

Software for sustainability – Green IT and Sustainable Computing

2023. 9. 27 – ADAC13 Symposium – Environmental Sustainability
Accelerated Data Analytics and Computing Institute (ADAC)

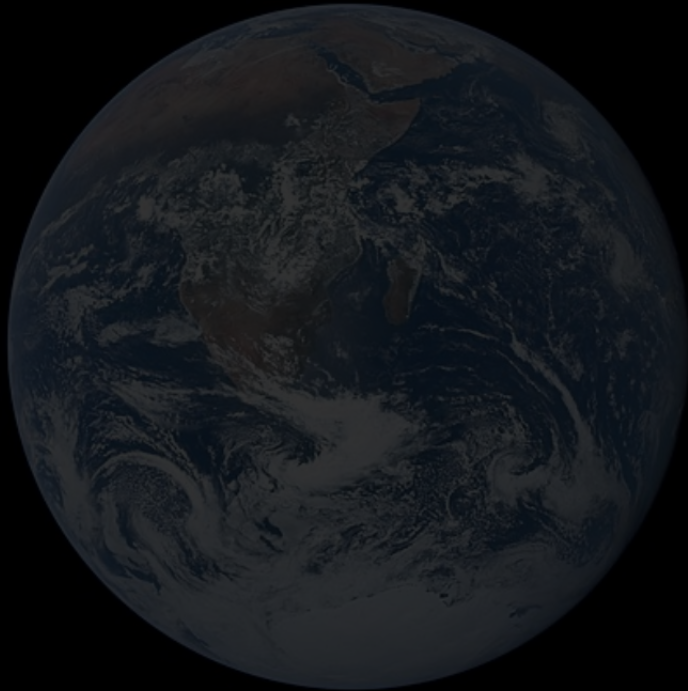
Woong Shin

Ph.D., HPC Systems Engineer
Analytics & AI Methods at Scale (AAIMS)
Advanced Technologies Section, National Center for Computational Sciences

Oak Ridge National Laboratory

ORNL is managed by UT-Battelle LLC for the US Department of Energy

This work was supported by, and used the resources of, the Oak Ridge Leadership Computing Facility, located in the National Center for Computational Sciences at ORNL, which is managed by UT Battelle, LLC for the U.S. DOE (under the contract No. DE- AC05-00OR22725).



Saving our earth...
Legacy to our next generation...

High Performance Computing (HPC) **Post-Exascale** Energy Efficiency

Energy Efficiency in
HPC has Many Faces

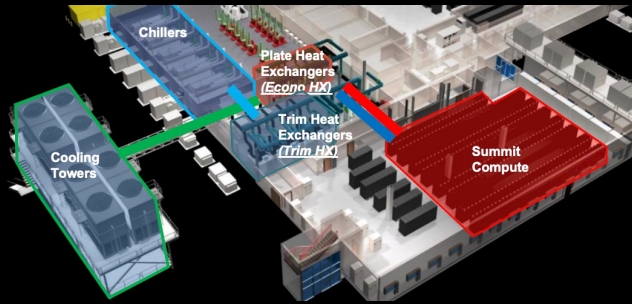
“There are still work to do!!”

Efforts in establishing user facility
support for energy efficiency

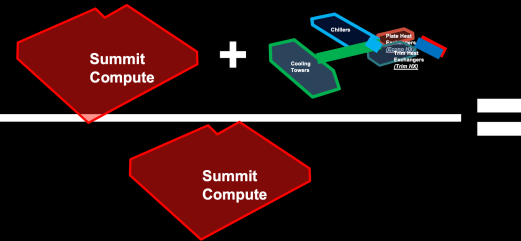
Towards “data-driven” application
level HPC energy efficiency

HPC Energy Efficiency

Which Energy Efficiency Are we Talking About?



PUE =
 Power Usage Effectiveness



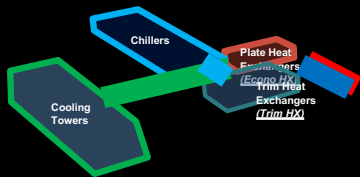
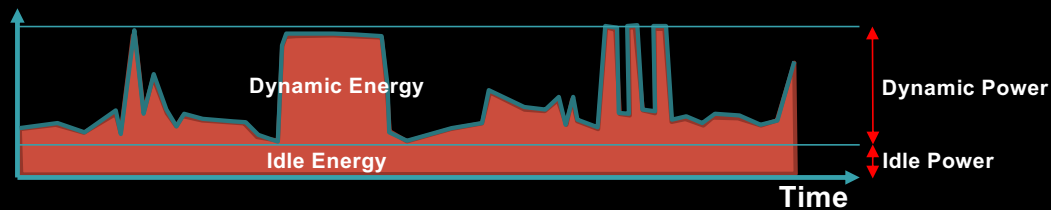
1.0+N
More Energy

Less Energy
1.0

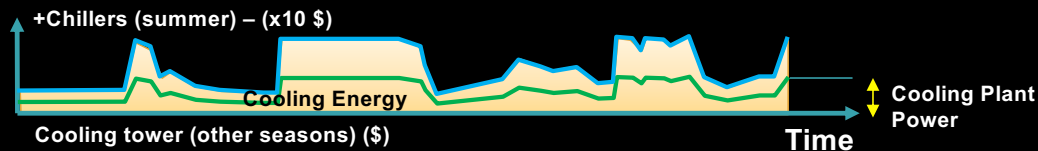
PUE = 1.03 ~ 1.1



Power For Compute



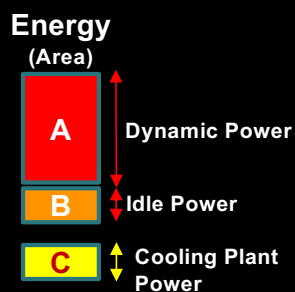
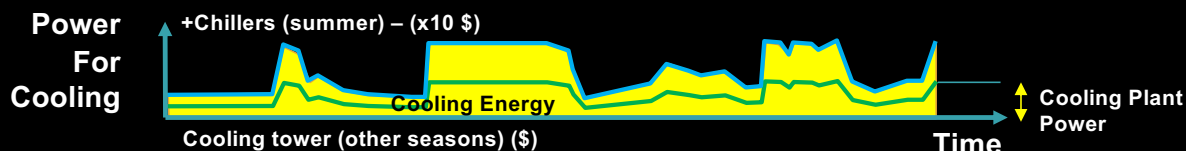
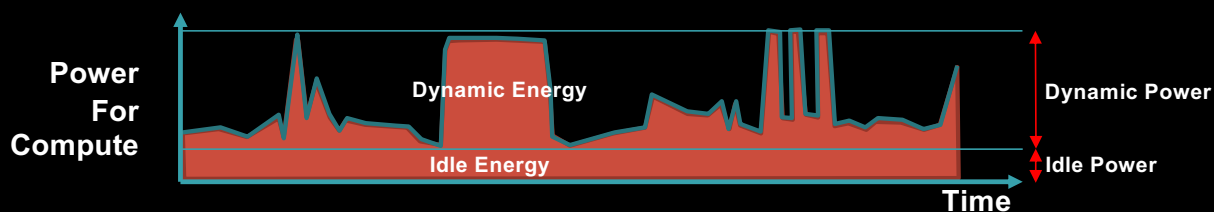
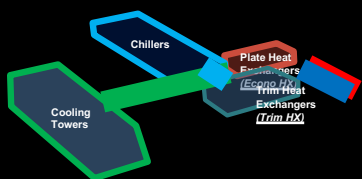
Power For Cooling



HPC Energy Efficiency

Which Energy Efficiency Are we Talking About?

$$\text{PUE} = 1.03 \sim 1.1$$



$$\text{HPC Energy Efficiency} =$$

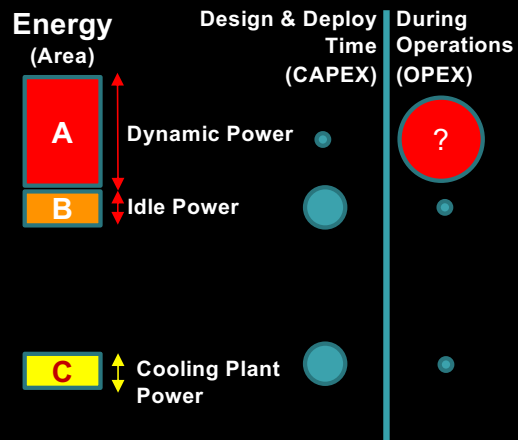
Useful Work

Total Energy

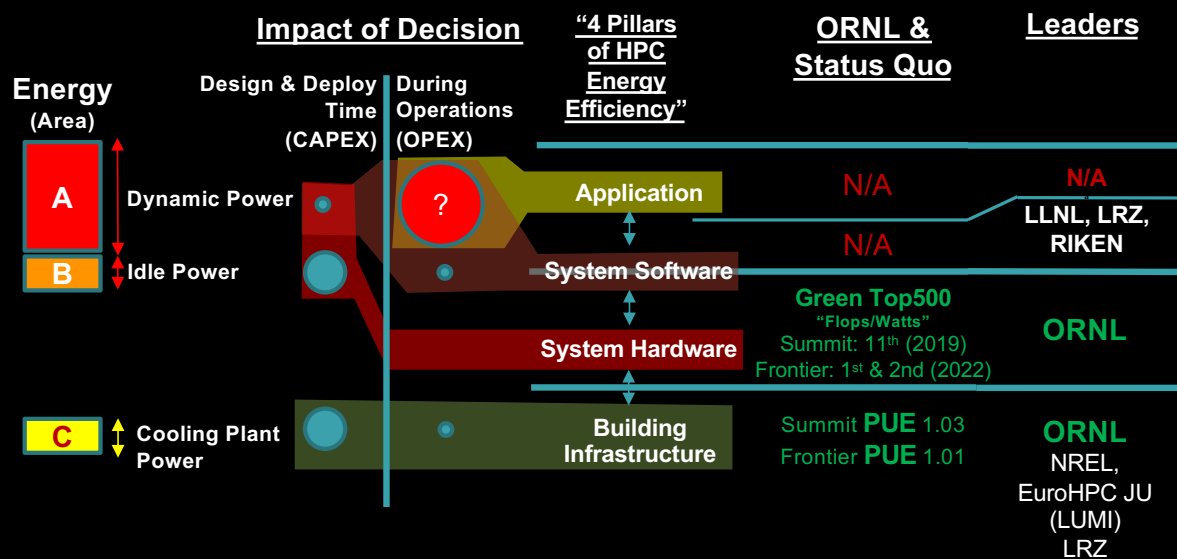
$$\text{Dynamic Power} + \text{Idle Power} + \text{Cooling Plant Power}$$

HPC Energy, When?

Impact of Decision

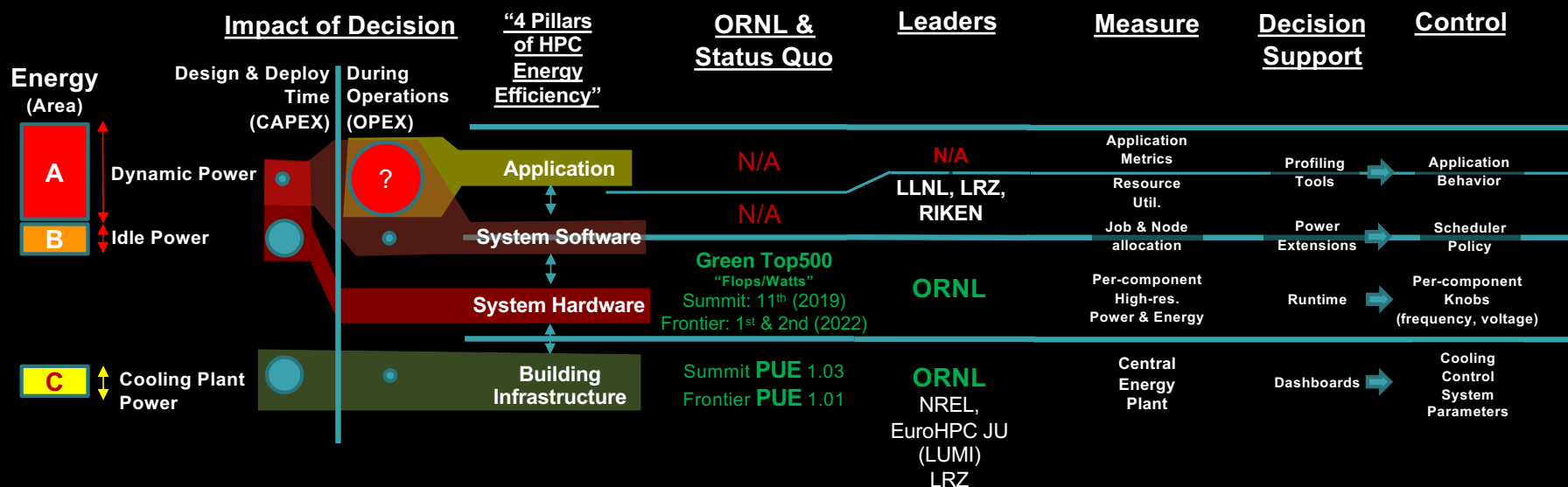


HPC Energy, When, Where?

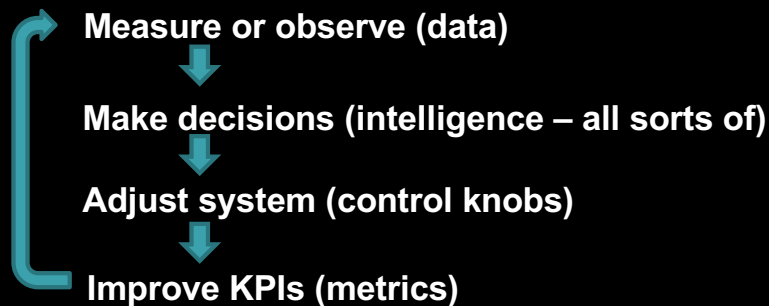


Optimizing during Operation?

Optimizing HPC Energy Efficiency

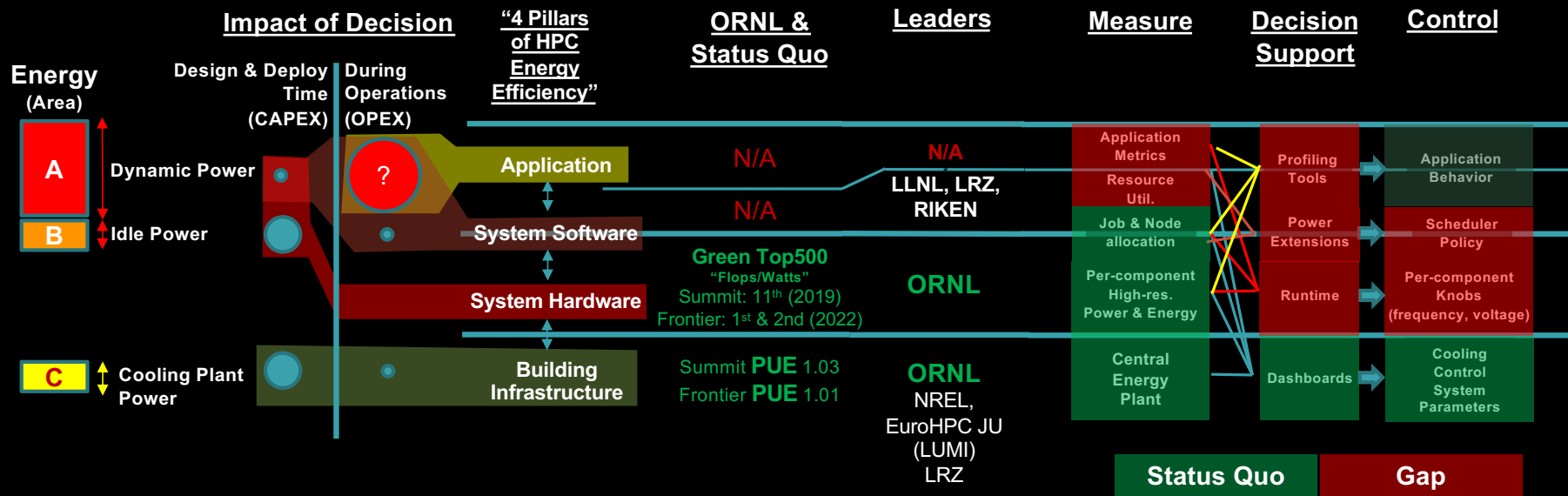


Continuous Improvement Loops

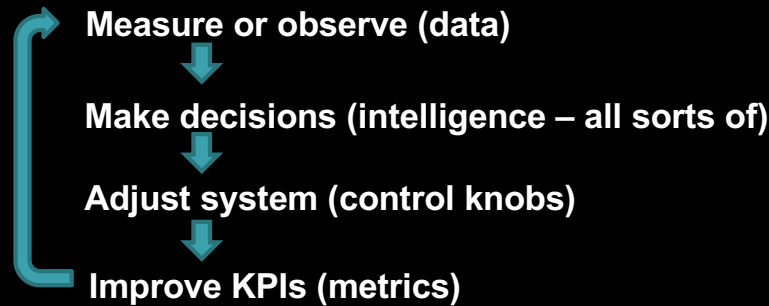
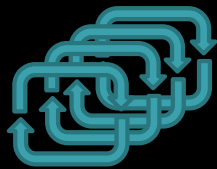


Gaps & Opportunities

Optimizing HPC Energy Efficiency



Continuous Improvement Loops



Post-Exascale Focus Areas towards HPC Energy Efficiency

"Preparation for the future – energy efficiency support"

Optimizing

HPC Energy Efficiency Measure

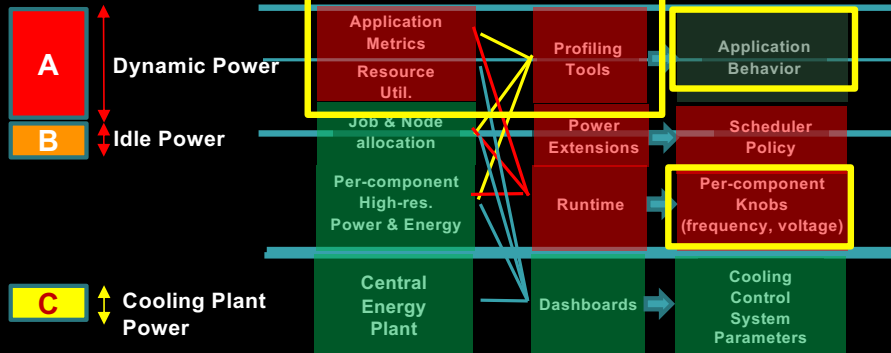
Decision Support

Control

Impact of Decision

"4 Pillars of HPC Energy Efficiency"

Energy (Area)



Potential new focus
"Application driven EE"

Support new focus
"Expose power control knobs"

Maintain course and improve
"New challenges"
"New methods"

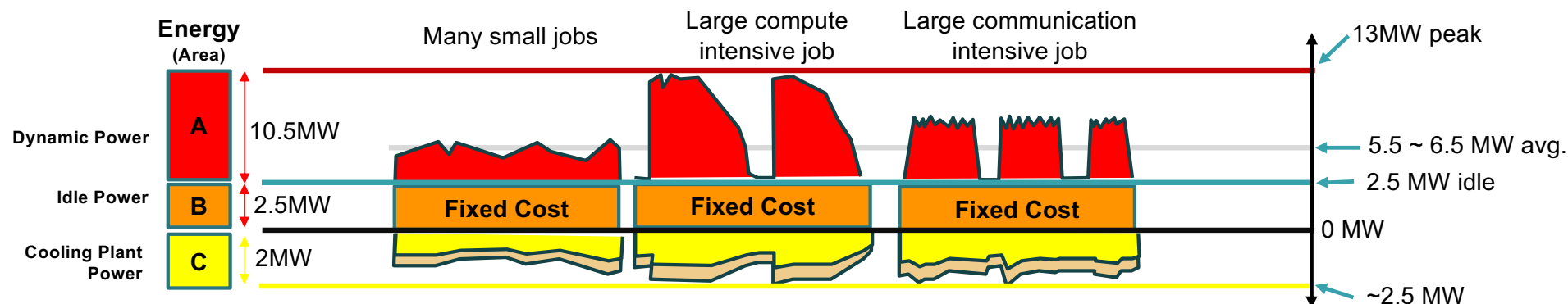
Taming the complexity of control

Energy Efficiency Support For Users



Research and Development
"User empowerment"

Can applications save energy? – Race to Halt (RTH)

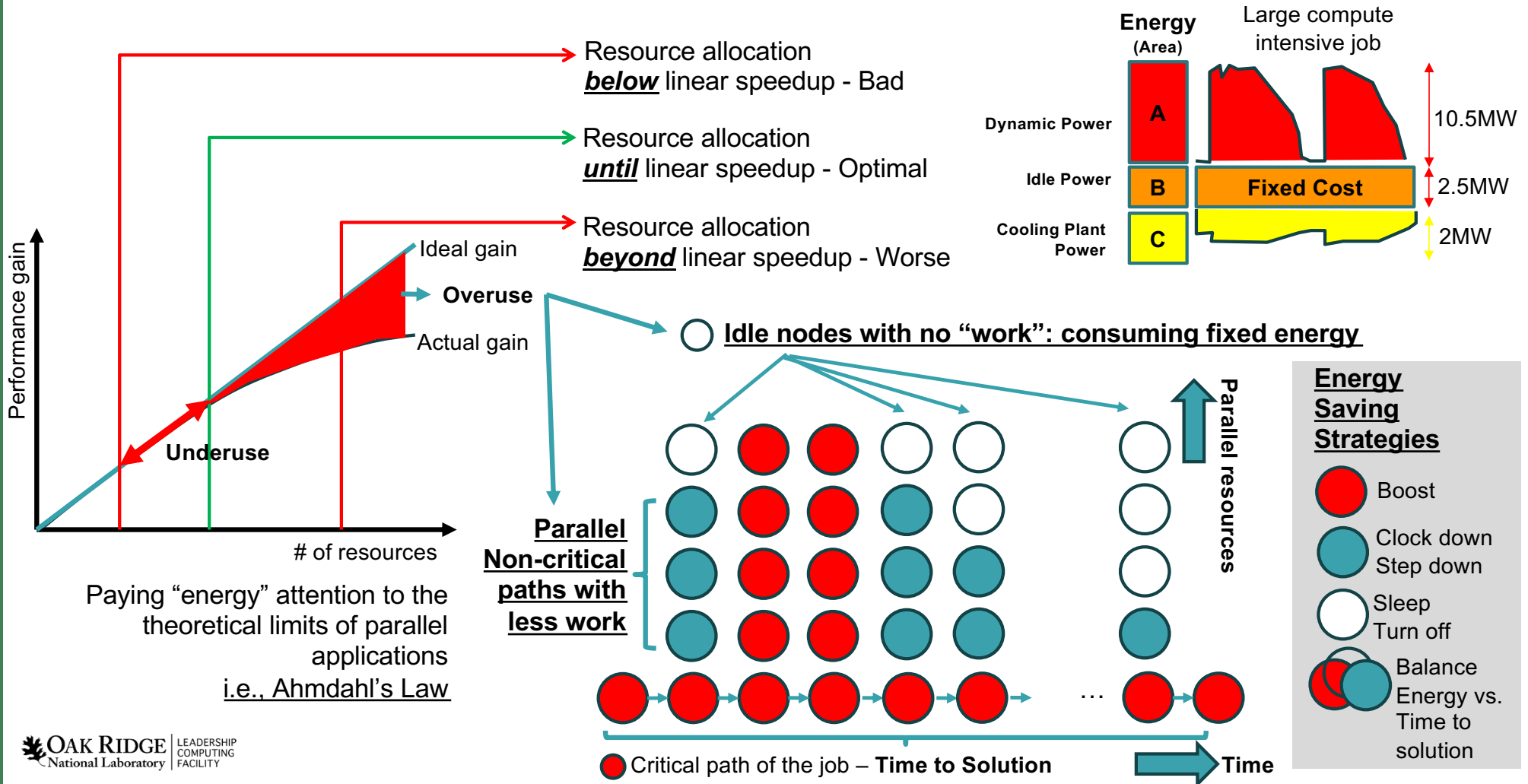


Finishing the compute early consumes the least amount of energy

Business as usual – performance driven optimizations
A good thing!!

But special attention towards energy is required!!

Energy Saving Opportunities on an RTH system - Basics



Behavioral Control Knobs – Strategies for an RTH system

Energy Saving Strategies



Boost



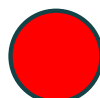
Clock down
Step down



Sleep
Turn off



Balance
Energy vs.
Time to
solution



Rush!!

- Optimize kernels for maximum performance
- Use accelerators to speed up where the app is the slowest



Trade Performance vs. Energy

- Within deadline use the slack to slow down and save energy



Turn off clock down or return unused resources

- Consolidate resources for higher per-node utilization
- Turn-off, sleep or deallocate unused resources

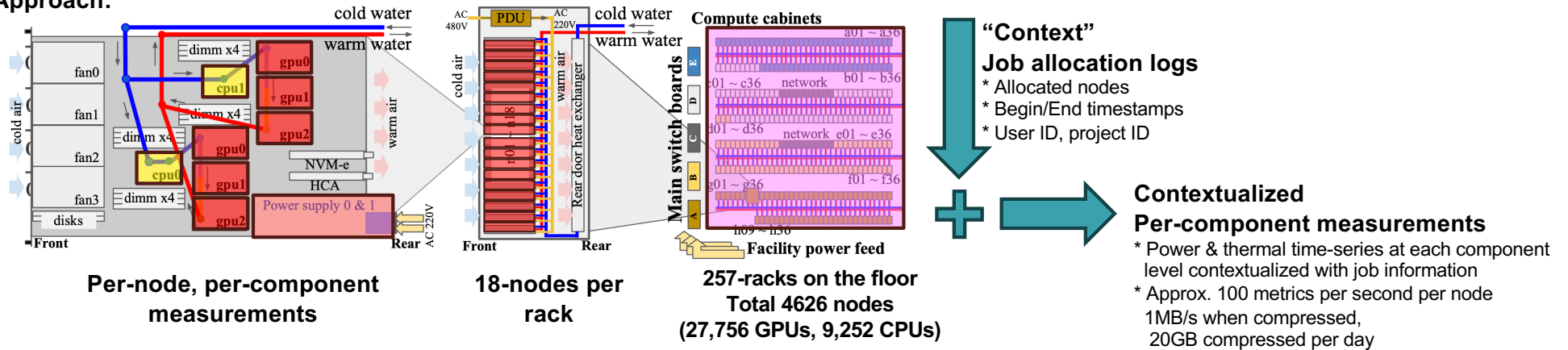


Red vs. others: balance

- Sometimes maybe less parallelism with higher utilization is better in terms of performance vs. energy

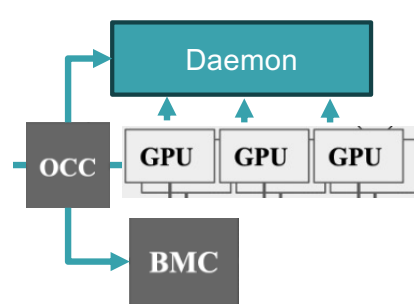
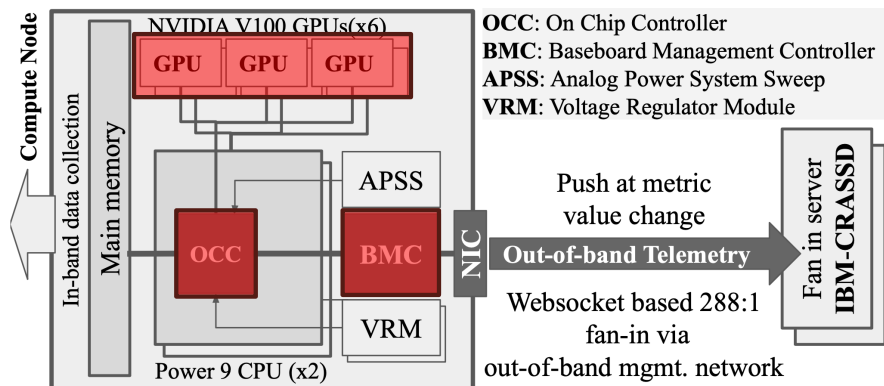
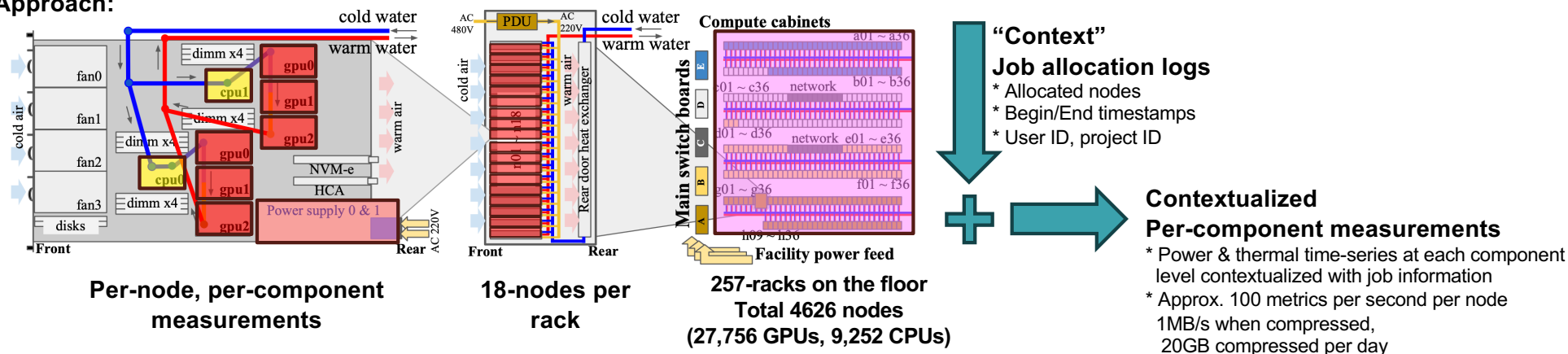
Instrumentation of the HPC compute system

Approach:



Instrumentation of the HPC compute system

Approach:



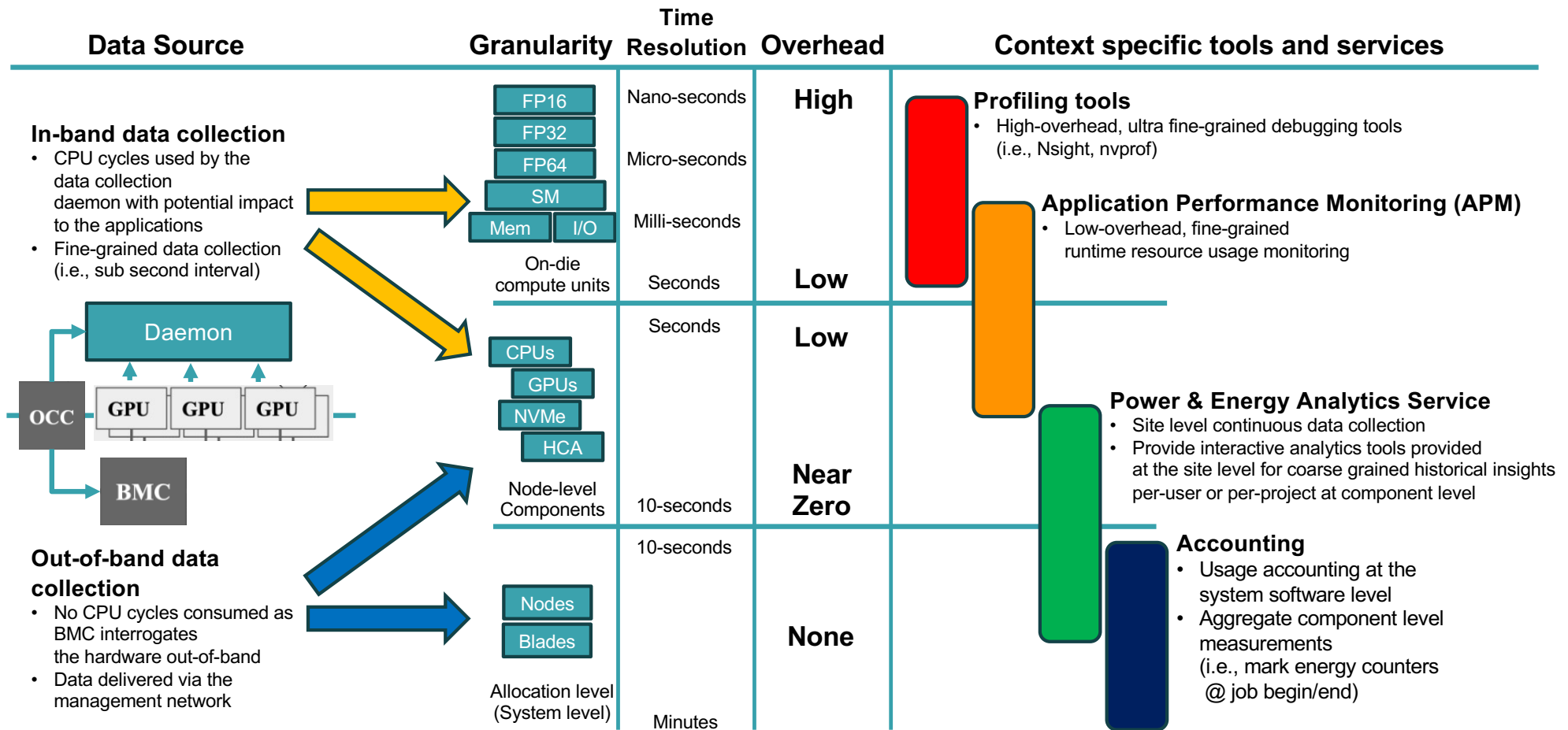
In-band data collection

- CPU cycles used by the data collection daemon with potential impact to the applications
- Fine-grained data collection (i.e., sub second interval)

Out-of-band data collection

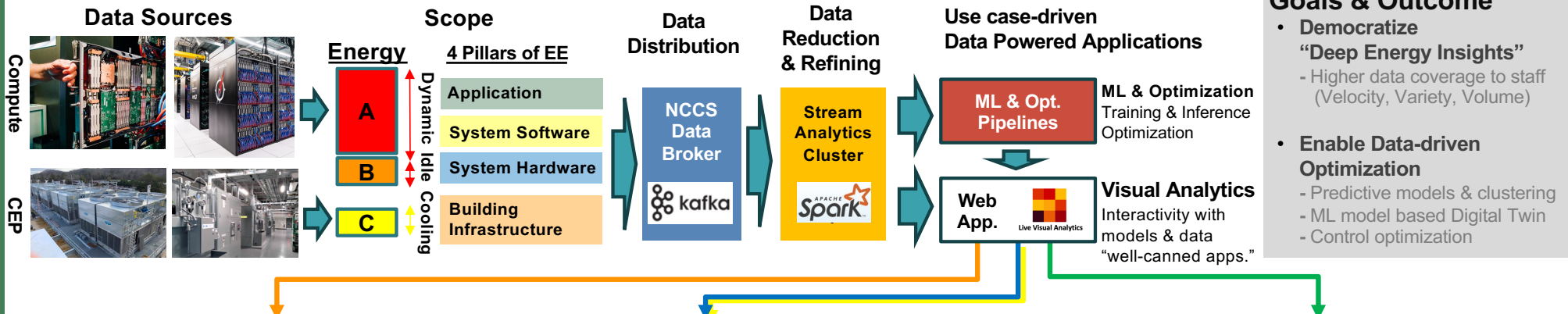
- No CPU cycles consumed as BMC interrogates the hardware out-of-band
- Data delivered via the management network

Tools and Services for Energy Efficiency

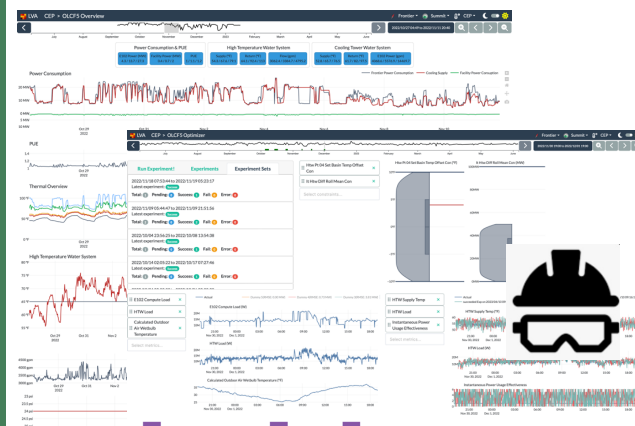


Power and Energy Analytics – “Increase Data & Model Usage”

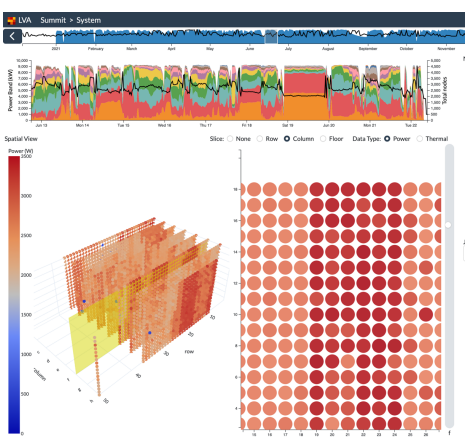
“Telemetry data, ML models, and interactive visual analytics tools for energy efficiency”



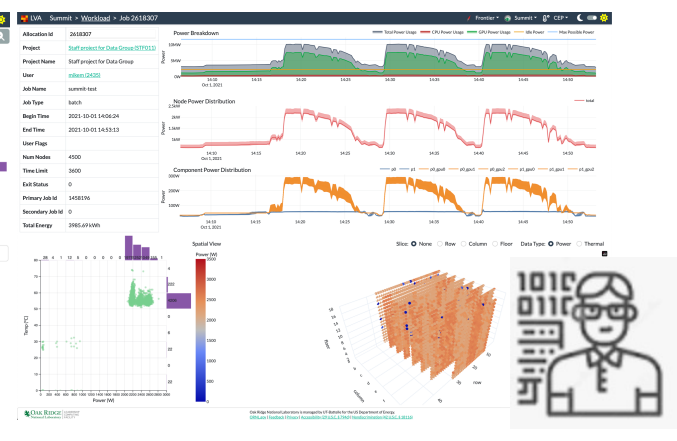
Central Energy Plant Control Analysis & ML Driven Optimization



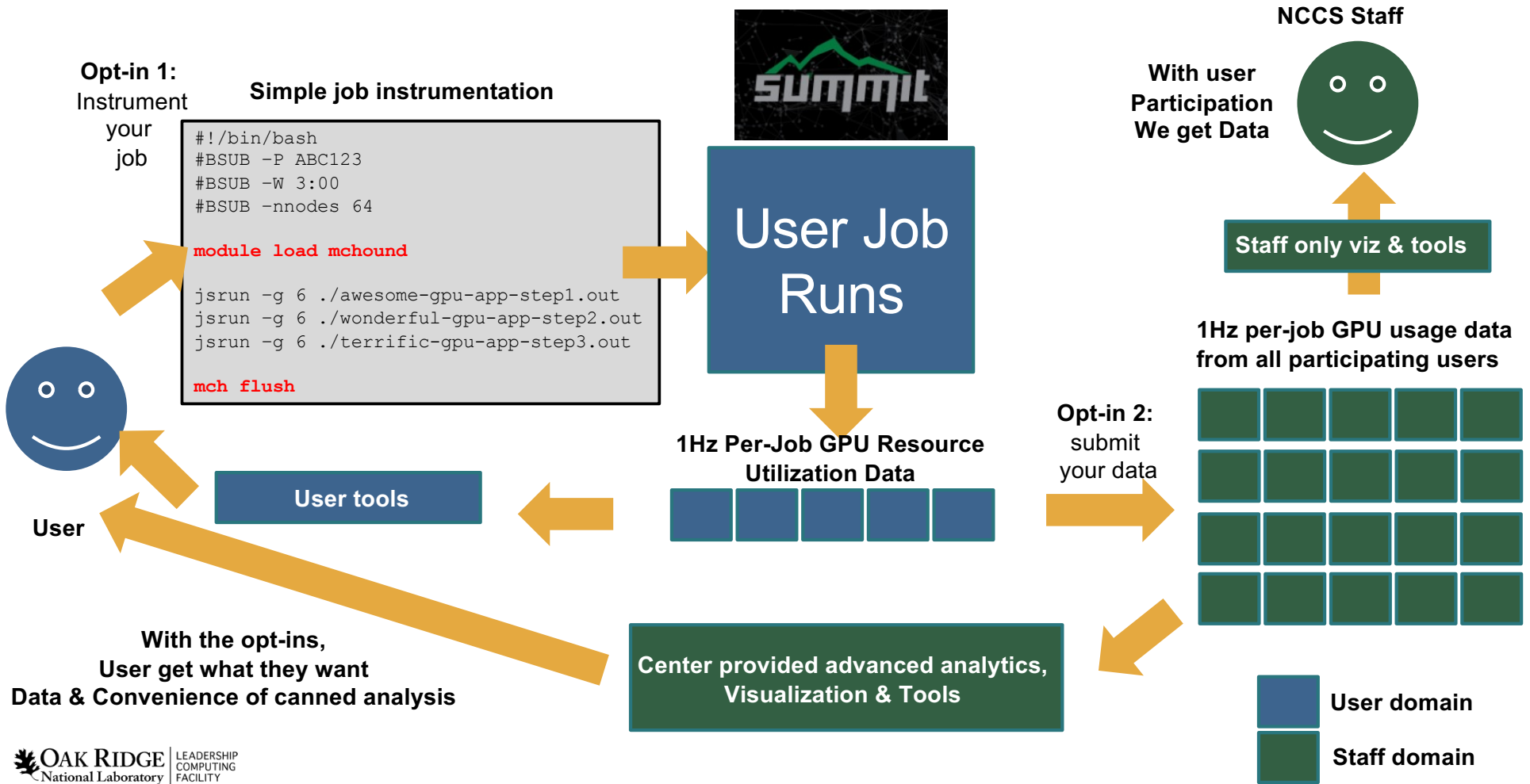
System Analysis



Per Application Profile

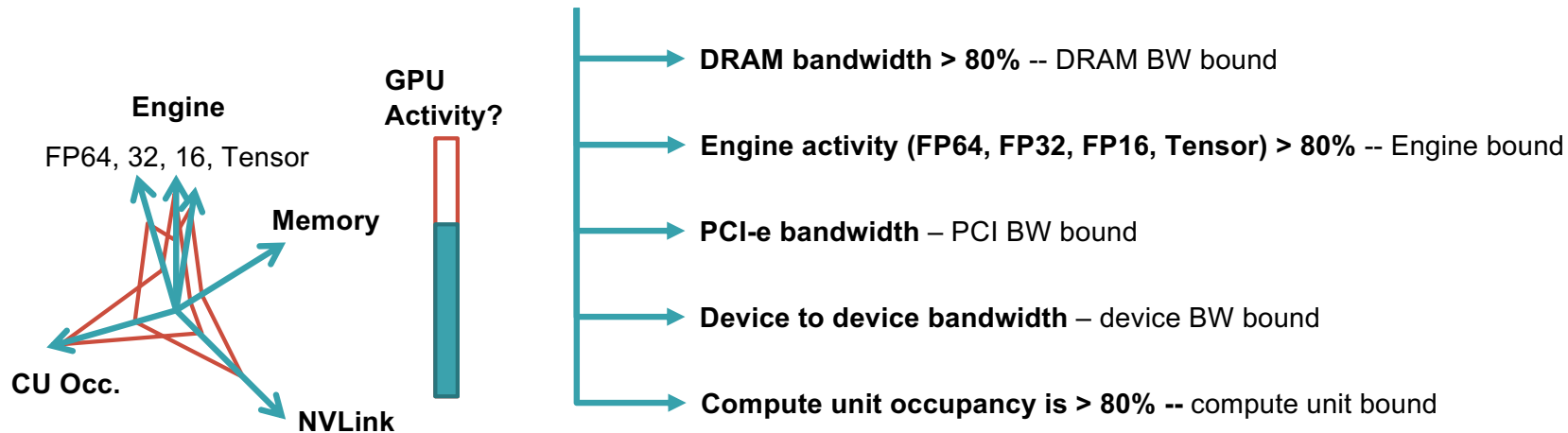


APM Concept – User opt-in driven data collection



Utilization of the compute device & Performance

Graphics activity < 100%?



Where is my application bound?

- **Power, energy & thermal**

- Power usage (watts), Total energy consumption (joules – counter)
- GPU Memory (HBM) & GPU core (SM/CU) temperature
- Pstate

- **Application auxiliary**

- GPU utilization (different)
- Framebuffer used (GPU memory used)
- SM, memory, video, sm_app, mem_app clock frequencies

Bringing everything together

Profiling tools

- High-overhead, ultra fine-grained debugging tools (i.e., Nsight, nvprof)

Application Performance Monitoring (APM)

- Low-overhead, fine-grained runtime resource usage monitoring

Power & Energy Analytics Service

- Site level continuous data collection
- Provide interactive analytics tools provided at the site level for coarse grained historical insights per-user or per-project at component level

Accounting

- Usage accounting at the system software level
- Aggregate component level measurements (i.e., mark energy counters @ job begin/end)

Low-level Debugging & Tuning

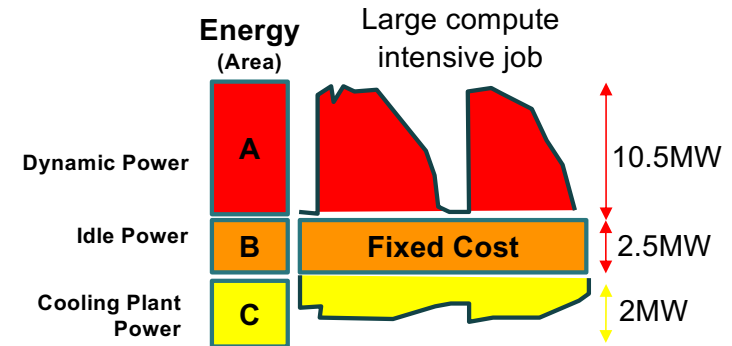
-- Active hammering and tinkering

Additional data in order to understand runtime resource usage and impact – where should I point my hammer?

Coarse grained first-degree indicator for power & energy issues – did I do everything else fine?

Final energy number

-- How good / bad did we do?



Energy Saving Strategies

- Boost
- Clock down Step down
- Sleep Turn off
- Balance Energy vs. Time to solution

Maximize Impact

Summary

- Post-exascale energy efficiency will require support for application level energy awareness
- There are low-hanging behavioral opportunities for applications
- We can start increasing energy awareness
 - First step is to kick-start the continuous improvement loop
- Filling in the workflow gap
 - Site level continuous data collection and interactive analytics
 - User opt-in based fine-grained runtime data collection
- Still a long way to go

Challenges, Future Work: A long way to go...

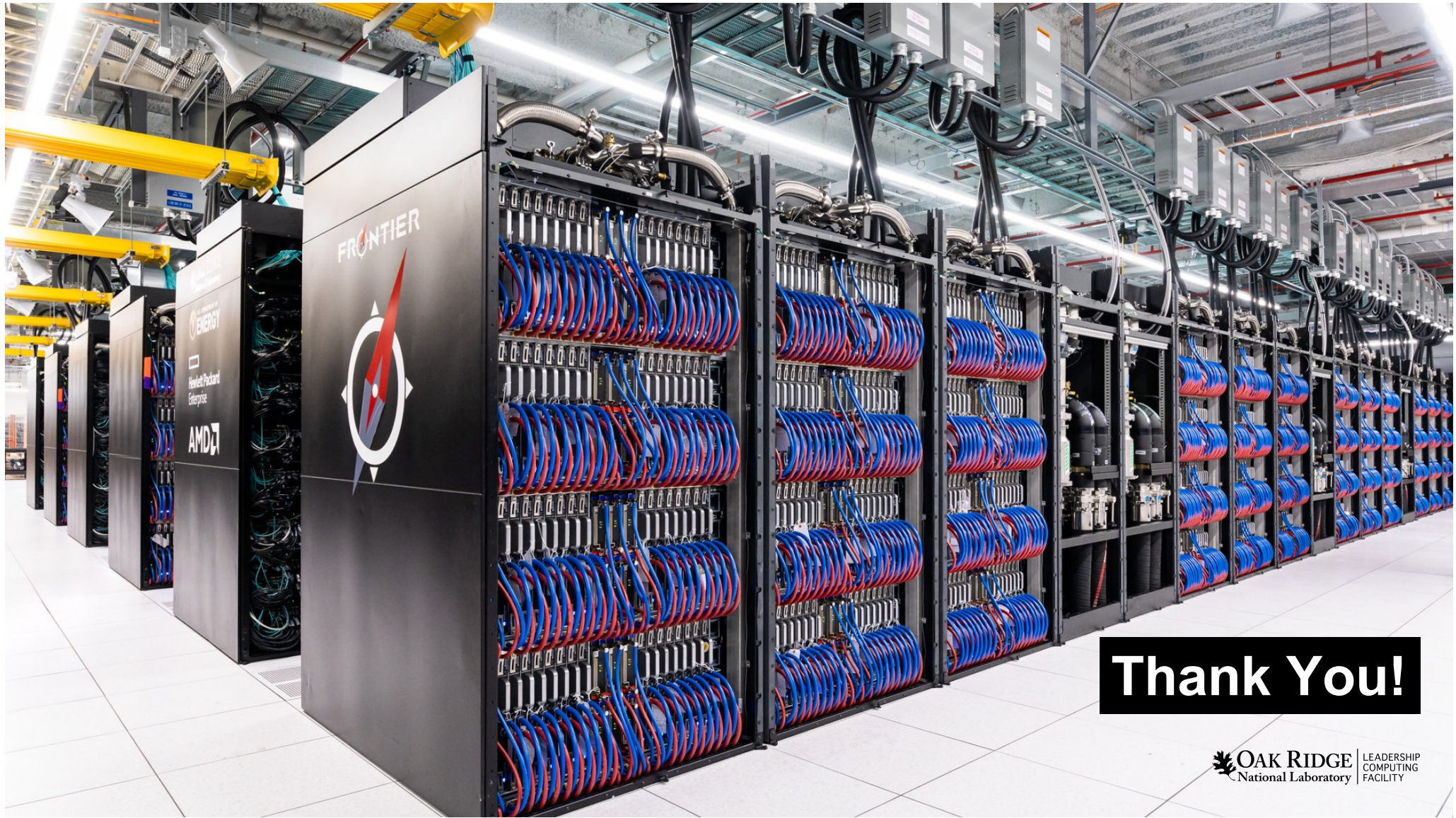
More Behavioral Options?

What is the incentive to put in E.E effort as a user?

Measuring Success? "Metrics"

Vendor Support? Telemetry, Tools, Capability

Can we automate the control?



Thank You!