



Going beyond CPUs: The Potential for Temperature-aware Data Centers

Justin Moore, Ratnesh Sharma, Rocky Shih
Jeff Chase, Chandrakant Patel
Partha Ranganathan



Motivation



asuwlink.uwoyo.edu/~jimkirk/

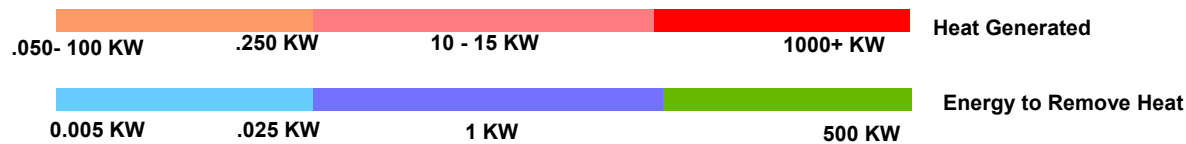
"TWO WORDS? ONE WORD!
STARTS WITH ... SOUNDS LIKE..."


- Several past studies on temperature-aware CPU designs
- BUT potential unexplored at higher-levels of system

The Temperature Problem in Data Centers

Cooling infrastructure

- At high-end, ~1W of cooling for every 1W of power!



- TCO costs: \$4-8 million cooling costs for 10MW data center
- Environmental costs:  11M GJ + 2M tons CO2 for US machines

Reliability and availability

- Mechanical parts – failure rates
- Thermal redlining if inlet exceeds 30C
- Lower operational efficiency at higher temperatures
 - 10-15C increase => server/disk failure rates up by 2X [Uptime, Cole]

Exacerbated by consolidation, overprovisioning & density trends

Addressing the Temperature Problem

- Conventional approaches at facilities level
 - New cooling approaches or better cooling delivery
- This work: **temperature-aware resource provisioning**
 - **Architecting a temperature-aware resource scheduler**
 - Characterizing the indirectly-controlled delayed-response metric
 - **Metrology**: Leverage thermo-dynamics-based air-flow equations
 - Combining IT level and facilities level (space and topology relations)
 - **Monitoring**: Deploy a location-aware knowledge plane [Splice]
 - Dealing with discrete power states
 - **Policies**: Algorithms for “zonal proximity”
 - **Preliminary results**
 - Significant cooling savings (within 94% of best-effort case)
 - Eliminate system failures caused by thermal emergencies

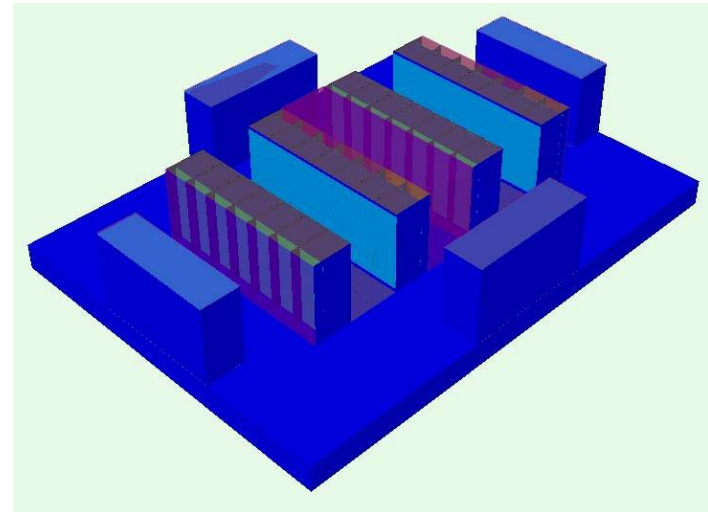
Outline for talk

- Motivation
- **Background and methodology**
- Temperature-aware resource scheduling
 - Metrology
 - Monitoring
 - Policy
- Summary and Future Work

Background and Methodology

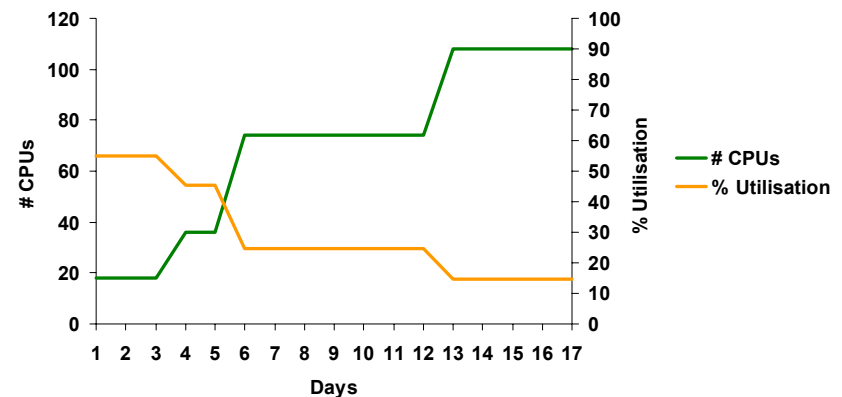
Conventional data center model

- 11.7mx8.5mx3.1m with 0.6m plenum
- 1120 servers
 - 4 rows x 7 racks x 40 1U servers
- 4 CRAC @ 86KW, hot/cold aisles
- Server-pair power states
 - 300W (idle), 580W (full)

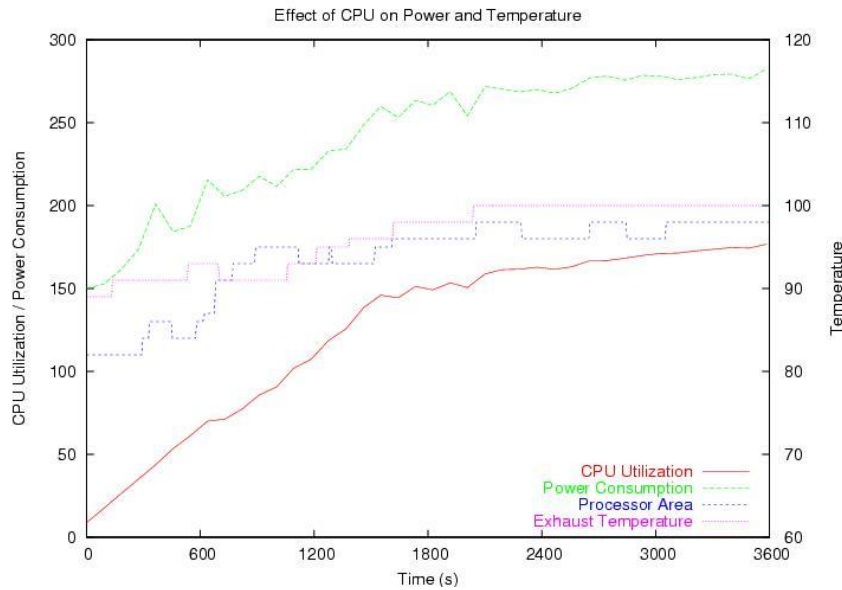


Scheduling media rendering workloads

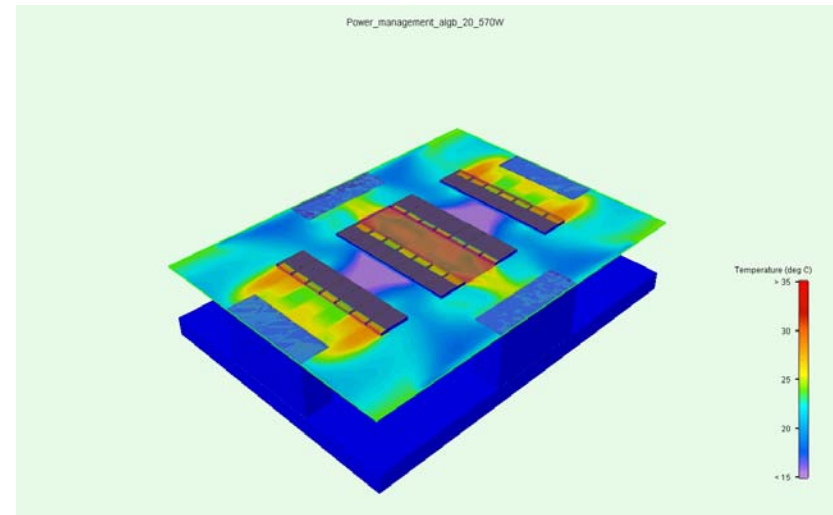
- @ utilizations of 25% and 50%



Defining the problem: temperature as an indirectly-controlled metric



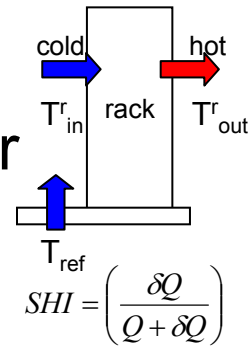
Individual server



Data Center

- Temperature as an indirectly-controlled metric
 - Non-intuitive correlations between system usage, power, temperature
 - Delayed response times
 - Need metrology to characterize these effects

Metrology to capture temperature variation effects [Sharma+2003]



- Thermodynamics-based proxies for thermal behavior
 - Model hot air infiltration into cold aisle and mixing
 - Model short circuiting (cold air directly to CRAC inlet)

$$\delta Q = \sum_j \sum_i m_{i,j}^r C_p \left((T_{in}^r)_{i,j} - T_{ref} \right)$$

- Thermal policies for heat distribution

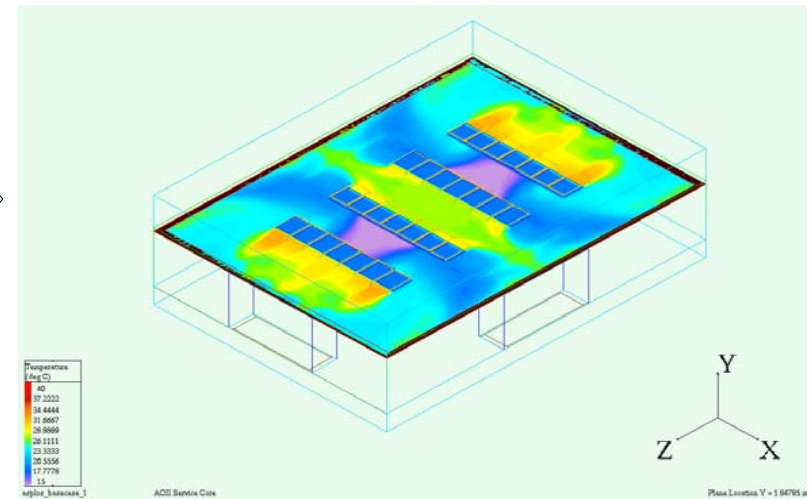
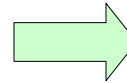
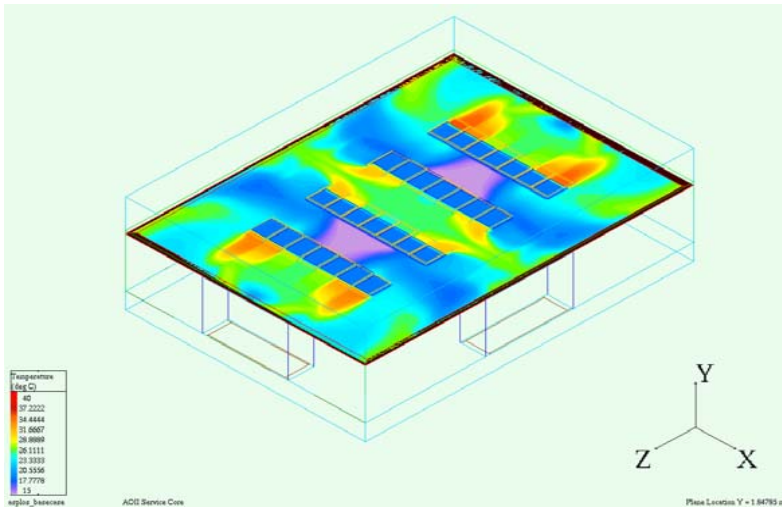
$$Q = \sum_j \sum_i m_{i,j}^r C_p \left((T_{out}^r)_{i,j} - (T_{in}^r)_{i,j} \right)$$

- $W = Q/COP$; $COP = f(T_{ref})$; $Q = mC_p(T_{return} - T_{ref})$
- Reducing T_{return} means T_{ref} can be increased correspondingly

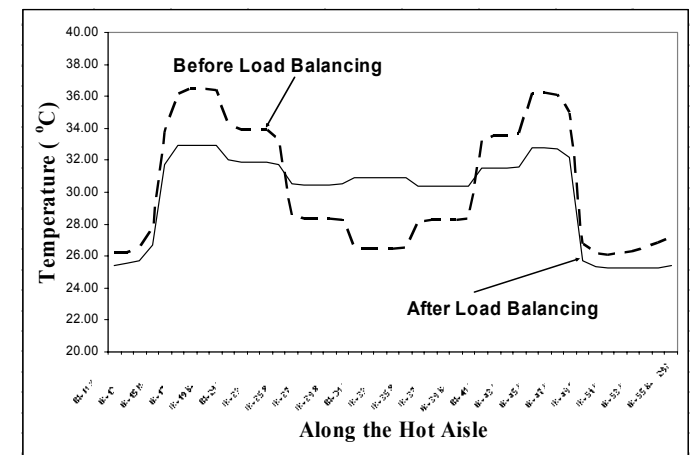
- This talk: Use exhaust temperature as first-order proxy
 - Make exhaust temp uniform to maximize inlet temperature (~25C)
 - Distribute heat inversely proportional to exhaust temp

- “Ideal distribution”
$$P_i = \left(\frac{T_i - T_{ref}}{T_{i=0} - T_{ref}} \right)^{-1} P_{i=0}$$

Savings from applying “ideal” policy

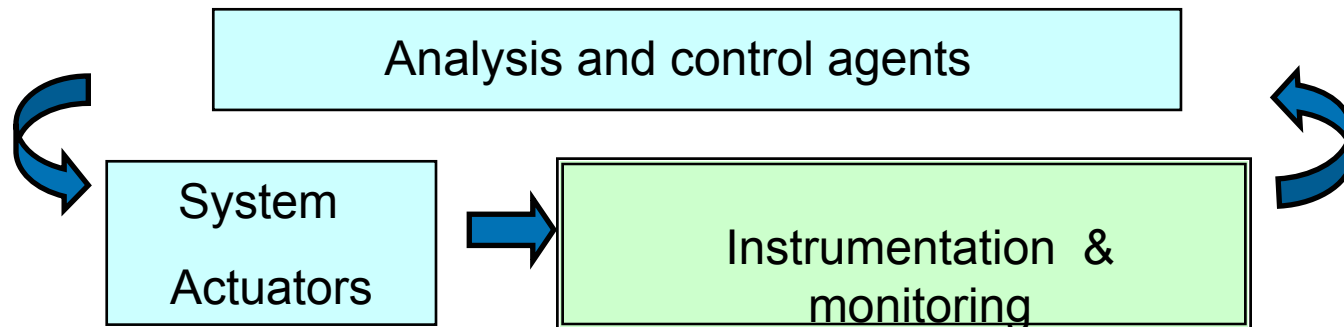


Smoothed exhaust temperature profile
 Higher CRAC efficiency + higher return temp
 Cooling energy savings 25%



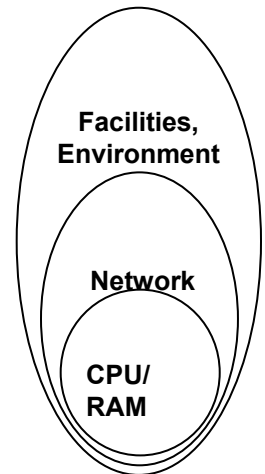
Implementation?

- Thermodynamics-based formulation of objective function and actuation impact
- BUT how do we implement this in a real system



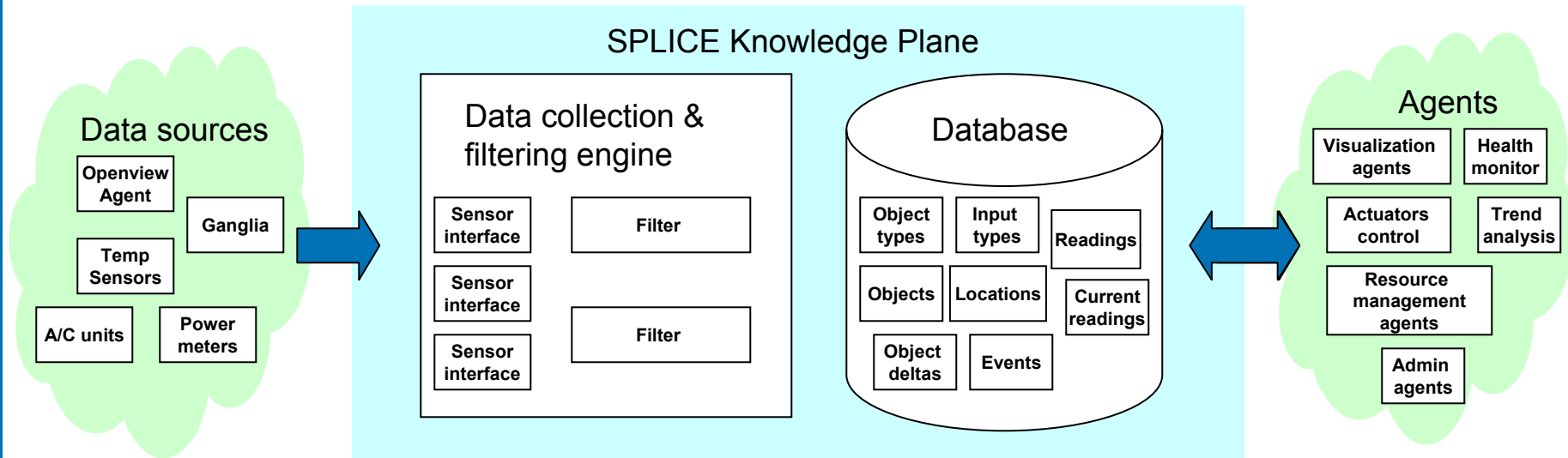
Instrumentation and Monitoring

- Current instrumentation approaches inadequate for temperature-aware resource scheduling
- Needs
 - Instrumentation across IT/facilities layers
 - “Expanded computing environment”
 - Conventional IT metrics (e.g., CPU, network, etc.)
 - Environmental sensors (power, temperature, humidity)
 - Proprietary and diverse “publish models” (e.g, OPC)
 - Synchronization
 - Data repository and access
 - Need for scalability to hundreds of sensors, millions of readings
 - Notion of higher-level and hierarchical object views
 - Speed of query access



A location-aware information plane

[Moore+2003]



- Instrumentation data sources
 - Unified correlated data collection and aggregation
- Data collection and filtering
 - Support for multiple interfaces
- Database schema
 - Enables higher-level object views, scalable, support for newer data types
- Analysis and control agents
 - SQL interface to database

Deployment

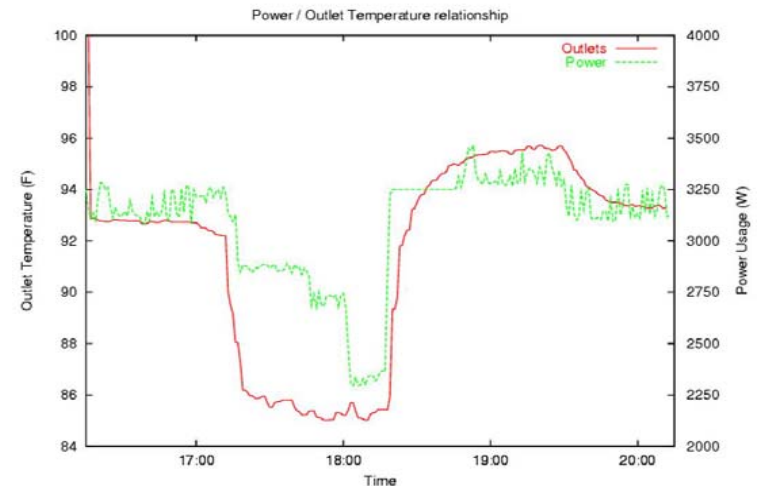
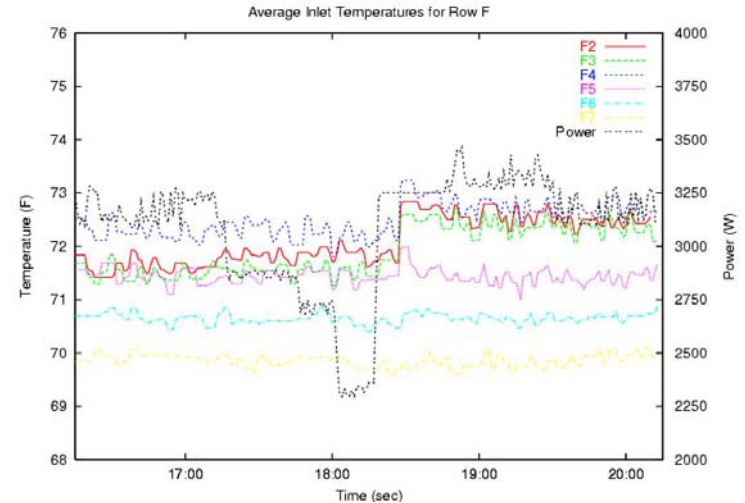
Splice deployed at HP Labs Utility Data Center (UDC)

- HP Openview for performance metrics and OPC interface for temperature and power sensors

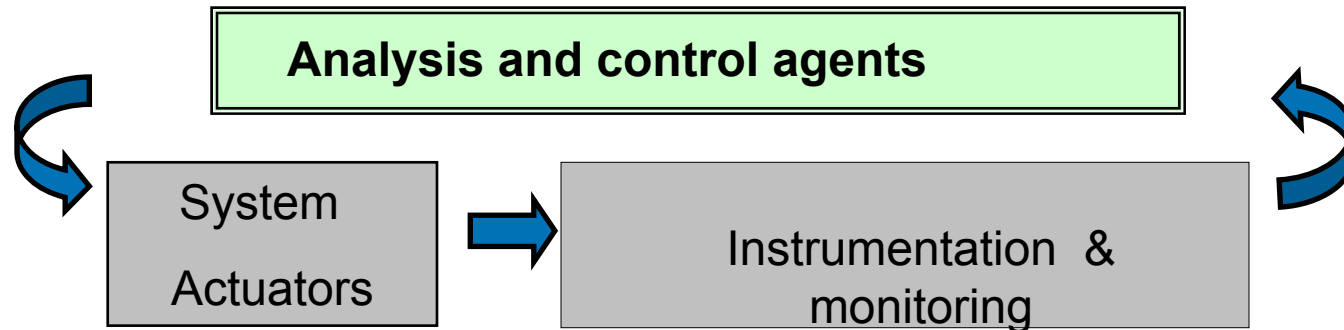
Use with temperature-aware scheduling

But also other IT-facilities- boundary optimizations

- E.g., operations automation (problem detection, cause-effect analysis, provisioning, ...)

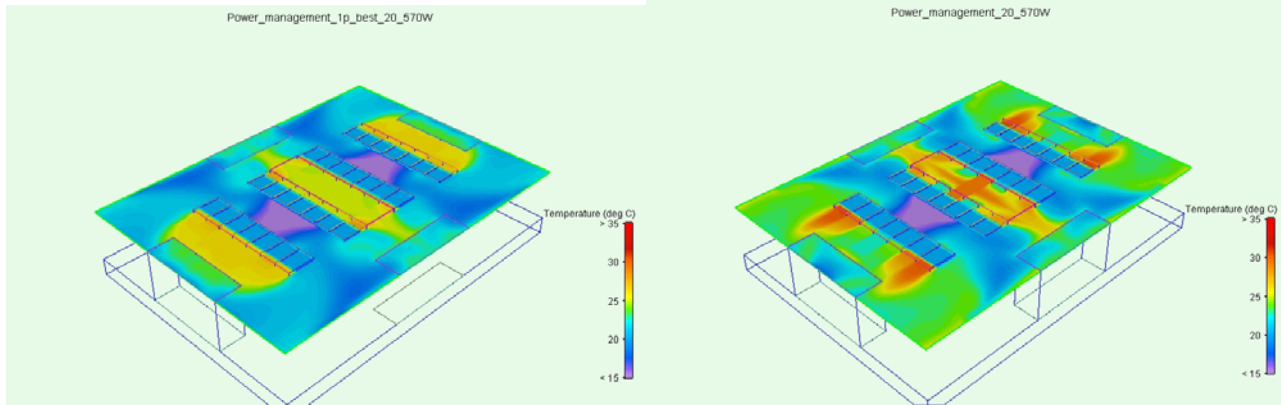


Policies



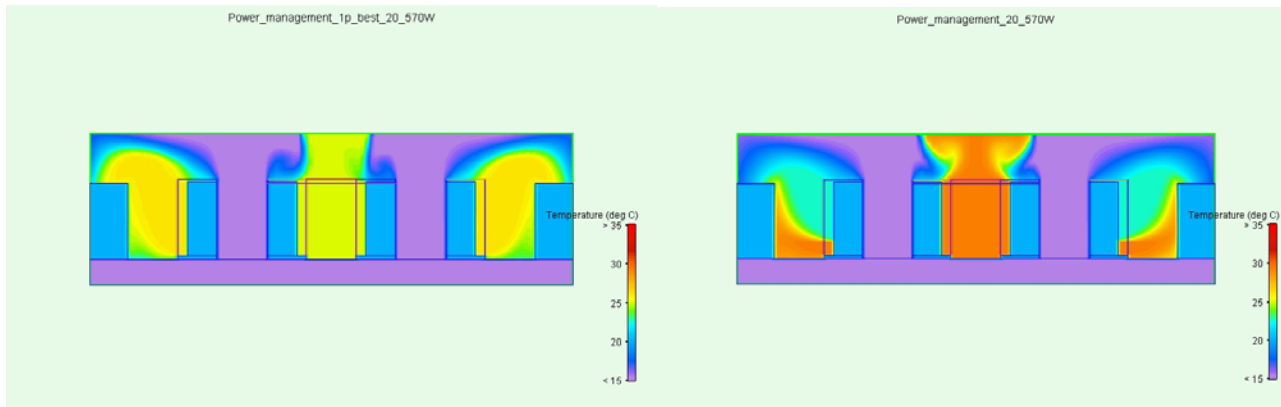
- “Ideal thermal policy” is analog
 - How do we discretize it for server power states and static task scheduling with no workload migration?
 - Simple heuristic based on thermal policy
 - Sort exhaust temperatures
 - Place hot loads on coolest spots
 - For our data center => interior middle racks

40% worse compared to ideal!



Ideal

ColdInlet

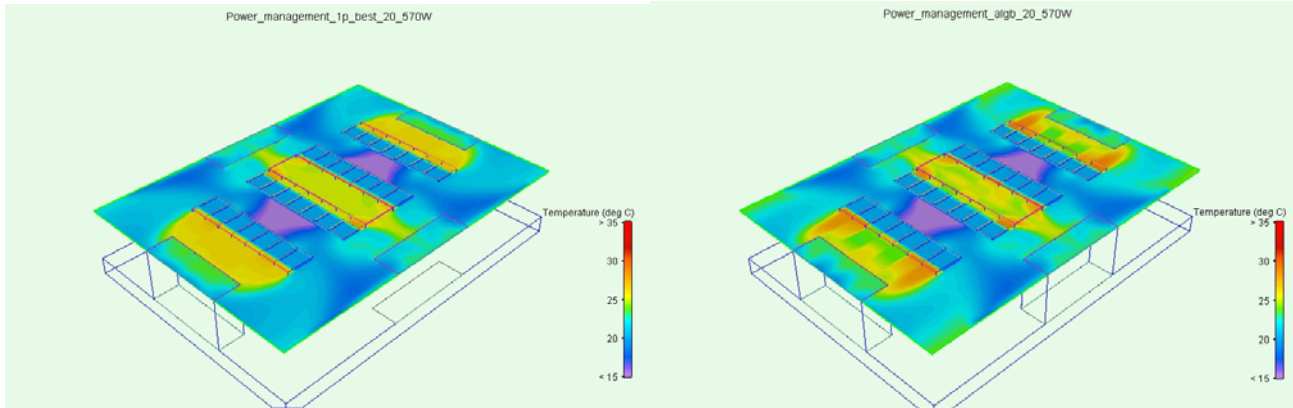


- Discreteness leads to imbalance and new hot spots
 - Increased energy to cool

Proximity-based algorithms

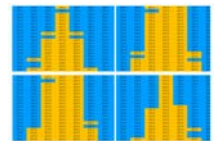
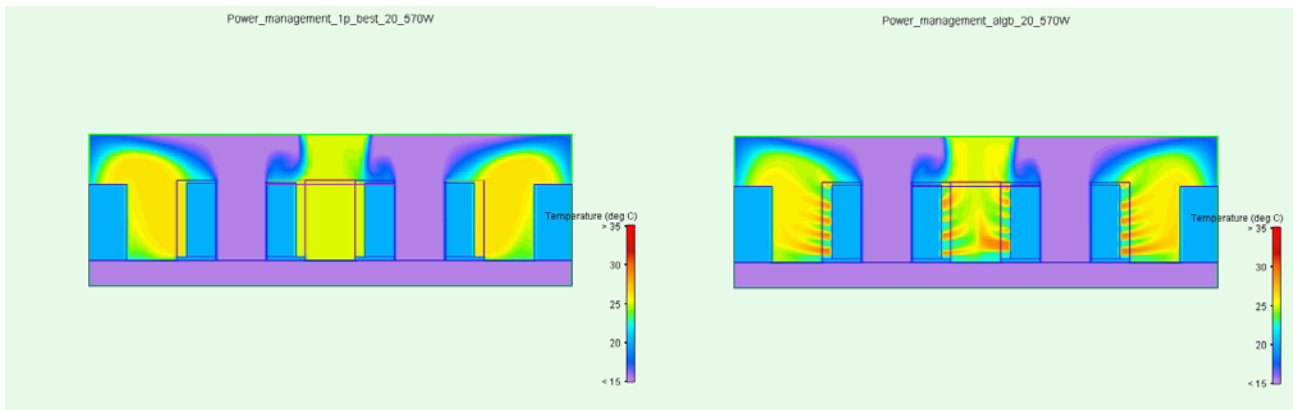
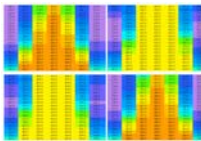
- “Two-pass Discretize”
 - Intra-row first-pass; inter-rack second-pass
 - Schedule per floor of analog allocation
 - Schedule excess with bias towards “median”
- “Proximity-based Poaching”
 - Single pass through three-dimensional space
 - Assign server load
 - derate adjacent servers for new analog allocation

Power savings close to ideal!



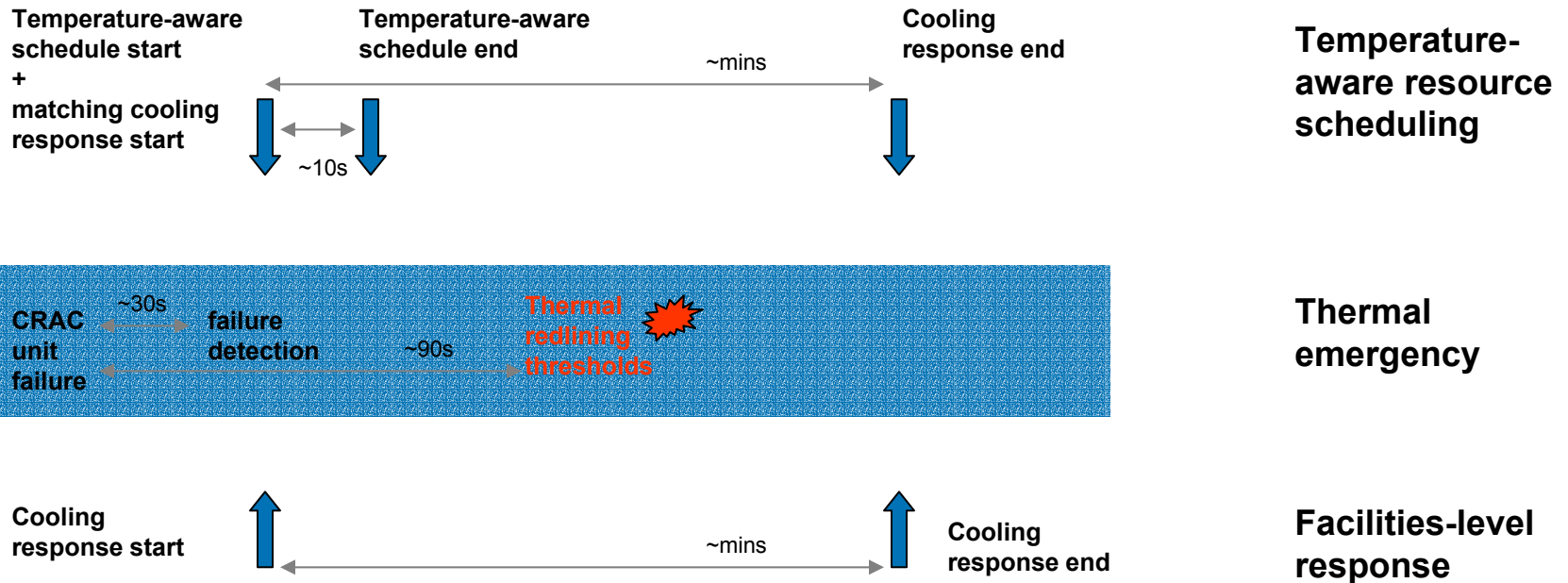
Ideal

Poaching



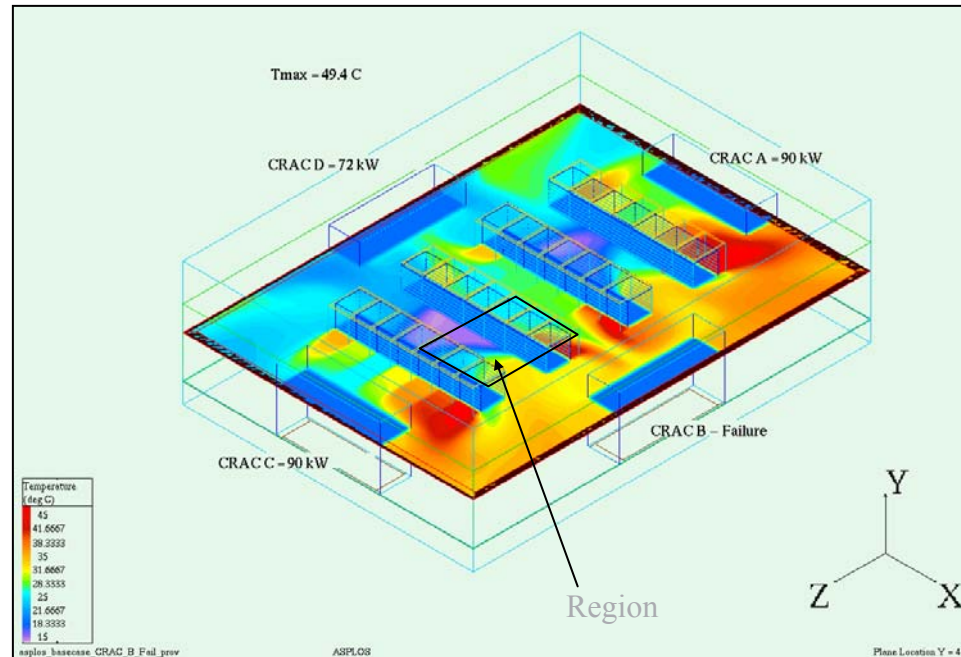
- Heat distribution matches at zonal level
 - Two-pass within 15%; Poaching within 6% of ideal
 - Poaching yields close to 25% energy savings w.r.t bad scheduling

Temperature-aware scheduling for thermal emergencies



- Faster response to thermal emergencies
 - Controlling heat source better than adjusting heat sinks
 - Same algorithms can be applied with emergency trigger

Temperature-aware scheduling for thermal emergencies



- Applying “proximity-based-poaching”
 - Reduce thermal redlining servers by 55% in first 30 sec
 - Potential to fully eliminate thermal redlining failures

Summary

- Temperature-aware provisioning valuable at data center level
 - Cooling costs reduction and increased reliability/availability
- This work: Architecting a temperature-aware resource scheduler
 - Characterizing the indirectly-controlled delayed-response metric
 - Metrology: Leverage thermo-dynamics-based air-flow equations
 - Combining IT level and facilities level (space and topology relations)
 - Monitoring: Deploy a location-aware knowledge plane [Splice]
 - Dealing with discrete power states
 - Policies: Algorithms for “zonal proximity”
 - Preliminary results
 - Significant cooling savings (within 94% of best-effort case)
 - Eliminate system failures caused by thermal emergencies
- Ongoing work
 - More elaborate thermal policies and coarser grain policies
 - More discrete power states (v/f scaling, virtual machines)
 - Control on CRAC air flow rates

Questions?



Related Work

- Traditional approaches
 - Facilities-level work on cooling systems [IPACK]
 - Costs, granularity of control and response, do not address heat
 - Power-aware IT resource scheduling [SOSP02, PACS02, WCOP01]
 - Focus on IT power, temperature can be improved or worsened
- Hybrid approach: Control at IT-facilities intersection
 - Workload migration proposed in Sharma et al [HPLTR03]
 - Focus on thermo-dynamic thermal policies in ideal scenario
- Our work: temperature-aware resource scheduling
 - Real-world constraints, architected solution

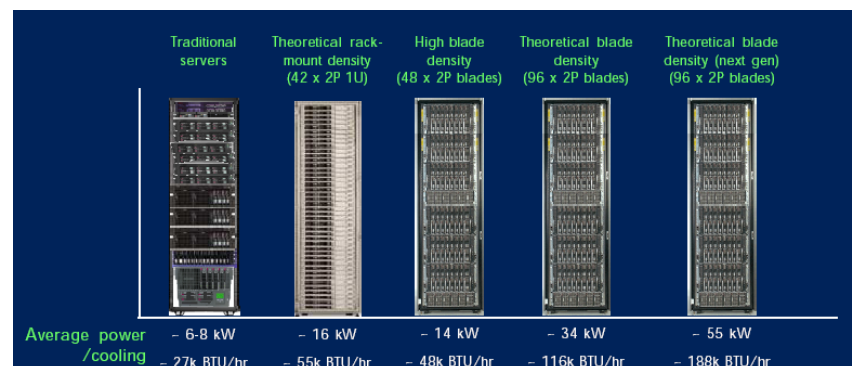
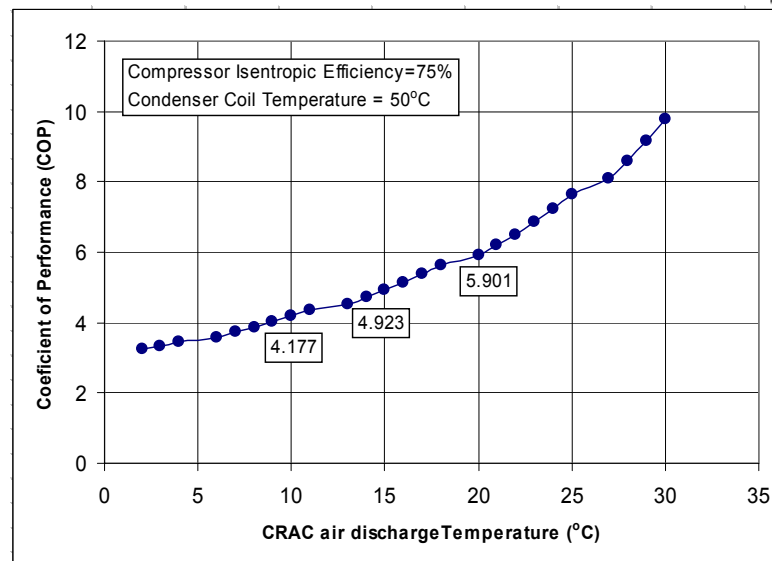
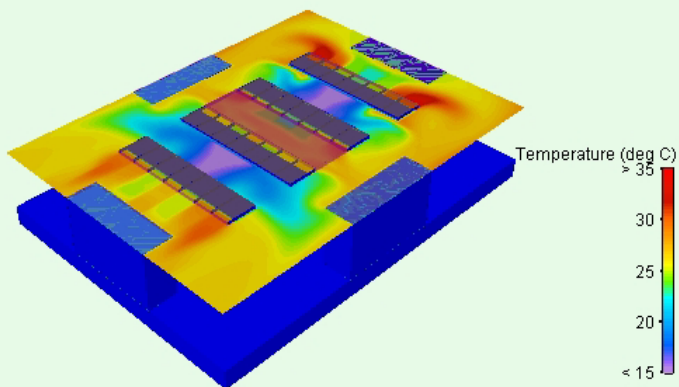
Temperature-aware scheduling: Challenges



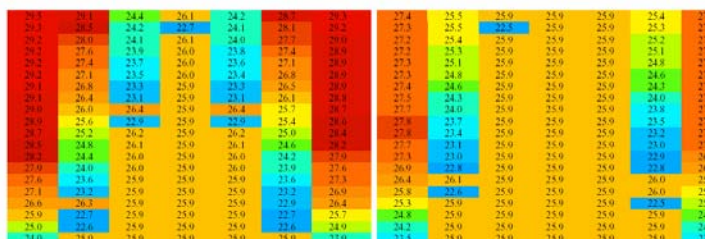
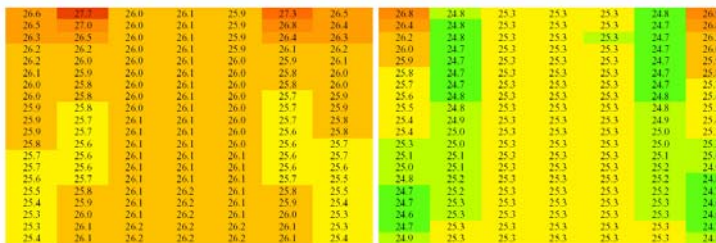
- Temperature as an indirectly-controlled metric
 - Non-intuitive correlations between system usage, power, temperature
 - Delayed response times
- Need for location-enhanced knowledge plane
 - Integrate IT-level metrics with facilities-level metrics
 - Capture spatial and topological relationships
- Discreteness in power states
 - Constraints on power modes in system
 - Constraints on workload migration modes

Backup

PM_20_failover_trans_steady

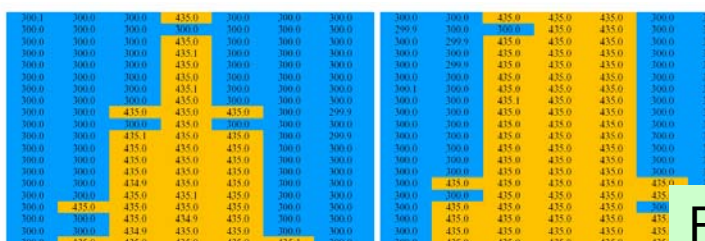


Proximity-based scheduling



Ideal

Poaching



Ideal

Poaching

