



# OCP SUMMIT

March 20-21  
2018  
San Jose, CA

**OPEN. FOR BUSINESS.**



# Power Capping and Scheduling on Racks with Flexible Power Supply

Justin Song, Chief Power Architect, Alibaba Cloud

**OPEN. FOR BUSINESS.**





# Agenda

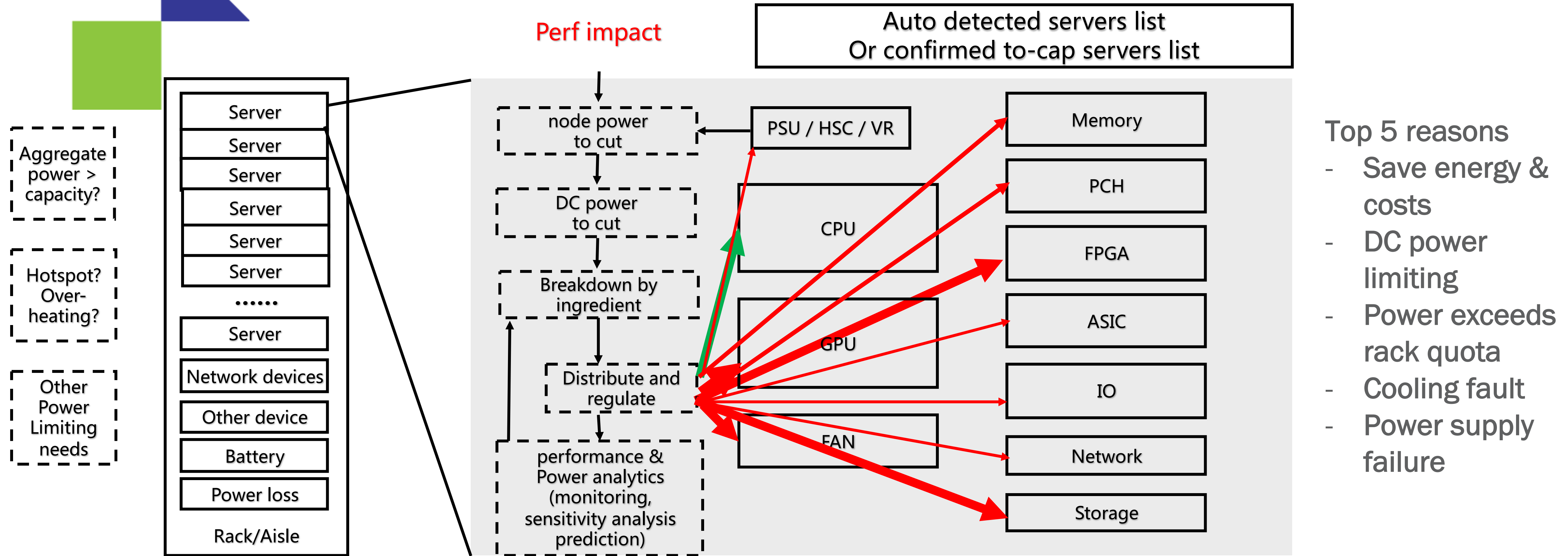
- Rack power capping
- Scheduling on racks with flexible power supply
- Summary & recommendation

**OPEN. FOR BUSINESS.**



# Rack Power Capping

green : h/w or f/w ; red : s/w



Conditions

Entities

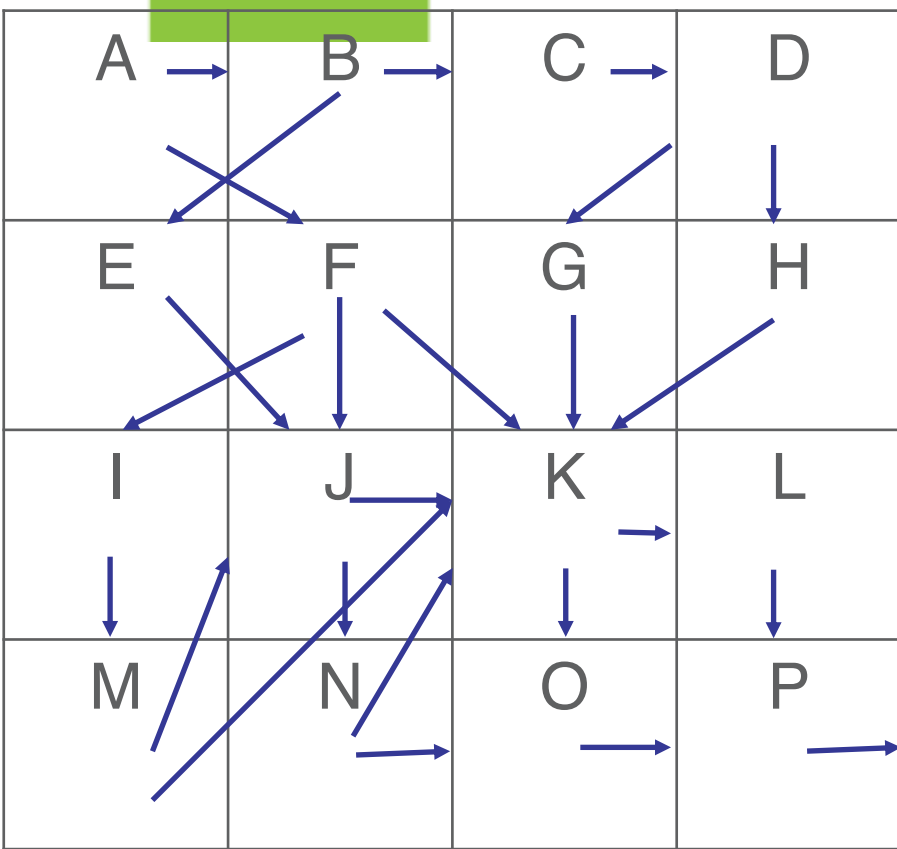
Actions

OPEN. FOR BUSINESS.

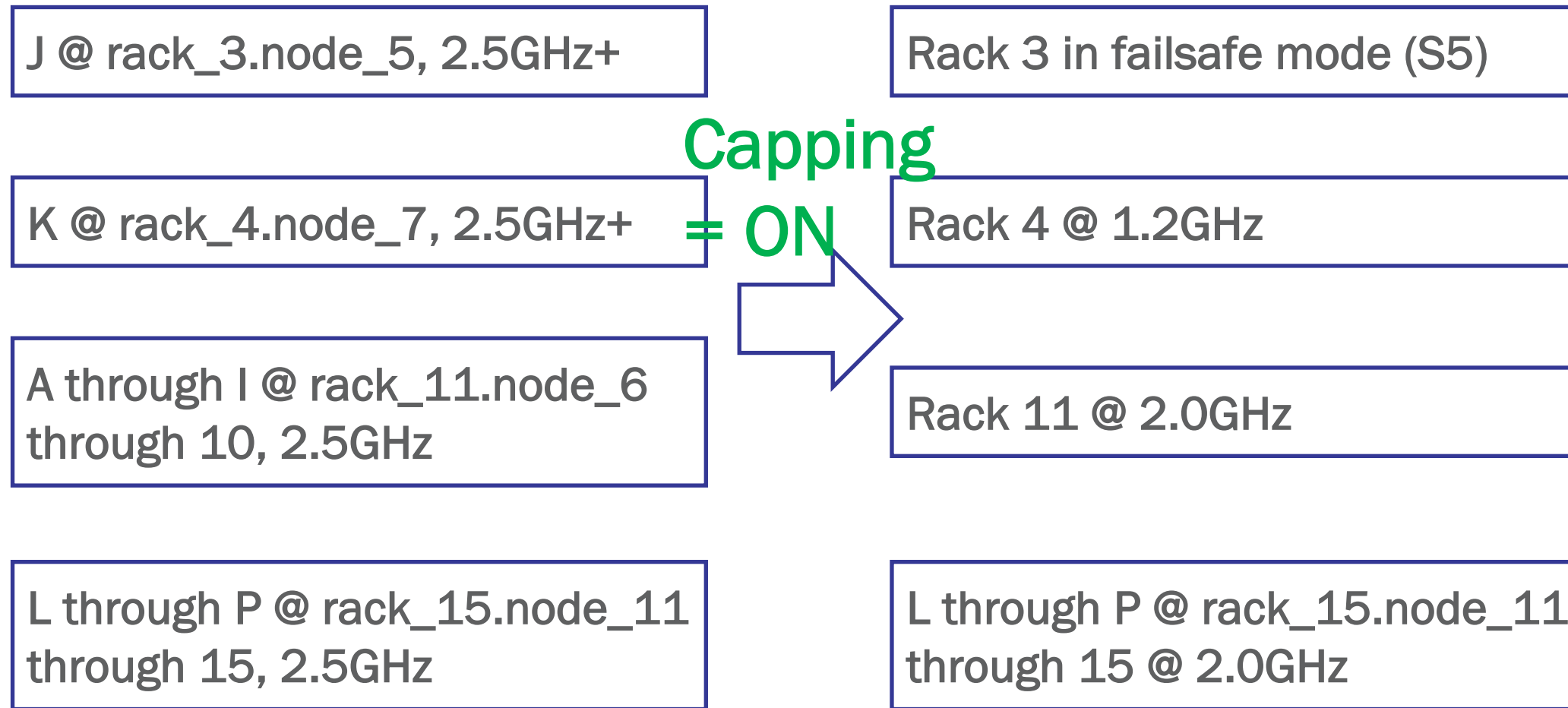


# Capping Helps and Hurts

16 instances and dependency

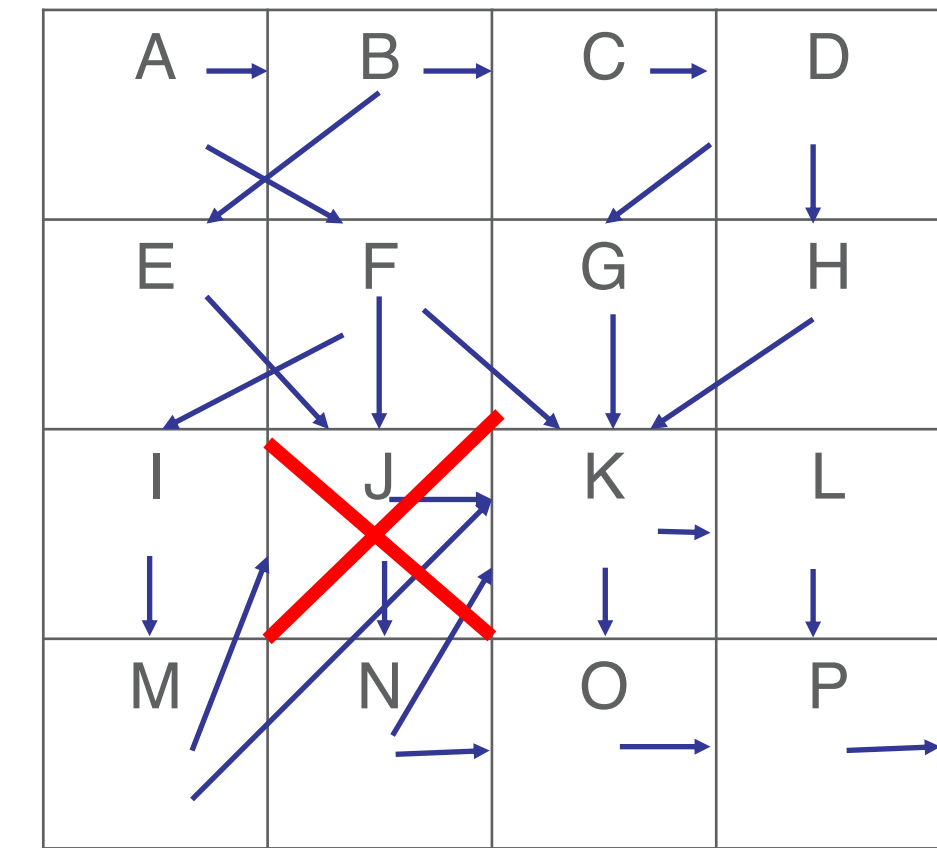


Mapped to 12 nodes @ 4 racks



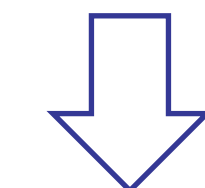
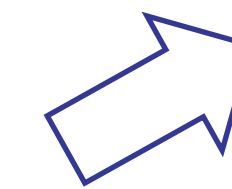
Power capping happens

Consequence T

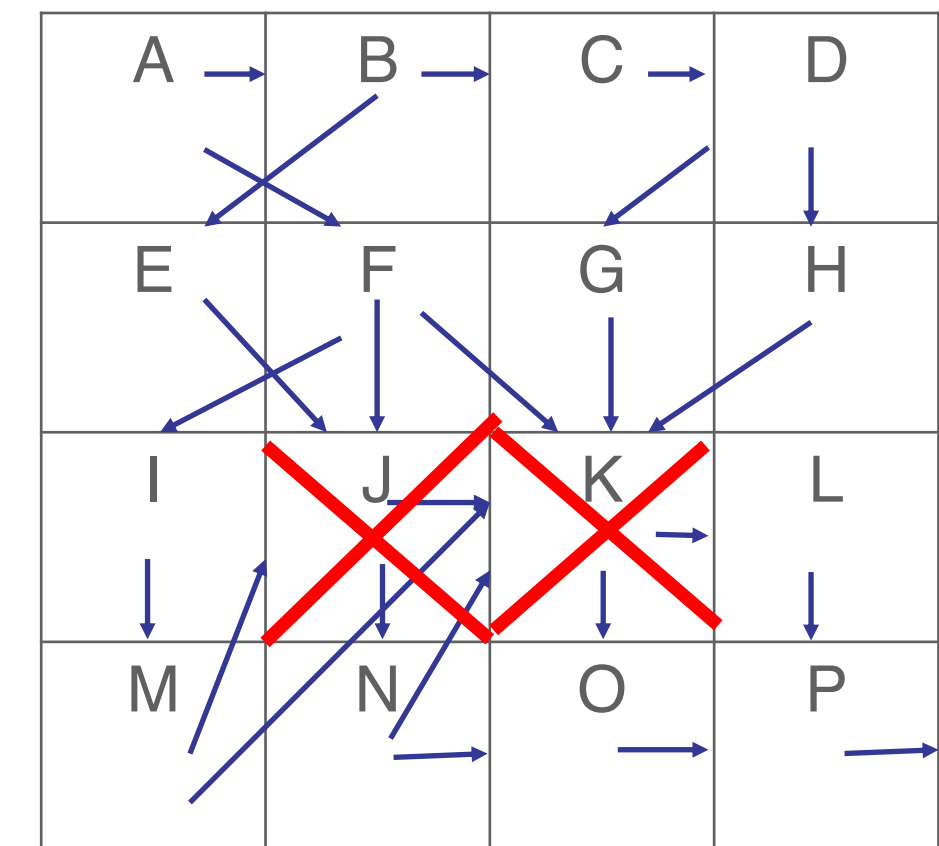


Capping = ON

Capping = OFF

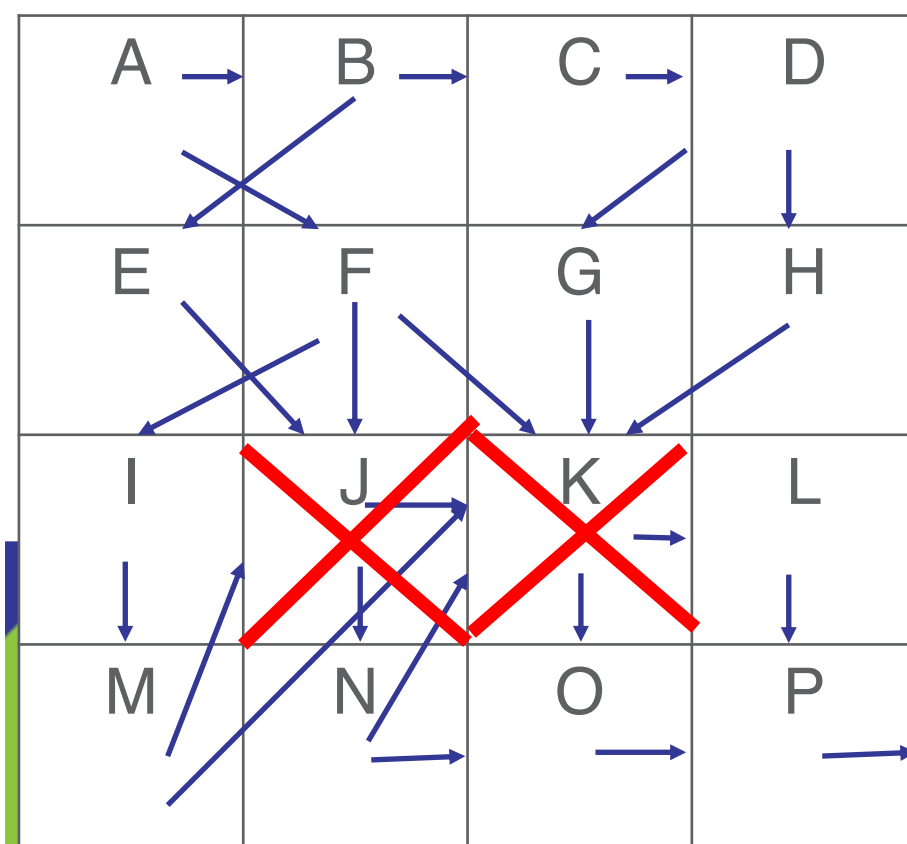


Consequence T+1

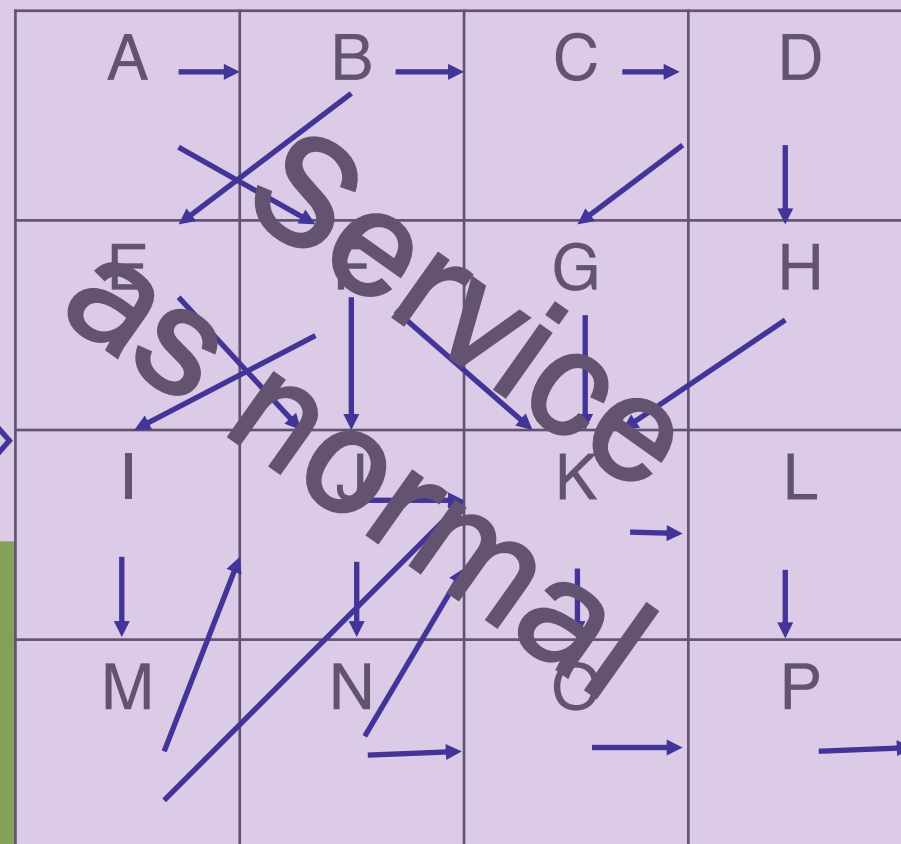


K's buffer blown

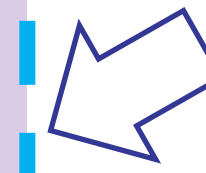
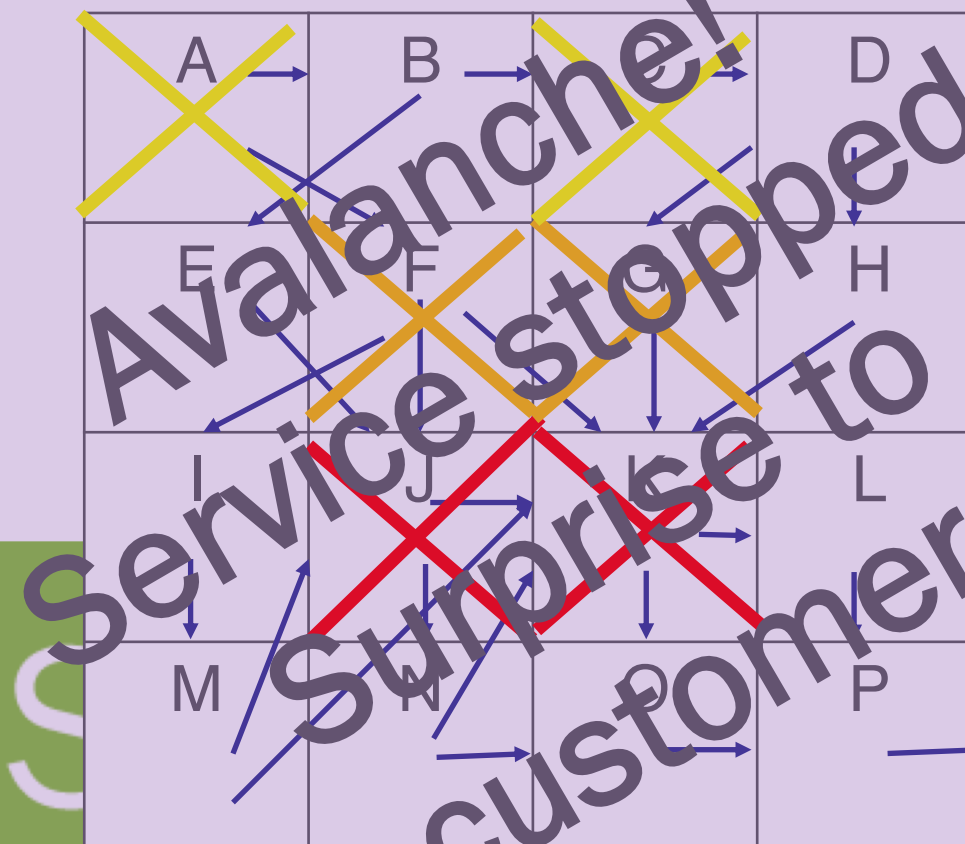
Consequence T+1



Consequence T+2

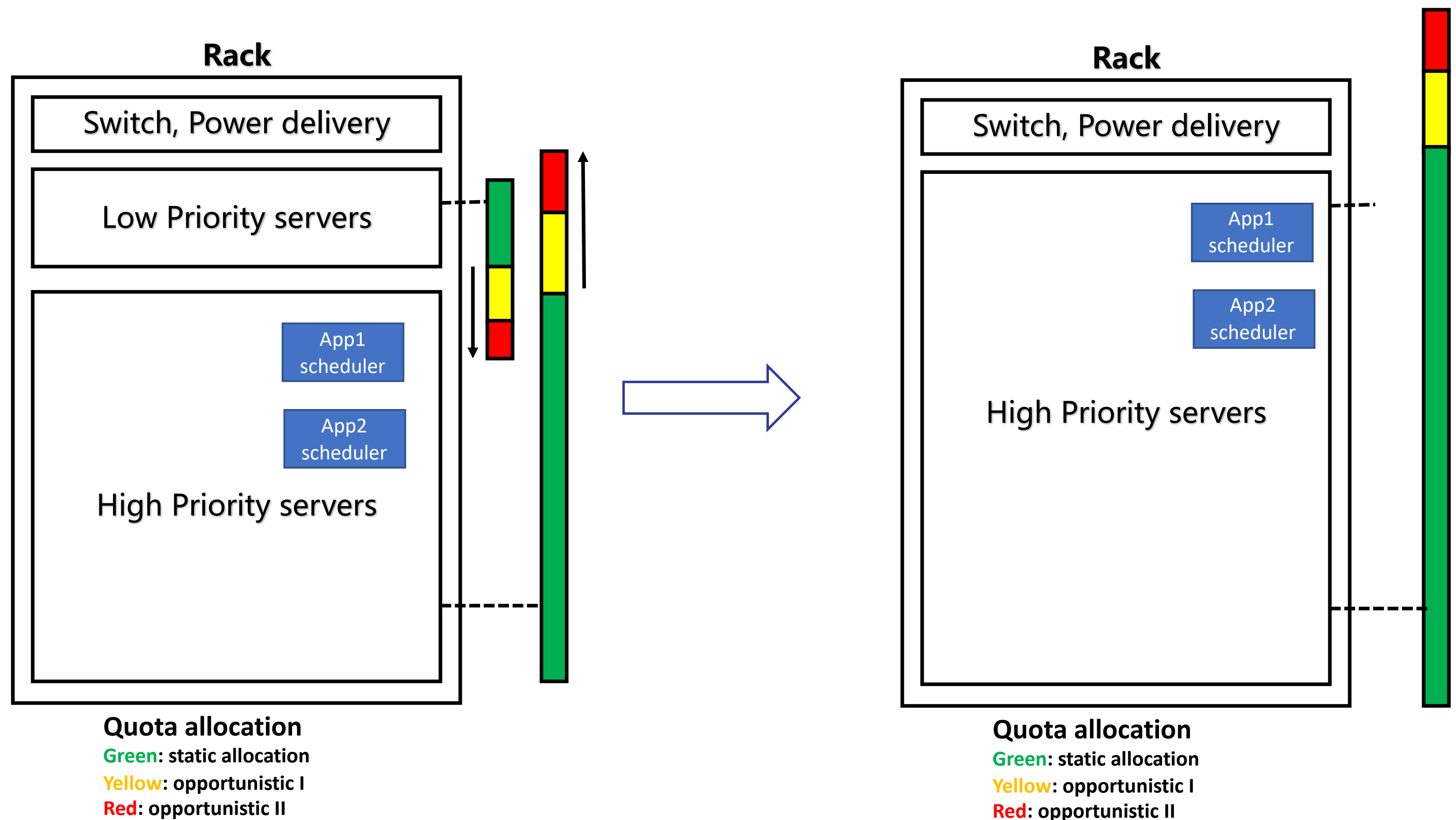


Consequence T+2



Pause, shutdown, migrate and resume

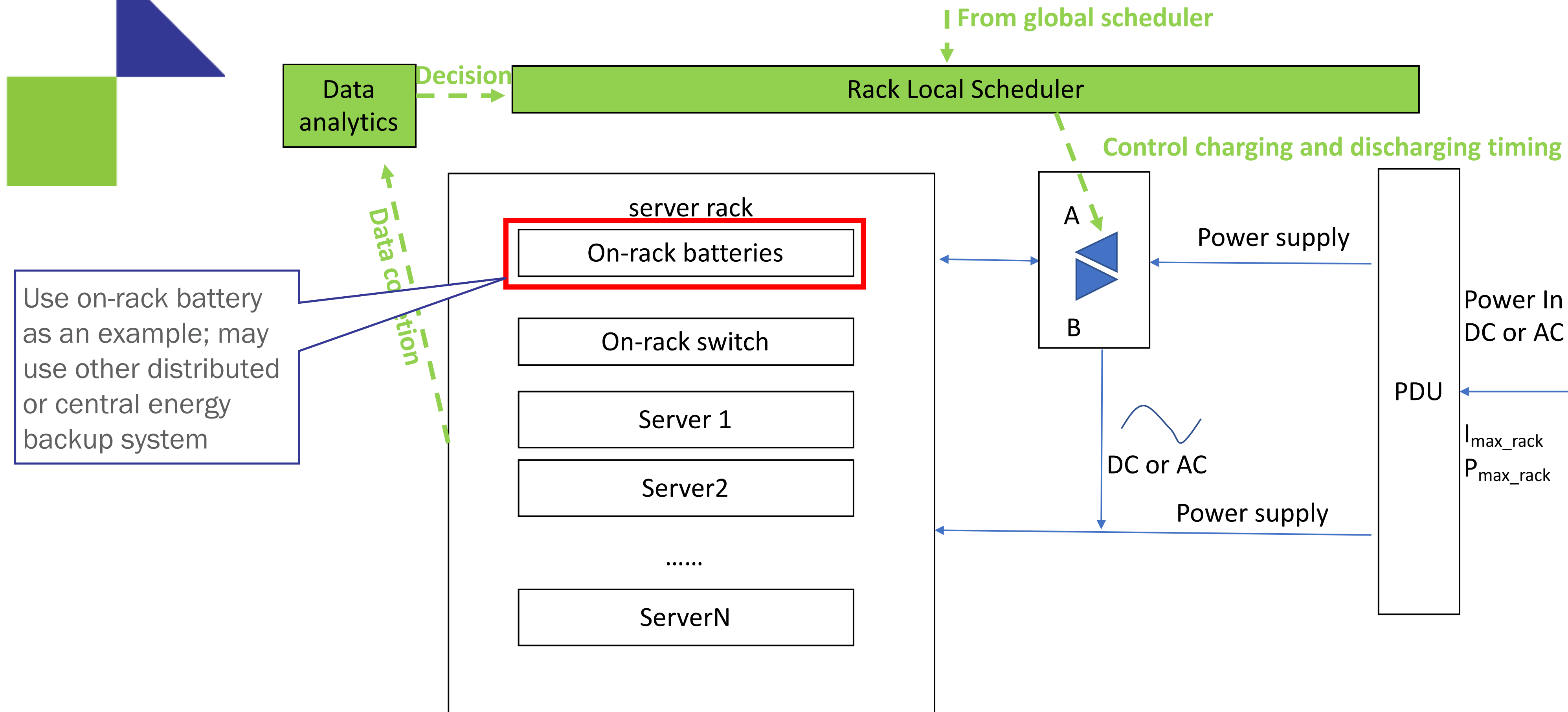
# Another Challenge: Apps Deployment



We need to increase density and simplify apps deployment without compromising SLA

OPEN. FOR BUSINESS.

# Scheduling with Flexible Power Supply



Scheduler + Data Analytics → Use Battery Smartly

OPEN. FOR BUSINESS.

# Manage Timing

## Charging

- Avoid peak usage of server loads
- Charge when energy is cheap

## Discharging

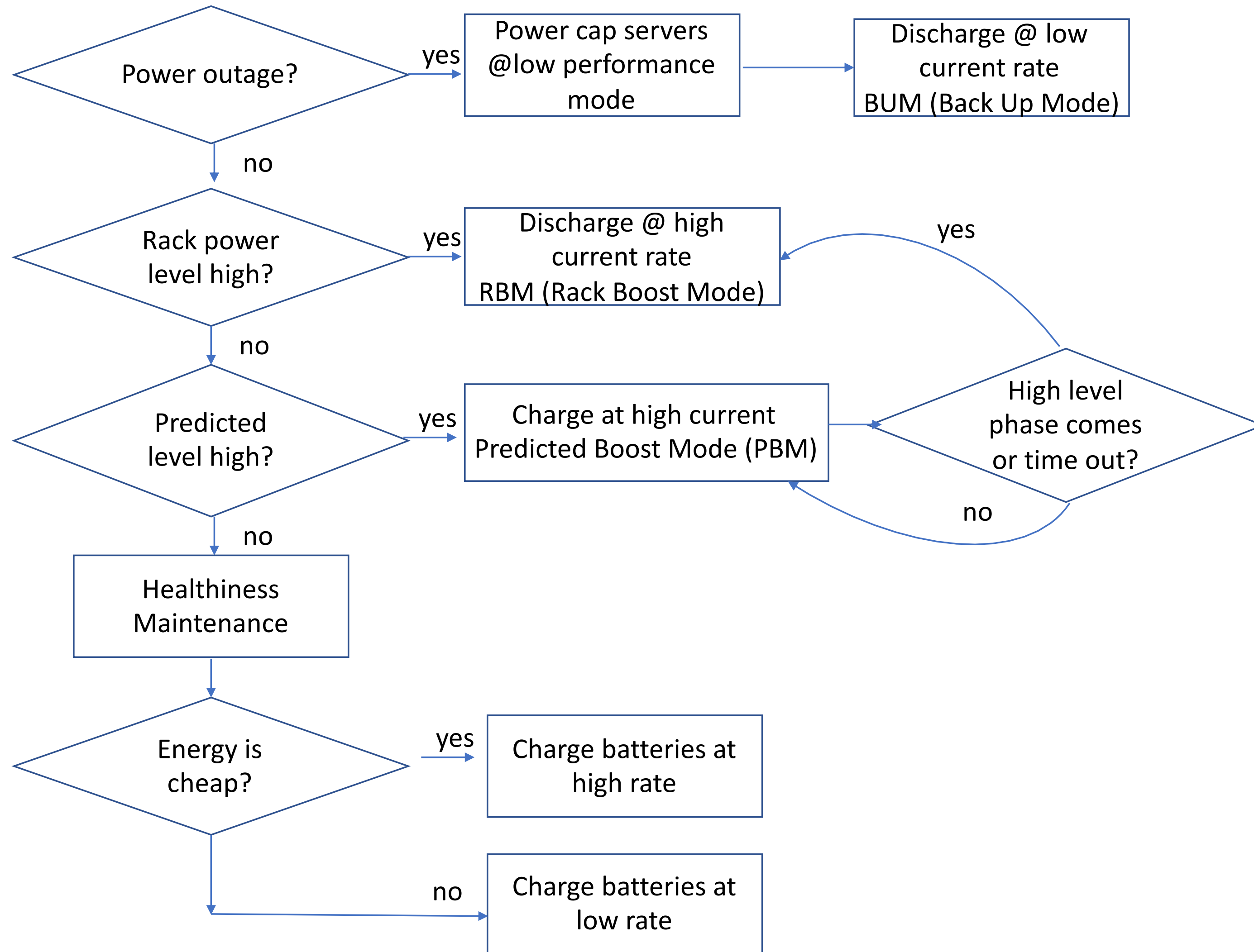
- To accommodate max usage (Rack Turbo)
- Use as a UPS (power shortage or outage)
- Predict duration of high demand

OPEN. FOR BUSINESS.





# Algorithm



OPEN. FOR BUSINESS.



# Summary and Recommendation


Use backup energy source to boost performance or enable higher deployment density

Nimble, flexible and cost effective deployment in data center

Software coordination, scheduling and application awareness is key

**OPEN. FOR BUSINESS.**





# Intel Rack Power Optimization Technology

Nishi Ahuja, Principal Engineer, Intel

**OPEN. FOR BUSINESS.**

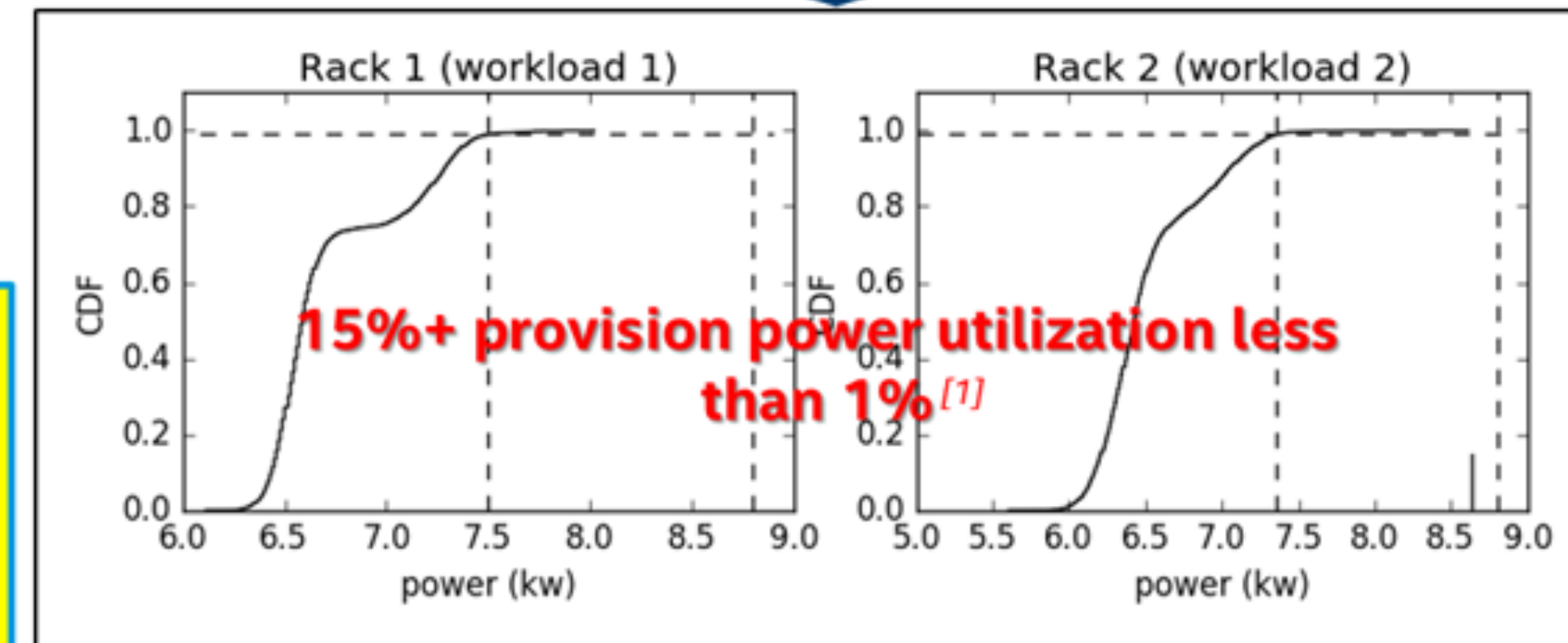
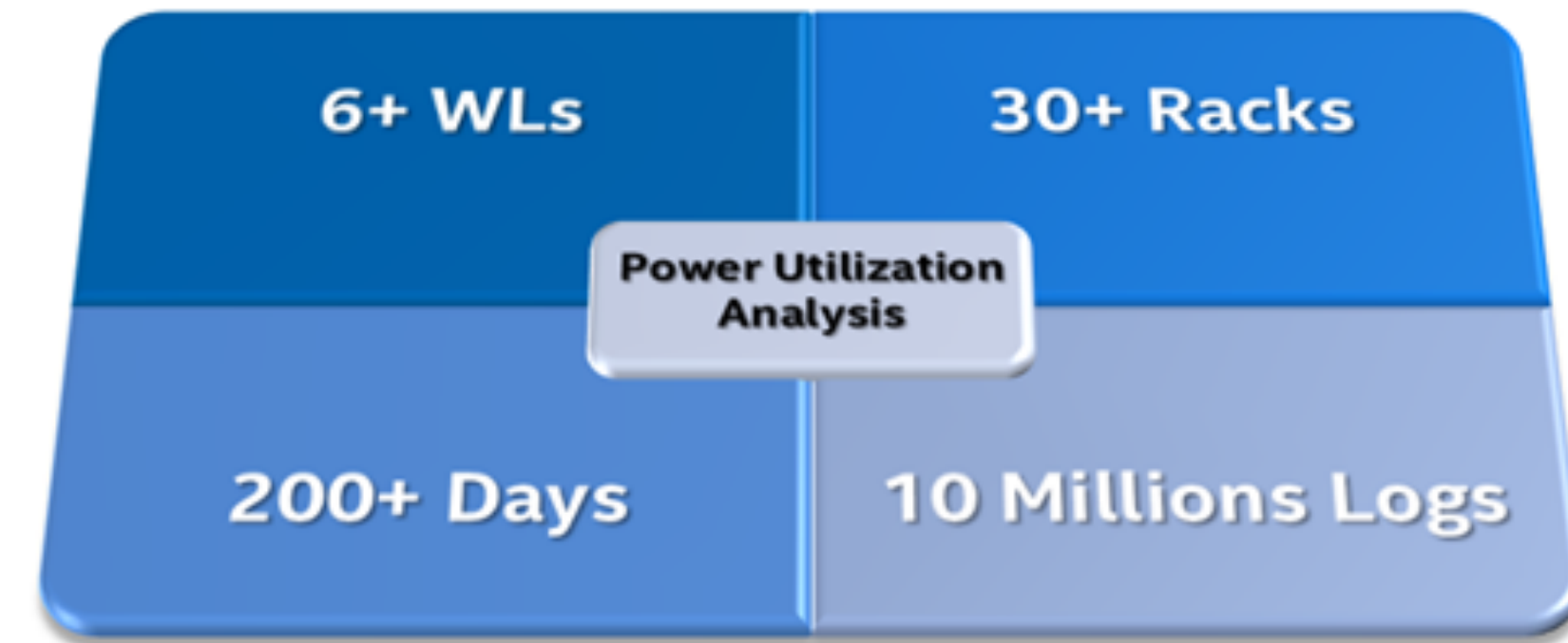


# Rack Density Decreasing vs. Rack Power Under-utilization



## Rack Density is limited by Power

Life Cycle	Stalled	Critical
<ul style="list-style-type: none"><li>• Rack lifecycle is 10~15 years, server 3~5 years in datacenter</li><li>• <b>3 Server generation per rack generation</b></li></ul>	<ul style="list-style-type: none"><li>• Rack Power Density can't be easily increased and <b>stalled at design target</b> (in PRC cloud datacenter, mostly 5-7KW)</li></ul>	<ul style="list-style-type: none"><li>• Server Node Density in Rack is <b>important TCO factor</b></li></ul>



**Problem Statement** – Rack density kept decreasing due to unmatched refresh cycle between infrastructure and IT equipment. The power usage analysis with top cloud service providers also showed the power utilization ratio is lower than then planned due to dynamic range of peak to average.

[1] The evaluation is based on specific CSP workload and Intel Xeon Processors.

OPEN. FOR BUSINESS.



# Node Manager Telemetry on Intel Xeon Scalable Processor

## Power Telemetry

- Total platform power
- Individual CPU, Memory and Xeon Phi power domains

## Thermal Telemetry

- Inlet & Outlet Airflow temperature
- Volumetric Airflow

## Utilization Telemetry

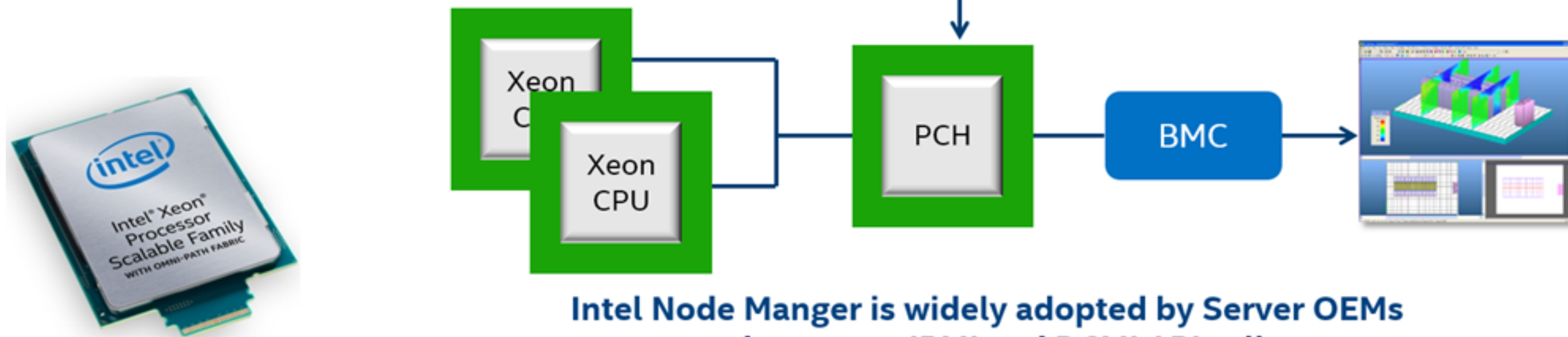
- Aggregate Compute Utilization per sec
- CPU, Memory and I/O Utilization Metrics

## Power Controls

- Power limiting during normal operation
- Power limit during boot

## Intel® Node Manager

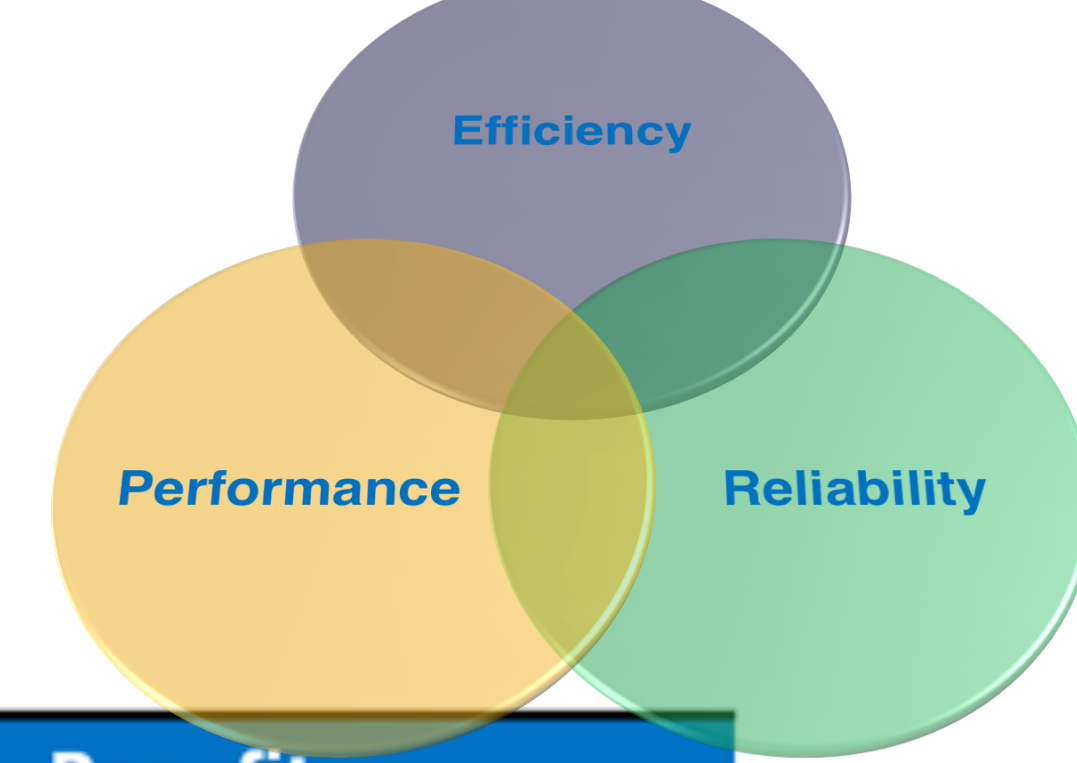
Firmware embedded in PCH's



Intel Node Manager is widely adopted by Server OEMs and supports IPMI and DCMI API calls

OPEN. FOR BUSINESS.

# Intel Rack Power Optimization Technology

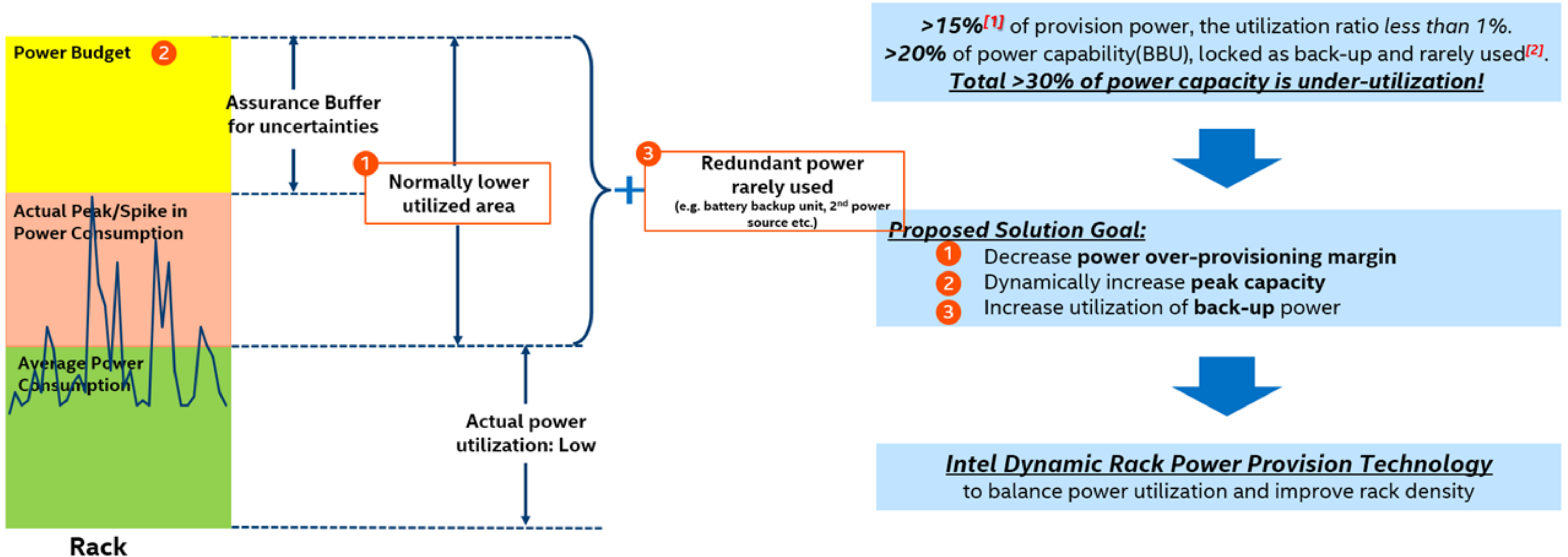


Options	Description	Benefits
<b>Static Power Capping</b>	Configure power budget for each server node to avoid power violations	Rack Power Safety
<b>Group Level Power Capping</b>	Dynamically adjust rack power allocation among servers to enforce rack wide power consumption is within power budget line	Rack Density and Power Utilization
<b>Dynamic Rack Power Provision (DRPP)</b>	Leverage rack battery backup system to shave those sporadic power spikes or peaks over budget line	Rack Density and Service Reliability
<b>Intelligent Orchestration – Power Awareness Scheduling</b>	Schedule workload according to telemetry intelligence, e.g. power awareness job scheduling according to real power demand as well as available power capacity in rack (or cluster)	Energy Efficiency and Service Reliability
<b>Dynamic Cluster Performance management</b>	Dynamically adjust platform performance states according to rack power demands and actual resource utilization.	Performance Per Watt Efficiency

OPEN. FOR BUSINESS.



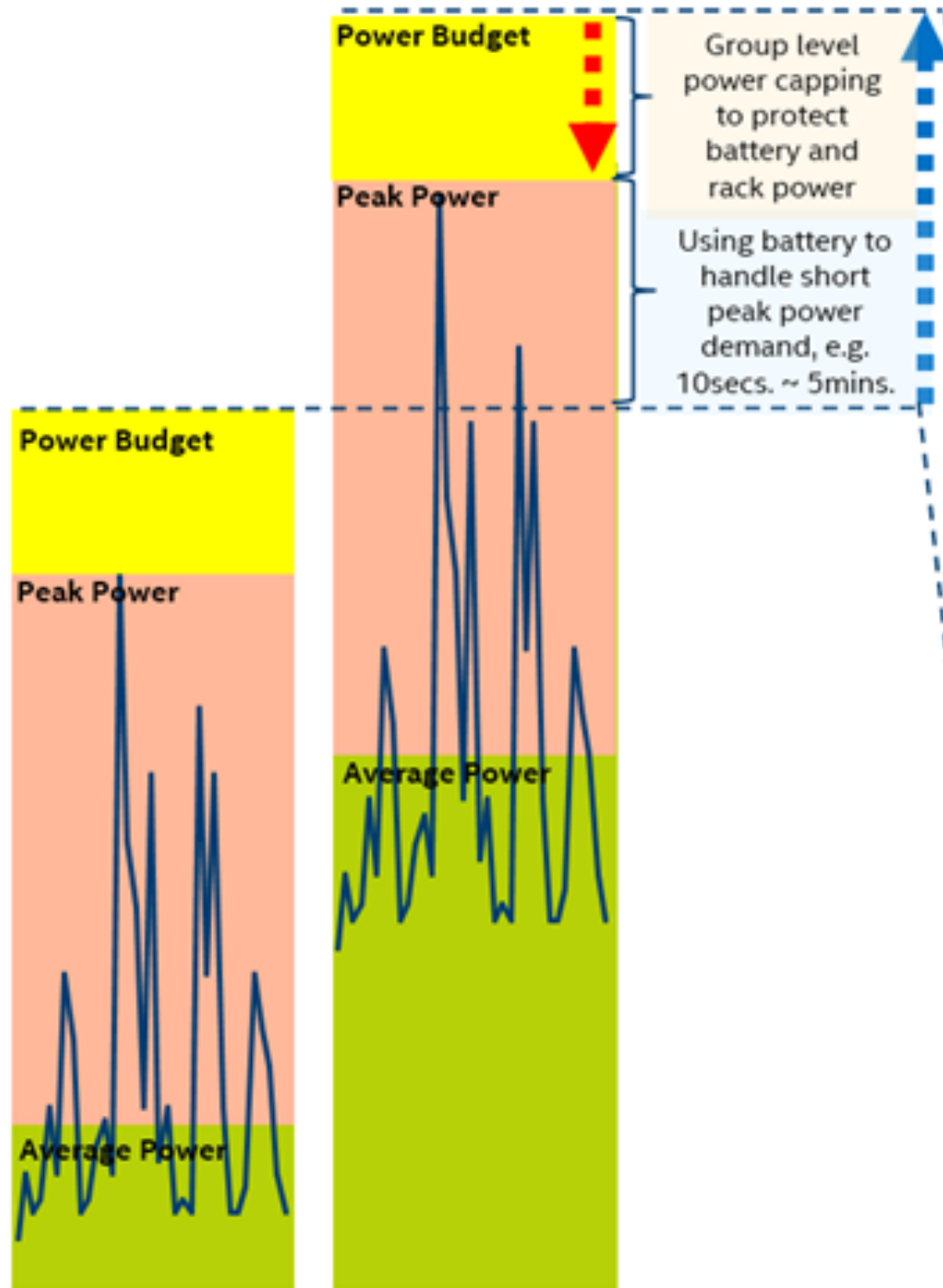
# DRPP - Dynamic Rack Power Provision



[1] - 15% is concluded from analysis of typical online rack system, and subject to change with environment change.  
[2] - The general design of backup power unit capacity is same as rack power capacity, and 20%~30% of backup power capacity is rarely used even under power failure condition.

OPEN. FOR BUSINESS.

# Turbo Rack - Innovative, Non-Disruptive Solution



Group level power capping to protect battery and rack power

Using battery to handle short peak power demand, e.g. 10secs. ~ 5mins.

**Software - Power Awareness Intelligence**

- Real-time Power Insights
- Group Power Capping and Dynamic Peak Capacity Provision
- Intelligent Job Scheduling/Migration

**Application Interface - RESTful**

- Contributed as Redfish API
- Easy Application Integration

**Hardware- Distributed Energy System**

- Open Interfaces for 3<sup>rd</sup> party battery system integration
- Unified architecture support, including OCP and conventional rack

**Firmware - Algorithm Optimization**

- Adaptive Power Re-balance Algorithm
- Configurable Power Limiting

Enabled by **Dynamic Rack Power Provision Technology** with improvement on power utilization and rack density **by 15%~25%**<sup>[1]</sup>



Power Shelf for OCP Rack [2]



Distributed Battery Backup System

**Rack**

[1] - 15%~25% is per evaluation analysis from specific workloads(WLs), and subject to change for different environments.

[2] - Refer to <https://datacenterfrontier.com/ocp-compute-rack-makers/>

OPEN. FOR BUSINESS.





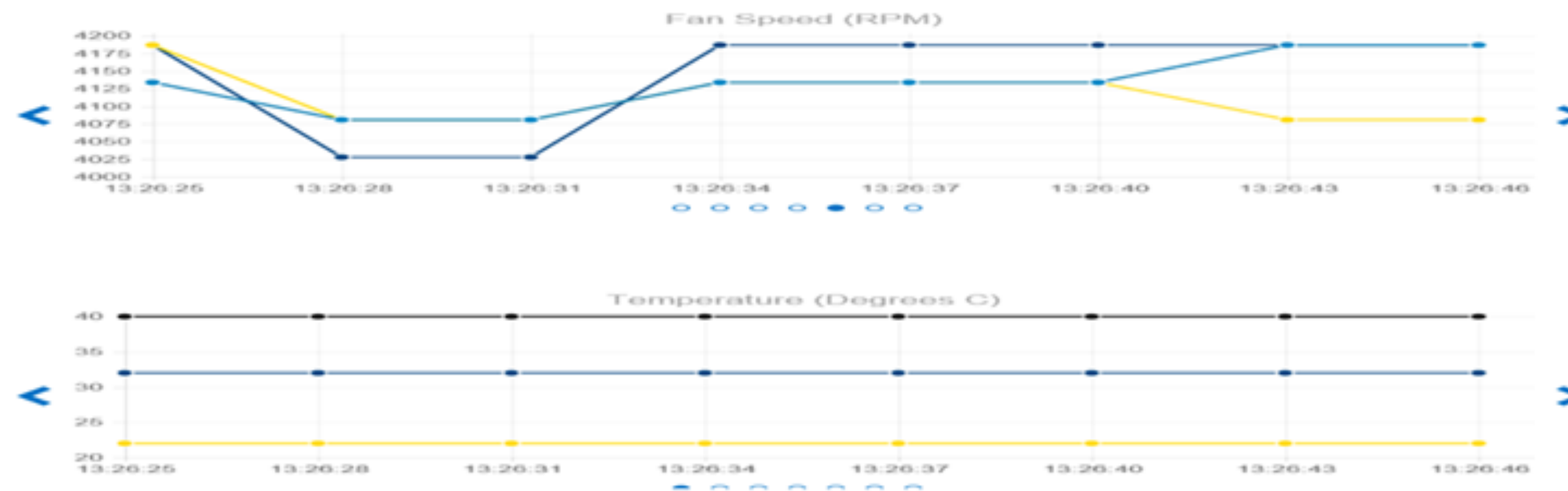
# Intelligent Orchestration based on HW Telemetry & Analytic

Strategy	Descriptions	Telemetry and Analytic
Power Awareness Scheduling	Schedule/Migrate workload according to power intelligence, e.g. energy efficiency, power budget, dynamic power capacity	System/rack power, power analysis.
Uniform airflow/Thermal Condition scheduling	Dynamic adjust workload placement to get one uniform thermal condition to avoid hot-spot in datacenter (with better PUE as well as)	Inlet/outlet temperature, Return Temperature Index(RTI), Rack Cooling Index(RCI)
Failure Event Triggered Migration	Migrate workload to other cluster or rack in case some critical failure events (e.g. power, or thermal) is identified	Failure events

Nodes in Cluster:



Power, Airflow, Compute Usage Per Second (CUPS) Data from Intel® Intelligent Power Node Manager 3.0 PTAS Feature:



OPEN. FOR BUSINESS.





# Turbo Rack Video (Intel)

2minutes 30 seconds

**OPEN. FOR BUSINESS.**





# OCP SUMMIT