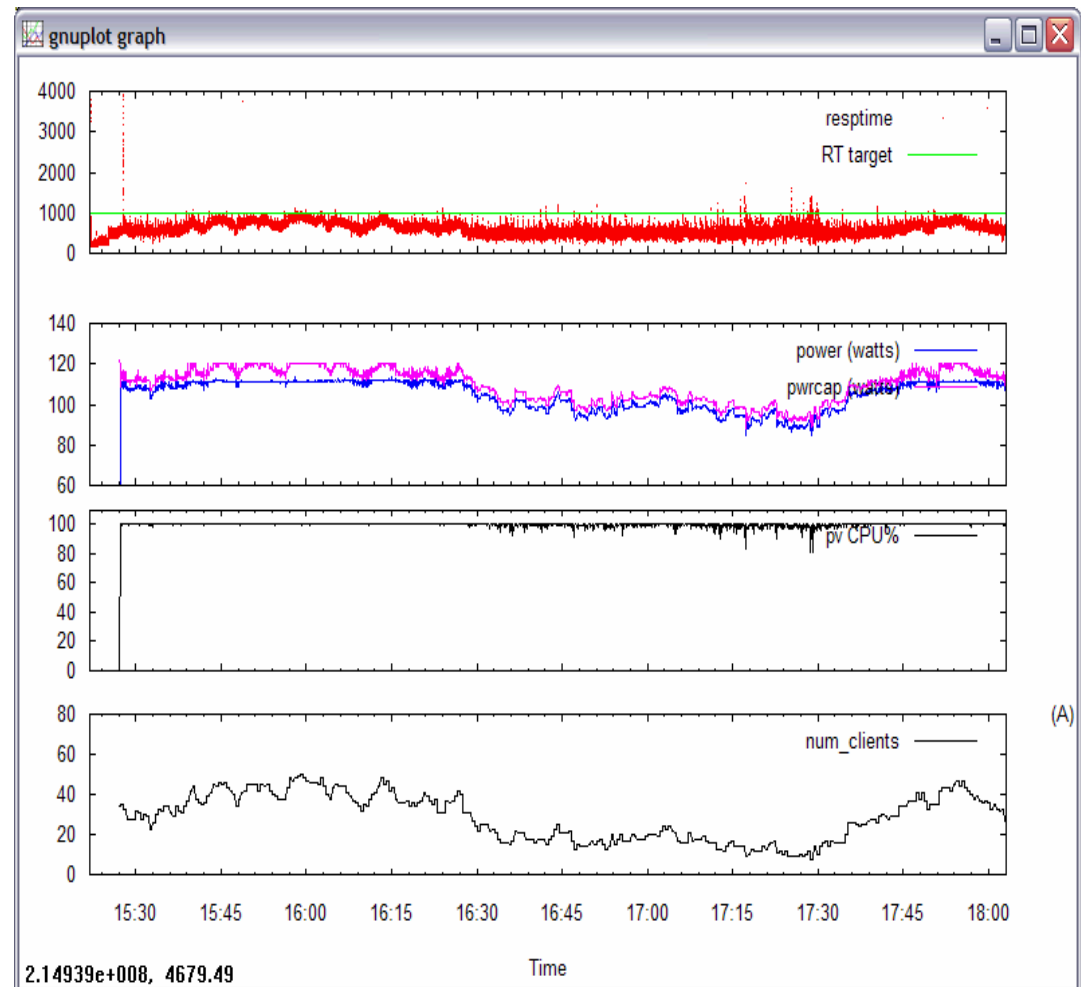
State = {power, performance metrics}

Action = {powercaps}

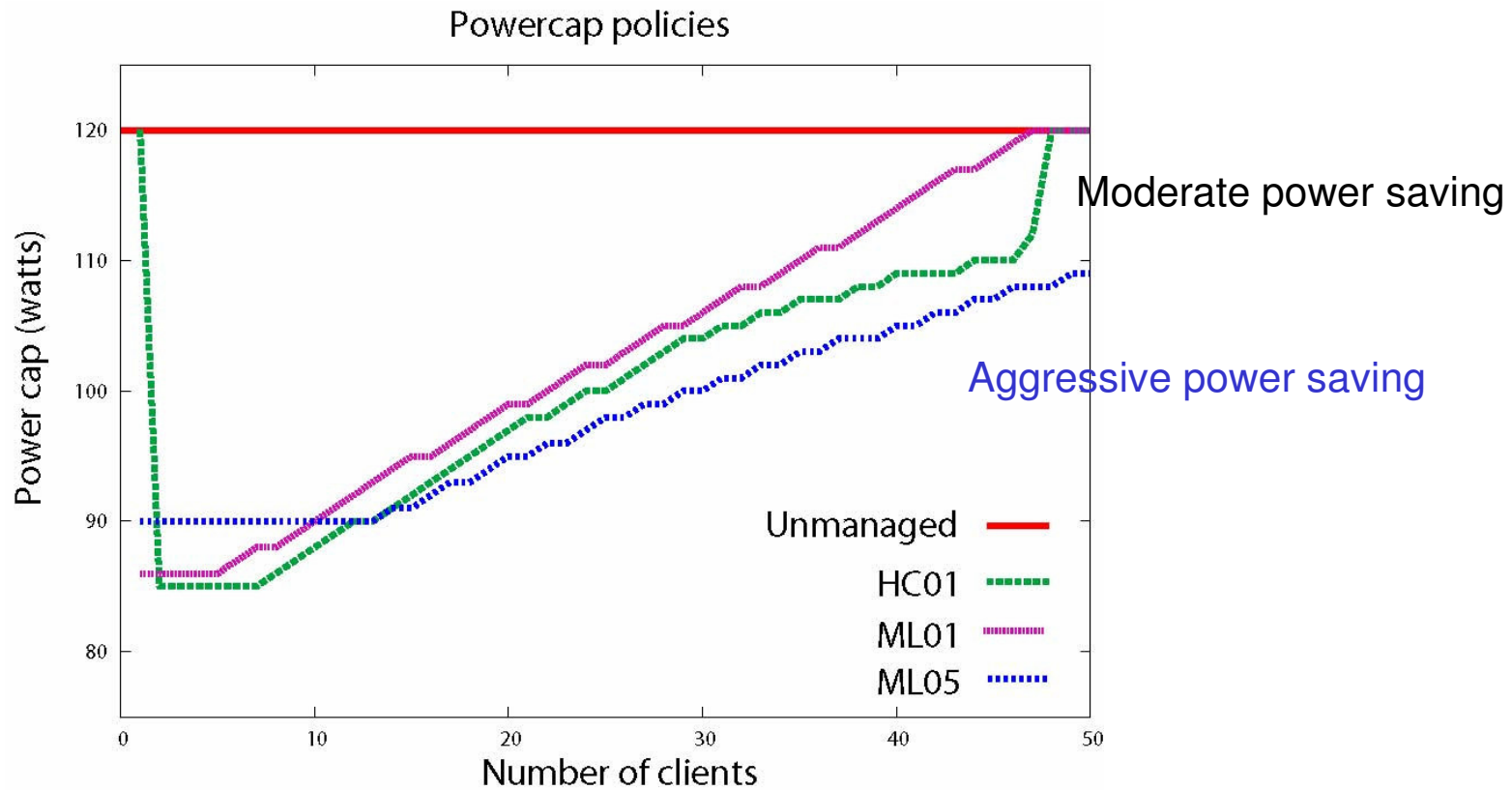
Reward = $U(\text{perf}, \text{pwr})$

Hybrid RL improves our initial power policy!

- **Good power savings**
 - 8.9% less power than for unmanaged case
 - Was 10.5% for hand-tuned policy
- **Reduced SLA violations**
 - 1.5% of response times exceed threshold
 - Was 21% for hand-tuned policy



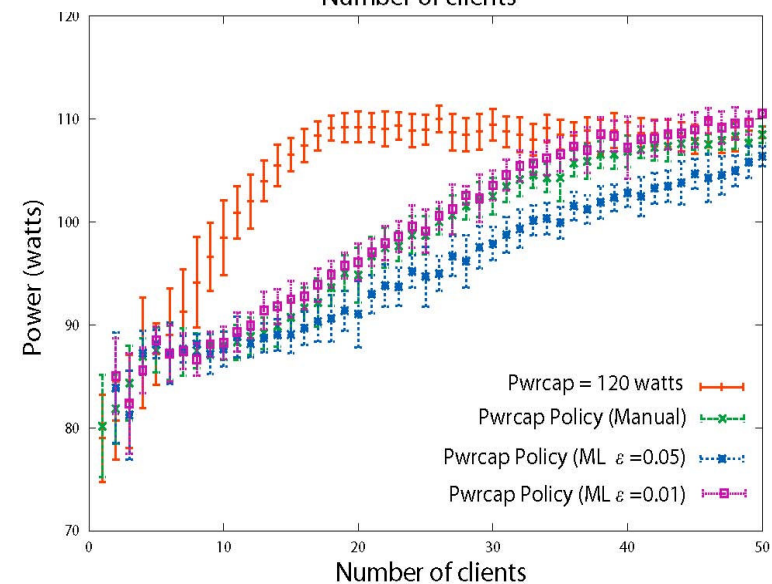
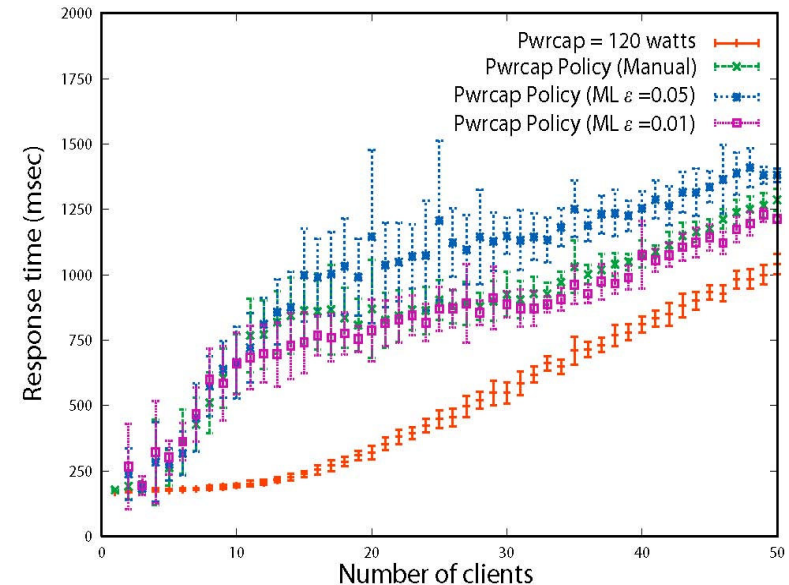
Powercap policies



Where is the power being saved?

- **Power is saved primarily when the number of clients is moderate**
 - For low workloads, power consumption is not constrained by powercaps
 - At high workloads, utility is maximized by setting high powercaps

Policy	Avg watts	Savings
Unmanaged	104.9	0%
HC01	95.3	9.2%
ML01	96.1	8.4%
ML05	92.7	11.6%

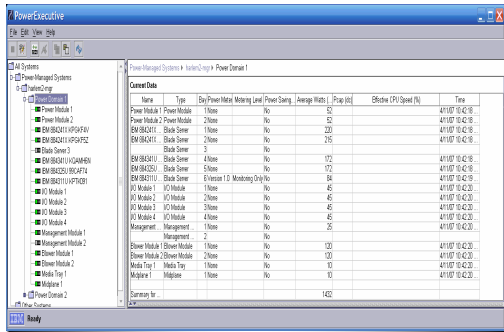


Agenda

- **Introduction**
- **Power-Performance Research**
 - Algorithms
 - Results
- **Commercialization**

ITM Power Agent

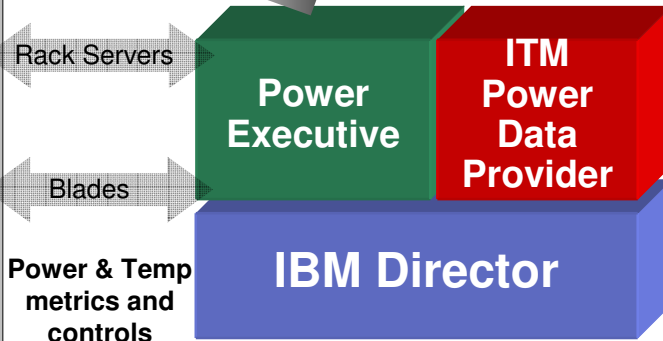
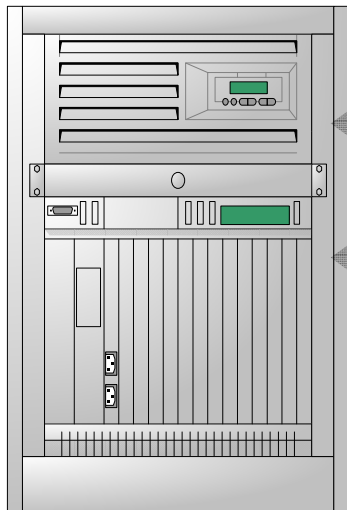
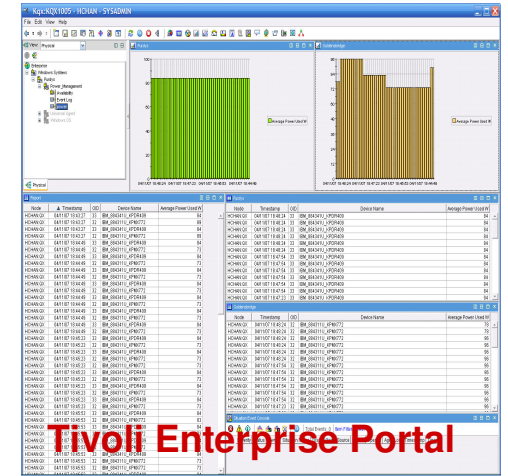
PowerExecutive is an IBM Director plug-in that interacts with the hardware management module to monitor power usage and temperature, and control power caps.



TUAM and other Tivoli managers can consume the historical aggregated performance and power data stored in **TDW**.



Admins can view and act upon live aggregated data and alerts on **TEP**.



ITM Power Data Provider is an IBM Director extension that feeds the **ITM Power Agent**.

Power, Temp



ITM Power Agent augments performance data traditionally collected from performance managers and the OS with power and temperature data. All of these data are aggregated for consumption by **TEP** and **TDW**.

Admin console (power and performance data)

The screenshot displays the IBM Admin console interface for 'ocean4SOCKdp:UAGENT00 - OCEAN4 - SYSADMIN'. The left-hand tree view shows the hierarchy: Enterprise > Windows Systems > OCEAN4 > Power_Management > dyn9002027240:POWER_MANAGEMENT > PERFORMANCE, POWER, and TEMPERATURE. Below the tree, there are summary tables for System Performance, Power, and Processor Temperature.

System Performance Table:

avRT	qLen	LocalTimeStamp
600	10261	04/20/07 15:05:11
595	10730	04/20/07 15:05:26
600	10370	04/20/07 15:05:44

Power Summary Table:

Machine Name	Power cap	Power used	LocalTi
goldensbridge	120	110	04/20/07

Processor Temperature Summary Table:

Machine name	Processor Temperature	LocalTi
goldensbridge	53	04/20/07
pleasantville	40	04/20/07

The main area contains four charts:

- System Response Time (ms):** A 3D bar chart showing response times over time. A red label 'Response Time' points to this chart. A specific data point is highlighted: 594 (04/20/07 15:10:56).
- Goldensbridge Power Cap and Power Use...:** A 3D bar chart comparing power cap (yellow) and power used (blue) for Goldensbridge. A red label 'Power' points to this chart.
- Pleasantville Power Cap and Power Used (Watt):** A 3D bar chart comparing power cap (yellow) and power used (blue) for Pleasantville.
- Goldensbridge Processor Temperature (C):** A 3D bar chart showing processor temperature for Goldensbridge. A red label 'Temperature' points to this chart.
- Pleasantville Processor Temperature (C):** A 3D bar chart showing processor temperature for Pleasantville.

At the bottom of the console, the status bar shows: Hub Time: Fri, 04/20/2007 03:30 PM, Server Available, and ocean4SOCKdp:UAGENT00 - OCEAN4 - SYSADMIN.

Power Management interface

Autonomic power mgmt

The screenshot shows the 'POWER - OCEAN4 - SYSADMIN' application window. A 'Power Control' dialog box is open, allowing configuration of power management settings. The 'Automatic Power Management Policy Control' option is selected and circled in red. A red arrow points from the text 'Autonomic power mgmt' to this option.

Select Power Mode

- Full Power: Default mode, set power cap to the maximum
- Manual Power Control: Set power cap to the desired level
- Automatic Power Management Policy Control: Set the policy for automatic power management control
- Local power control: Power control by each server locally

Select Policy

Optimized policy version 3 4/2/2006

Select Power Level and Machine

95W Goldensbridge

Status

Manual power control mode..

Summary and Future Directions

- **Thus far, we have achieved coordination across**
 - Multiple levels (from chip to application)
 - Two management disciplines (performance, power)
 - Approach
 - Express joint objectives in terms of utility functions
 - Combine modeling, optimization and state-of-art ML technologies
 - ML can save time and yield better policy
- **We are currently exploring**
 - Dynamic *voltage* and frequency scaling
 - Dynamic power-off of servers, exploiting virtualization and load balancing
 - **> 30% power savings in initial tests**
 - An additional management discipline: availability
- **Opportunities to extend to data center level**