

# Introduction to Variorum

Module 1 of 2, ECP Lecture Series

6 August 2021 11:30AM-1PM ET

Tapasya Patki, Aniruddha Marathe,  
Stephanie Brink, and Barry Rountree



# Module 1 Agenda

---

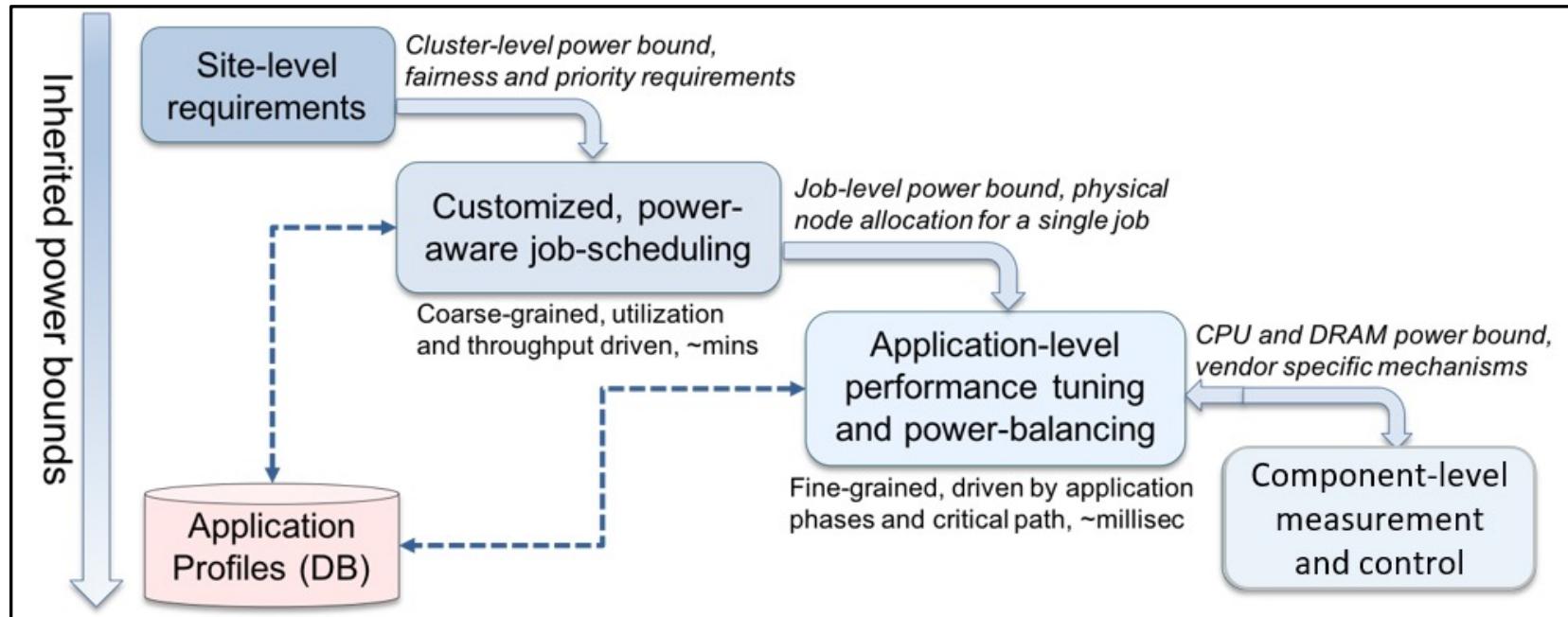
- Challenges in Power Management and the HPC Power Stack (20 min)
- Understanding Power Management Knobs on Intel, IBM, NVIDIA, ARM, and AMD Platforms (20 min)
- Variorum Library (40 min)
- Wrap Up (10 min)

# Welcome: HPC Systems and Power Management



- Compute Nodes, I/O Nodes, Network
- Static and **Dynamic** Power Management
- DOE's ongoing project: ECP Argo
  - Other synergistic projects: Flux, Caliper, GEOPM
  - Workflows: MuMMI, E3SM

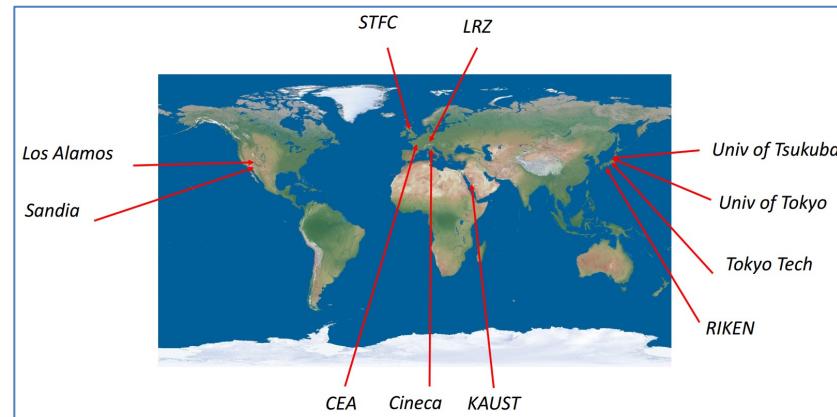
# HPC PowerStack: Community Effort on System-wide, dynamic power management



<https://hpcpowerstack.github.io/>

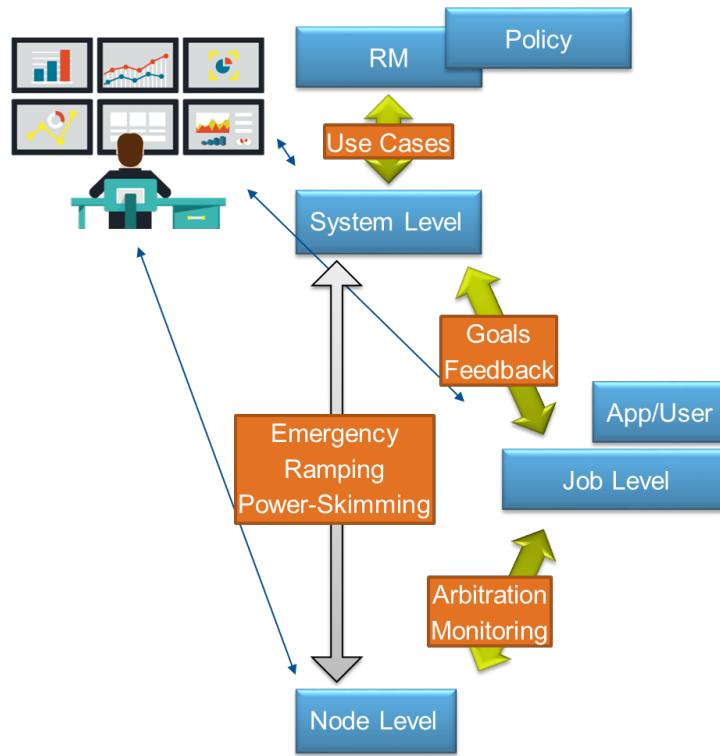
# PowerStack: Stakeholders

- Current industry collaborators: Intel, IBM, AMD, ARM, NVIDIA, Cray/HPE, Fujitsu, Altair, ATOS/Bull, and PowerAPI community standard
- Multiple academic and research collaborators across Europe, Asia, US
- Three working groups established
- Dynamic power management at all levels, along with prioritization of the critical path, application performance and throughput
- One of the prototypes developed as part of ECP using SLURM, GEOPM, Variorum/msr-safe (close collaboration with Intel)
- Additional software with Flux and Variorum underway



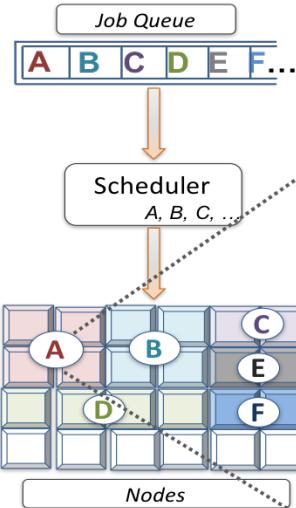
*EEHPC-WG's insight into sites investing in Energy- and Power-aware Job Scheduling and Resource Management (EPA-JSRM)*

# PowerStack: Layers



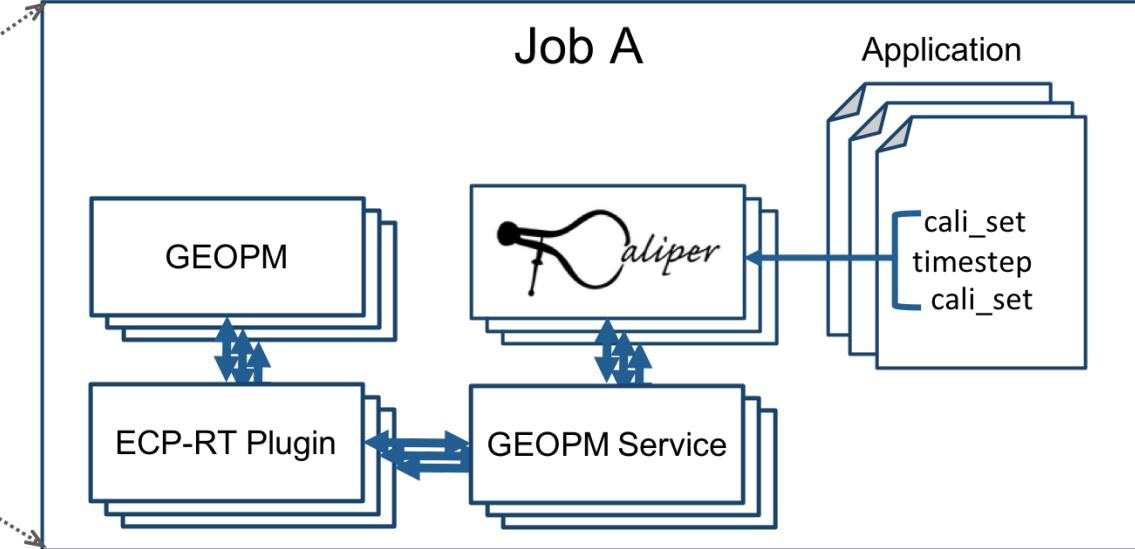
- Site-level
  - Policies, objectives
- System-level
  - Scheduling, runtimes, application-awareness
- Platform-level
  - Hardware interaction, PowerAPI
- Idle and Emergency Management

# ECP Argo uses SLURM/Flux, GEOPM and Variorum as vehicles for power management at Exascale

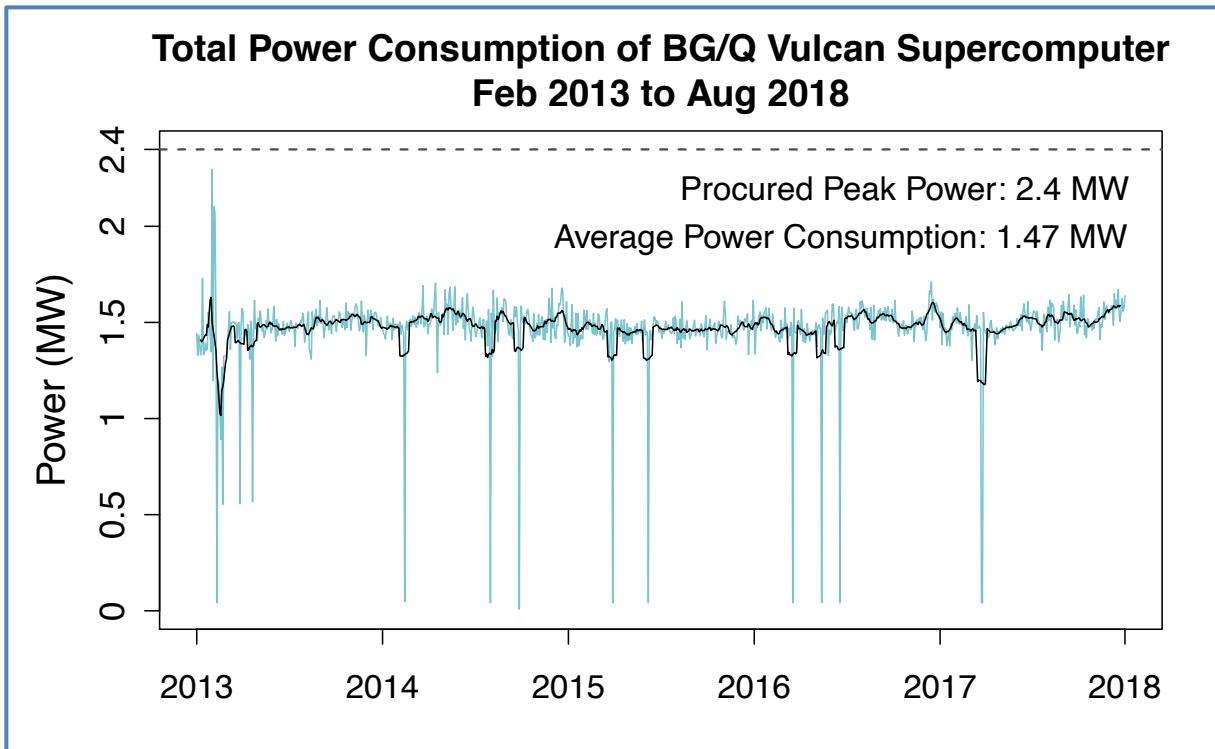


**ECP Argo, ECP Flux:**  
Power-aware scheduler such  
as Flux, SLURM or PowSched  
(low-overhead static decisions)

**ECP Argo:** Job-level runtime system for fine-grained, dynamic tuning and balancing. Supports diverse architectures, workloads and programming models

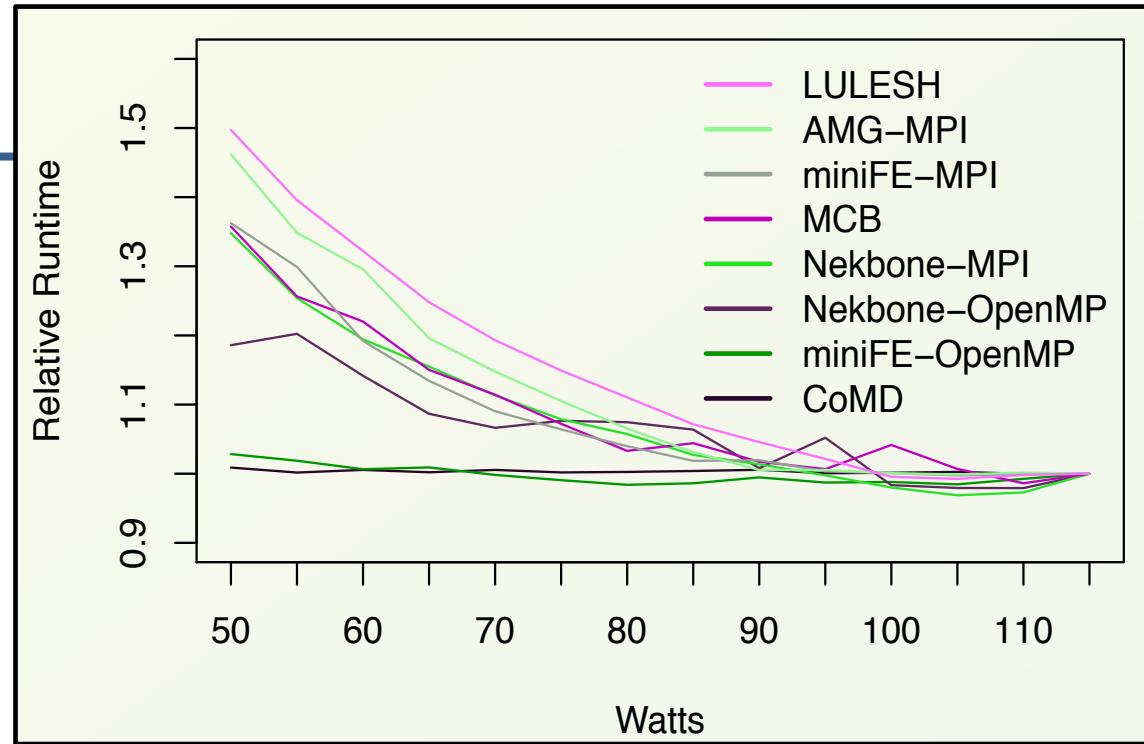


# Unused Power: 40%



# Why is it so?

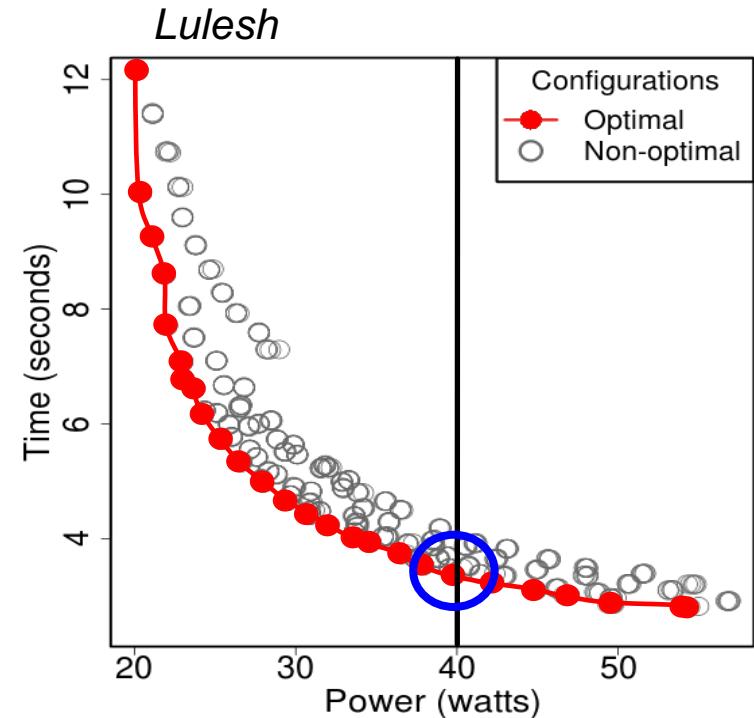
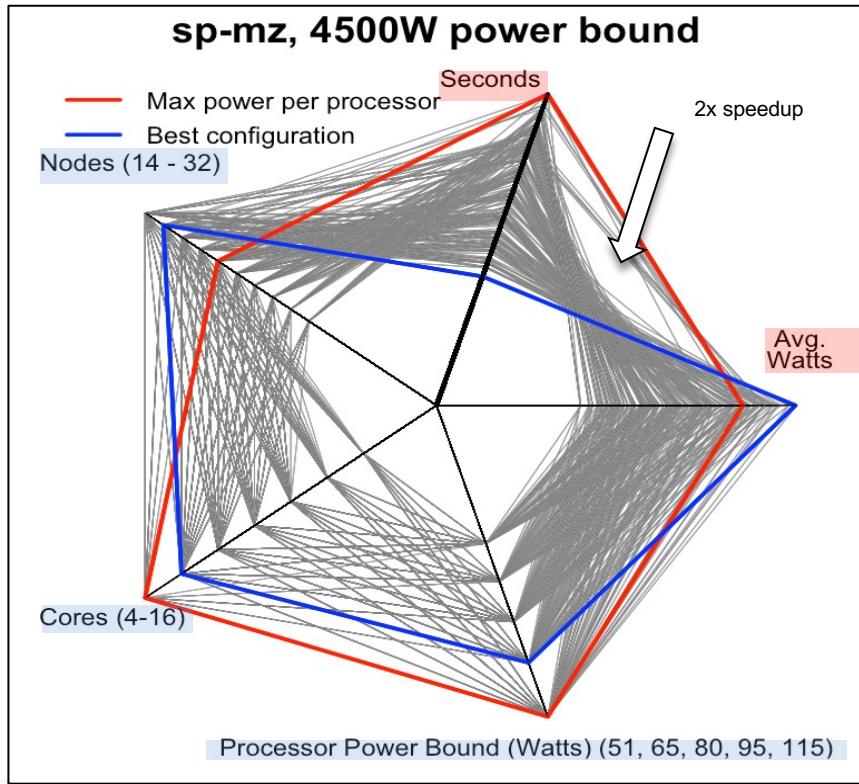
- Applications have different memory, communication, I/O requirements and phase behaviors
- Applications don't utilize all of the allocated power, thus **allocating more power doesn't always improve performance**



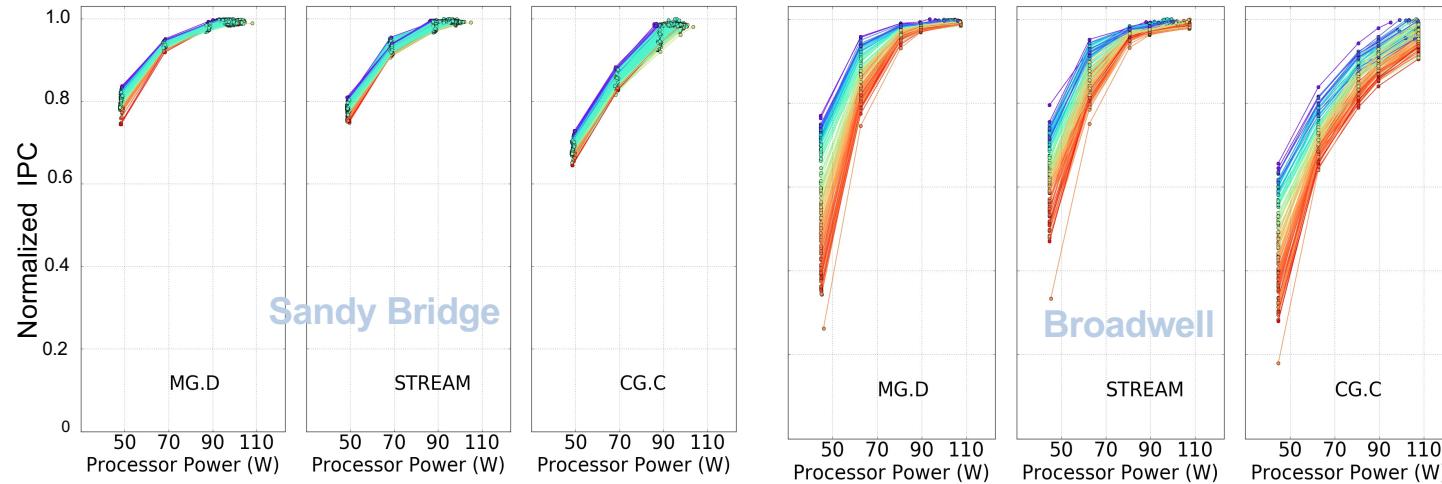
Intel Sandy Bridge, 8 nodes, (2 sockets, 8 cores)

Min: 51 W, Max: 115 W

# Automatic configuration selection is crucial for performance and energy efficiency



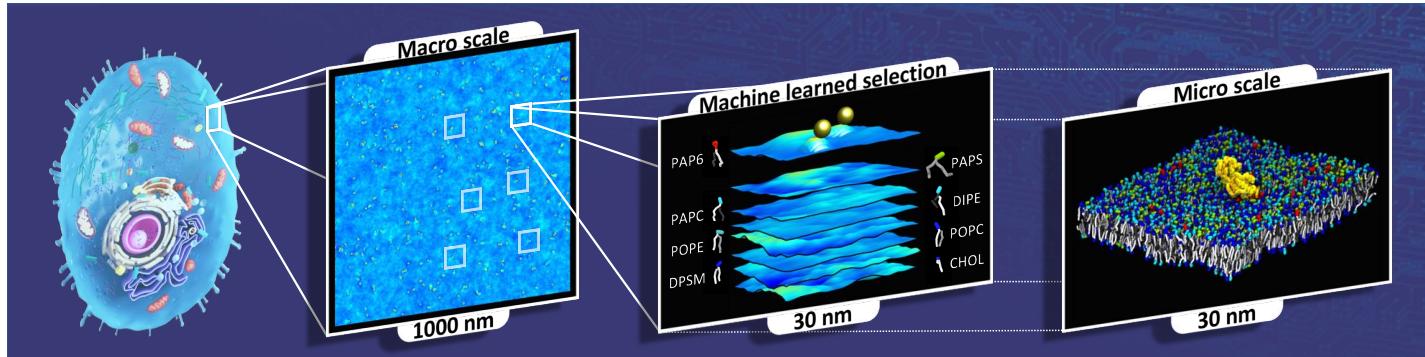
# Significant performance variability can result from processor manufacturing process



- Manufacturing variability in hardware continues to increase, and will worsen with heterogeneous systems
- 4x difference between two generations of Intel processors, needs advanced runtime options for mitigation

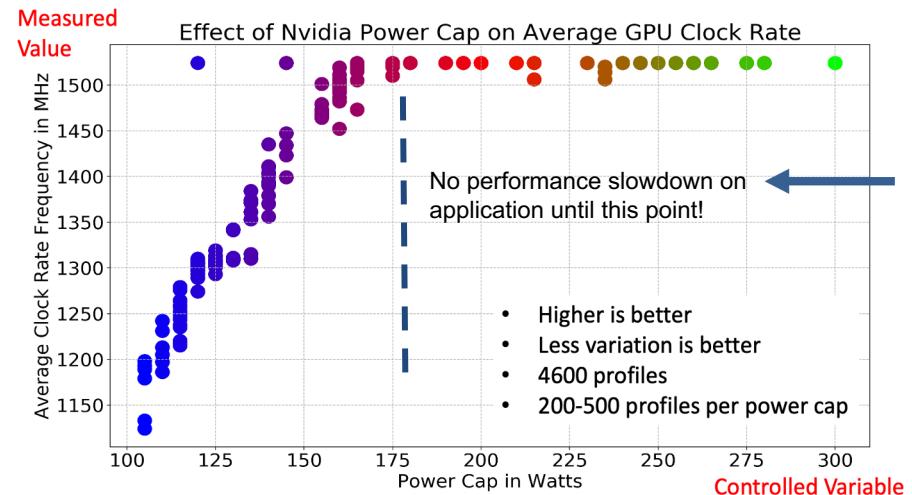
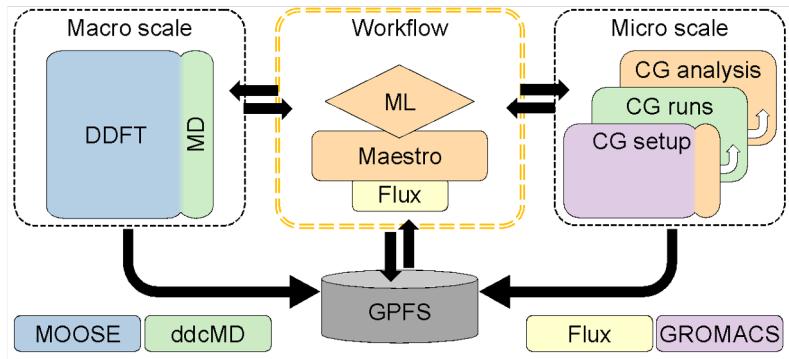
# Demo of Application Performance and Processor Variability

# Power Management of workflows is necessary: Example of the Multiscale Machine-Learned Modeling Infrastructure (MuMMI)



- MuMMI is designed to understand cancer-signaling mechanisms using multi-scale machine learning.
- Based on the mutated RAS protein, which is responsible for **a third of all human cancers**
- Understanding the role of RAS requires exploration at multiple length- and time-scales
- The MuMMI workflow makes significant use of GPUs on the Lassen/Sierra supercomputers

# Power Management Workflow Example: Significant Cost Savings with Active GPU Power Management on MuMMI Workflow



- ddcMD is one of the main compute and GPU-oriented components of MuMMI Workflow
- Cluster wide savings of 382kW with no performance slowdown with GPU power management!
- 254.6 kWh energy savings, **\$43k cost savings per day** @ 7c per kWh

*Patki et al., Comparing GPU Power and Frequency Capping: A Case Study with the MuMMI Workflow, WORKS'19 (SC)*

# Facilities Perspective: Mitigating Power Swings on Sierra/Lassen with in-depth application analysis with Variorum



Example: LBANN on Sierra at full scale has significant fluctuations impacting LLNL's electrical grid -- workload swings expected to worsen at exascale

- Livermore Big Artificial Neural Network toolkit (LBANN) -- infrastructure used for deep learning in HPC
- LBANN utilizes all 4 GPUs per node
- Data shows 3 minute samples over 6 hours on Sierra with >200 KW swings
- Other workflows have similar trends with power fluctuations at scale
- Mitigation of power fluctuations is required to avoid electrical supply disruption
- Variorum + Flux can dynamically analyze applications and prevent future fluctuations

# So, how do we even manage power? (measurement vs control)

---

Two primary mechanisms in hardware for control:

- Dynamic voltage and frequency scaling (ARM)
- Power capping, such as RAPL (Intel, AMD) or OPAL (IBM) or NVML (NVIDIA)
- Automatic tuning through Intel Turbo Boost or IBM UltraTurbo
- ACPI defines P-states, which are voltage/frequency pairs
- DVFS: `cpufreq`, ‘userspace’ governor
- RAPL:
  - Low-level register programming, supported with `msr` kernel module along with `msr-tools`
  - LLNL’s `msr-safe` and `variorum` provide safe and clean access from user space across architectures
- OPAL:
  - Firmware layer providing in band monitoring with sensors and out of band power capping

# Introduction to Intel's RAPL

---

- Intel provides programmable, machine-specific registers for power, energy and thermal management (power capping)
- MSR Domains for server architectures:
  - **Package** – represents processor cores, caches, and other things on the socket
  - **DRAM** – represents memory components
  - Other **uncore** registers also exist

*Intel SDM: Vol 3, Chapter 14.9, <https://software.intel.com/en-us/articles/intel-sdm>*

# Deep dive into the MSR\_PKG\_POWER\_LIMIT, (0x610h, rw)

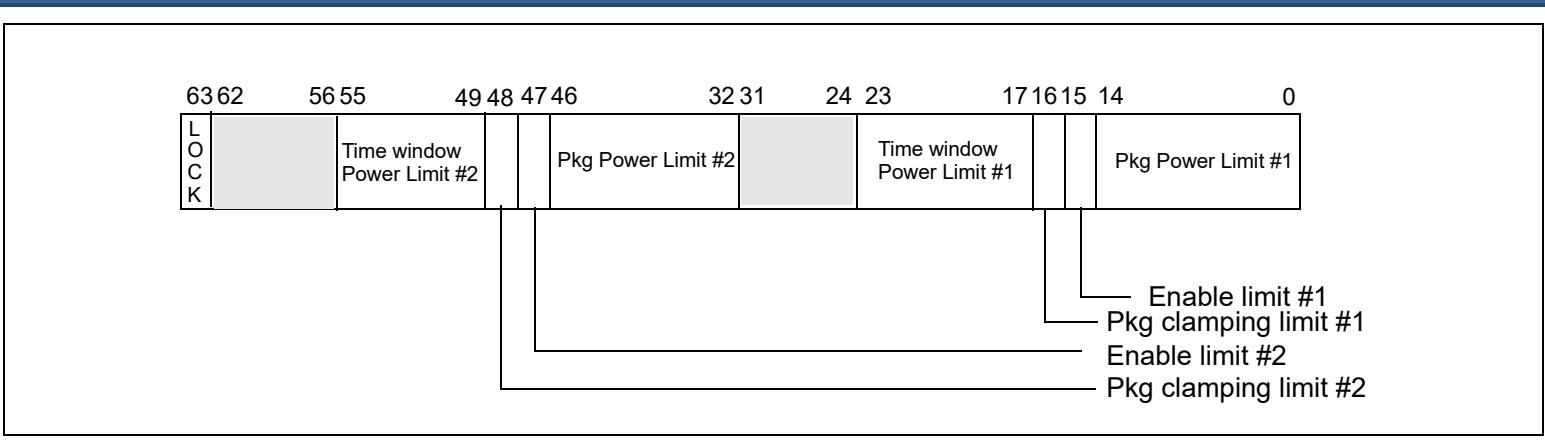


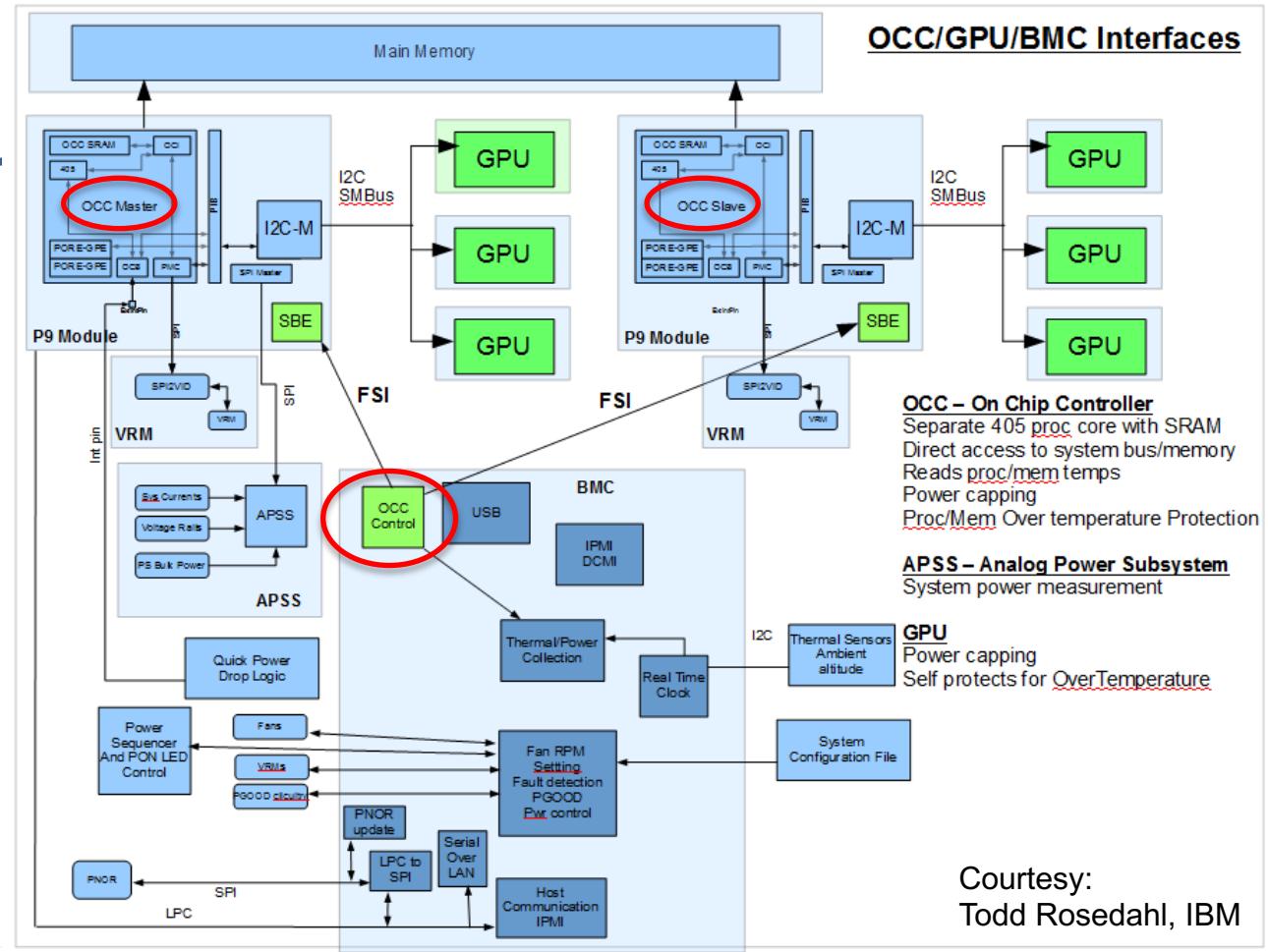
Figure 14-32. MSR\_PKG\_POWER\_LIMIT Register

- Pkg Power Limit#1, bits 14:0, sets average power limit corresponding for window 1, unit 0.125 W
- Enable Limit#1: bit 15, 0 = disabled
- Pkg Clamping Limit#1, bit 16, allow going below OS requested P/T state
- Time Window Limit#1, bits 23:17, minimum is typically 1 millisec
- Pkg Limit #2 is typically not used

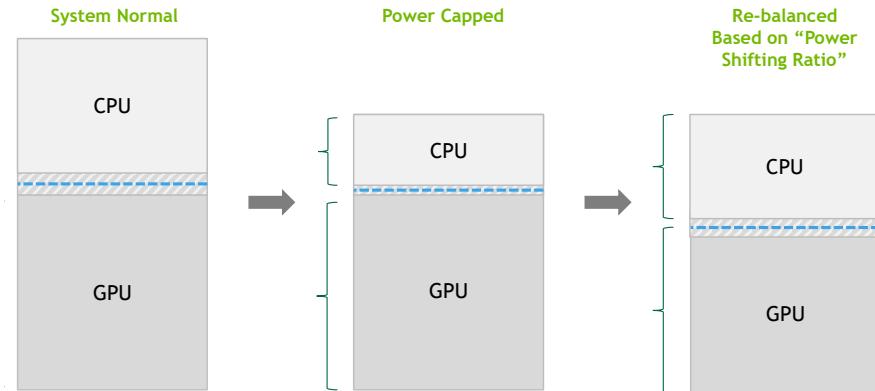
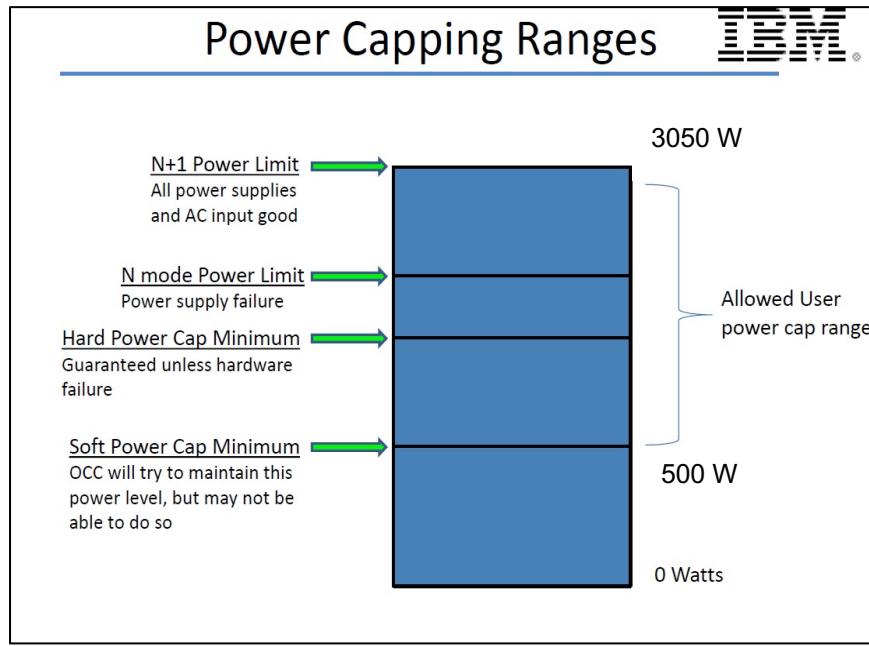
# Demo of `msr-tools` and introduction to `msr-safe`, `libmsr`

# Power9 OCCs have a Primary-Secondary Design

- OCC collects thermal data every 250us
- Power knobs are exposed with the Open Power Abstraction Layer (OPAL)
- When power button is pressed, BMC selects master chip and releases SBEs, OPAL firmware is loaded



# Capping Node and GPU power



- Power-shifting ratio determines distribution between CPU and GPU
- Per-socket cap and memory cap cannot be specified, master OCC has control

# OPAL interfaces for in-band sensors and power capping

---

Requires OPAL firmware v6.0.18+ – includes fix for a firmware bug we found based on soft power cap

Read-only access:

- `/sys/firmware/opal/exports/occ_inband_sensors`
  - Over 336 sensors reported on alehouse
  - Including power, temperature, frequencies for CPU, Memory, GPU (1ms granularity)

Read/write access:

- `/sys/firmware/opal/powercap/system-powercap/powercap-current`
- `/sys/firmware/opal/psr/cpu_to_gpu_*` (per socket)

# AMD power management can be accessed with E-SMI Library, Energy Driver, and Host System Management Port (HSMP) module

- EPYC™ System Management Interface In-band Library (E-SMI library) is available at:  
[https://github.com/amd/esmi\\_ib\\_library](https://github.com/amd/esmi_ib_library)
- AMD Energy Driver allows access for core and socket energy counters through MSRs and RAPL (with hwmon), available at: [https://github.com/amd/amd\\_energy](https://github.com/amd/amd_energy)
  - Power, Energy and Time Units **MSR\_RAPL\_POWER\_UNIT**/ C001\_0299: shared with all cores in the socket
  - Energy consumed by each Core **MSR\_CORE\_ENERGY\_STATUS**/ C001\_029A: 32-bitRO, Accumulator, core-level power reporting
  - Energy consumed by Socket **MSR\_PACKAGE\_ENERGY\_STATUS**/ C001\_029B: 32-bitRO, Accumulator, socket-level power reporting, shared with all cores in socket
  - These registers are updated every 1 ms and cleared on reset of the system.
- HSMP driver for power metrics is available at: ([https://github.com/amd/amd\\_hsmp](https://github.com/amd/amd_hsmp))
  - Allows for power capping, setting of boost limits and PCIe access.
  - Internal mechanism uses a mailbox approach for register access
- Structure of SysFS interface (/sys/devices/system/cpu/) includes CPU, Socket and general system management (does not include GPU management, which is provided through rocm-smi)

# AMD Power Control Knobs are exposed through the Host System Management Port (HSMP) kernel module

- SysFS structure:

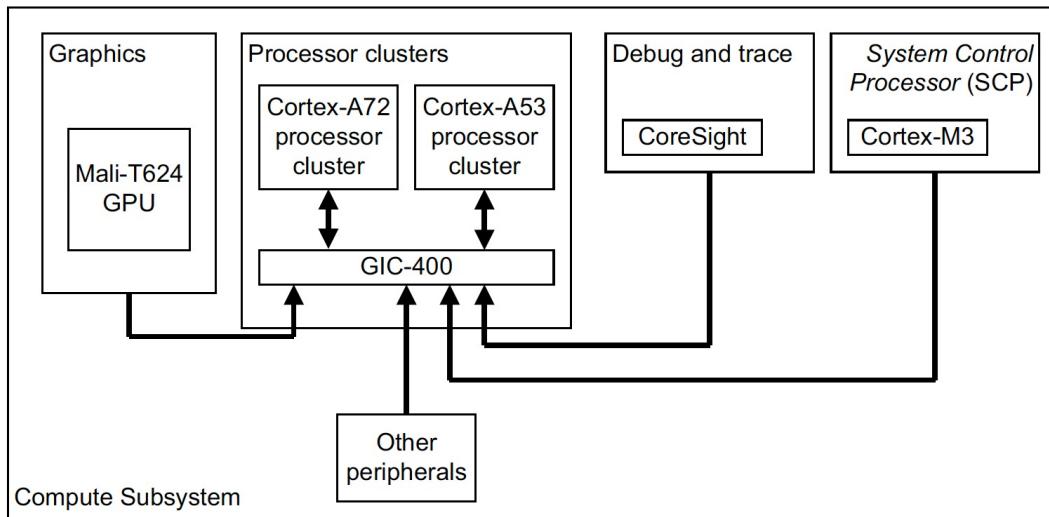
`amd_hsmp/cpuX/`  
    `boost_limit`

Directory for each possible CPU  
(RW) HSMP boost limit for the core in MHz

`amd_hsmp/socketX/`  
    `boost_limit`  
    `c0_residency`  
    `cclk_limit`  
    `fabric_clocks`  
    `fabric_pstate`  
    `power`  
    `power_limit`  
    `power_limit_max`  
    `proc_hot`  
    `tctl`

Directory for each possible socket  
(WO) Set HSMP boost limit for the socket in MHz  
(RO) Average % all cores are in C0 state  
(RO) Most restrictive core clock (CCLK) limit in MHz  
(RO) Data fabric (FCLK) and memory (MCLK) in MHz  
(WO) Set data fabric P-state, -1 for autonomous  
(RO) Average socket power in milliwatts  
(RW) Socket power limit in milliwatts  
(RO) Maximum possible value for power limit in mW  
(RO) Socket PROC\_HOT status (1 = active, 0 = inactive)  
(RO) Thermal Control value (not temperature)

# ARM Juno r2 : SoC Architecture



— Compute Subsystem internal peripheral interrupts →

- Cortex-A72 MP2 cluster (r0p0eac)
  - Dual cluster, SMP configuration
  - Overdrive 1.2GHz speed
- Cortex-A53 MP4 cluster (r0p3)
  - Quad cluster, SMP configuration
  - Overdrive 950MHz speed
- Quad Core MALI T624 r1p0
  - Nominal 600MHz operating speed
  - Caches: L2 128KB
- Control and telemetry
  - DVFS and power gating via SCP
  - 4 energy meters
  - Temperature, clocks telemetry

# ARM Juno r2 : System Interface

- Monitoring and control through Sysfs interface<sup>1</sup>
  - Exposed through the Linux HWMon interface

## Power telemetry (mW)

Location: /sys/class/hwmon/hwmon0/

SYS\_POW\_SYS : power1\_input  
SYS\_POW\_A72 : power2\_input  
SYS\_POW\_A53 : power3\_input  
SYS\_POW\_GPU : power4\_input

## Clocks telemetry (kHz)

Location: /sys/devices/system/cpu/cpufreq/

big clocks: policy0/scaling\_cur\_freq  
LITTLE clocks: policy1/scaling\_cur\_freq

## Thermal telemetry (C)

Location: /sys/class/hwmon/hwmon0/

SoC temperature: temp1\_input  
big temperature: temp2\_input  
LITTLE temperature: temp3\_input  
GPU temperature: temp4\_input

## Frequency control (kHz)

Location: /sys/devices/system/cpu/cpufreq/

big clocks: policy0/scaling\_setspeed  
LITTLE clocks: policy1/scaling\_setspeed

# NVML<sup>1</sup>: Nvidia Measurement Interface

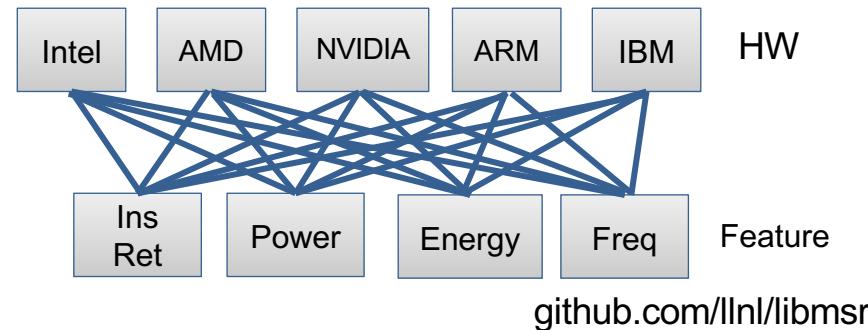
---

- **Power usage:** nvmlDeviceGetPowerUsage (device ID, &power)
  - Returns power usage of the target GPU device in watts
- **Power limit:** nvmlDeviceGetPowerManagementLimit(device ID, &powerlimit)
  - Retrieves the power management limit associated with this device in watts
- **Temperature:** nvmlDeviceGetTemperature (device ID, NVML\_GPU, &temp)
  - Returns temperature of the target GPU device in degree Celsius
- **Clocks:** nvmlDeviceGetClock (device ID, NVML\_CLOCK\_SM, &clocks)
  - Returns clock speed of the target GPU device in MHz
- **GPU utilization:** nvmlDeviceGetUtilizationRates(device ID, &utilization)
  - Instantaneous utilization rate for the target GPU device

# Power management capabilities differ across vendors

- libmsr (~2012) has been quite successful in the community
  - Provided simpler monitor/control interfaces translating bit fields into 64-bit MSRs
  - But, was Intel-specific
- Interfaces, domains, latency, capabilities
  - But, may also vary across generations within the same vendor
- Goals of variorum:
  - Target 95% of users (friendly APIs)
  - More devices
    - i.e., CPUs, Accelerators, IPMI, PCIe (CSRs, MMIO)
  - More platforms
    - i.e., Intel Skylake, IBM P9, NVIDIA Volta, AMD Epyc, ARM
  - More hardware knobs and controls

Feature	Ivy Bridge	Haswell
CPU Freq	Per-socket implementation	Per-core implementation
Energy Status	Uses energy unit fused into MSR_POWER_UNIT for conversion	Uses standard 15.3 micro-Joules for conversion
Power Limit	Same implementation	

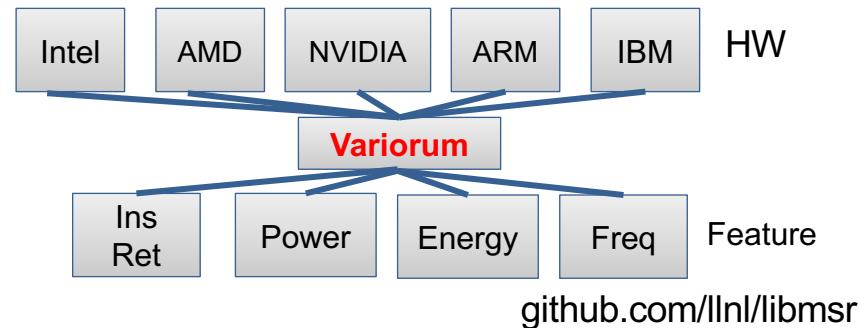


[github.com/llnl/libmsr](https://github.com/llnl/libmsr)

# Power management capabilities differ across vendors

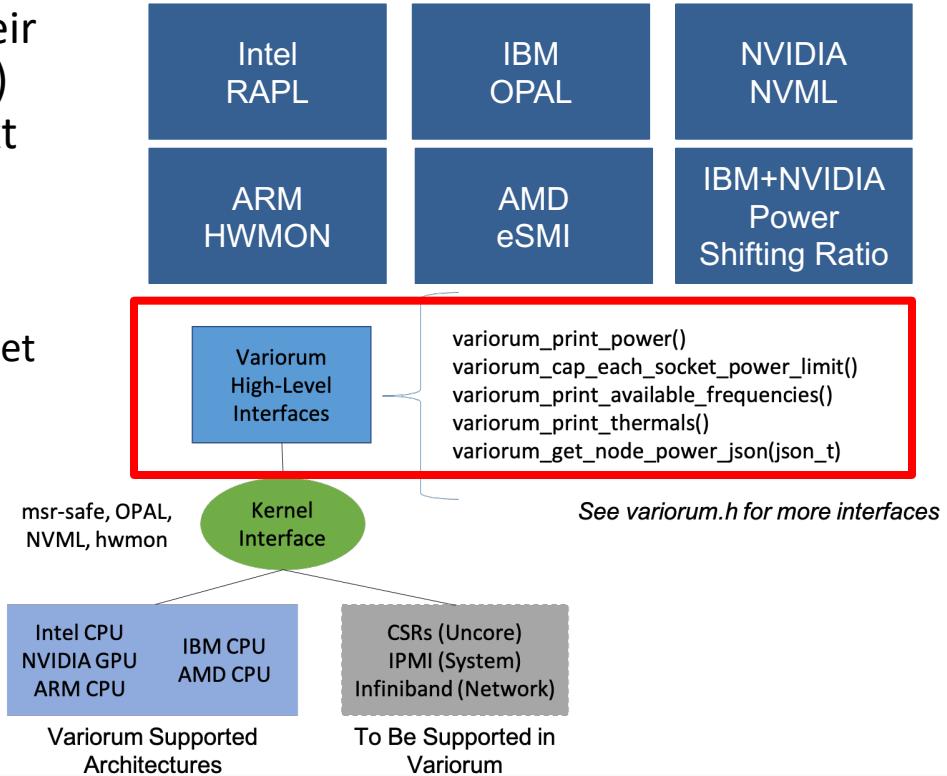
- libmsr (~2012) has been quite successful in the community
  - Provided simpler monitor/control interfaces translating bit fields into 64-bit MSRs
  - But, was Intel-specific
- Interfaces, domains, latency, capabilities
  - But, may also vary across generations within the same vendor
- Goals of variorum:
  - Target 95% of users (friendly APIs)
  - More devices
    - i.e., CPUs, Accelerators, IPMI, PCIe (CSRs, MMIO)
  - More platforms
    - i.e., Intel Skylake, IBM P9, NVIDIA Volta, AMD Epyc, ARM
  - More hardware knobs and controls

Feature	Ivy Bridge	Haswell
CPU Freq	Per-socket implementation	Per-core implementation
Energy Status	Uses energy unit fused into MSR_POWER_UNIT for conversion	Uses standard 15.3 micro-Joules for conversion
Power Limit	Same implementation	



# Variorum: Vendor-neutral user space library for power management

- Power management capabilities (and their interfaces, domains, latency, capabilities) widely differ from one vendor to the next
- Variorum: Platform-agnostic vendor-neutral, simple front-facing APIs
  - Evolved from *libmsr*, and designed to target several platforms and architectures
  - Abstract away tedious and chaotic details of low-level knobs
  - Implemented in C, with function pointers to specific target architecture
  - Integration with higher-level power management software through JSON



# Demo of Variorum Build and APIs

# Variorum Examples: Source code annotations as well as non-intrusive monitoring are supported

```
[brink2 octomore:~/variorum/build/examples (dev)]$ ./variorum-print-power-example | grep PACKAGE  
_PACKAGE_ENERGY_STATUS Offset Host Socket Bits Energy_J Power_W Elapsed_sec Timestamp_sec  
_PACKAGE_ENERGY_STATUS 0x611 octomore 0 0x1b684b74 28065.178955 0.000000 0.000000 0.000000  
_PACKAGE_ENERGY_STATUS 0x611 octomore 1 0x3ef4d257 64467.286560 0.000000 0.000000 0.000006  
_PACKAGE_ENERGY_STATUS 0x611 octomore 0 0x1b684cf9 28065.202698 32.128113 0.000739 0.000708  
_PACKAGE_ENERGY_STATUS 0x611 octomore 1 0x3ef4d257 64467.286560 0.000000 0.000739 0.000708
```

```
[brink2 octomore:~/variorum/build/examples (dev)]$ ./variorum-print-power-example | grep DRAM  
_DRAM_ENERGY_STATUS Offset Host Socket Bits Energy_J Power_W Elapsed_sec Timestamp_sec  
_DRAM_ENERGY_STATUS 0x619 octomore 0 0x5290a443 21136.641647 0.000000 0.000000 0.000006  
_DRAM_ENERGY_STATUS 0x619 octomore 1 0x547ee184 21630.880920 0.000000 0.000000 0.000006  
_DRAM_ENERGY_STATUS 0x619 octomore 0 0x5290a614 21136.648743 9.335970 0.000760 0.000732  
_DRAM_ENERGY_STATUS 0x619 octomore 1 0x547ee357 21630.888046 9.376124 0.000760 0.000732
```

```
[patki1@lassen36:examples]$ ./variorum-print-power-example  
_IBMPOWER Host Socket PWRSYS_W PWRPROC_W PWRMEM_W PWRGPU_W Timestamp_sec  
_IBMPOWER lassen36 0 443 116 21 72 0.000019  
_IBMPOWER lassen36 1 0 104 20 68 0.000082  
_IBMPOWER lassen36 0 443 116 21 72 0.000128  
_IBMPOWER lassen36 1 0 104 20 68 0.000158
```

```
[marathe1@lassen31 (dev) 0]$ examples/variorum-print-power-example  
_GPU_POWER_USAGE Host Socket DeviceID Power_W  
_GPU_POWER_USAGE lassen31 0 0 37.233002  
_GPU_POWER_USAGE lassen31 0 1 35.787002  
_GPU_POWER_USAGE lassen31 1 2 36.744002  
_GPU_POWER_USAGE lassen31 1 3 37.723002
```

```
[genericarmv8:/mnt/ssd/am/variorum/build$ examples/variorum-print-power-example  
_ARM_POWER Host Sys_mW Big_mW Little_mW GPU_mW  
_ARM_POWER genericarmv8 773.00 668.93 48.98 71.67  
_ARM_POWER genericarmv8 795.67 320.31 87.14 69.33
```

Intel print power:  
Read energy  
register and  
convert to power

IBM print  
power

NVIDIA  
print power

ARM print  
power

# Variorum Examples: Source code annotations as well as non-intrusive monitoring are supported

Intel cap power: Cap power at the socket level

```
[brink2 octomore ~/variorum/build/examples (dev)]$ ./variorum-cap-socket-power-limits-example -l 100
Capping each socket to 100W.
```

NVIDIA cap power: Feature is available in NVML, but we haven't implemented yet

```
[marathe1@lassen31 (dev) ~]$ examples/variorum-cap-socket-power-limits-example -l 120
```

```
Capping each socket to 120W.
```

```
lassen31:/g/g92/marathe1/myworkspace/variorum-tutorial/variorum/src/variorum/variorum.c:variorum_cap_each_socket_power_limit():383: _ERROR_VARIORUM_FEATURE_NOT_IMPLEMENTED: Feature not yet implemented or is not supported
```

ARM cap power: Feature is not supported, only DVFS

```
[genericarmv8:/mnt/ssd/am/variorum/build$ examples/variorum-cap-socket-power-limits-example -l 120
```

```
Capping each socket to 120W.
```

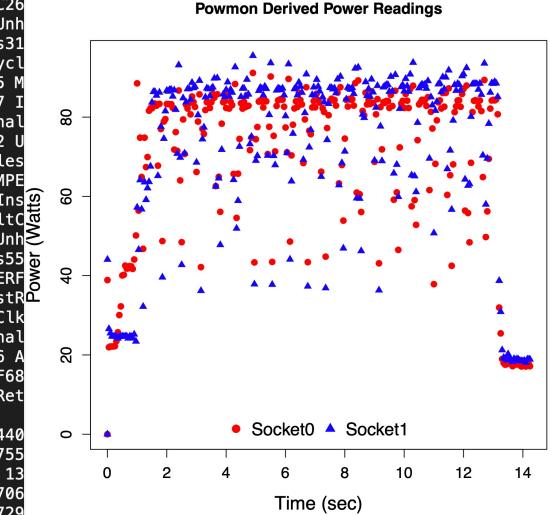
```
(null):/mnt/ssd/am/variorum/src/variorum/variorum.c:variorum_cap_each_socket_power_limit():383: _ERROR_VARIORUM_FEATURE_NOT_IMPLEMENTED: Feature not yet implemented or is not supported
```

# Demo of Variorum Non-Intrusive Monitoring: Powmon

# Powmon: Non-intrusive monitoring tool across architectures

```
_POWMON time pkg0_joules pkg0_lim1watts pkg0_lim2watts dram0_joules dram0_limwatts pkg1_joules pkg1_lim1watts pkg1_lim2watts dram1_joules dram1_limwatts InstRet0 UnhaltClkCycles0 UnhaltRefCycles0 APERF0 MPERF0 TSC0 InstRet1 UnhaltClkCycles1 UnhaltRefCycles1 APERF1 MPERF1 TSC1 InstRet2 UnhaltClkCycles2 UnhaltRefCycles2 APERF2 MPERF2 TSC2 InstRet3 UnhaltClkCycles3 UnhaltRefCycles3 APERF3 MPERF3 TSC3 InstRet4 UnhaltClkCycles4 UnhaltRefCycles4 APERF4 MPERF4 TSC4 InstRet5 UnhaltClkCycles5 UnhaltRefCycles5 APERF5 MPERF5 TSC5 InstRet6 UnhaltClkCycles6 UnhaltRefCycles6 APERF6 MPERF6 TSC6 InstRet7 UnhaltClkCycles7 UnhaltRefCycles7 APERF7 MPERF7 TSC7 InstRet8 UnhaltClkCycles8 UnhaltRefCycles8 APERF8 MPERF8 TSC8 InstRet9 UnhaltClkCycles9 UnhaltRefCycles9 APERF9 MPERF9 TSC9 InstRet10 UnhaltClkCycles10 UnhaltRefCycles10 APERF10 MPERF10 TSC10 InstRet11 UnhaltClkCycles11 UnhaltRefCycles11 APERF11 MPERF11 TSC11 InstRet12 UnhaltClkCycles12 UnhaltRefCycles12 APERF12 MPERF12 TSC12 InstRet13 UnhaltClkCycles13 UnhaltRefCycles13 APERF13 MPERF13 TSC13 InstRet14 UnhaltClkCycles14 UnhaltRefCycles14 APERF14 MPERF14 TSC14 InstRet15 UnhaltClkCycles15 UnhaltRefCycles15 APERF15 MPERF15 TSC15 InstRet16 UnhaltClkCycles16 UnhaltRefCycles16 APERF16 MPERF16 TSC16 InstRet17 UnhaltClkCycles17 UnhaltRefCycles17 APERF17 MPERF17 TSC17 InstRet18 UnhaltClkCycles18 UnhaltRefCycles18 APERF18 MPERF18 TSC18 InstRet19 UnhaltClkCycles19 UnhaltRefCycles19 APERF19 MPERF19 TSC19 InstRet20 UnhaltClkCycles20 UnhaltRefCycles20 APERF20 MPERF20 TSC20 InstRet21 UnhaltClkCycles21 UnhaltRefCycles21 APERF21 MPERF21 TSC21 InstRet22 UnhaltClkCycles22 UnhaltRefCycles22 APERF22 MPERF22 TSC22 InstRet23 UnhaltClkCycles23 UnhaltRefCycles23 APERF23 MPERF23 TSC23 InstRet24 UnhaltClkCycles24 UnhaltRefCycles24 APERF24 MPERF24 TSC24 InstRet25 UnhaltClkCycles25 UnhaltRefCycles25 APERF25 MPERF25 TSC25 InstRet26 UnhaltClkCycles26 UnhaltRefCycles26 APERF26 MPERF26 TSC26 InstRet27 UnhaltClkCycles27 APERF27 MPERF27 TSC27 InstRet28 UnhaltClkCycles28 UnhaltRefCycles28 APERF28 MPERF28 TSC28 InstRet29 UnhaltRefCycles29 APERF29 MPERF29 TSC29 InstRet30 UnhaltClkCycles30 UnhaltRefCycles30 APERF30 MPERF30 TSC30 InstRet31 UnhaltClkCycles31 APERF31 MPERF31 TSC31 InstRet32 UnhaltClkCycles32 UnhaltRefCycles32 APERF32 MPERF32 TSC32 InstRet33 UnhaltClkCycles33 UnhaltRefCycles33 APERF33 MPERF33 TSC33 InstRet34 UnhaltClkCycles34 UnhaltRefCycles34 APERF34 MPERF34 TSC34 InstRet35 UnhaltClkCycles35 UnhaltRefCycles35 APERF35 MPERF35 TSC35 InstRet36 UnhaltClkCycles36 UnhaltRefCycles36 APERF36 MPERF36 TSC36 InstRet37 UnhaltClkCycles37 UnhaltRefCycles37 APERF37 MPERF37 TSC37 InstRet38 UnhaltRefCycles38 APERF38 MPERF38 TSC38 InstRet39 UnhaltClkCycles39 UnhaltRefCycles39 APERF39 MPERF39 TSC39 InstRet40 UnhaltRefCycles40 APERF40 MPERF40 TSC40 InstRet41 UnhaltClkCycles41 UnhaltRefCycles41 APERF41 MPERF41 TSC41 InstRet42 UnhaltClkCycles42 APERF42 MPERF42 TSC42 InstRet43 UnhaltClkCycles43 UnhaltRefCycles43 APERF43 MPERF43 TSC43 InstRet44 UnhaltClkCycles44 UnhaltRefCycles44 APERF44 MPERF44 InstRet45 UnhaltClkCycles45 UnhaltRefCycles45 APERF45 MPERF45 TSC45 InstRet46 UnhaltClkCycles46 UnhaltRefCycles46 APERF46 MPERF46 InstRet47 UnhaltClkCycles47 UnhaltRefCycles47 APERF47 MPERF47 TSC47 InstRet48 UnhaltClkCycles48 UnhaltRefCycles48 APERF48 MPERF48 TSC48 InstRet49 UnhaltRefCycles49 APERF49 MPERF49 TSC49 InstRet50 UnhaltClkCycles50 UnhaltRefCycles50 APERF50 MPERF50 TSC50 InstRet51 UnhaltClkCycles51 APERF51 MPERF51 TSC51 InstRet52 UnhaltClkCycles52 UnhaltRefCycles52 APERF52 MPERF52 TSC52 InstRet53 UnhaltClkCycles53 UnhaltRefCycles53 APERF53 MPERF53 TSC53 InstRet54 UnhaltClkCycles54 UnhaltRefCycles54 APERF54 MPERF54 TSC54 InstRet55 UnhaltClkCycles55 UnhaltRefCycles55 InstRet56 UnhaltClkCycles56 UnhaltRefCycles56 APERF56 MPERF56 TSC56 InstRet57 UnhaltClkCycles57 UnhaltRefCycles57 APERF57 MPERF57 InstRet58 UnhaltClkCycles58 UnhaltRefCycles58 APERF58 MPERF58 TSC58 InstRet59 UnhaltClkCycles59 UnhaltRefCycles59 APERF59 MPERF59 TSC59 InstRet60 UnhaltClkCycles60 APERF60 MPERF60 TSC60 InstRet61 UnhaltClkCycles61 UnhaltRefCycles61 APERF61 MPERF61 TSC61 InstRet62 UnhaltClkCycles62 APERF62 MPERF62 TSC62 InstRet63 UnhaltClkCycles63 UnhaltRefCycles63 APERF63 MPERF63 TSC63 InstRet64 UnhaltClkCycles64 UnhaltRefCycles64 APERF64 MPERF64 TSC64 InstRet65 UnhaltClkCycles65 UnhaltRefCycles65 APERF65 MPERF65 TSC65 InstRet66 UnhaltClkCycles66 UnhaltRefCycles66 APERF66 MPERF66 TSC66 InstRet67 UnhaltClkCycles67 UnhaltRefCycles67 APERF67 MPERF67 TSC67 InstRet68 UnhaltClkCycles68 UnhaltRefCycles68 APERF68 MPERF68 TSC68 InstRet69 UnhaltClkCycles69 APERF69 MPERF69 TSC69 InstRet70 UnhaltClkCycles70 UnhaltRefCycles70 APERF70 MPERF70 TSC70 InstRet71 UnhaltClkCycles71 APERF71 MPERF71 TSC71 _POWMON 1628208953868 0.000000 100.000000 144.000000 0.000000 0.000000 100.000000 144.000000 0.000000 0.000000 895574401226890161 632743034114 100.000000 144.000000 0.000000 0.000000 100.000000 144.000000 0.000000 0.000000 8955744037 392256130260 454468949151 982249453375 1399677666440 8755940531408343 159197386780 381119499471 426251214186 949790167949 13410975 136462196306 347052449930 392652822009 924160986472 1301406531220 8755940531415175 132748021525 332227876178 35739148170632438 8755940531401518 130324684934 333077226054 371584667727 779232433354 1058055591365 8755940531413649 245780832714 4315867294594953 1342673242309 8755940531412543 221747875933 389400867040 416948368095 775187689543 1002596799690 8755940531414755 214158673281 423898284103 507596248893 1342673242309 8755940531412543 221747875933 389400867040 416948368095 775187689543 1002596799690 8755940531414755 214158673281 423898284103 5087596248893 887008004746 1232053319407 8755940531417765 206962096267 384042638043 38944395628 735122632857 923189400909 8755940531414916 247293457672
```

Intel



# Variorum Current Support (as of v0.4.1)

---

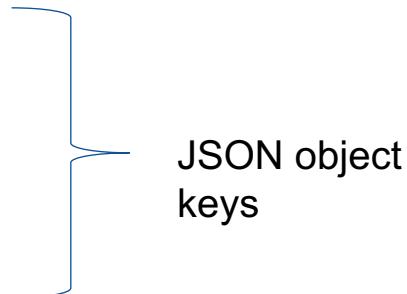
- Initial v0.1.0 released Nov 11, 2019
  - Platforms and microarchitectures supported:
    - Intel: Kaby Lake, Skylake, Broadwell, Haswell, Ivy Bridge, Sandy Bridge
    - IBM: Power9
- Current release (April 2021), v0.4.1:
  - Platforms and microarchitectures supported:
    - Nvidia: Volta
    - ARM: Juno
    - AMD (under review)
  - JSON API to integrate with external tools (e.g., Kokkos, Flux)



<https://github.com/llnl/variorum>

# Adding a vendor-neutral JSON interface

- Many of Variorum's APIs are printing output to stdout for user to parse
  - While nice for providing a friendly interface to understanding the hardware-level metrics, this **limits ability for Variorum to provide these metrics to an external tool**
- Added int variorum\_get\_node\_power\_json(json\_t \*) to integrate variorum with other tools (*e.g.*, Flux and Kokkos)
  - { "hostname": (string),
  - "timestamp": (int),
  - "power\_node": (int),
  - "power\_cpu\_socket\_<id>": (int)
  - "power\_mem\_socket\_<id>": (int)
  - "power\_gpu\_socket\_<id>": (int) }
- Example: Reporting end-to-end power usage for Kokkos loops
- Example: Provide power-awareness to Flux scheduling model enabling resources to be assigned based on available power



# Implementing a vendor-neutral node-level power cap interface

- What if the device does not provide an explicit power knob for the node?
- IBM Power9 provides a node-level power knob through OPAL (part of IBM firmware)
  - But, not all platforms provide this same interface
  - Intel provides a socket-level power knob
- Intel provides CPU- and DRAM-level power knobs through MSRs, while IBM provides a power shifting ratio to determine the power limit of the CPU and GPU components
- Variorum's current implementation:

```
int variorum_cap_best_effort_node_power_limit(int power_lim)
```

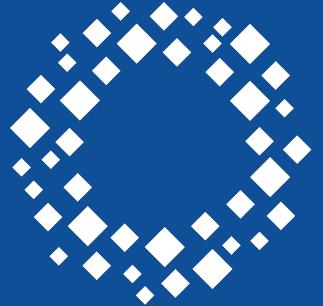
  - IBM: caps the power limit to the node as expected
  - Intel: uniformly distribute the power limit across all sockets in the node (ignore memory and uncore power)

# Join Us for Module 2 on Aug 13, 11:30AM-1:00PM ET

## “Integrating Variorum with System Software and Tools”

- The HPC Power Stack revisited: need for power management at various levels
- GEOPM: job-level power management
- Flux and SLURM (Research Extensions): system-level power management
- Kokkos and Caliper: application and workflow power management
- Upcoming Features in Variorum
- The HPC Power Stack Roadmap

Both modules will be repeated  
on Aug 20 and Aug 23



# CASC

Center for Applied  
Scientific Computing

 Lawrence Livermore  
National Laboratory

#### **Disclaimer**

This document was prepared as an account of work sponsored by an agency of the United States government. Neither the United States government nor Lawrence Livermore National Security, LLC, nor any of their employees makes any warranty, expressed or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States government or Lawrence Livermore National Security, LLC. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States government or Lawrence Livermore National Security, LLC, and shall not be used for advertising or product endorsement purposes.