



ZeroSum:

*User Space Monitoring of Resource Utilization and Contention
on Heterogeneous HPC Systems*

Kevin A. Huck

Oregon Advanced Computing Institute for Science and Society (OACISS)



UNIVERSITY OF
OREGON

Performance Analysis Perspective on Monitoring

Three main classes of performance optimizations:

1. Algorithmic replacement (usually a high level of difficulty)
 - E.g. replace $O(n^2)$ with $O(n \log n)$
 - Can involve data structure changes, new dependencies, major rewrites
2. Code optimization (usually a medium difficulty)
 - Improve cache reuse, reduce stalls (branching, instructions, I/O, etc)
3. Optimized launch configuration (low difficulty, high embarrassment potential)
 - Misconfiguration
 - Wrong assumptions from another system/application
 - Changes to system policies/defaults (e.g. reserved cores)

Behavioral Economics & Computer Science

“People aren’t dumb – the world is hard.”

– Richard Thaler, behavioral economist and co-author (with Cass Sunstein) of *Nudge: Improving Decisions About Health, Wealth, and Happiness*, 2008

People make bad choices for good reasons...

If you want to change behavior, change the defaults.

(not always possible – “one size fits most”)

Motivation: Why do users (want to) monitor at runtime?

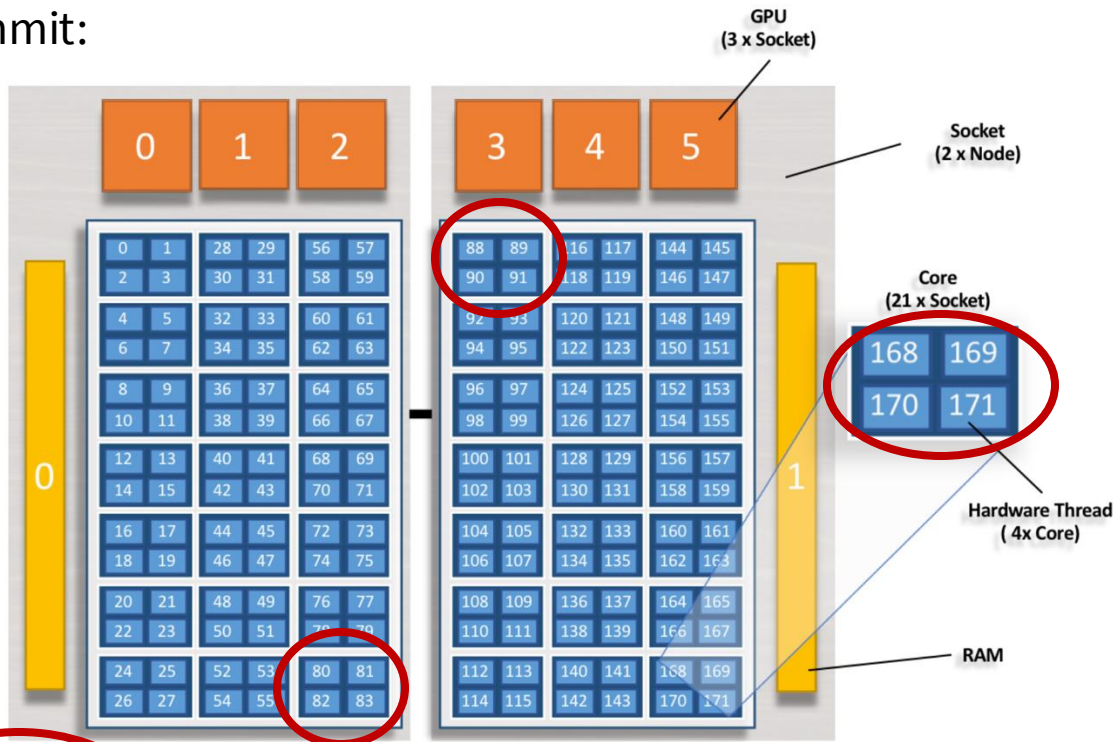
- **Sanity check** / curiosity / impatience (logging, essentially)
- **Check for misconfiguration**
- **Check for efficient utilization**
- Confirmation of expected hardware / operating system behavior
- **Identify cause of failure – deadlock, crash, stalls, etc.**
- Adaptation / computational steering / feedback & control
- ~~Identify system failures (out of scope)~~

(Mis)Configuration – what could *possibly* go wrong?

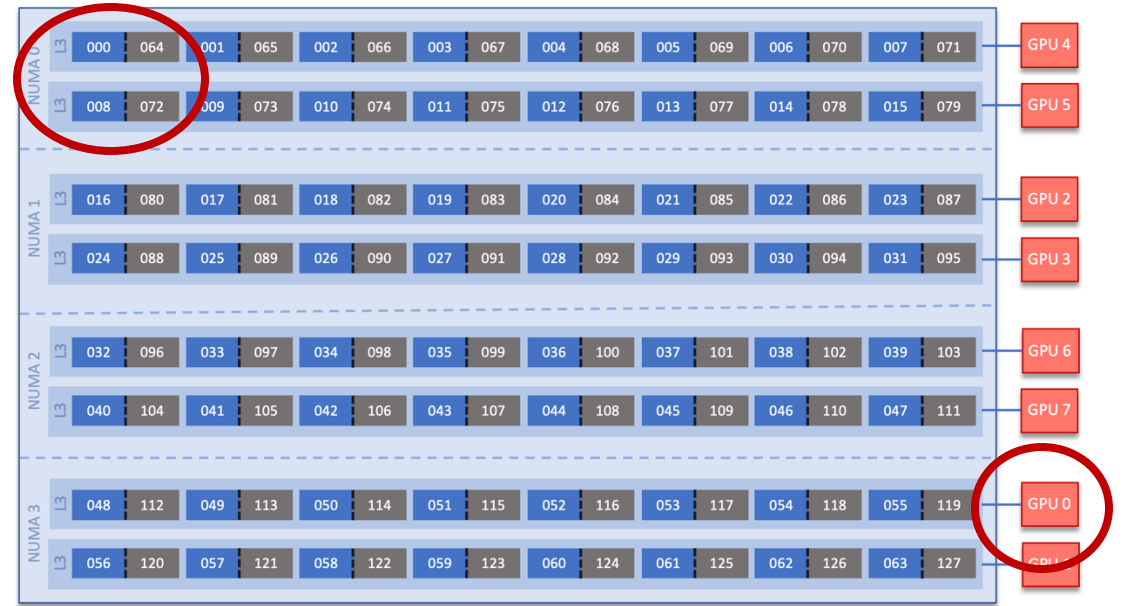
- Process placement:
 - Logical/physical mappings, resources assigned/constrained to each process
 - Are there reserved core(s) for system? GPU mapping?
- Thread placement: Socket, NUMA domain, core, thread (HWT)
- Undersubscribing: Wasted hardware, energy, time (under-utilization)
 - Can provide better performance in some situations (memory-bound code)
- Oversubscribing: Increased contention with no realized benefit
- Imbalances
 - What is the communication frequency/volume between pairwise MPI ranks?
- Slurm/PBS/Alps/Torque/Flux are *complicated to use*
 - ...especially when combined with MPI, OpenMP, GPUs, or other model settings

Although accurate, system documentation can be confusing

Summit:



Frontier:



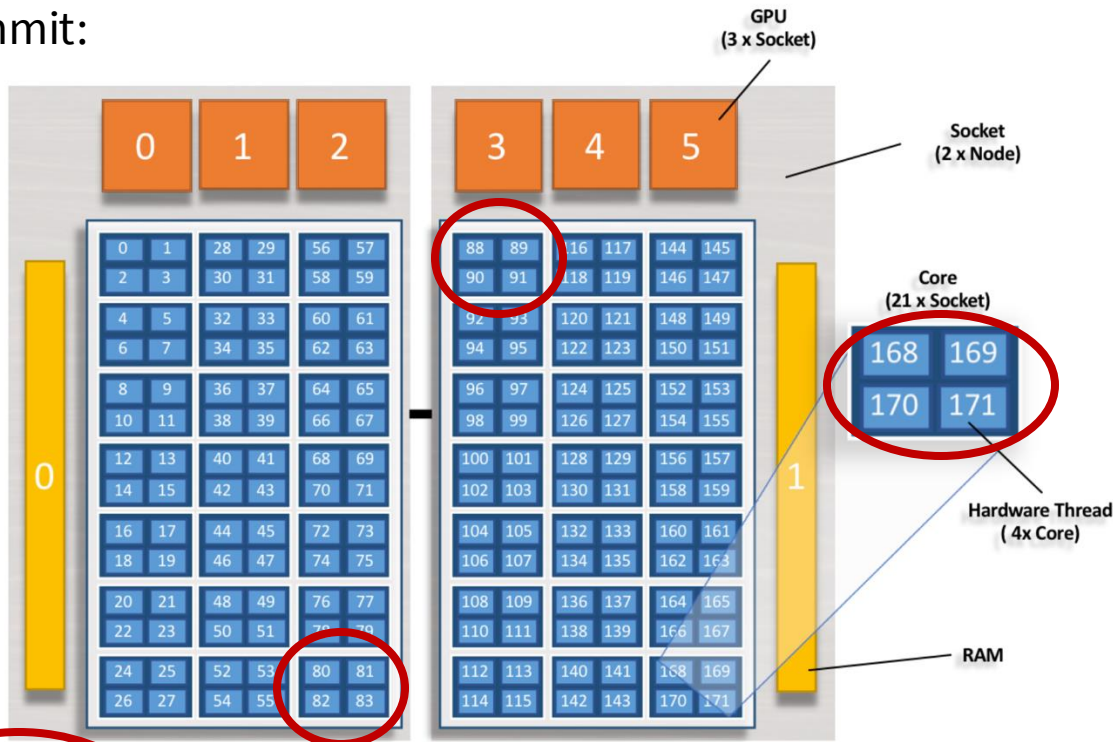
- 1 node
- 2 sockets (grey)
- 42 physical cores* (dark blue)
- 168 hardware cores (light blue)
- 6 GPUs (orange)
- 2 Memory blocks (yellow)

Core indexing, HWT indexing, GPU mapping not universal (or even intuitive)

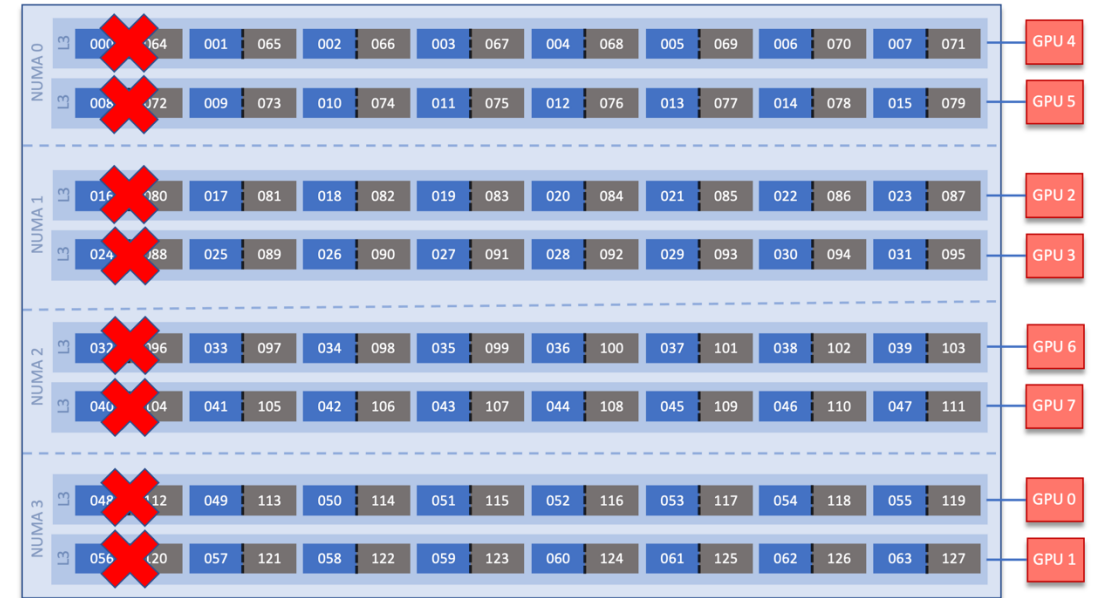
*Core Isolation: 1 core on each socket has been set aside for overhead and is not available for allocation through jsrun. The core has been omitted and is not shown in the above image.

Although accurate, system documentation can be confusing

Summit:



Frontier:



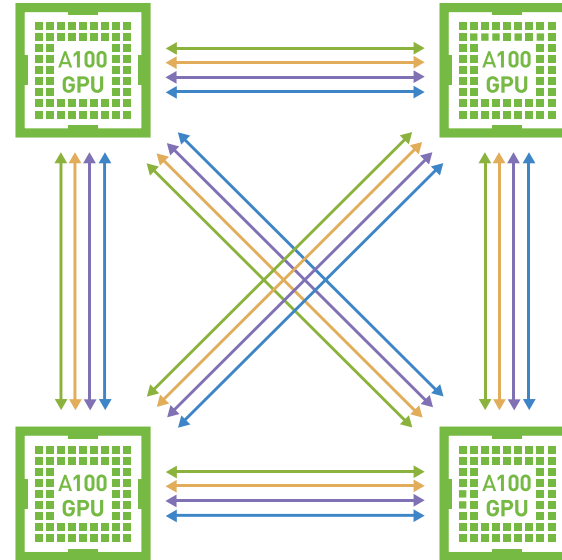
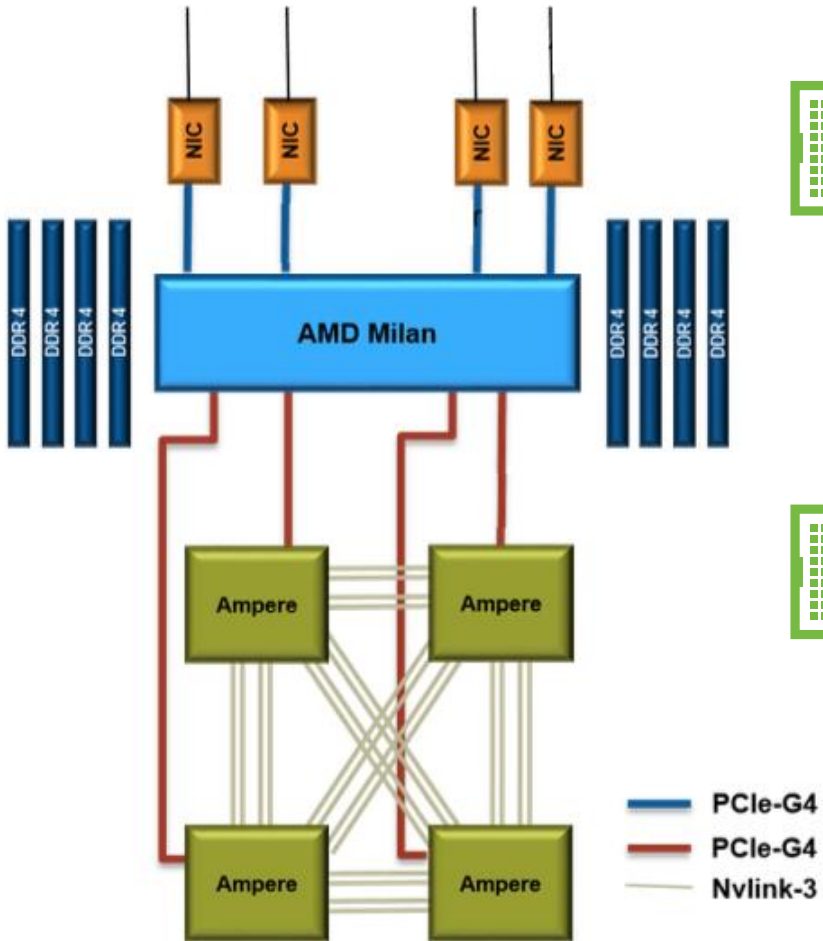
- 1 node
- 2 sockets (grey)
- 42 physical cores* (dark blue)
- 168 hardware cores (light blue)
- 6 GPUs (orange)
- 2 Memory blocks (yellow)

***Core Isolation:** 1 core on each socket has been set aside for overhead and is not available for allocation through jsrun. The core has been omitted and is not shown in the above image.

Modified default: system reserves 8 cores...but user controllable

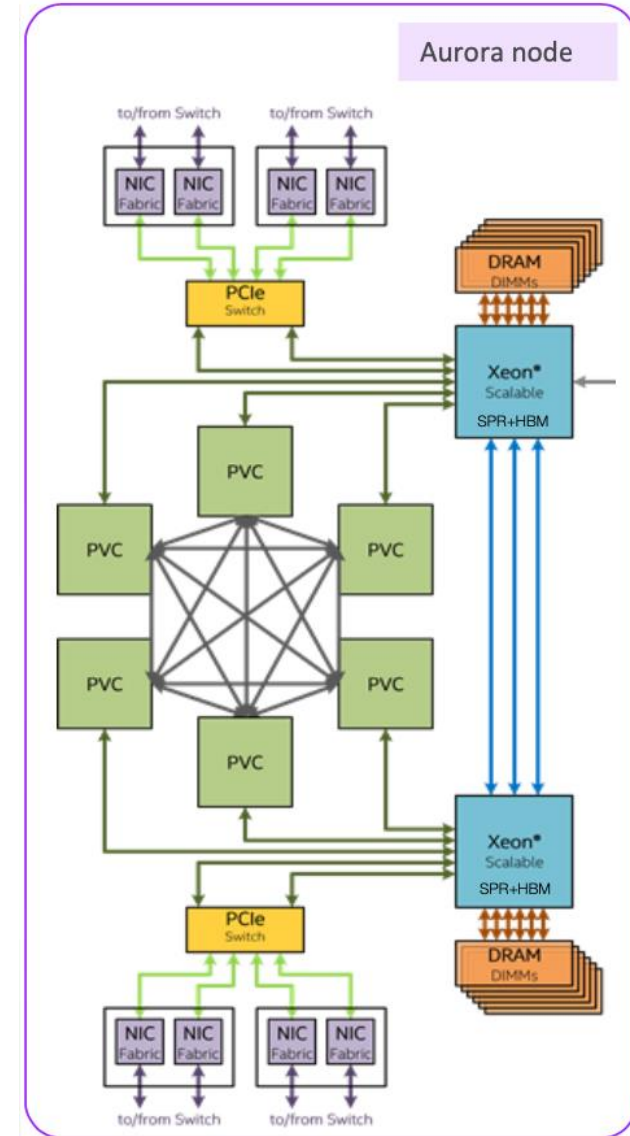
More systems - information missing

Perlmutter:



No physical layout information available in the documentation, no core or GPU indexing...

Aurora:



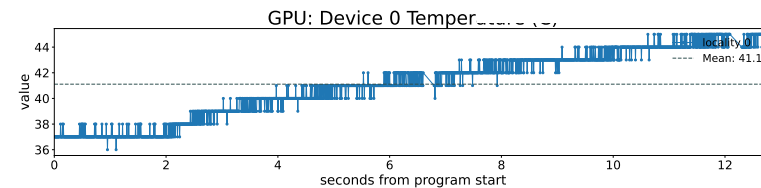
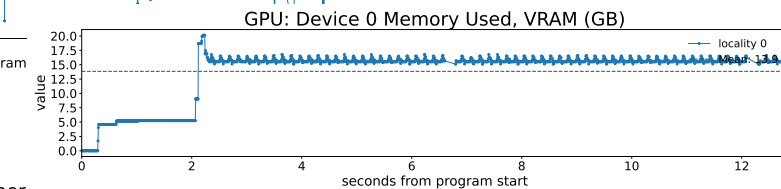
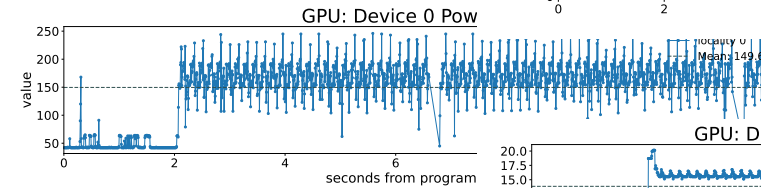
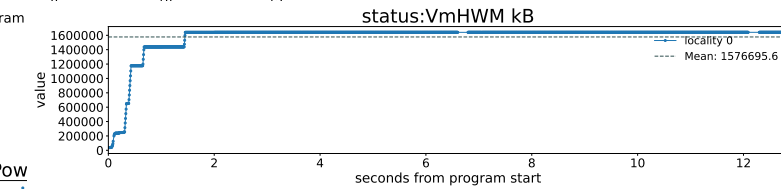
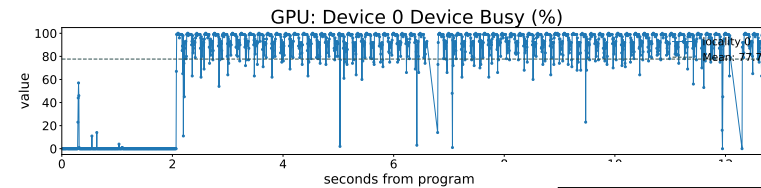
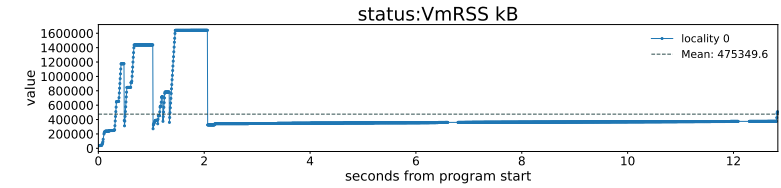
Utilization

- Resource Monitor, Activity Monitor, menumeters, top, htop, NVML, ROCm-SMI, etc.
 - We (users) LOVE these!
- Why can't we have that for HPC?
 - LDMS, Ganglia, Puppet Console, TACC Stats, JobStats, etc.
- ...for the user?

```

khuck ~ ssh - ssh sever - zsh - 100x20
1  [|||||] 100.0% 5  [|||||] 100.0%
2  [|||||] 100.0% 6  [|||||] 100.0%
3  [|||||] 100.0% 7  [|||||] 100.0%
4  [|||||] 100.0% 8  [|||||] 100.0%
Mem [|||||] 2.03G/31.1G Tasks: 106, 264 thr; 8 running
Swp [|||||] 0K/75.2G Load average: 1.54 0.57 0.20
Uptime: 34 days, 18:10:58

  PID USER   PRI  NI  VIRT   RES   SHR  S  CPU% MEM%  TIME+  Command
1522757 khuck  20   0 1461M 363M 194M R 799. 1.1 0:26.34 ./build/bin/lu-decomp
1522759 khuck  20   0 1461M 363M 194M R 100. 1.1 0:03.29 ./build/bin/lu-decomp
1522760 khuck  20   0 1461M 363M 194M R 100. 1.1 0:03.29 ./build/bin/lu-decomp
1522764 khuck  20   0 1461M 363M 194M R 100. 1.1 0:03.27 ./build/bin/lu-decomp
1522765 khuck  20   0 1461M 363M 194M R 100. 1.1 0:03.26 ./build/bin/lu-decomp
1522762 khuck  20   0 1461M 363M 194M R 99.5 1.1 0:03.27 ./build/bin/lu-decomp
1522763 khuck  20   0 1461M 363M 194M R 99.5 1.1 0:03.27 ./build/bin/lu-decomp
1522761 khuck  20   0 1461M 363M 194M R 99.5 1.1 0:03.26 ./build/bin/lu-decomp
1522758 khuck  20   0 1461M 363M 194M S 0.7 1.1 0:00.06 ./build/bin/lu-decomp
F1Help F2Setup F3Search F4Filter F5Tree F6SortBy F7Nice F8Nice +F9Kill F10Quit
    
```



ZeroSum: User Space Monitoring of Resource Utilization and Contention

- Inspired by the *hello_jsub* program from Tom Papatheodore:
https://code.ornl.gov/t4p/Hello_jsrun
- Why “ZeroSum”? – An advantage for one side results in an equivalent loss for the other
 - the fixed number of resources available in an allocation – no elasticity
 - the need to periodically use *some* resources to monitor
- Available on GitHub: <https://github.com/UO-OACISS/zerosum>
- Monitors application threads (LWP), CPU hardware (HWT), Memory, and GPU hardware for all processes, all nodes in the allocation

ZeroSum Functionality

- ✓ **Detect the initial/changing configuration of the application**
- ❑ Evaluate the configuration to automatically detect misconfigurations
- ✓ Provide runtime feedback to the user that the program is progressing
- ✓ **Provide a report of how effectively the hardware was utilized**
- ✓ **Provide a report of how much contention was identified in the execution**
- ❑ Provide a way to expose the observed data to other tools that can perform computational steering / runtime optimization / reconfiguration

- ✓ Implemented
- ❑ Future work

How to use ZeroSum

- Wrapper script(s) – `zerosum` and `zerosum-mpi`
 - Periodicity – default 1 second
 - Detailed/verbose output – true/false, default false
 - Heartbeat (memory consumption) – true/false, default false
 - Register signal handler – true/false, default false
 - Run in debugger – true/false, default false (also specify executable name)
 - Detect deadlocks – true/false, default false
 - Deadlock period – how many periods of “inactivity” is considered “deadlock”
 - HWT/Core for async thread – defaults to last core/HWT in affinity list for process
- Preloads the library, wraps `__libc_start_main` or creates global static constructor/destructor functions

What does it do? ...Configuration Detection

- Query `/proc/[self|pid]/status` to get the allowable cores
- Query `/proc/meminfo` to get total memory available
- Query MPI rank, size, hostname (after `MPI_Initialized()` returns `true`)
- If available, use `hwloc` to query hardware topology (`lstopo`)
- Asynchronous background thread is started, it periodically queries:
 - `/proc/[self|pid]/status` to get process thread count, memory usage
 - `/proc/[self|pid]/task` directory to get all thread IDs
 - For each thread, query the affinity list for that thread (it may change!), utilization, state, core/HWT it's running on, context switches, other metrics
 - `/proc/stat` to query core utilization of all cores
- OMP-Tools (v5.0+) callback used to identify OpenMP threads at creation*
- NVML/ROCm-SMI/SYCL libraries used to query GPU(s)

Utilization Report

- Rank 0 writes a summary report to the screen
- All ranks write a report to a log file – including full time series data as CSV
- All observed threads (LWP) are reported – user/system/idle, context switches, affinity list
- All assigned cores (HWT) are reported – user/system/idle
- All assigned GPUs stats are reported

Example Output – miniQMC (OpenMP target offload) on Frontier (OLCF)

Duration of execution: 210.878 s

Process Summary:

Process Summary

MPI 000 - PID 51334 - Node frontier09085 - CPUs allowed: [1,2,3,4,5,6,7]

LWP (thread) Summary:

LWP 51334: Main,OpenMP - stime: 12.48, utime: 63.94, nv_ctx: 4, ctx: 365488, CPUs: [1]
LWP 51343: ZeroSum - stime: 0.15, utime: 0.26, nv_ctx: 9, ctx: 679, CPUs: [7]
LWP 51374: Other - stime: 0.00, utime: 0.00, nv_ctx: 0, ctx: 6, CPUs:
[1-7,9-15,17-23,25-31,33-39,41-47,49-55,57-63,65-71,73-79,81-87,89-95,97-103,
105-111,113-119,121-127]
LWP 51384: OpenMP - stime: 12.60, utime: 64.00, nv_ctx: 3, ctx: 365742, CPUs: [3]
LWP 51385: OpenMP - stime: 12.63, utime: 64.27, nv_ctx: 2, ctx: 352574, CPUs: [5]
LWP 51386: OpenMP - stime: 12.74, utime: 63.76, nv_ctx: 473, ctx: 368585, CPUs: [7]

LWP (thread) Summary

Note: times are in jiffies

Example Output – miniQMC (OpenMP target offload) on Frontier (OLCF)

HWT (core/thread) Summary

Hardware Summary:

```
CPU 001 - idle: 22.70, system: 12.42, user: 64.52
CPU 002 - idle: 99.82, system: 0.00, user: 0.00
CPU 003 - idle: 23.08, system: 12.60, user: 63.97
CPU 004 - idle: 99.83, system: 0.00, user: 0.00
CPU 005 - idle: 22.79, system: 12.62, user: 64.23
CPU 006 - idle: 99.83, system: 0.00, user: 0.00
CPU 007 - idle: 22.94, system: 12.89, user: 63.81
```

Note: times are in jiffies

GPU Summary

```
GPU 0 - (metric: min avg max)
Clock Frequency, GLX (MHz): 800.000000 1614.691943 1700.000000
Clock Frequency, SOC (MHz): 1090.000000 1090.000000 1090.000000
Device Busy %: 0.000000 14.616114 52.000000
Energy Average (J): 0.000000 8.328571 10.000000
GFX Activity: 0.000000 17223.704762 38443.000000
GFX Activity %: 0.000000 13.706161 41.000000
Memory Activity: 0.000000 623.623810 1536.000000
Memory Busy %: 0.000000 0.355450 3.000000
Memory Controller Activity: 0.000000 0.303318 2.000000
Power Average (W): 90.000000 126.483412 138.000000
Temperature (C): 35.000000 37.909953 39.000000
UVD|VCN Activity: 0.000000 0.000000 0.000000
Used GTT Bytes: 11624448.000000 11624448.000000 11624448.000000
Used VRAM Bytes: 15044608.000000 4743346651.601895 4839596032.000000
Used Visible VRAM Bytes: 15044608.000000 4743346884.549763 4839596032.000000
Voltage (mV): 806.000000 891.848341 906.000000
```

Logs have full time series of all samples

Contention Detection Support

- Voluntary/non-voluntary context switches (analysis todo)
- Minor/major page faults, pages swapped (analysis todo)
- System time analysis (analysis todo)
- Comparing affinity lists – across threads *and* across processes (todo)
- Memory consumption (analysis todo)
- GPU memory consumption (analysis todo)
- However, **can detect deadlocks** – both **active** (i.e. spinning at MPI collective) and **passive** (i.e. waiting for mutex)

Evaluation / Example usage

MPI+OpenMP version of miniQMC on Frontier, 8 processes, 7 threads (64 cores, 1 thread per core, 8 cores reserved)

LWP	Type	stime	utime	nvctx	ctx	CPUs
18351	Main [†]	1.54	15.17	332905	1838	1
18356	ZeroSum	0.42	1.10	194	1007	1
18385	Other	0.00	0.00	0	41	1-127 [‡]
18405	OpenMP	0.31	13.09	232689	5	1
18407	OpenMP	0.44	12.93	353365	11	1
18408	OpenMP	0.21	13.22	92528	3	1
18409	OpenMP	0.47	12.93	394014	10	1
18410	OpenMP	0.37	13.03	302371	7	1
18411	OpenMP	0.41	12.97	348829	10	1

Table 1: Frontier results, default configuration. [†]indicates that the main thread is also an OpenMP thread. [‡]indicates that the first core of each L3 region was set aside for system processes, not all threads in the sequence 1-127 are allowed but summarized for brevity in the table (see LWP 51274 in Listing 2).

LWP	Type	stime	utime	nvctx	ctx	CPUs
18552	Main [†]	3.13	88.40	5	704	1-7
18561	ZeroSum	0.79	2.64	2	2790	7
18588	Other	0.00	0.00	0	41	1-127 [‡]
18589	OpenMP	1.10	90.00	9	716	1-7
18590	OpenMP	1.10	93.00	8	724	1-7
18591	OpenMP	1.07	90.52	9	692	1-7
18592	OpenMP	1.10	89.83	14	766	1-7
18593	OpenMP	1.10	90.48	7	728	1-7
18594	OpenMP	1.10	91.93	300	849	1-7

Table 2: Frontier results, configuration requesting 7 cores per process. [†]indicates that the main thread is also an OpenMP thread. [‡]indicates that the first core of each L3 region was set aside for system processes, not all threads in the sequence 1-127 are allowed but summarized for brevity in the table (see LWP 51274 in Listing 2).

LWP	Type	stime	utime	nvctx	ctx	CPUs
18948	Main [†]	3.07	88.57	2	386	1
18954	ZeroSum	0.71	2.57	2	291	7
18981	Other	0.00	0.00	0	41	1-127 [‡]
18992	OpenMP	1.18	96.36	0	422	2
18993	OpenMP	1.14	96.50	1	391	3
18994	OpenMP	1.18	96.46	0	381	4
18995	OpenMP	1.11	93.89	0	324	5
18996	OpenMP	1.14	93.29	0	370	6
18997	OpenMP	1.14	95.54	208	358	7

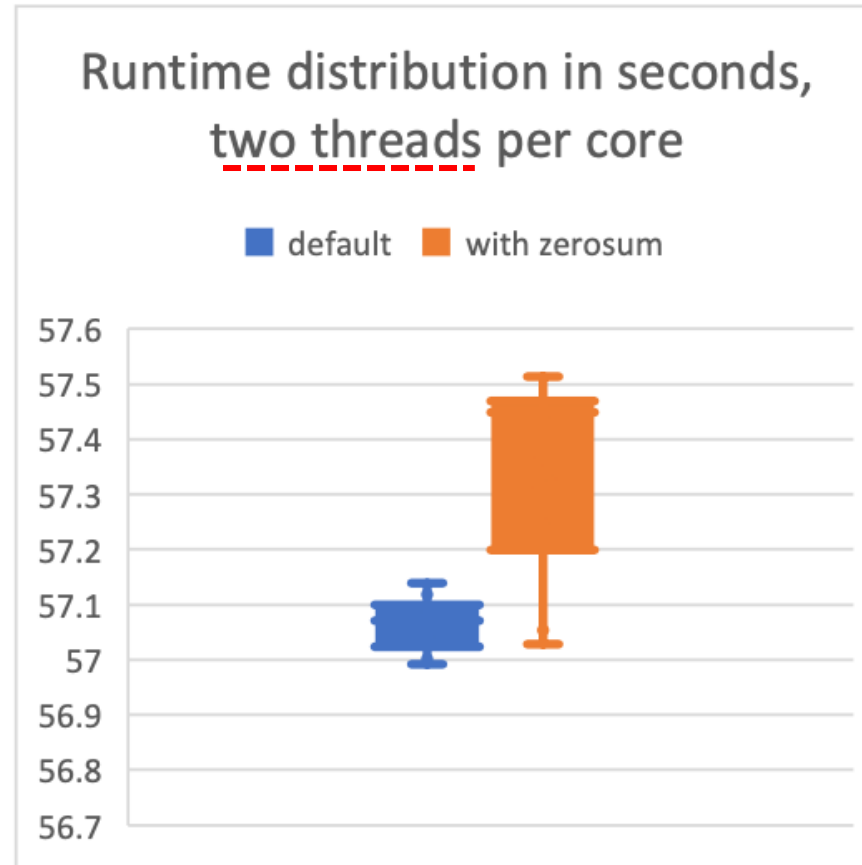
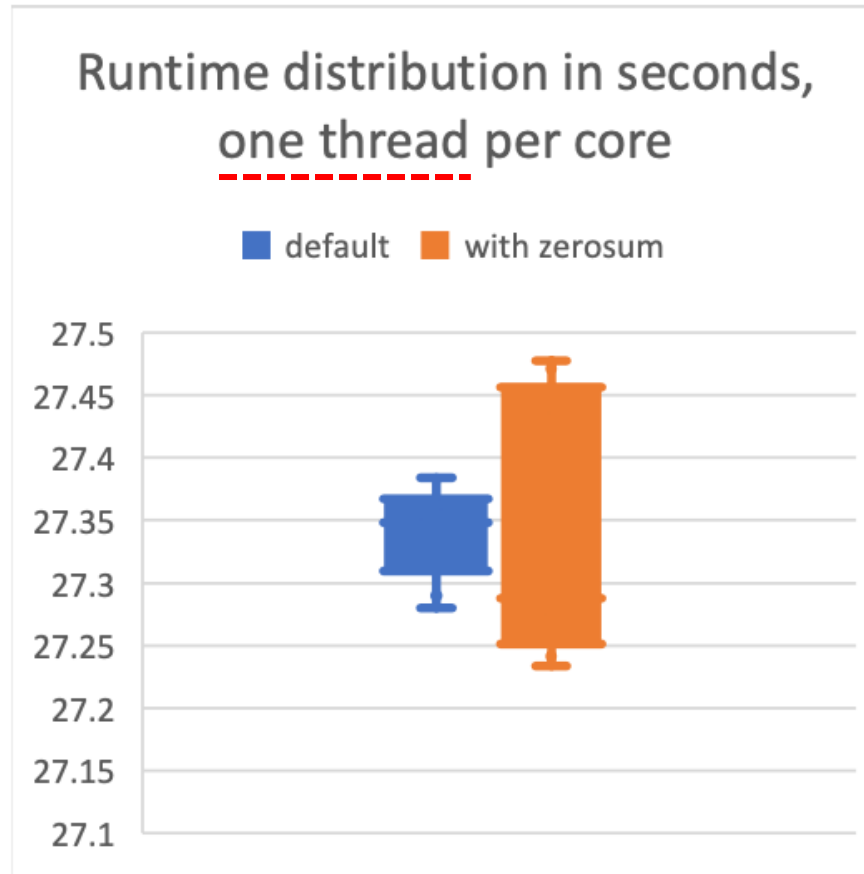
Table 3: Frontier results, configuration requesting 7 cores per process and binding OpenMP threads to cores. [†]indicates that the main thread is also an OpenMP thread. [‡]indicates that the first core of each L3 region was set aside for system processes, not all threads in the sequence 1-127 are allowed but summarized for brevity in the table (see LWP 51274 in Listing 2).

```
export OMP_NUM_THREADS=7
srun -n8 zerosum-mpi miniqmc
```

```
export OMP_NUM_THREADS=7
srun -n8 -c7 zerosum-mpi miniqmc
```

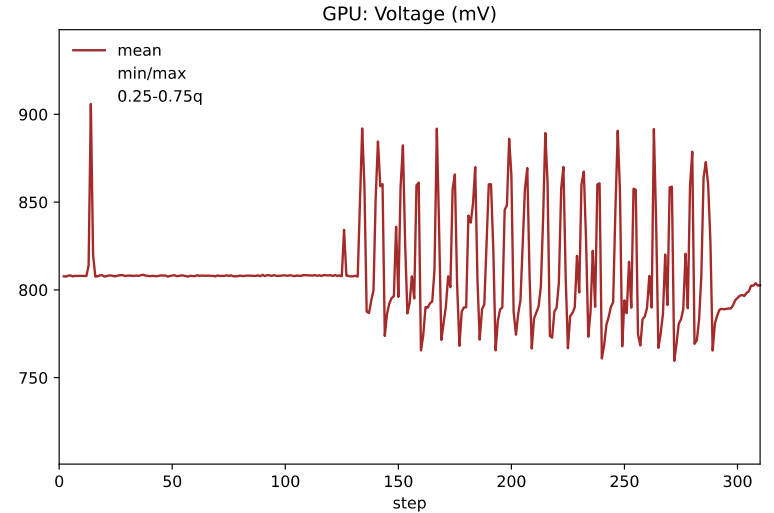
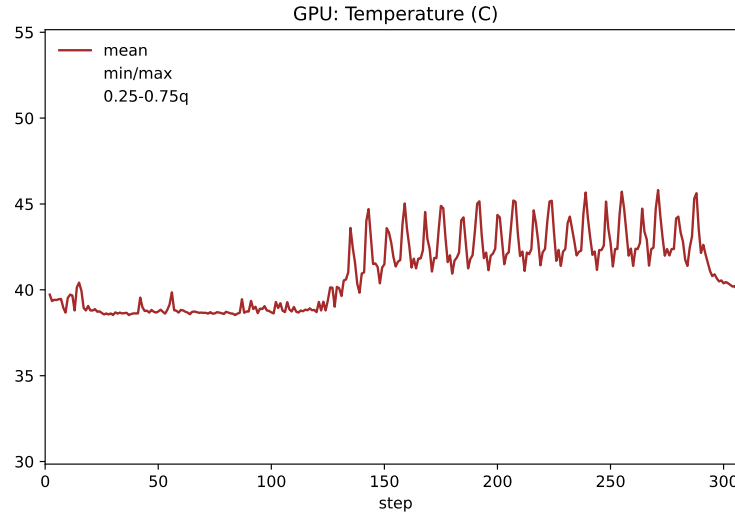
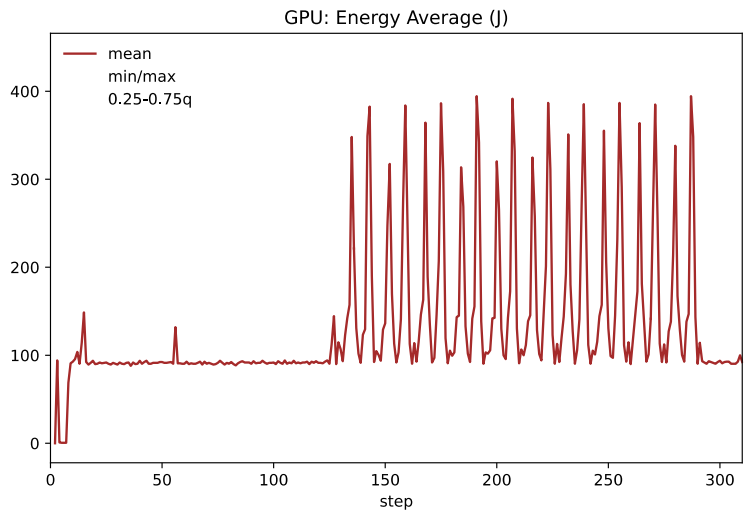
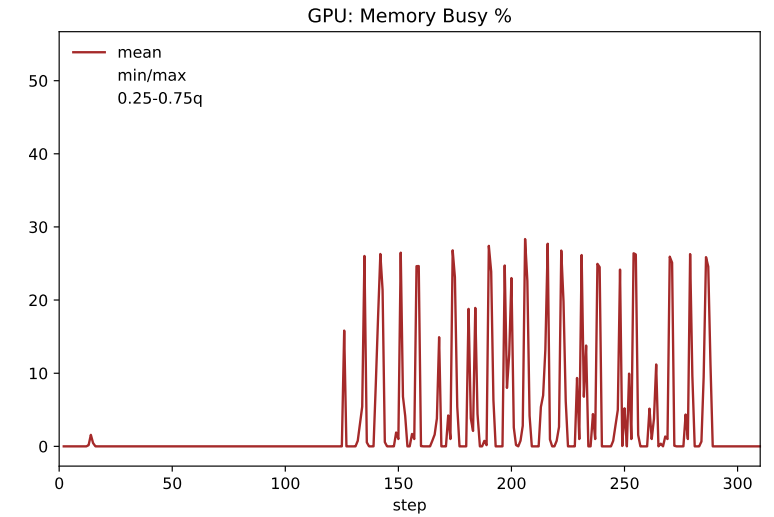
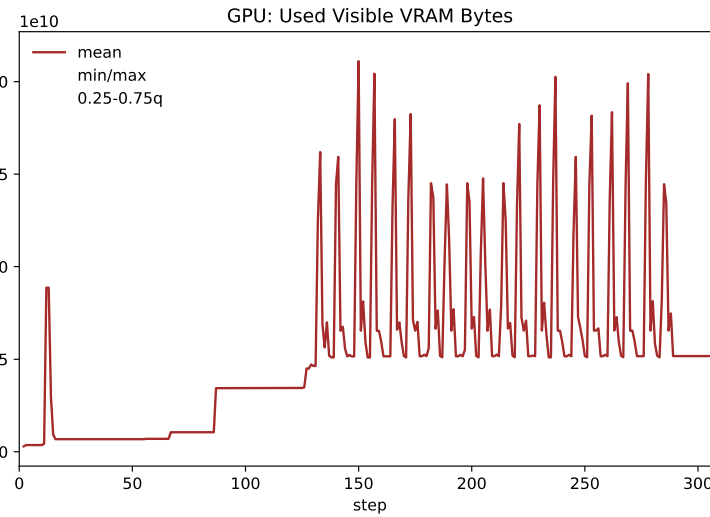
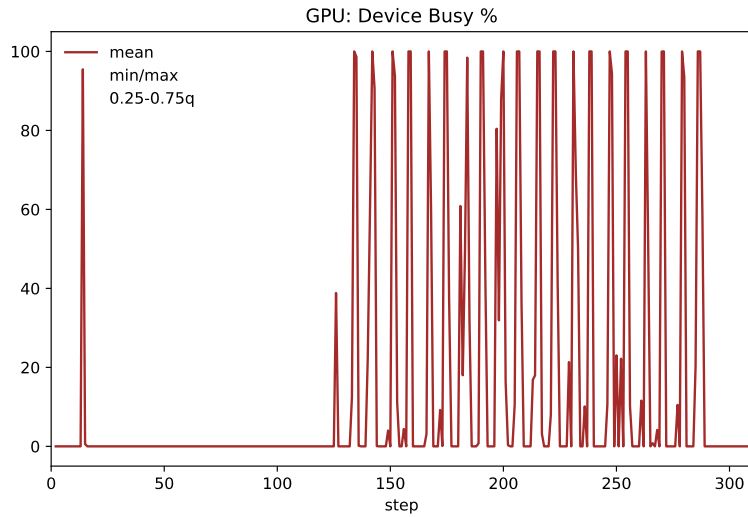
```
export OMP_NUM_THREADS=7
export OMP_PROC_BIND=spread
export OMP_PLACES=cores
srun -n8 -c7 zerosum-mpi miniqmc
```

Overhead – less than 0.5% in resource constrained example



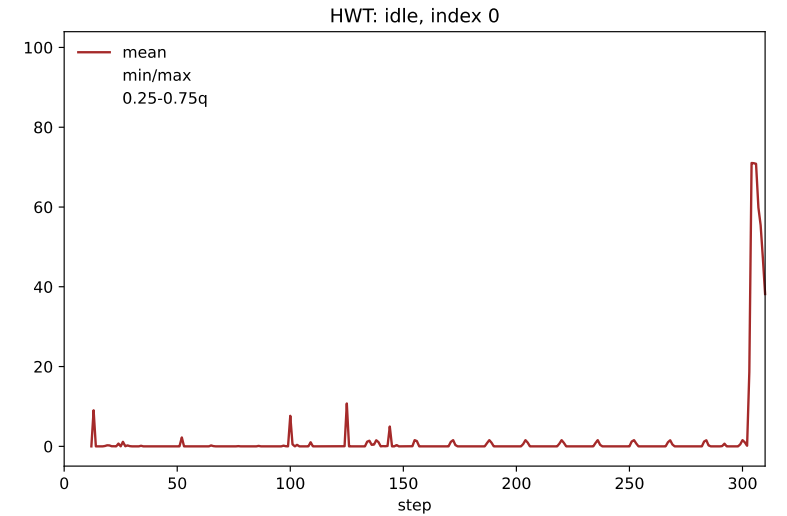
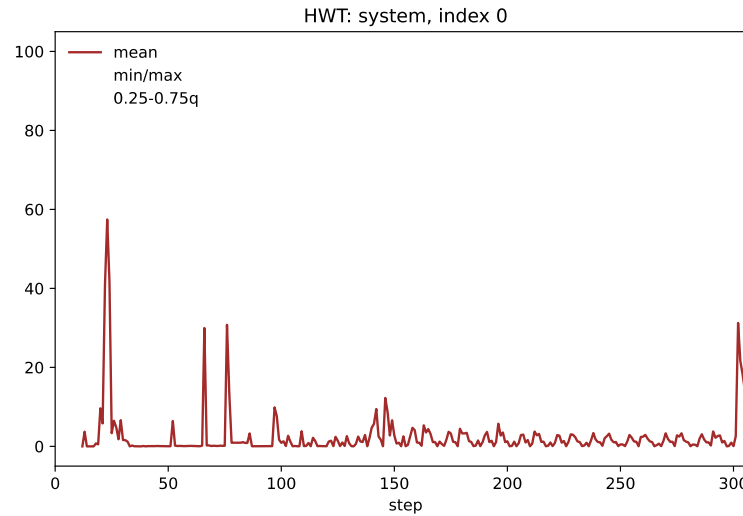
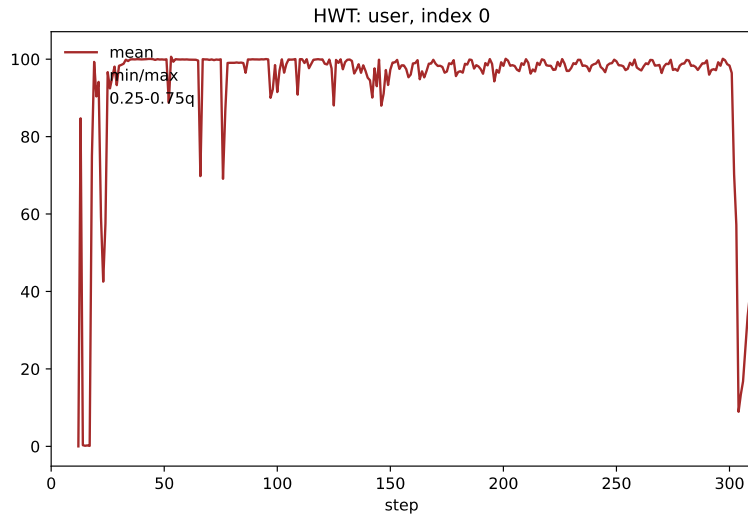
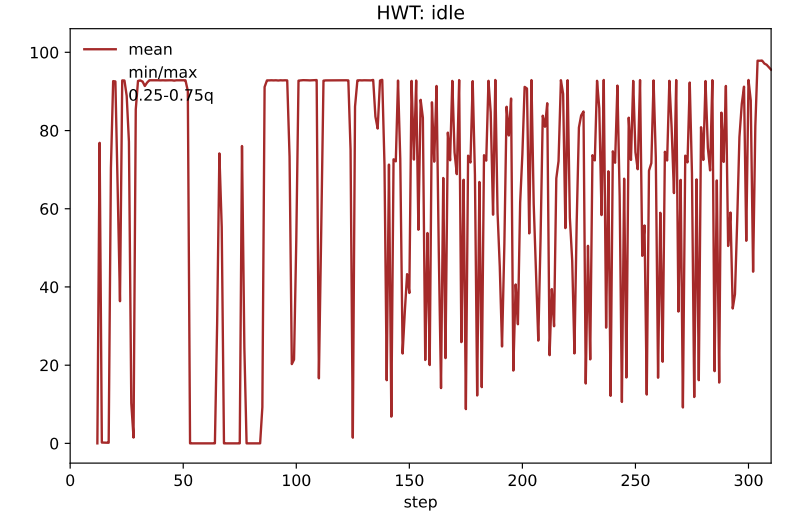
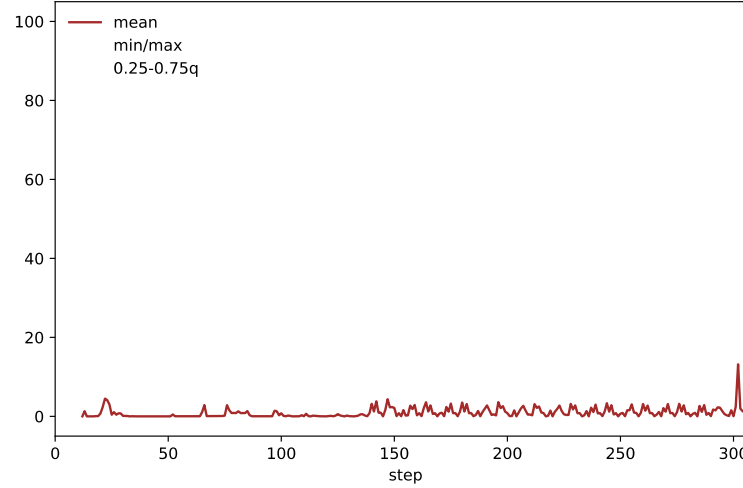
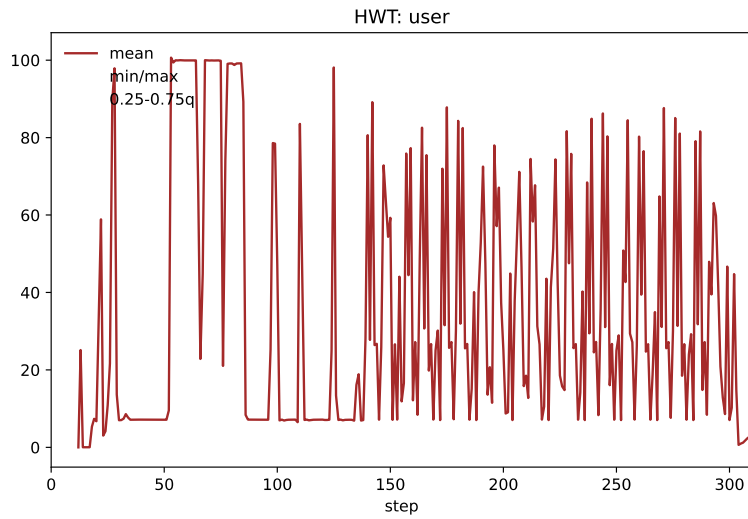
miniQMC time distributions executed 10 times using one OpenMP thread per core (left). In this comparison, the distribution of times with ZeroSum is noisier, but there is no significant observation of measurable overhead. The right figure shows the time distributions using two OpenMP threads per core. In this comparison, the distribution of times with ZeroSum is both noisier and longer tailed, and does show an observation of overhead, averaging about 0.2752 seconds, or 0.5%.

Summary views of collected data – 64 GCDs, XGC running on Frontier



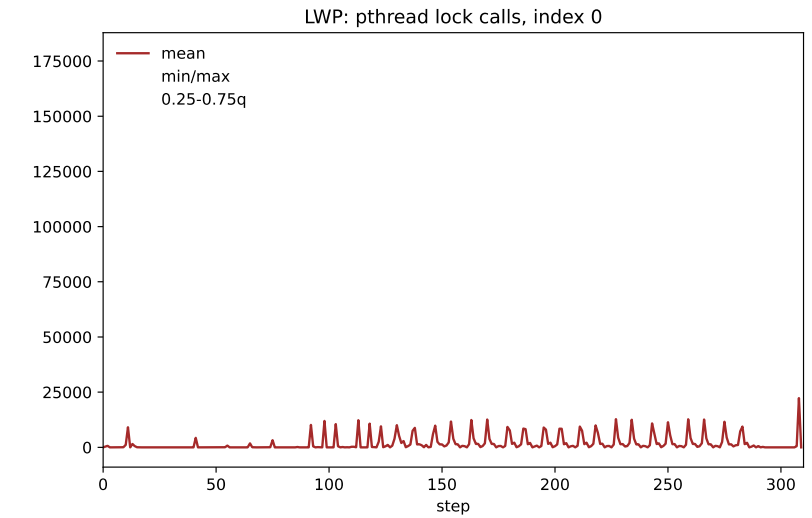
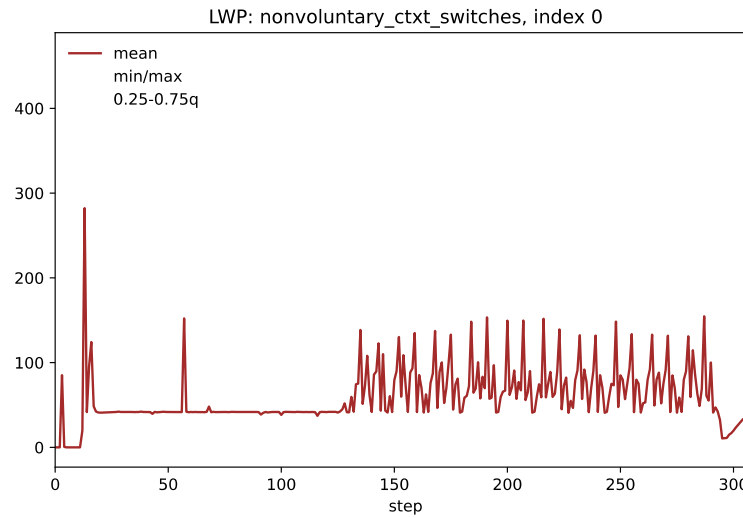
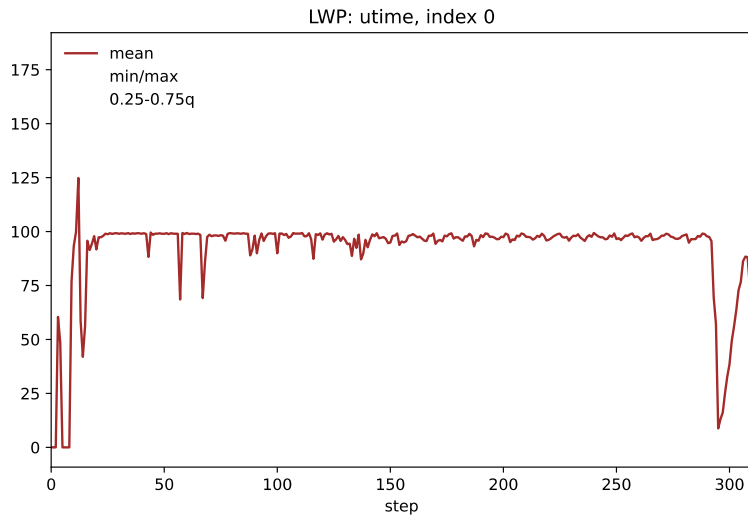
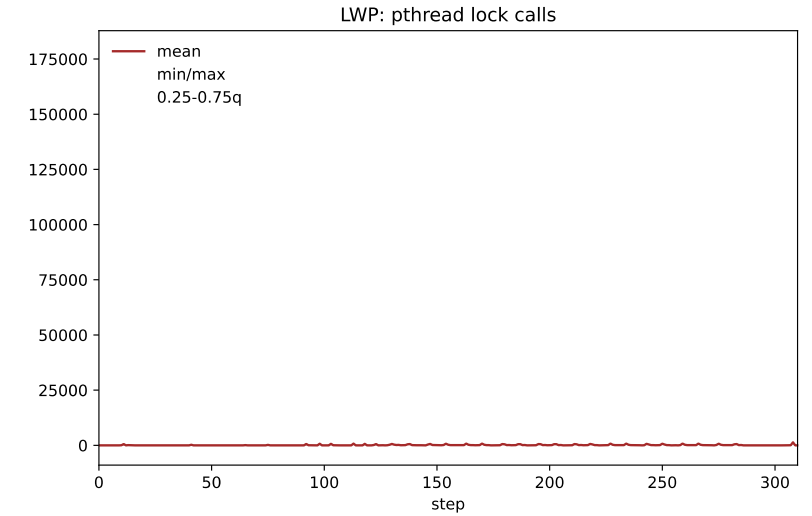
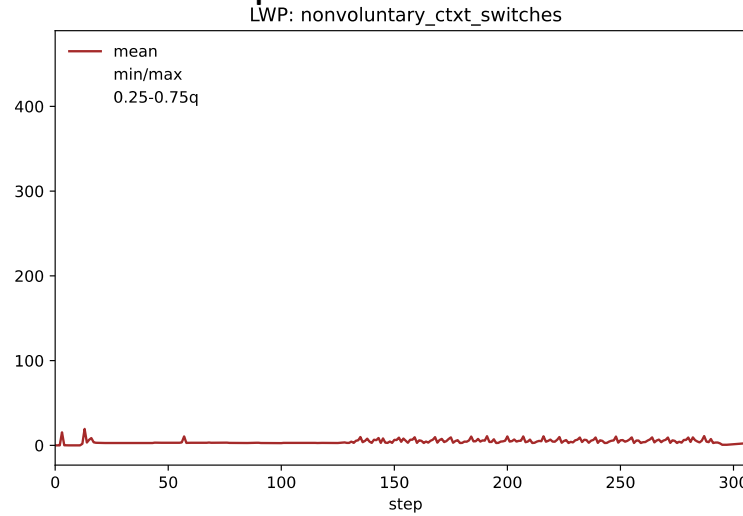
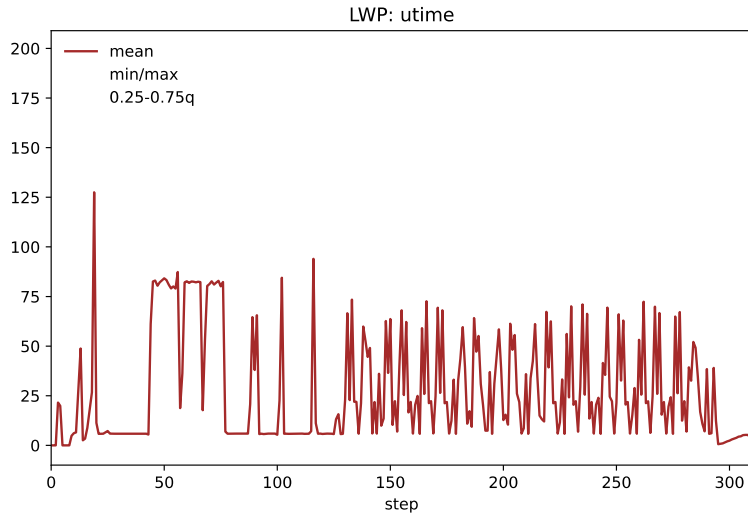
Summary views of collected data – 896 HWTs, XGC running on Frontier

↙ All hardware threads from all ranks ↘



Summary views of collected data – 1088 LWPs, XGC running on Frontier

↙ All process threads ↘

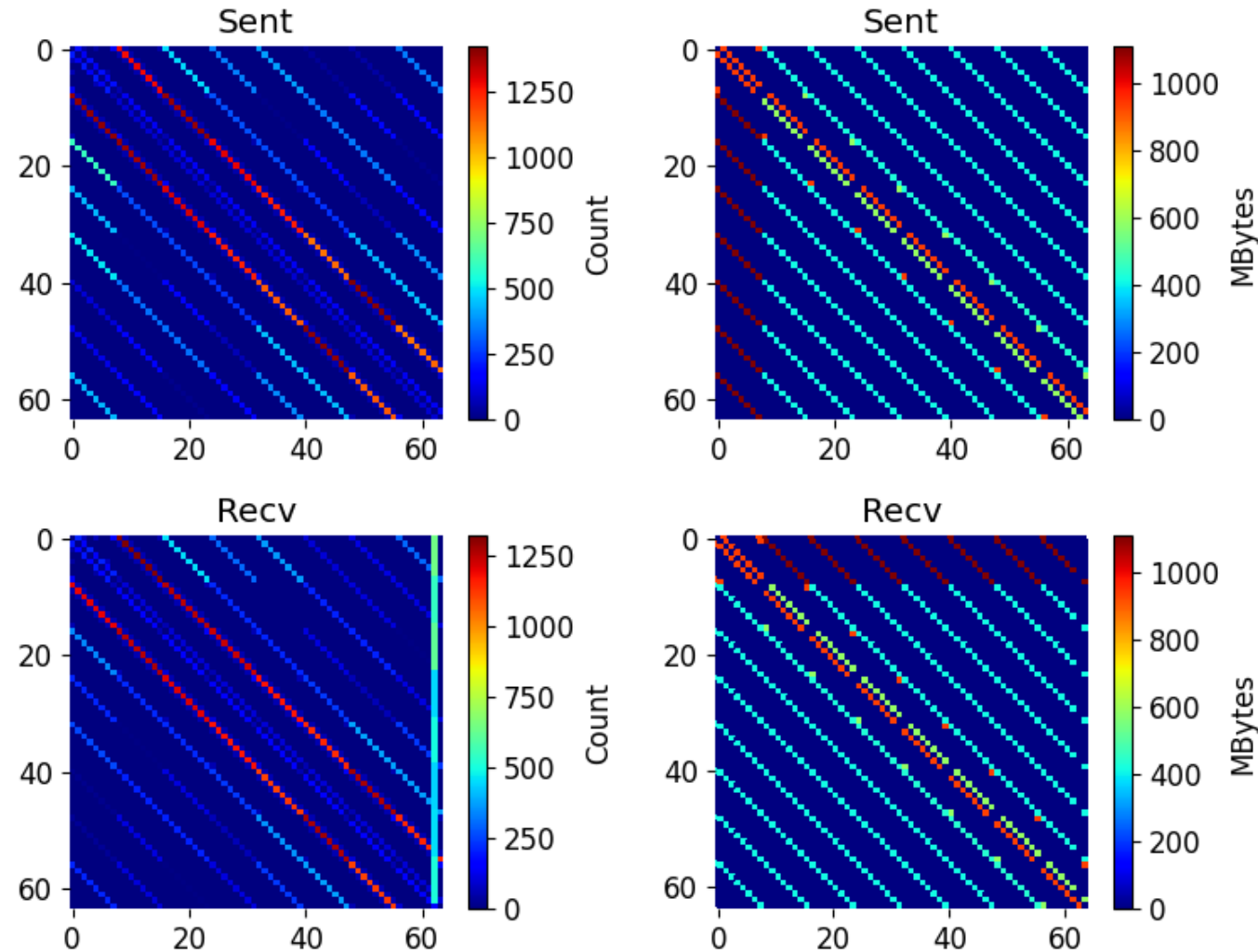


↙ Only main process thread ↗

MPI P2P data

- MPI implementation includes wrappers for `MPI_Recv/Irecv`, `MPI_Sendrecv`, `MPI_Send/Bsend/Isend/Rsend/Ssend` to capture P2P frequency, volume
- Obviously can be done by other/better performance measurement tools, but for a quick view, can be useful
- Apply analysis from “Optimizing Process-to-Core Mappings for Application Level Multi-dimensional MPI Communications”, Karlsson et al. 2012.

<https://doi.org/10.1109/CLUSTER.2012.47>



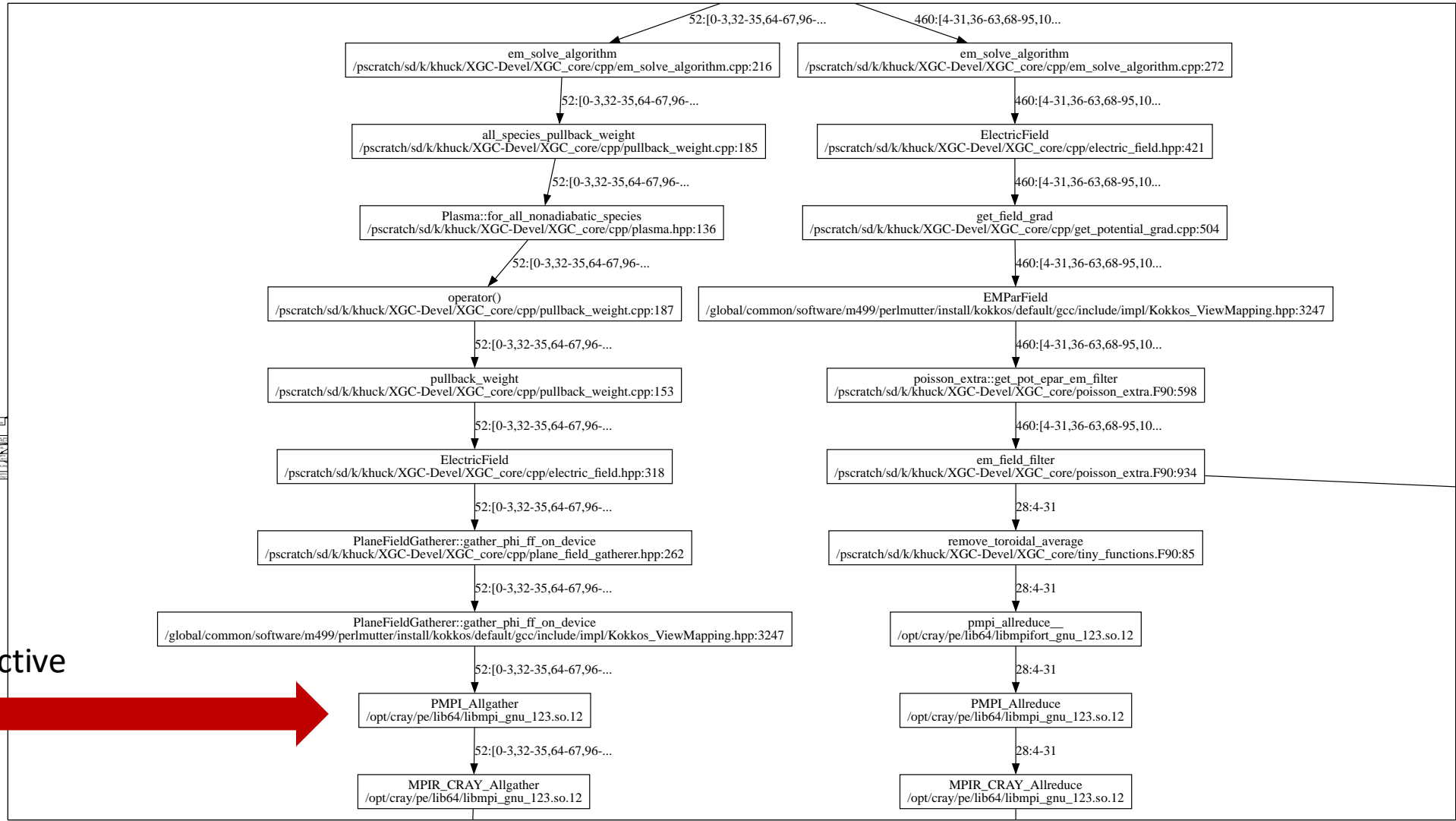
Successful Use Cases

- Octo-Tiger (HPX) thread placement on Fugaku
 - Detected unusual core indexing, needs to be mapped using hwloc information
- Argobots example
 - Detected that the user was running client & server on same node, with overlapping affinity lists
- Grid (<https://github.com/paboyle/Grid>) on Frontier (lattice QCD library)
 - Helped find/understand ideal resource mapping with custom bind arguments
- XGC on Frontier
 - Empirically demonstrated that 2 threads per core provided better performance than 1 thread per core
 - Helped detect deadlock in GPU-aware libfabric/MPI implementation
 - Helped detect algorithmic deadlock in MPI (application bug)

Deadlock detection – at scale

- Runs application in gdb/cuda-gdb/rocgdb/gdb-oneapi in “batch” mode – sequence of commands are scripted
- ZeroSum monitors all threads during execution:
 - If all sleeping for X seconds (X is user-configurable):
 - pthread_kill(SIGQUIT) the main thread
 - If all but 1 sleeping for X seconds, *and* that thread doesn’t have at least one minor page fault during that time, it’s *probably* deadlocked at an MPI collective:
 - pthread_kill(SIGQUIT) the main thread
- Gdb will intercept the signal and is scripted to get a backtrace from all threads
- Post-process the threads and make a tree (just like STAT/Cray-STAT)
 - *However*, ZeroSum is automated - queue times on Perlmutter were around 8 hours, user didn’t want to babysit the 128 node slurm request

Example: XGC deadlocked on Perlmutter with 512 ranks



Only 52 of 64 in communicator participated in collective



Conclusions, Future Work

- ZeroSum addresses the *configuration optimization* problem
 - <https://github.com/UO-OACISS/zerosum>
- Still need automated misconfiguration/contention detection
- Output data could be better – use ADIOS2 BP5? HDF5? SQLite?
 - Depends on analysis needs/wants
 - Current log file “format” means analysis process is manual/ad hoc (via Python)
- Streaming data to Mochi (SOMA) or other service (LDMS)?
 - Enables robust monitoring approach – but likely redundant in many cases
- Integration with performance tools (TAU, APEX, etc)
 - Analysis of application performance data in context of system monitoring data
- Input for automated feedback/control (APEX, Argobots, SOMA, application)

Acknowledgements

Parts of this research was supported by the Exascale Computing Project (17-SC-20-SC), a joint project of the U.S. Department of Energy's Office of Science and National Nuclear Security Administration, responsible for delivering a capable exascale ecosystem, including software, applications, and hardware technology, to support the nation's exascale computing imperative.

This research used resources of the Oak Ridge Leadership Computing Facility at the Oak Ridge National Laboratory, which is supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC05-00OR22725.



U.S. DEPARTMENT OF
ENERGY

Office of
Science



EXASCALE COMPUTING PROJECT



Thanks! Questions?

Extra slides

Users need help...

NERSC slack:

can users get fine-grained statistics of their usage? Let's say I run a GPU job and I'd like to have a look retrospectively on things like memory usage, how much the GPUs are consuming in Watts or whatever, and other statistics like that. In a previous institution I was associated with, there was like a weird dashboard where one could get all sorts of interesting info, and it was pretty helpful to track stuff

I was hoping for something more like nvidia-smi...

Hi everybody, I am fairly new to the perlmutter system. I am wondering if there is some tool to **monitor** node performances like the jupyter online one.

My ideal use case is: launch a job, ssh to node running -> **monitor** performances [on cpu systems I would generally use htop].

This is not thought for continuous use, just to make sure that the settings are exploiting all node resources that they can. Thank you very much!

ALCF slack:

Are there any tools on Aurora to monitor GPU usage similar to **nvidia-smi** or **rocm-smi** ?

Quick quiz:

- Where does MPI rank 0 of the following get allocated/scheduled on 1 Frontier node?

`salloc -N 1 --threads-per-core=1 ...`

```
Process Summary:  
MPI 000 - PID 64341 - Node frontier00765 - CPUs allowed: [1,2,3,4,5,6,7]
```

`srun -n1 -c7 --gpus-per-task=1 --gpu-bind=closest`

```
Process Summary:  
MPI 000 - PID 64384 - Node frontier00765 - CPUs allowed: [49,50,51,52,53,54,55]  
MPI 000 - Node frontier00765 - RT_GPU_ID 0 - GPU_ID 0 - Bus_ID 0000:c1:00.0
```

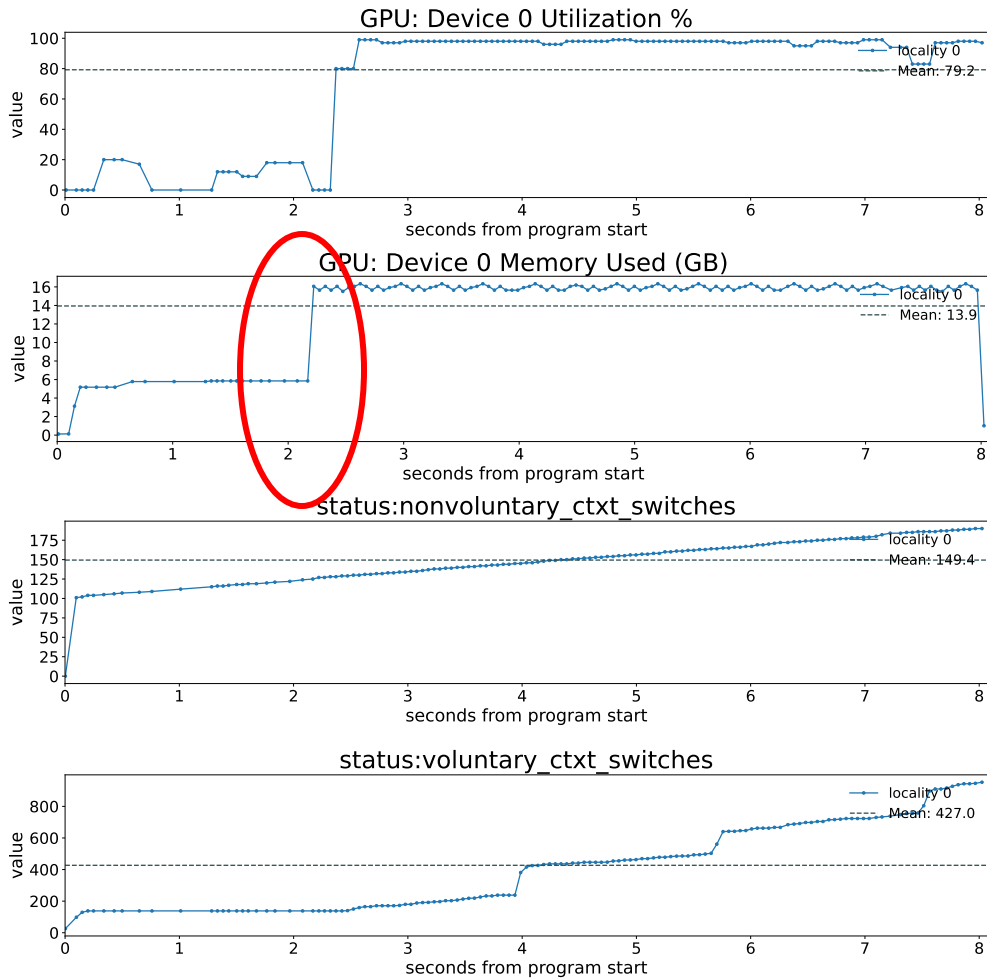
`srun -n8 -c7 --gpus-per-task=1 --gpu-bind=closest`

```
Process Summary:  
MPI 000 - PID 64516 - Node frontier00765 - CPUs allowed: [1,2,3,4,5,6,7]  
MPI 000 - Node frontier00765 - RT_GPU_ID 0 - GPU_ID 4 - Bus_ID 0000:d1:00.0
```

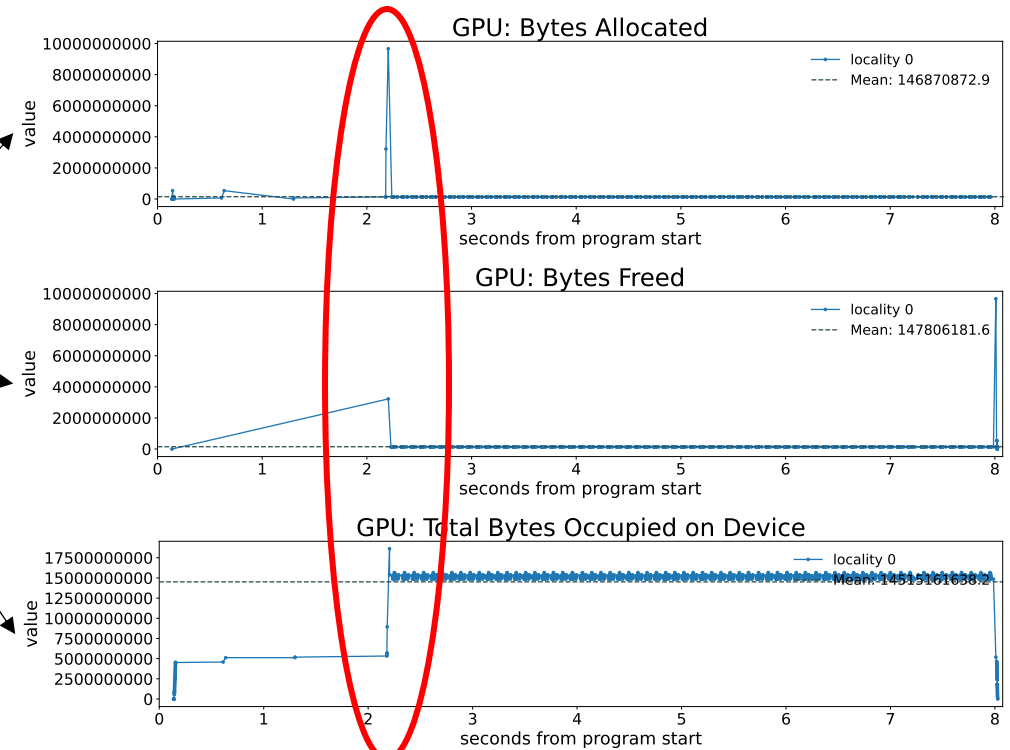


Caveat: Monitoring != Explicit Measurement

- Monitoring Data (periodic)



- Progress Data (events)



Example: APEX data from Lulesh Kokkos with CUDA back end, visualized with Python – clear need for GPU HWM

What is the hardware / operating system doing?

- Kernel monitoring is popular and useful (strace, ptrace, dtrace, dtruss, ftrace, KUtrace, kprobe, system-tap, KTAU, STaKTAU, eBPF, bpftrace, BCC, ...)
- Nataraj, Morris, Malony, Sottile, and Beckman. "The Ghost in the Machine: Observing the Effects of Kernel Operation on Parallel Application Performance." In *SC2007*, pp. 1-12. 2007.
- Subsystem monitoring (Darshan)
- Did I ask for the right thing?
- Will I run out of a limited resource?
- Are there alternatives?
- Can I get this information *without being root? Without spawning another process to monitor my process? Without kernel patches/support? Low/no overhead?*

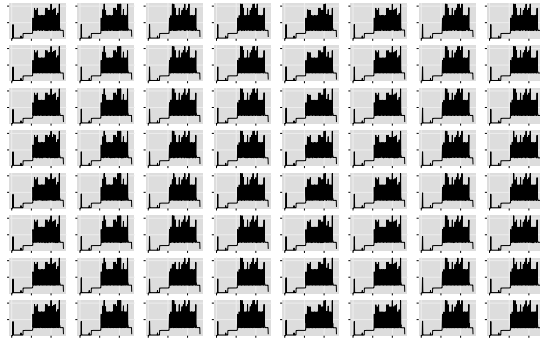
What (we think) users DON'T want:

(BTW, I am happy to ask for some help visualizing/simplifying the data...sparklines aren't going to cut it)

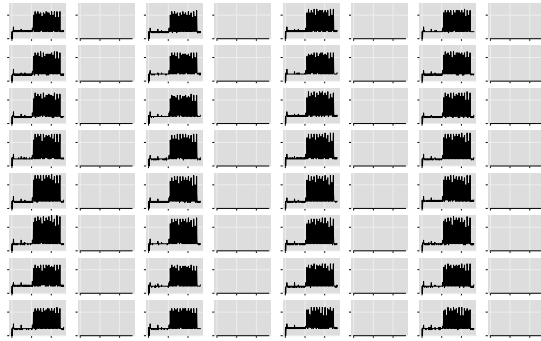


Sparklines: 64 MPI Ranks on Frontier – need better

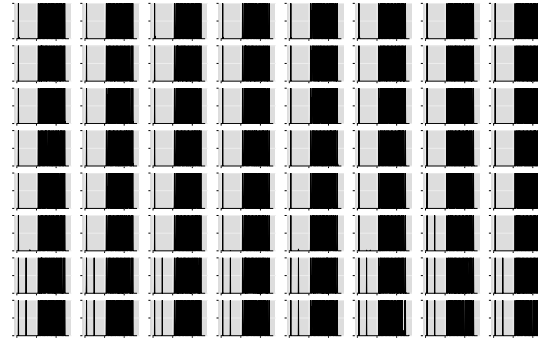
Used VRAM Bytes, range: [295,063,552, 21,668,823,040]



Energy Average (J), range: [0, 30]



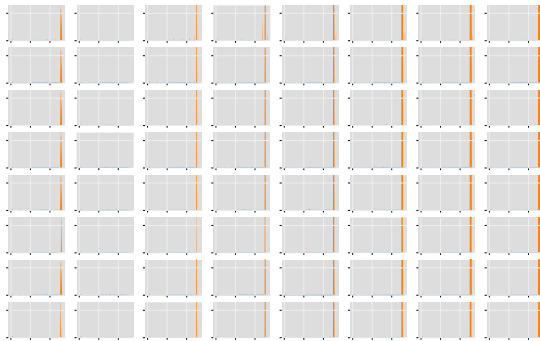
Device Busy %, range: [0, 100]



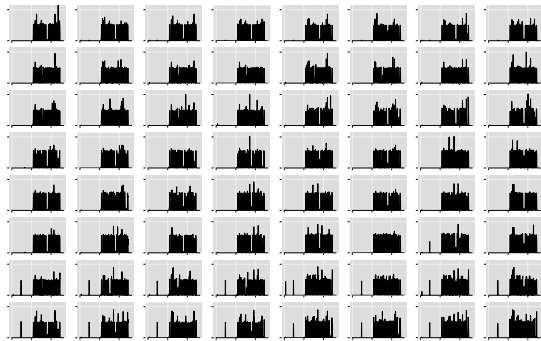
Voltage (mV), range: [712.0, 943.0]



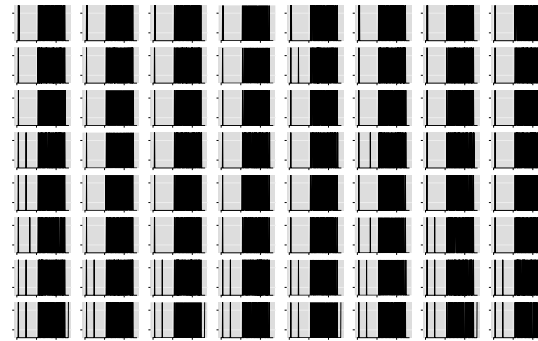
Context Switches, range: [0, 6,452]



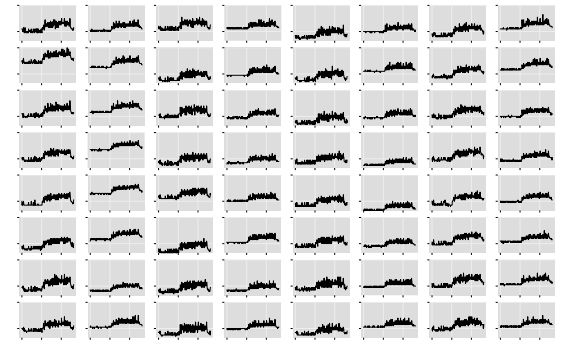
Memory Busy %, range: [0, 57]



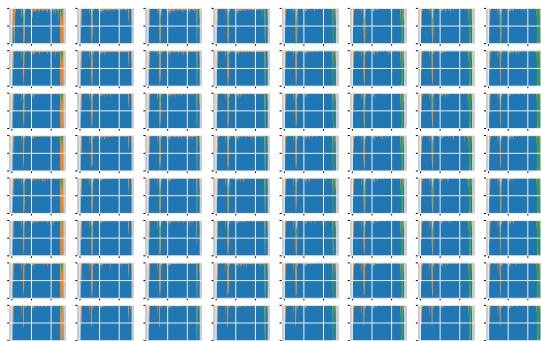
Clock Frequency, GLX (MHz), range: [800, 1,700]



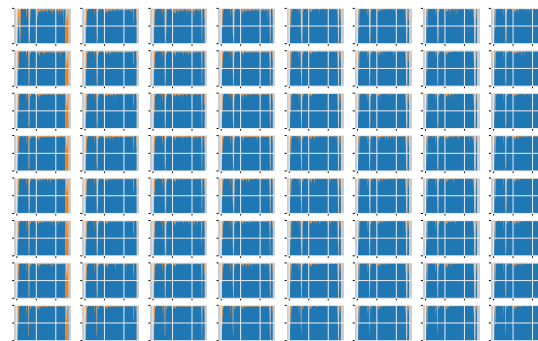
Temperature (C), range: [33.0, 60.0]



HWT utilization

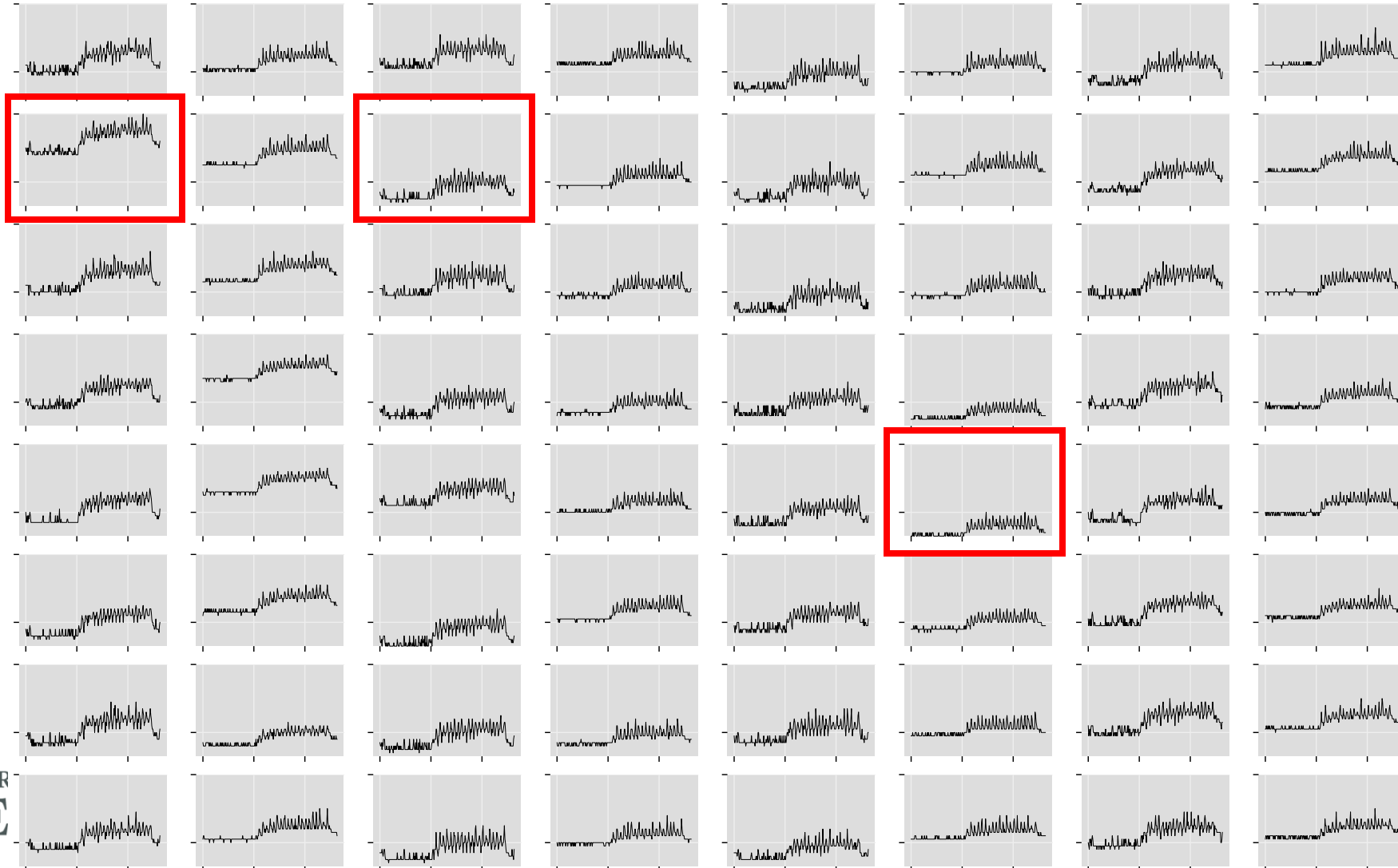


LWP utilization



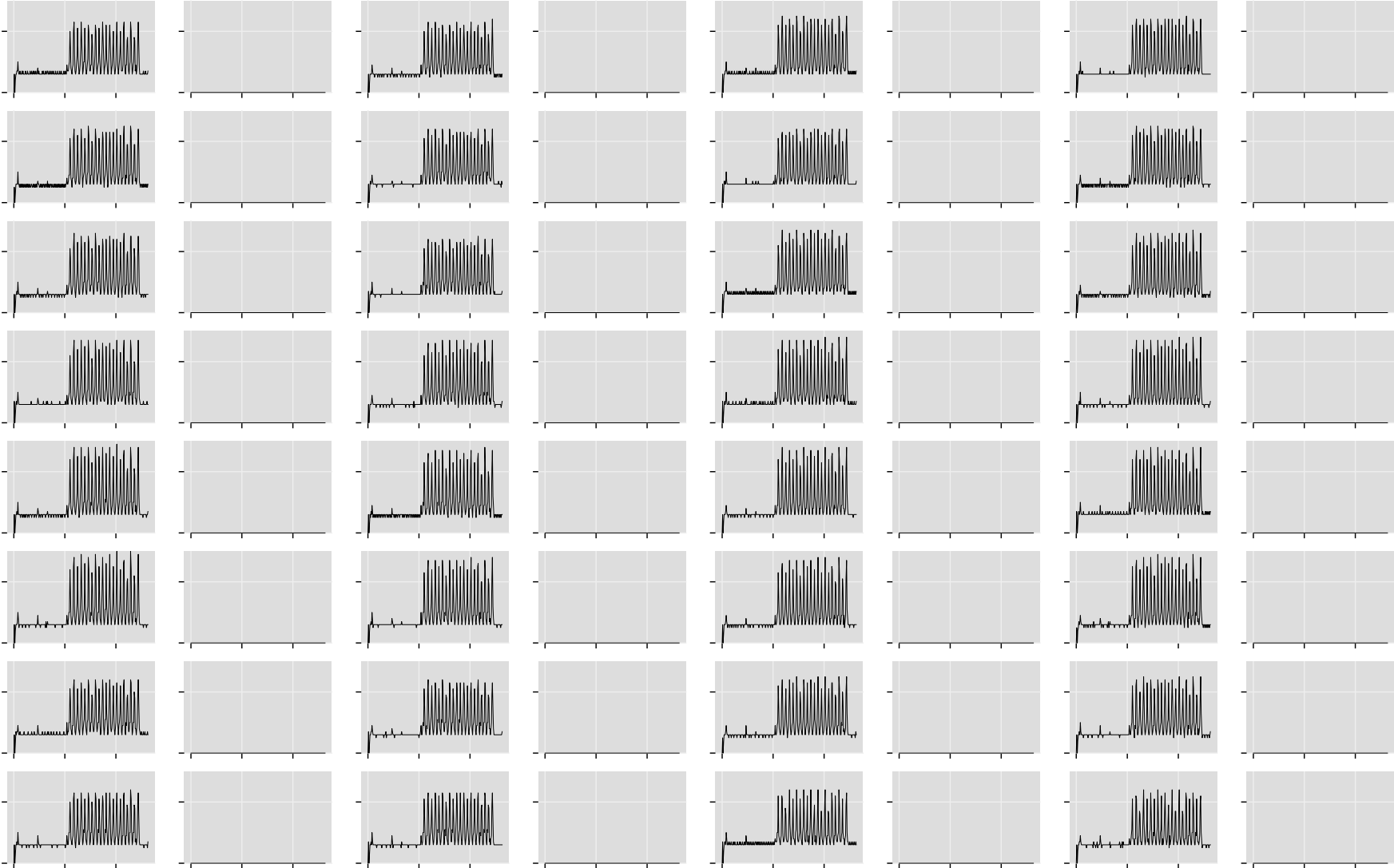
Sparklines: 64 MPI Ranks on Frontier

Temperature (C), range: [33.0, 60.0]



Sparklines: 64 MPI Ranks on Frontier

Energy Average (J), range: [0, 30]



Sparklines: 64 MPI Ranks on Frontier

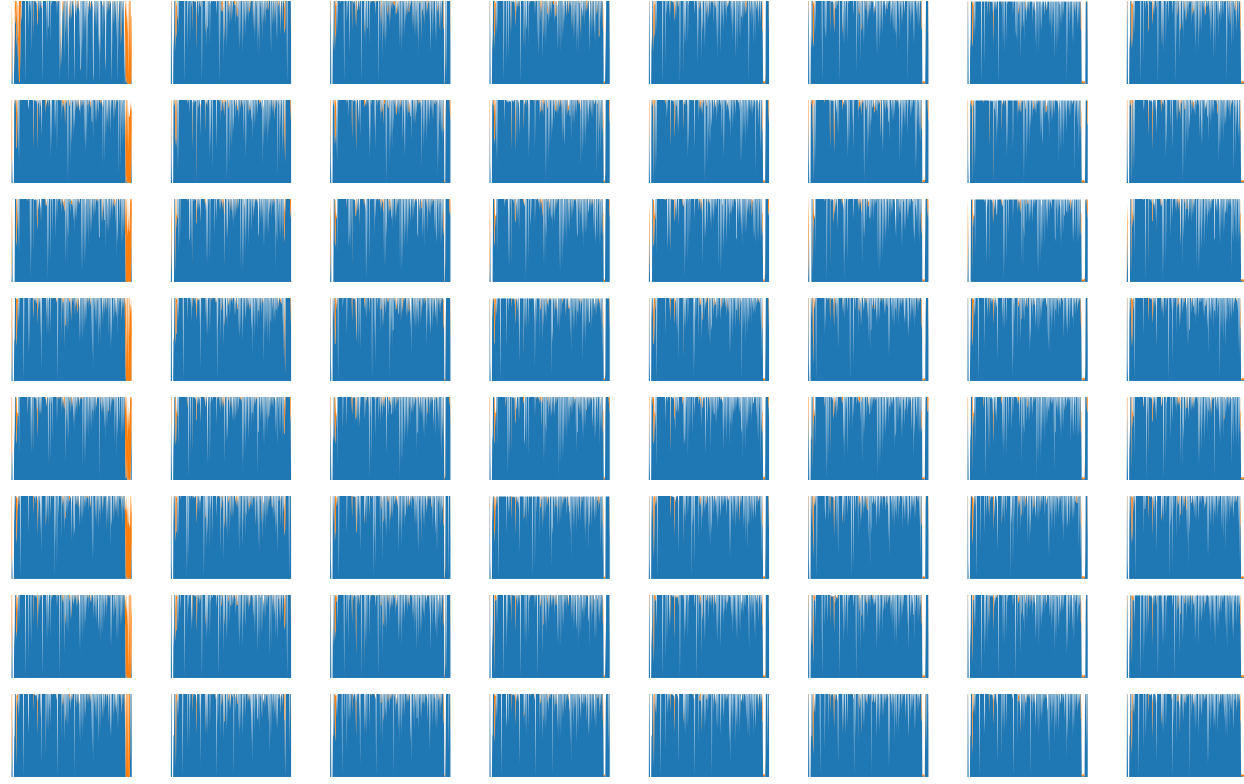
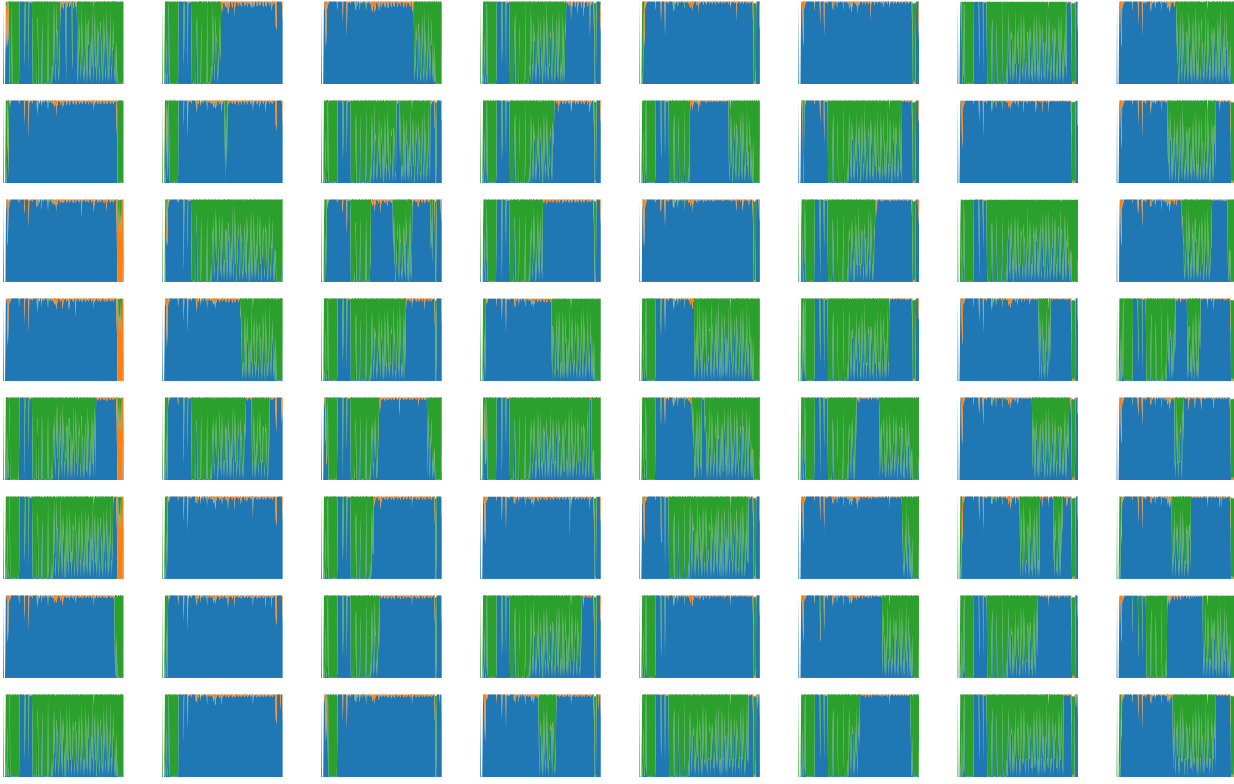
Used VRAM Bytes, range: [295,063,552, 21,668,823,040]



Sparklines: 64 MPI Ranks on Frontier – OMP_BIND=cores

HWT utilization

LWP utilization



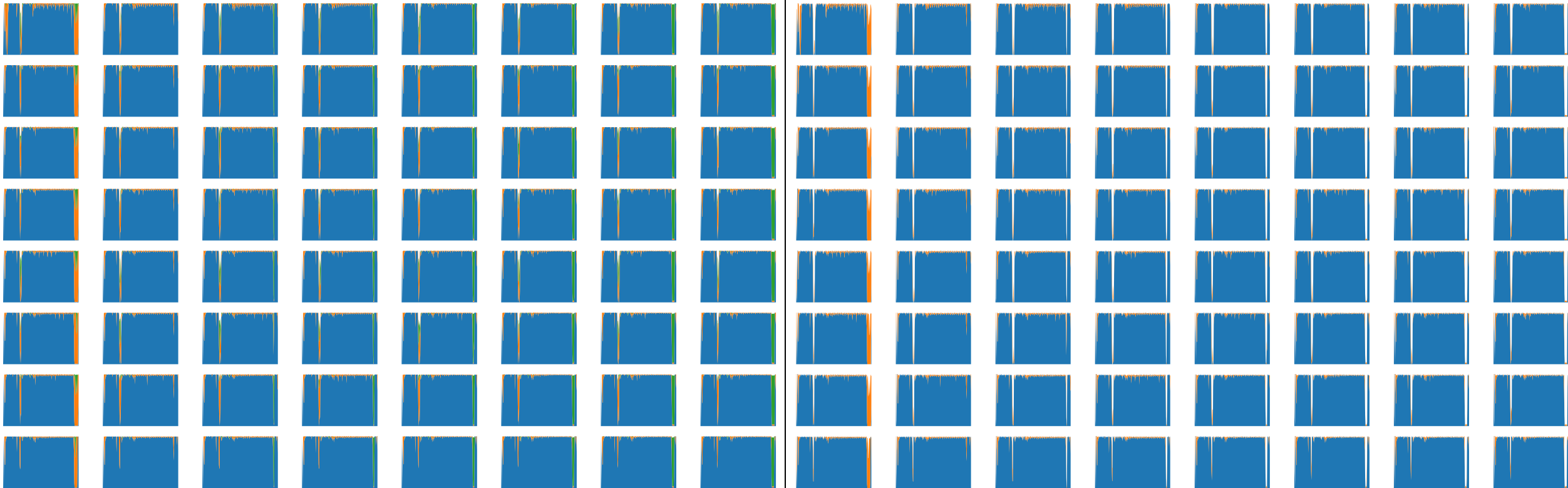
■ User
■ System
■ Idle

■ User
■ System

Sparklines: 64 MPI Ranks on Frontier – OMP_BIND=threads

HWT utilization

LWP utilization



■ User
■ System
■ Idle

■ User
■ System