

Faraz Ahmad ASPLOS 10

Joint Optimization of Idle and Cooling Power in Data Centers While
Maintaining Response Time

Data Center Power

- ✦ Cooling is a large fraction of total power
- ✦ Idle power in servers is wasted and needs more cooling
- ✦ Shutting off idle servers concentrates heat in those that are active, requiring more cooling
- ✦ Handling surges in load is slowed by having servers shut off or in standby

PowerTrade

- ✦ Divide data center into thermal zones
- ✦ Calibrate load against temperature
- ✦ Put hot zone servers in standby when possible
- ✦ Spread load among cool zones to balance server power and temperature
 - ✦ Running more processors reduces hot spots, but increases idle power -- need to optimize

Static vs. Dynamic

- ✦ One set of measurements is better than none, but doesn't match varying loads
- ✦ Can dynamically measure power as load changes
 - ✦ Found that groups of 10 servers, tuned every 20 minutes works well

SurgeGuard

- Optimum power may result in too many processors that are slow to respond to surges in load
- Keep some servers in reserve to handle load (balance once per hour - but can run out)
- Add servers to replenish reserves at finer granularity (5 min)
- Deactivate excess servers at hour intervals only

Simulation on Traces

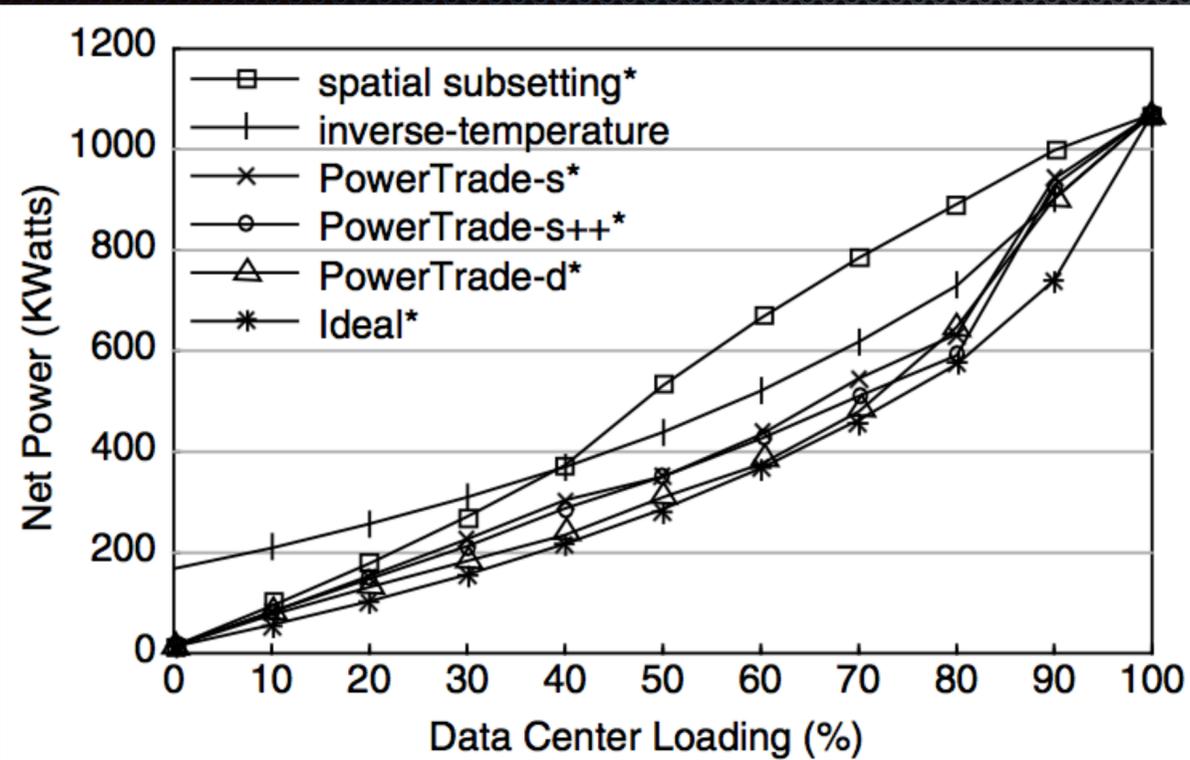


Figure 11: Net Power (* includes SurgeGuard)

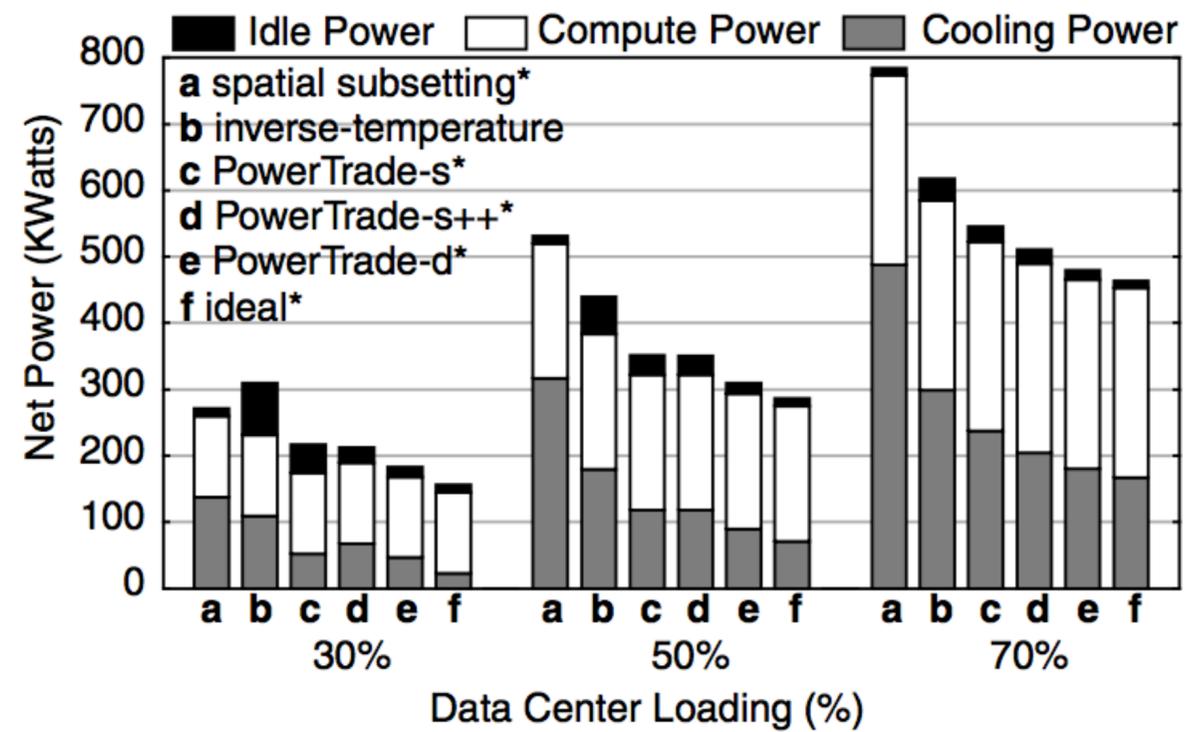


Figure 12: Net Power Breakdown (* includes SurgeGuard)

Discussion

Matt Skach ISCA 15

Thermal Time Shifting: Leveraging Phase Change Materials to Reduce Cooling Costs in Warehouse Scale Computers

Thermal Load

- Warehouse computers vary in load over time
- Need to be provisioned for cooling under peak load
- Can be wasteful if peak is limited in time
- Use material that can capture some of exhaust heat and store it until lower load, then release it
- Phase change exploits latent heat
- Tested multiple materials, wax was chosen

Wax Box in Exhaust

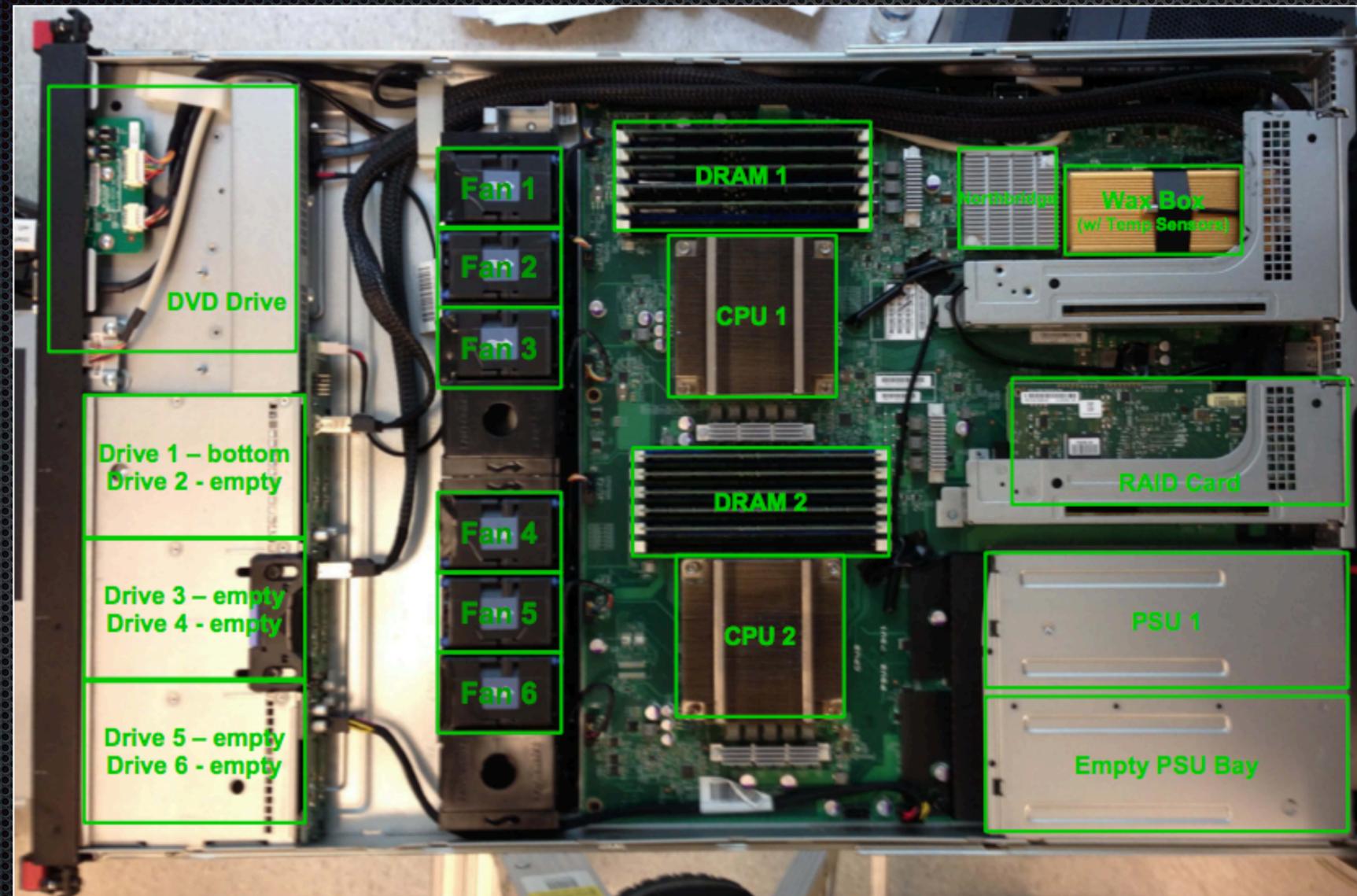


Figure 3: RD330 Server with major components labeled.

Checking the Model

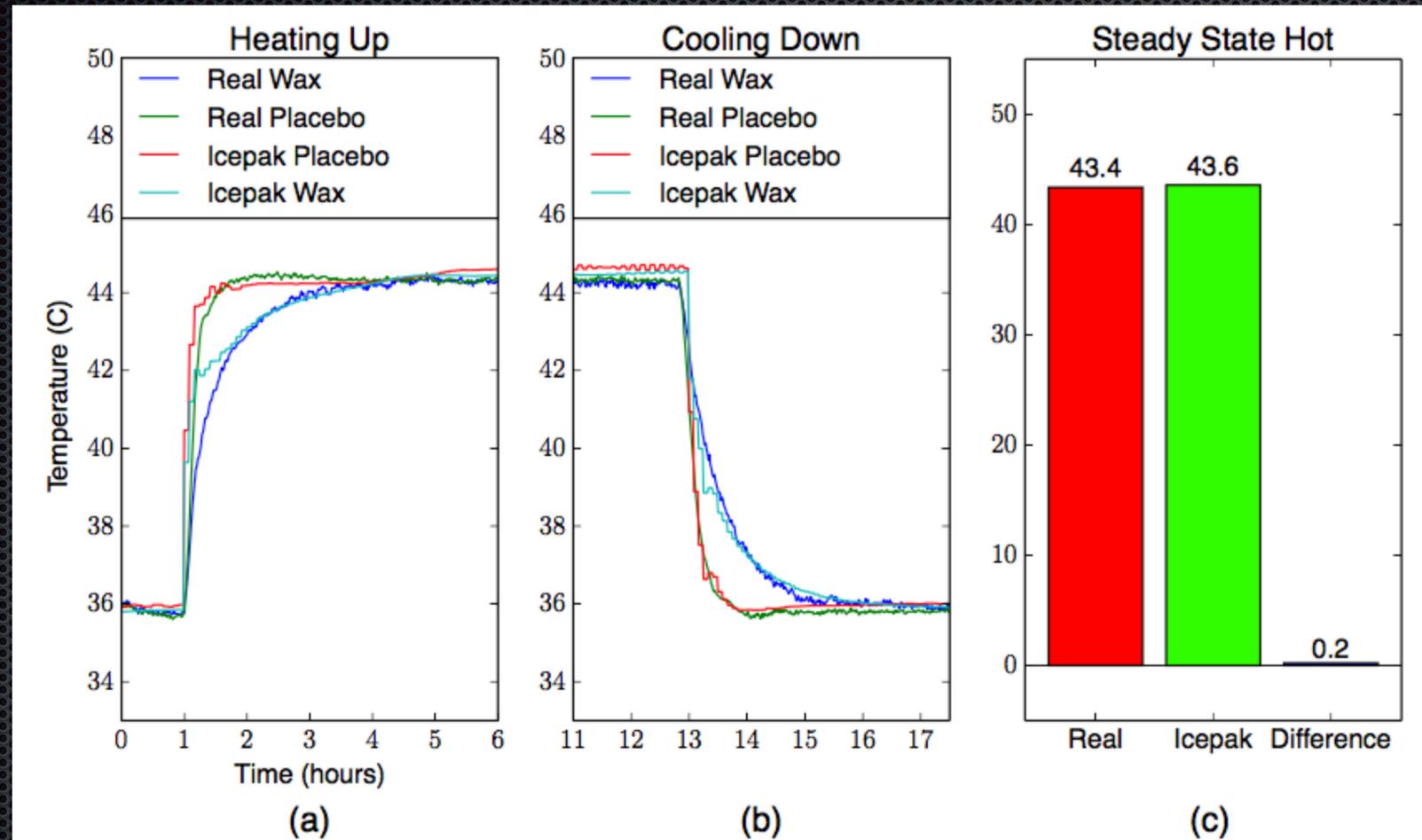


Figure 4: Model Validation. Transient traces while heating up (a) and cooling off (b), and steady state while hot (c) comparison of temperatures around the wax in the real server and our Icepak model.

Results

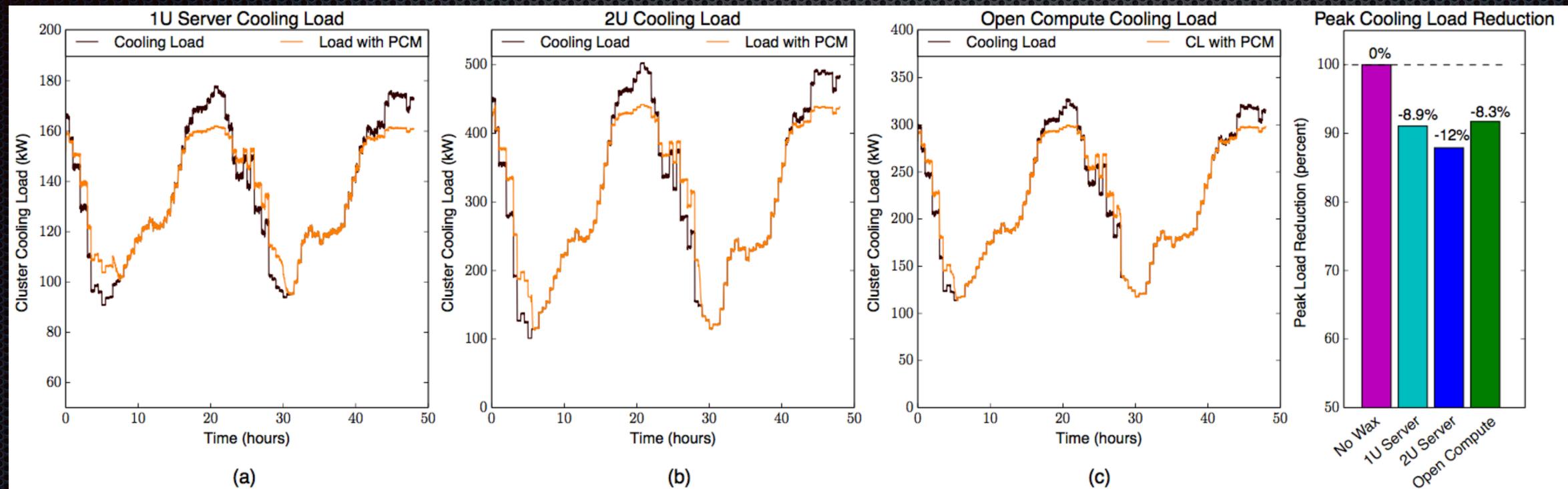
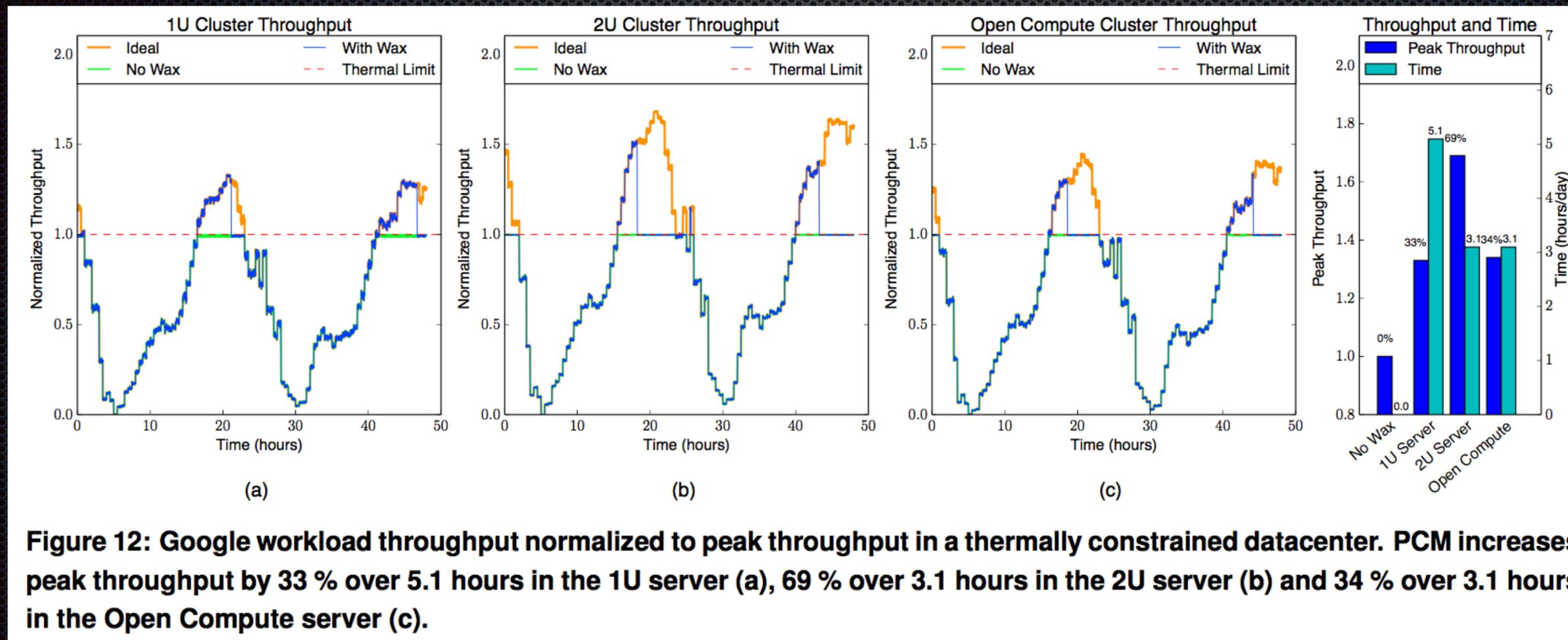


Figure 11: Cooling load per cluster over a two day Google trace in a datacenter with a fully subscribed cooling system. PCM reduces peak cooling load by 8.9 % in a cluster of low power 1U servers (a), 12 % in a cluster of 2U high throughput commodity servers (b), and by 8.3 % in a cluster of high density Open Compute servers (c).

Throughput Results



Matt Skach ISCA 18

Virtual Melting Temperature: Managing Server Load to Minimize Cooling Overhead with Phase Change Materials

Phase Change

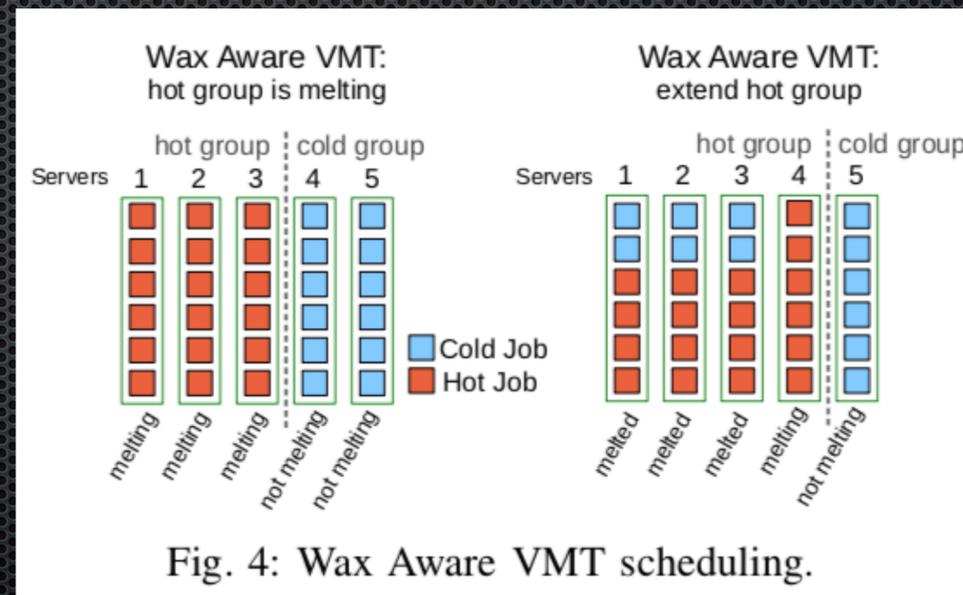
- ✦ Wax absorbs heat at the highest rate when it melts
 - ✦ Combination of sensible and latent heat
- ✦ Thermal Time Shifting doesn't necessarily reach melting point
 - ✦ Sensible energy storage only
- ✦ Introduce Virtual Melting Point scheduling to concentrate loads to melt wax

Thermal-Aware Placement

- Determine thermal profile of workload
- Gather hot processes into a subgroup designated as hot
- Keep other processes in cold group
- Hot group will melt wax and absorb more heat

Wax Aware Placement

- Measure wax temperature
- Once wax is melted, keep it melted (so it doesn't release heat)
- As all the wax in the hot group becomes melted, add from the cold group



Workloads for Evaluation

TABLE I: Workloads considered for scaleout study (power is normalized to a single 8 core Xeon E7-4809 v4 CPU; each server contains four CPUs).

Workload	CPU Power	VMT Class
WebSearch	37.2 W	hot
DataCaching	13.5 W	cold
VideoEncoding	60.9 W	hot
VirusScan	3.4 W	cold
Clustering	59.5 W	hot

Reliability Concerns

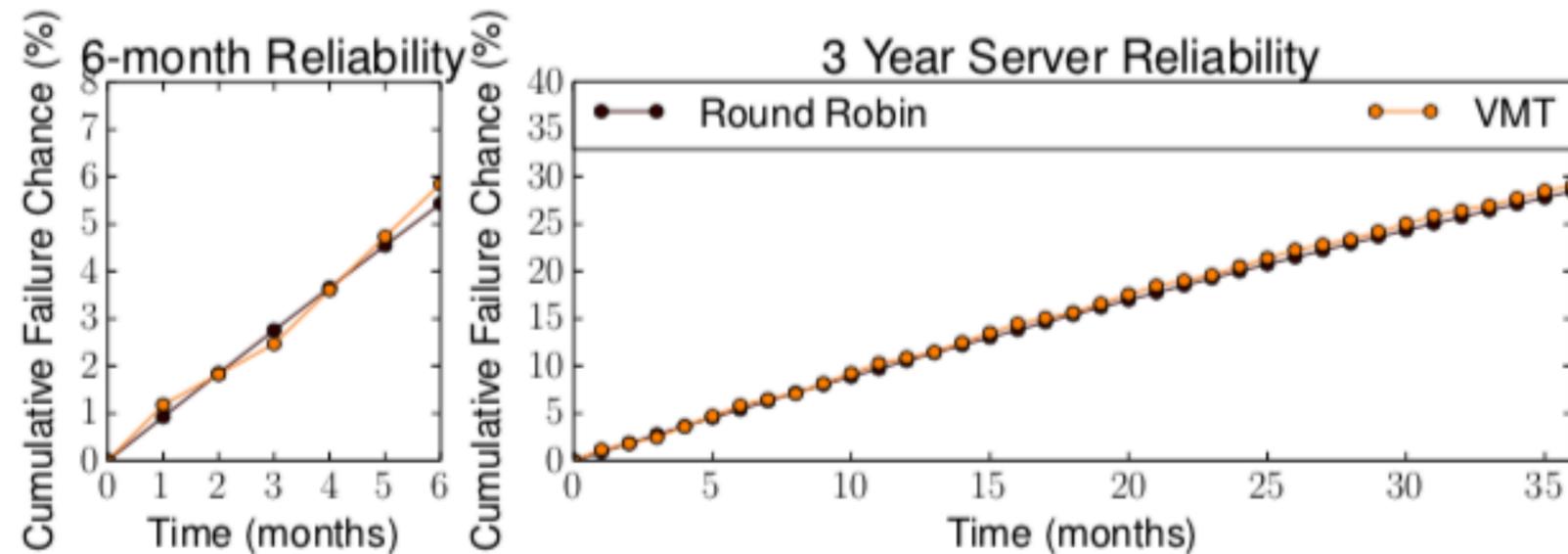
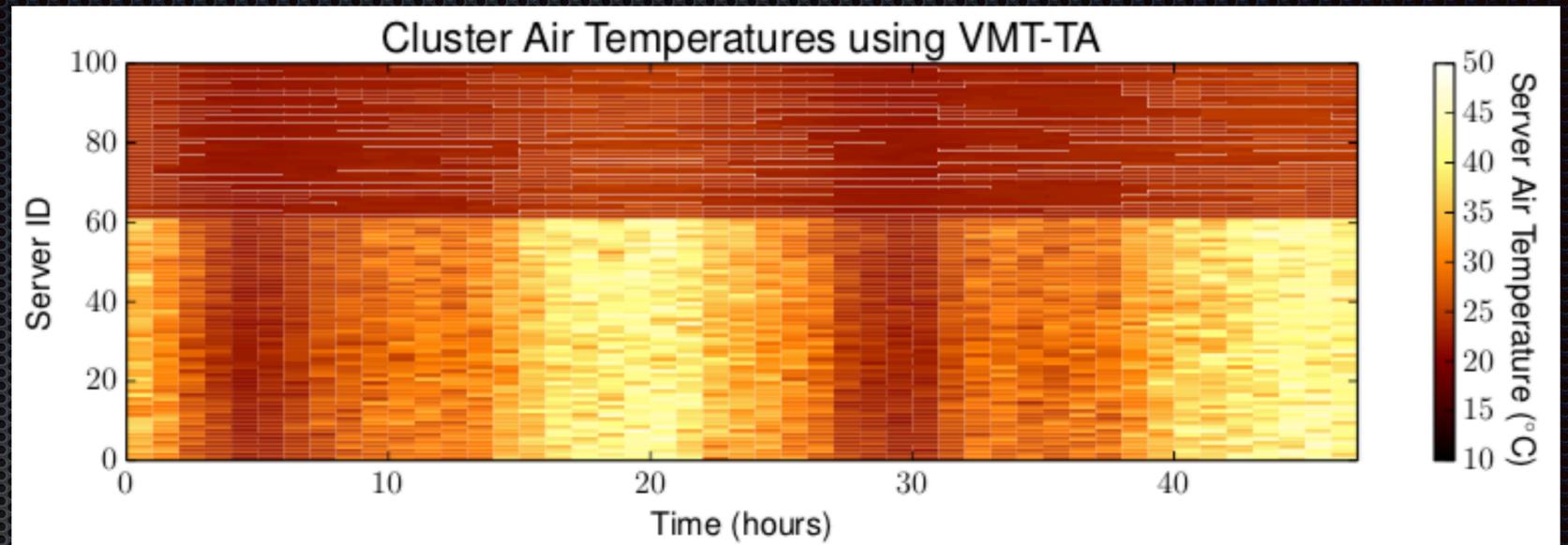
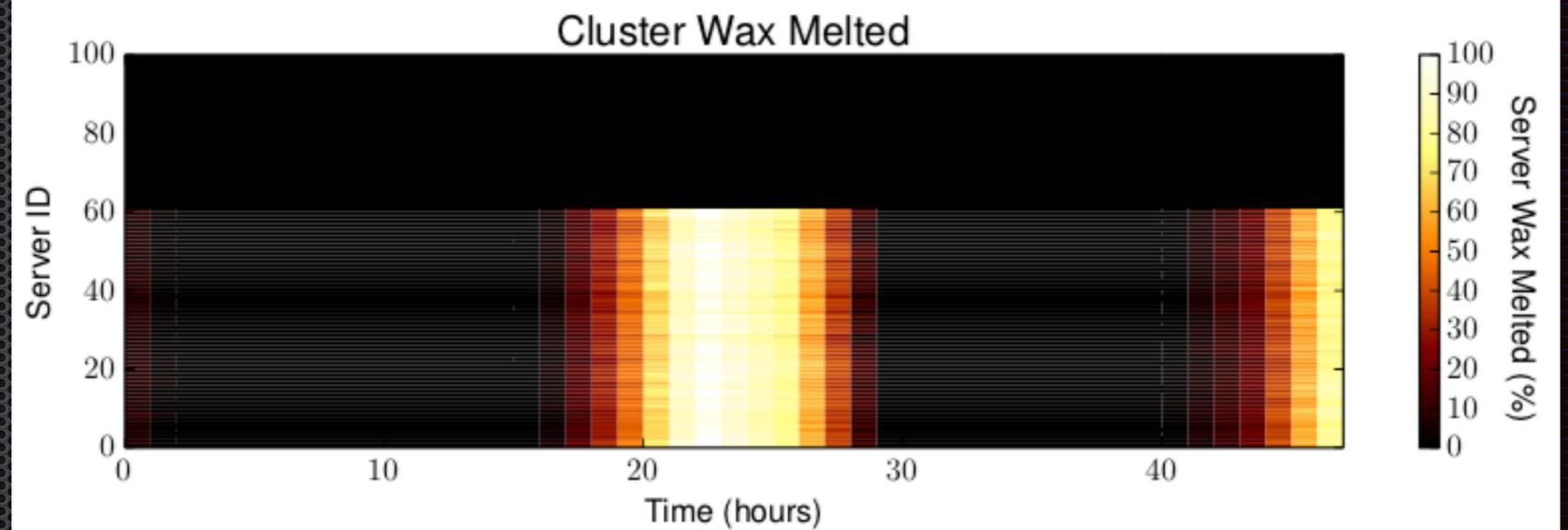


Fig. 7: Server reliability for round robin versus VMT-WA when 20% of servers are rotated each month (3 months in the hot group, 2 months in the cold group). After 3 years, the cumulative failure rate for VMT-WA is 0.4% higher than for RR.

Thermal Profile



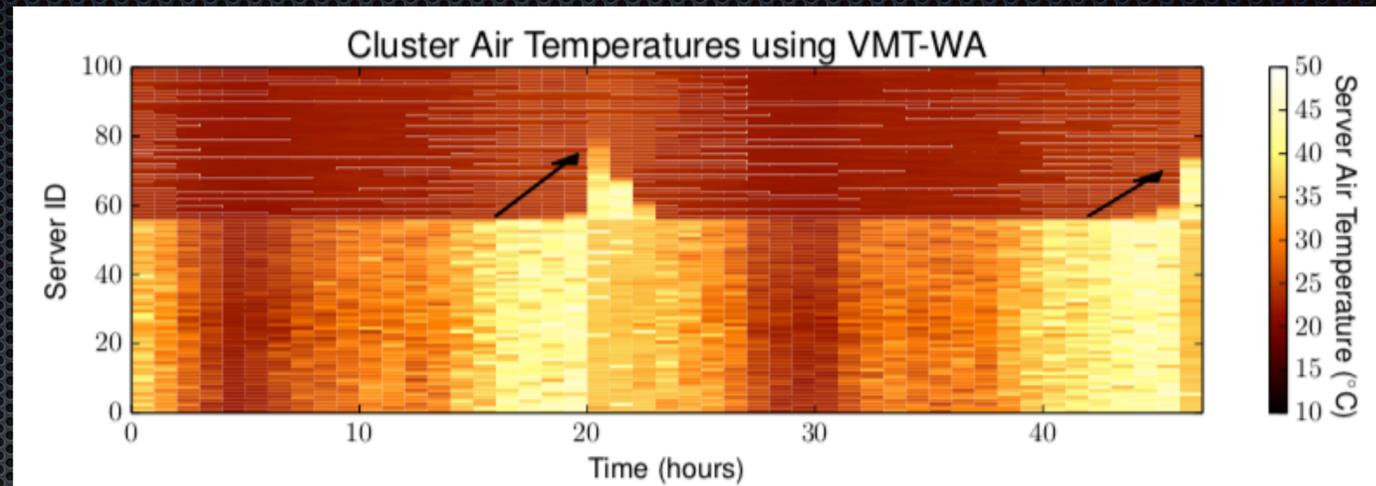
(a) Air temperatures at the wax using VMT-TA.



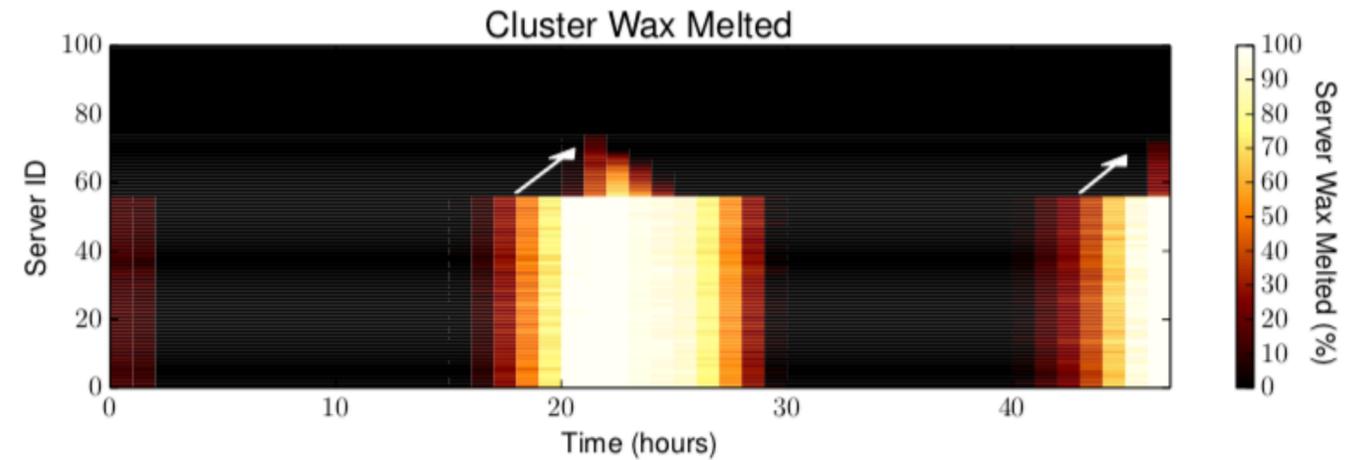
(b) Wax melted using VMT-TA.

Fig. 11: Air temperatures and wax melting for 100 servers using VMT-TA with $GV=22$.

VMT-WA Profile



(a) Air temperatures at the wax using VMT-WA.



(b) Wax melted using VMT-WA.

Fig. 14: Heat map of Air temperature at the wax, and wax melted, for a cluster of 100 servers using VMT-WA scheduling ($GV=20$). The hot group servers (bottom) have a consistently higher temperature than the cold group servers (top). Note the expansion of the hot group around 20 and 45 hours correspond with peak load and wax in the hot group reaching the wax threshold.

Cooling Load Reduction

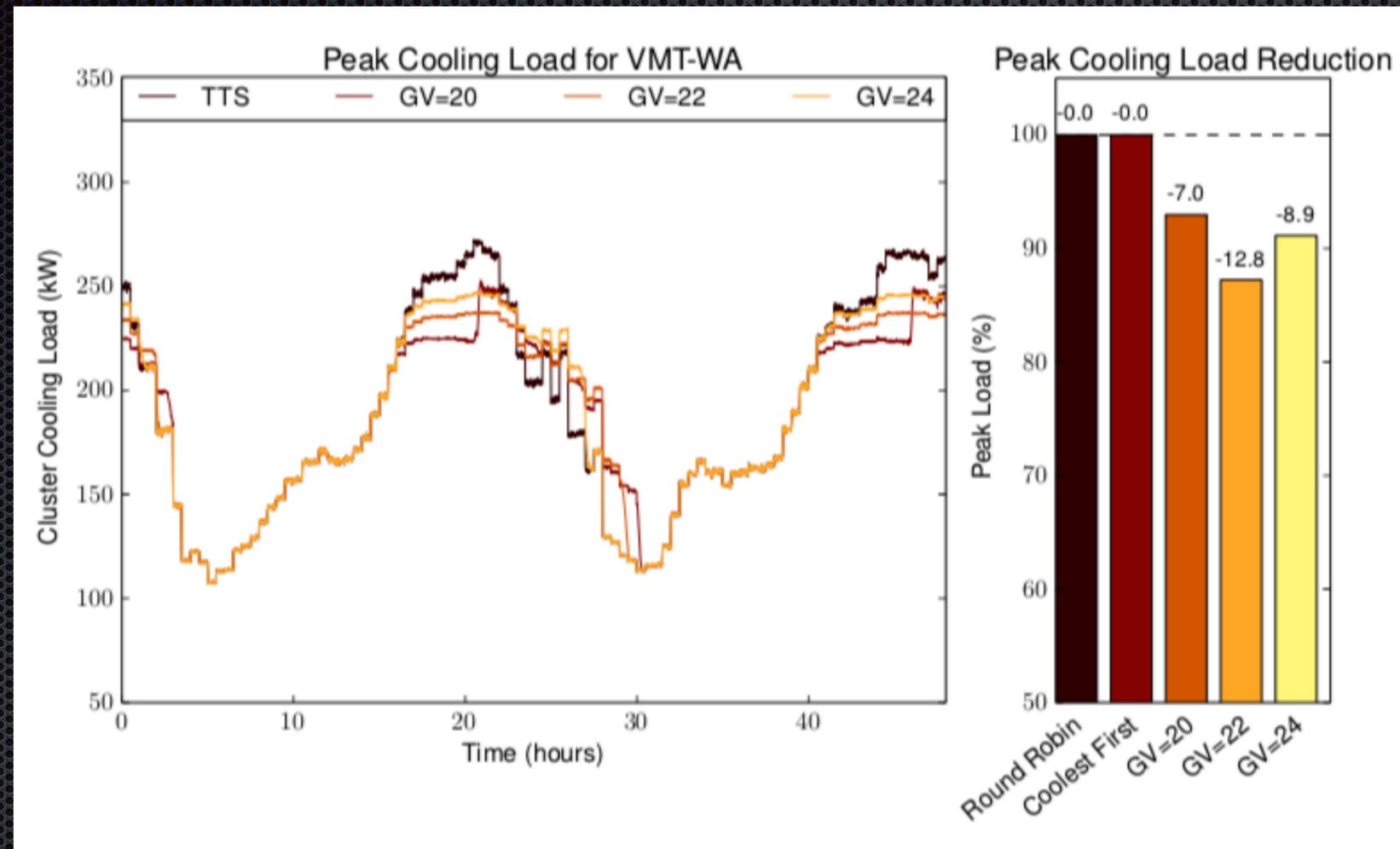


Fig. 16: Cooling load reduction with VMT-WA at 3 different GV levels for a cluster of 1000 servers. For GV=20 when the hot group becomes fully melted, VMT-WA adds more servers to the hot group to and rebalanced load to continue melting wax.

Discussion

Qiang Wu ISCA 16

Dynamo: Facebook's Data Center Wide Power Management System

Power Supply

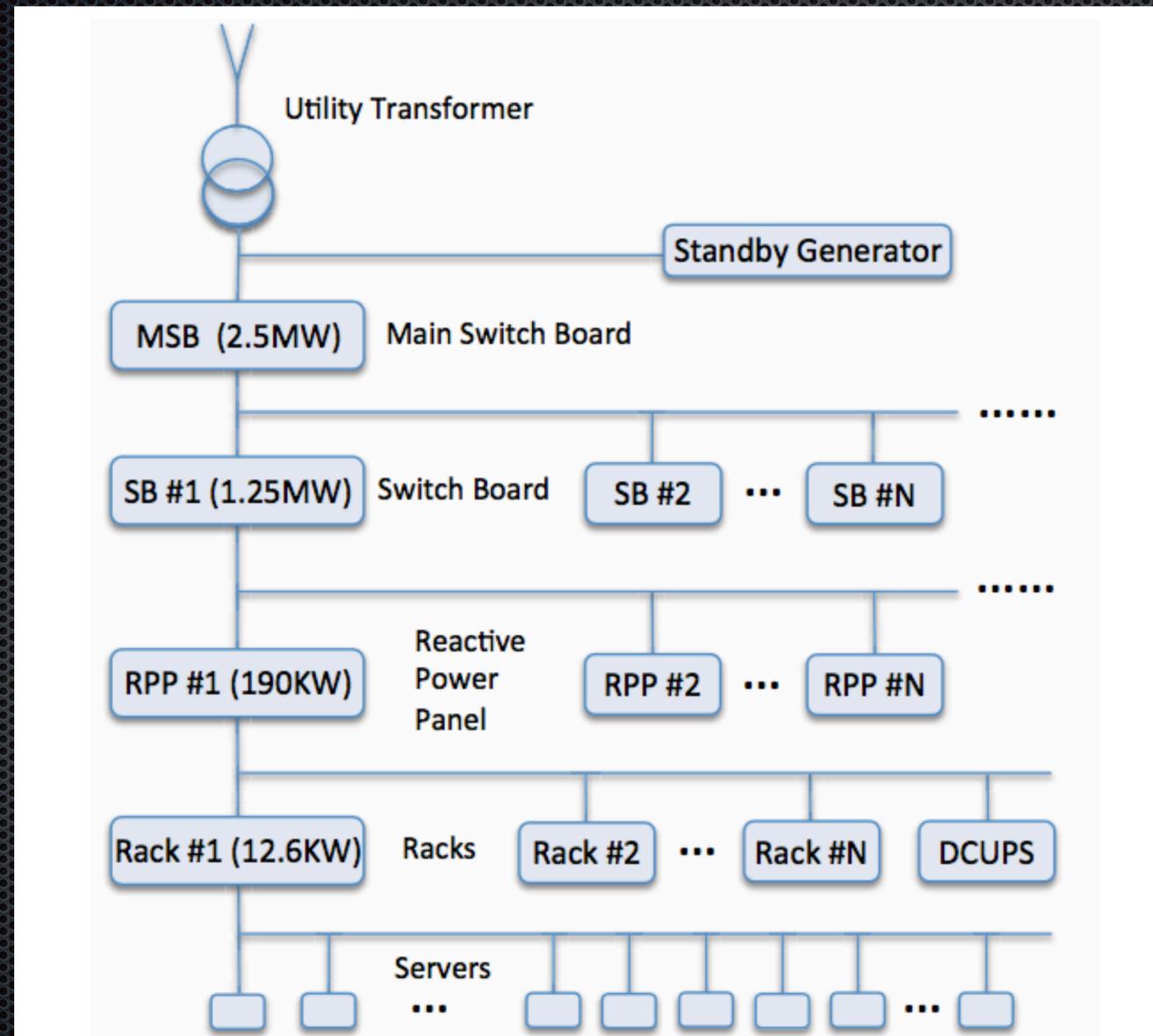


Figure 2. Typical Facebook data center power delivery infrastructure [14].

Power Variation

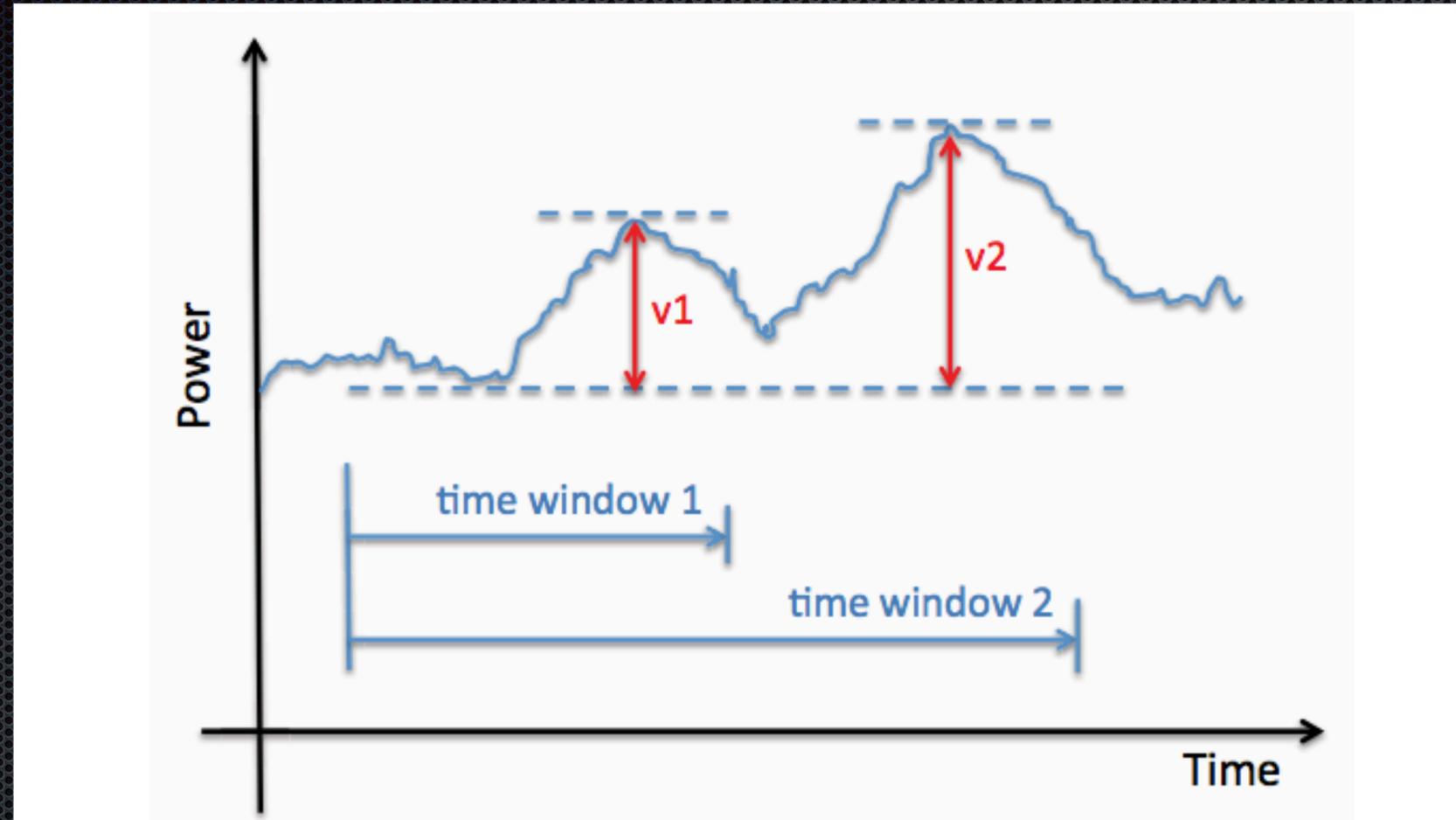


Figure 4. An illustration of calculated power variation for a time window. The maximum power variation is the difference between the maximum and minimum power values in the time window. Here, v_1 and v_2 are the maximum power variations for time windows 1 and 2, respectively.

Services Vary

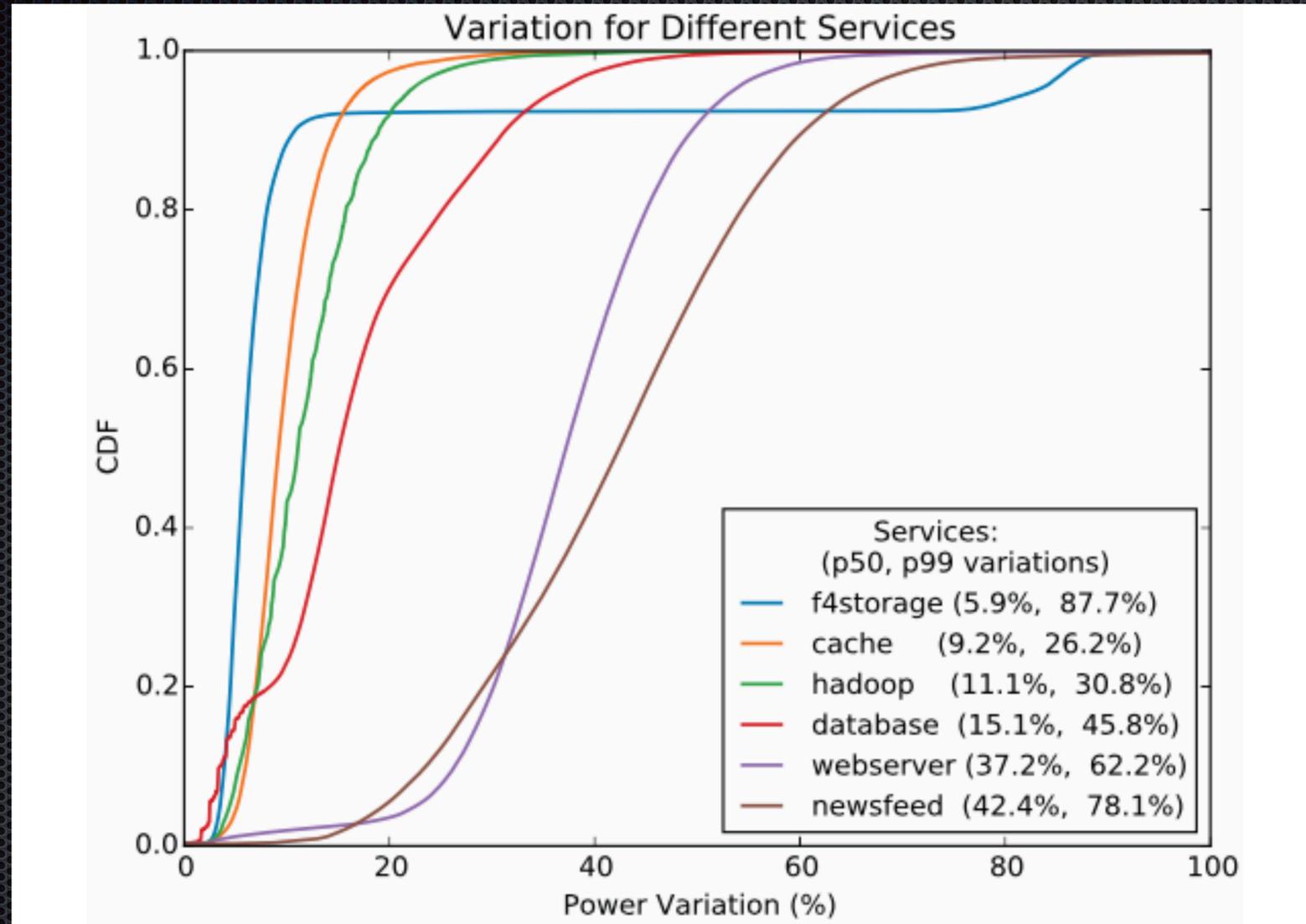


Figure 6. The measured power variations for several services in Facebook at the server level. The default time window is 60 s. For convenience, we also show the median (50th percentile, or p50) and p99 power variation values for each case.

Dynamo Coordinates Levels

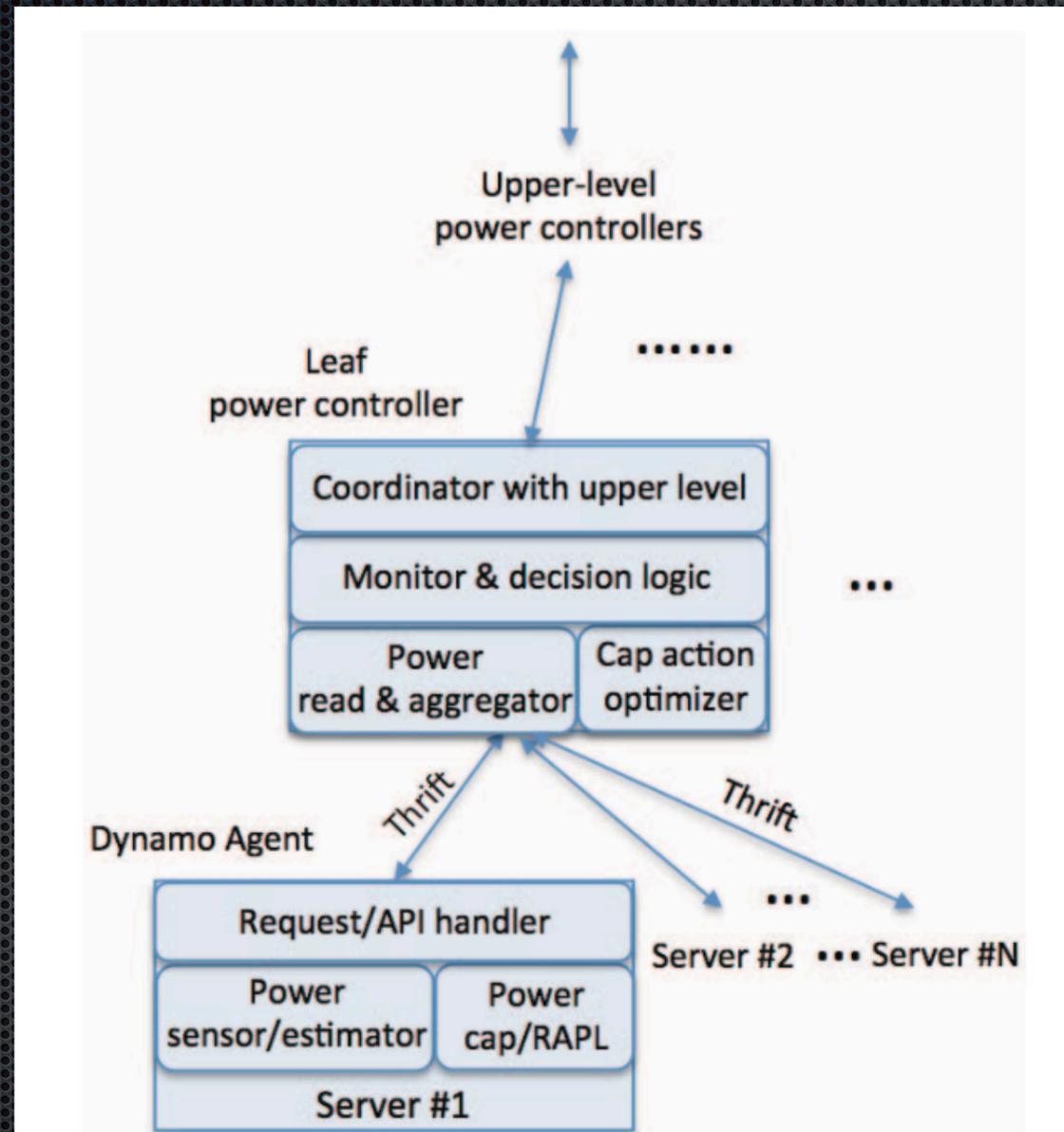


Figure 7. An illustration of how Dynamo's major components interact with each other. Dynamo has two major components: the agent, which runs on every server at Facebook; and the controller, which monitors the power of each device in the power delivery hierarchy.

Power Capping (single)

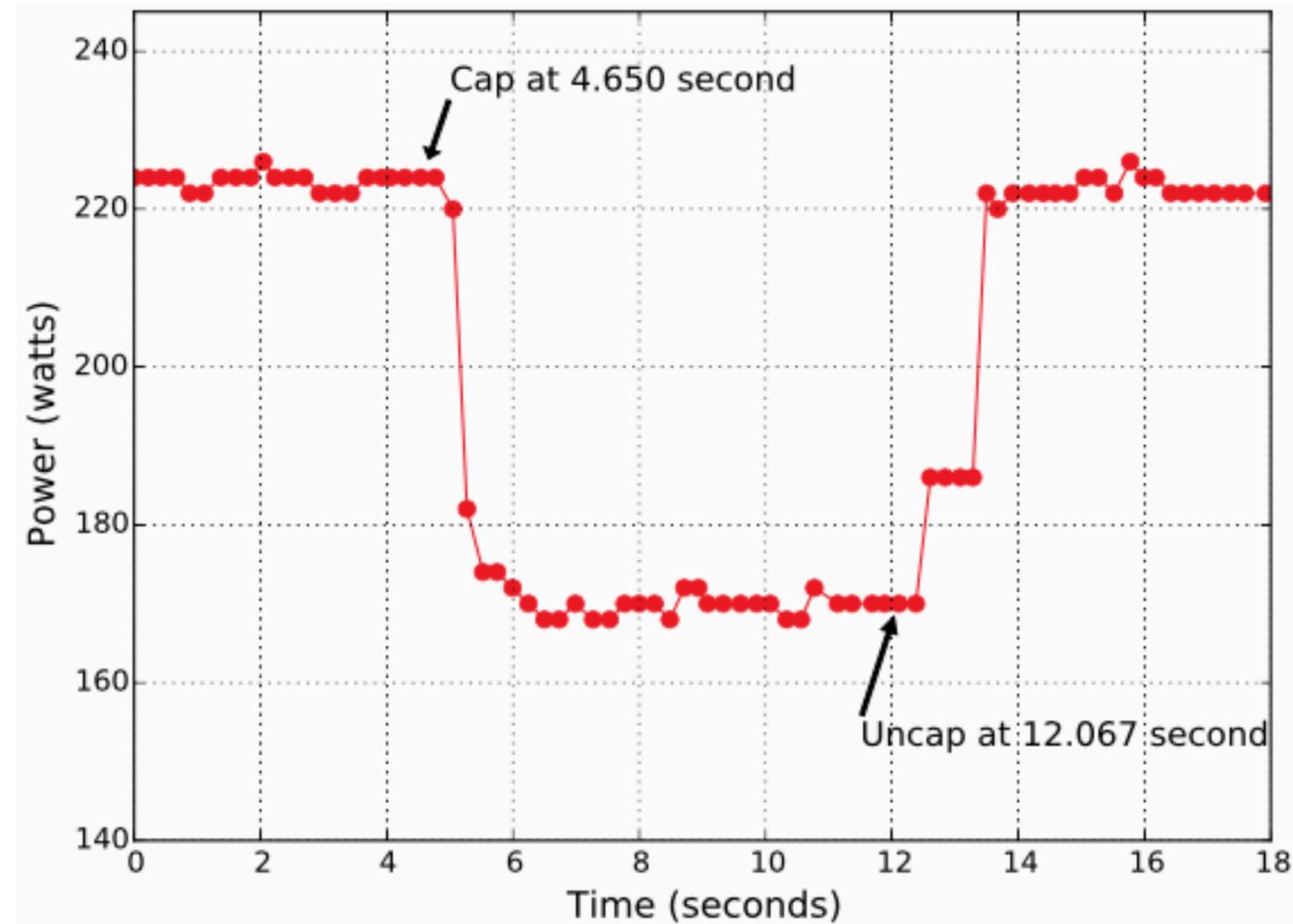


Figure 9. The measured single-server power capping and uncapping test results using Dynamo agent and RAPL.

Example

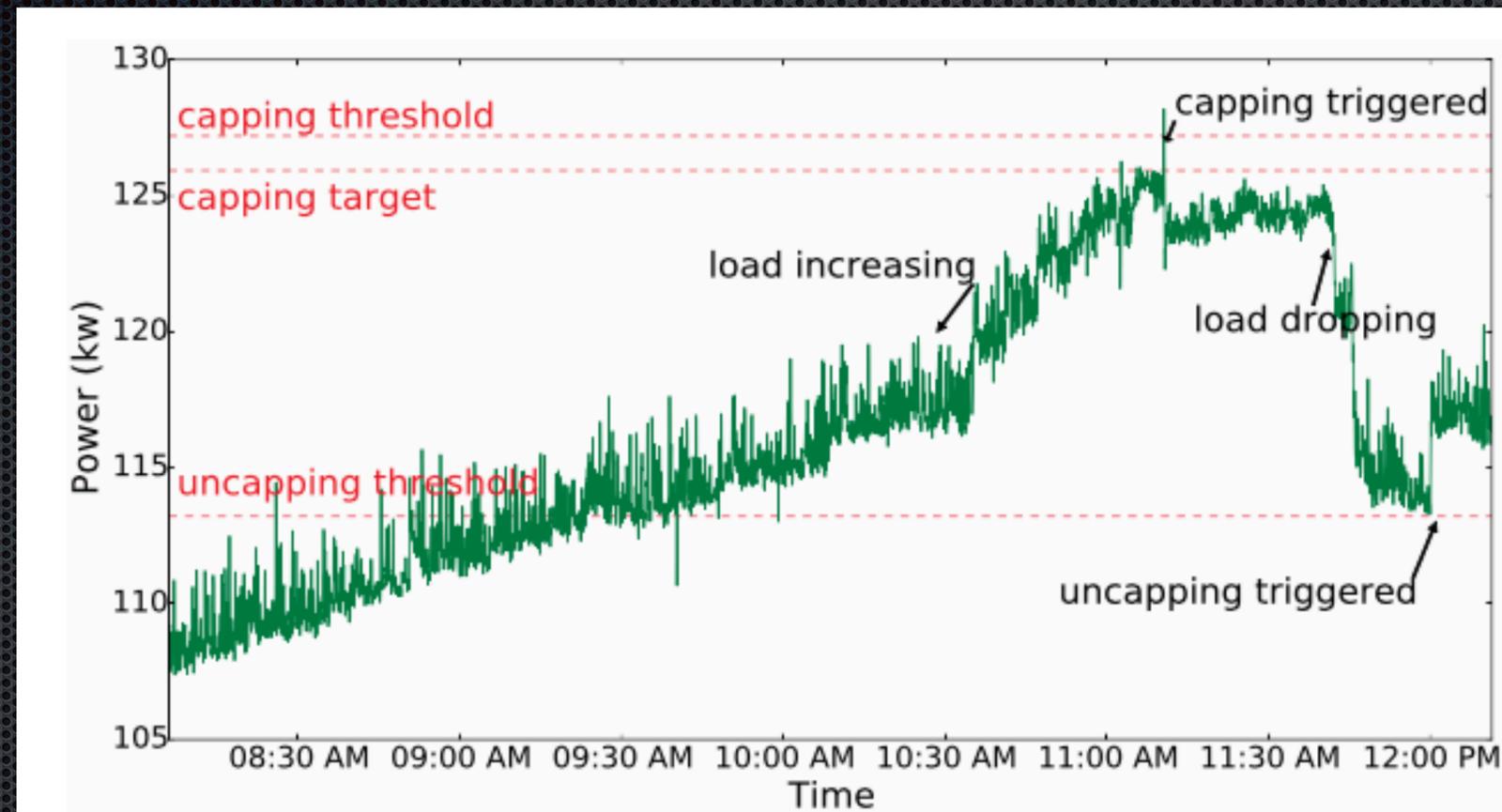


Figure 11. Power measurement showing power capping and uncapping for a row of servers in one front-end cluster located in Ashburn, Virginia.

Surge Protection

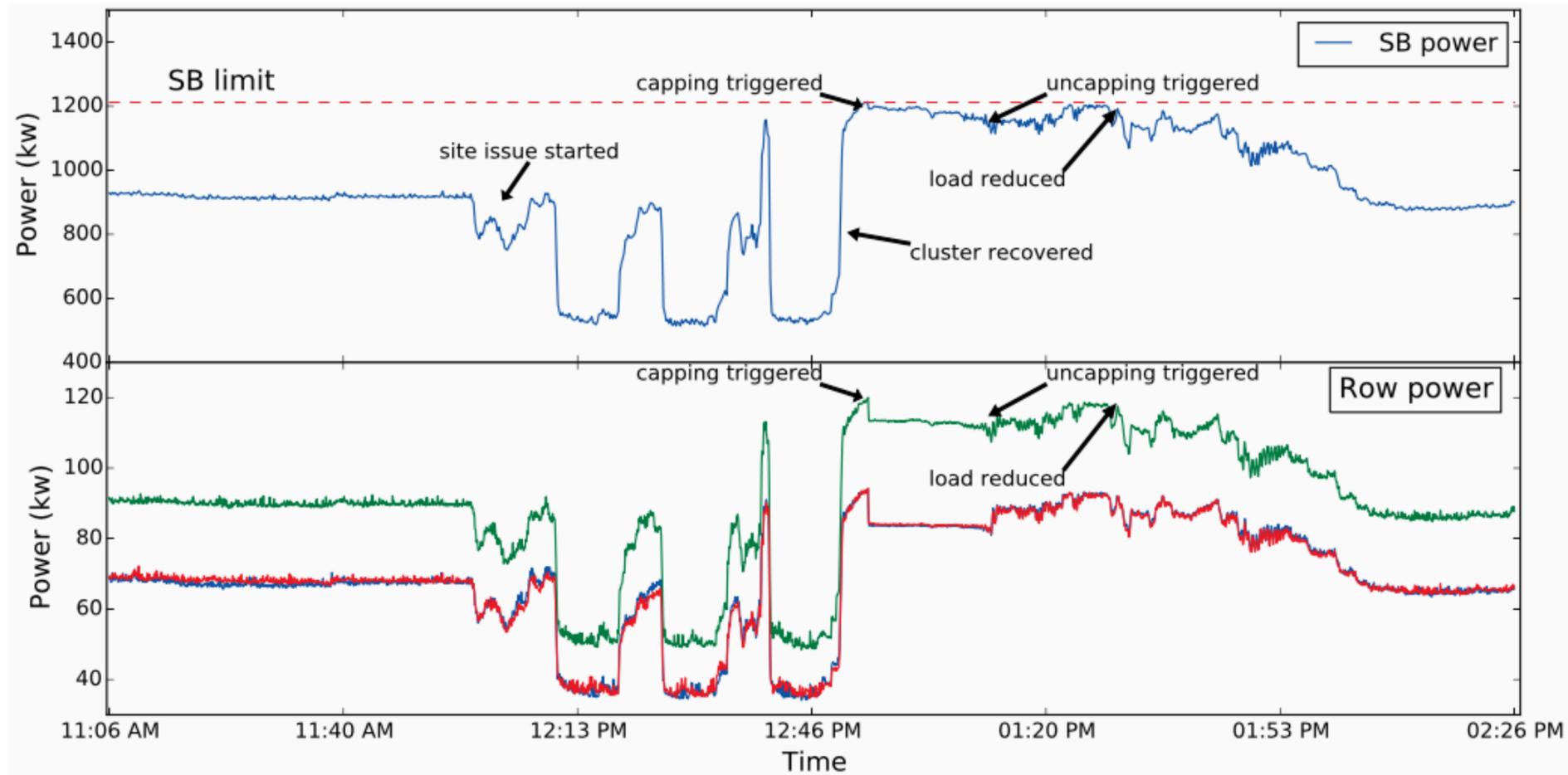


Figure 12. A real-world case study of how Dynamo prevented a potential power outage. A power surge occurred during recovery from an unplanned site issue and led one SB in Facebook's Altoona, Iowa data center to exceed its power limit. An upper-level power controller kicked in around 12:48 PM and three offender rows/RPPs got capped. The upper graph is the power consumption of the SB, while the lower graph is the power consumption for the three rows/RPPs.

Overclock Enabling

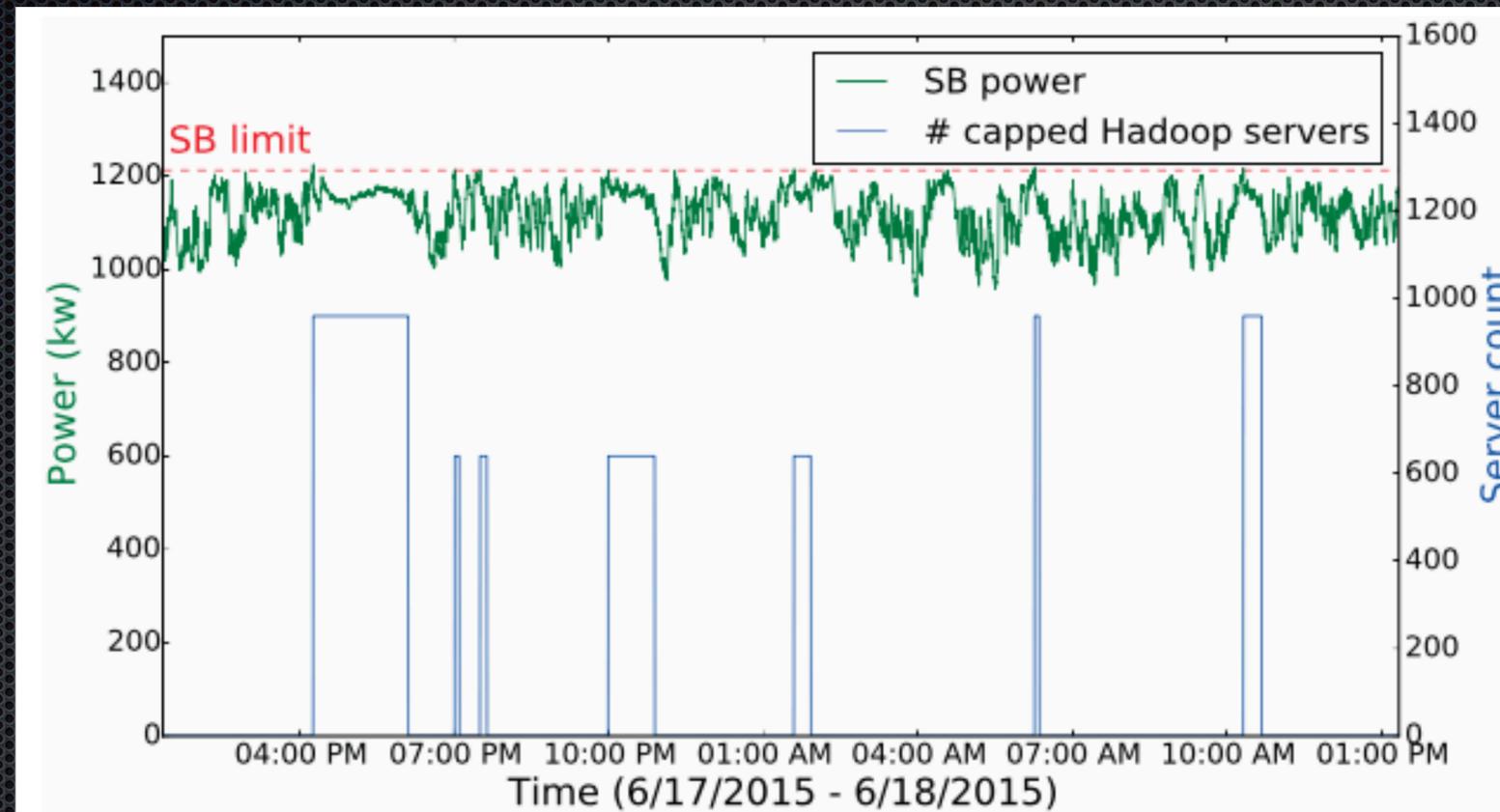


Figure 14. An example showing Dynamo-enabled performance boosting in a production Hadoop cluster in our Prineville, Oregon data center. The top line is the SB power consumption with Turbo Boost enabled for Hadoop servers. The bottom line is the number of servers being capped by Dynamo. The data spans 24 hours.

Conclusions

- ✦ Monitoring is as important as capping
- ✦ Service-aware design simplifies capping testing
- ✦ Design capping systems to be hardware agnostic
- ✦ Use accurate simulation for missing power information
- ✦ Keep it simple to be reliable at scale

Discussion