

HPCToolkit: Experiences at Exascale



Jonathon Anderson

Rice University

STW 2023

June 19, 2023



U.S. DEPARTMENT OF
ENERGY

Office of
Science

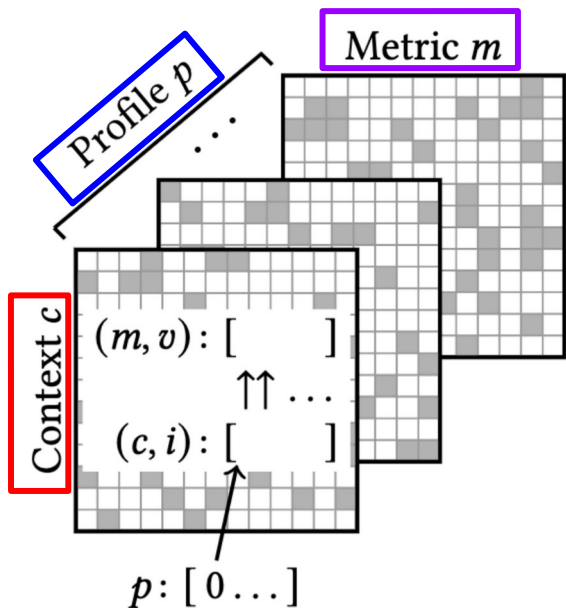
The Team

- PI: John Mellor-Crummey
- Staff
 - Mark Krentel, Laksono Adhianto, Marty Itzkowitz, Wil Phan, Matt Barnett
- Students
 - Jonathon Anderson, Yumeng Liu, Dejan Grubisic, Dragana Grbic

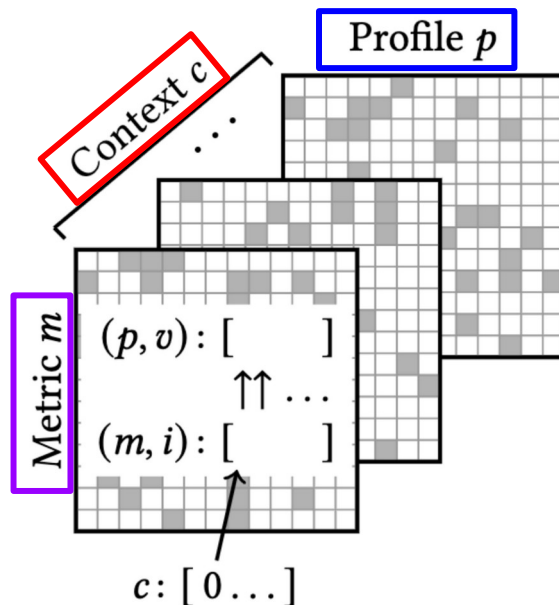
Recap: Preparing HPCToolkit for Exascale

- No other tool provides **detailed** performance analysis **at scale**!
- The problem: performance data at scale is **huge** and **slow** to process!
- Two-pronged solution, next slides

Recap: Efficient Sparse Formats for Performance Data

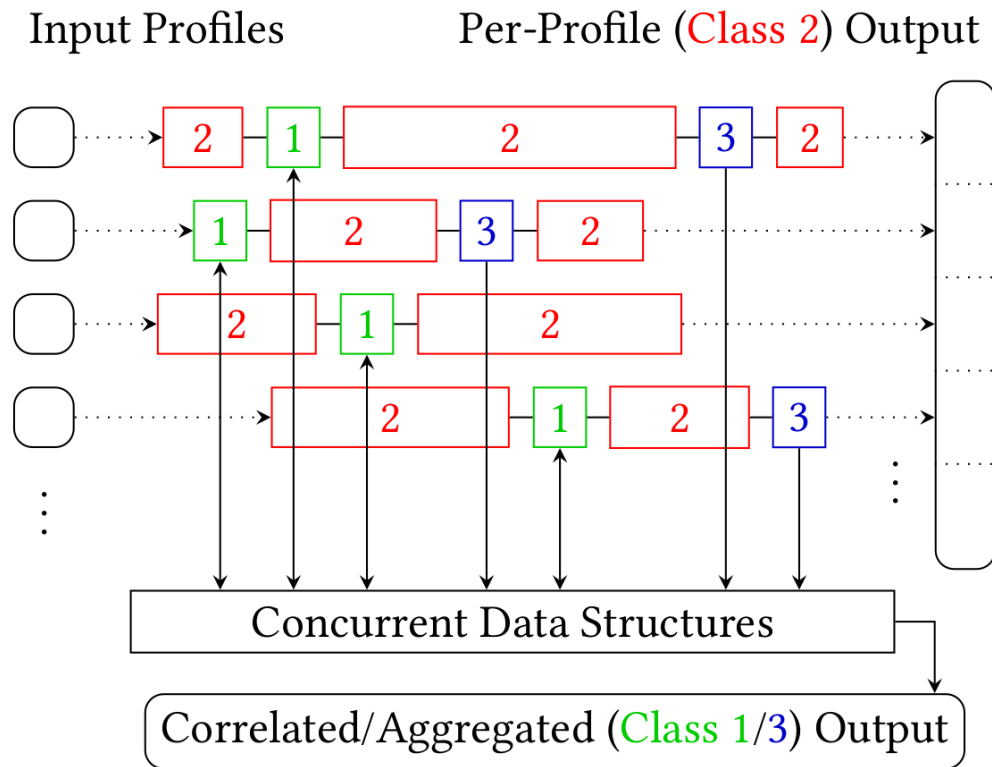


Profile-Major-Sparse (PMS):
data for each profile is contiguous



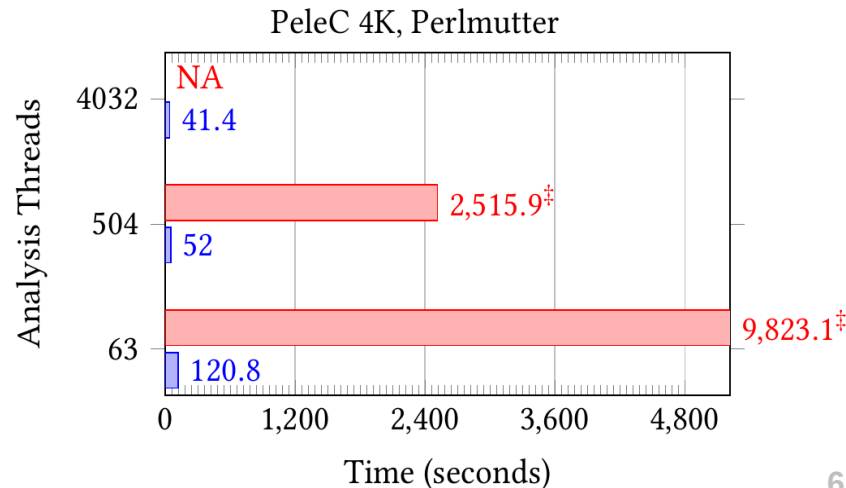
Context-Major-Sparse (CMS):
data for each context is contiguous

Recap: Highly-parallel Multithreading for Performance Analysis



Recap: Preparing HPCToolkit for Exascale

- End result: **Stellar improvements for GPU-accelerated codes!**
 - Running PeleC on 512 nodes of Perlmutter, ~17.1 PFLOPs
 - Size reduction from sparsity: 14.3 TB → **11.4 GB (1254x)**
 - Min. node reduction from threading: 16 nodes → **1 node**
 - Overall time reduction: 2.7 hr → **2 min (81.3x)**
- Also significant improvements for CPU-only cases
 - For details, see our paper in ICS'22



Last year, the only question asked was

Have you used this at Exascale?

Today, the answer to that question is

Yes, we have.

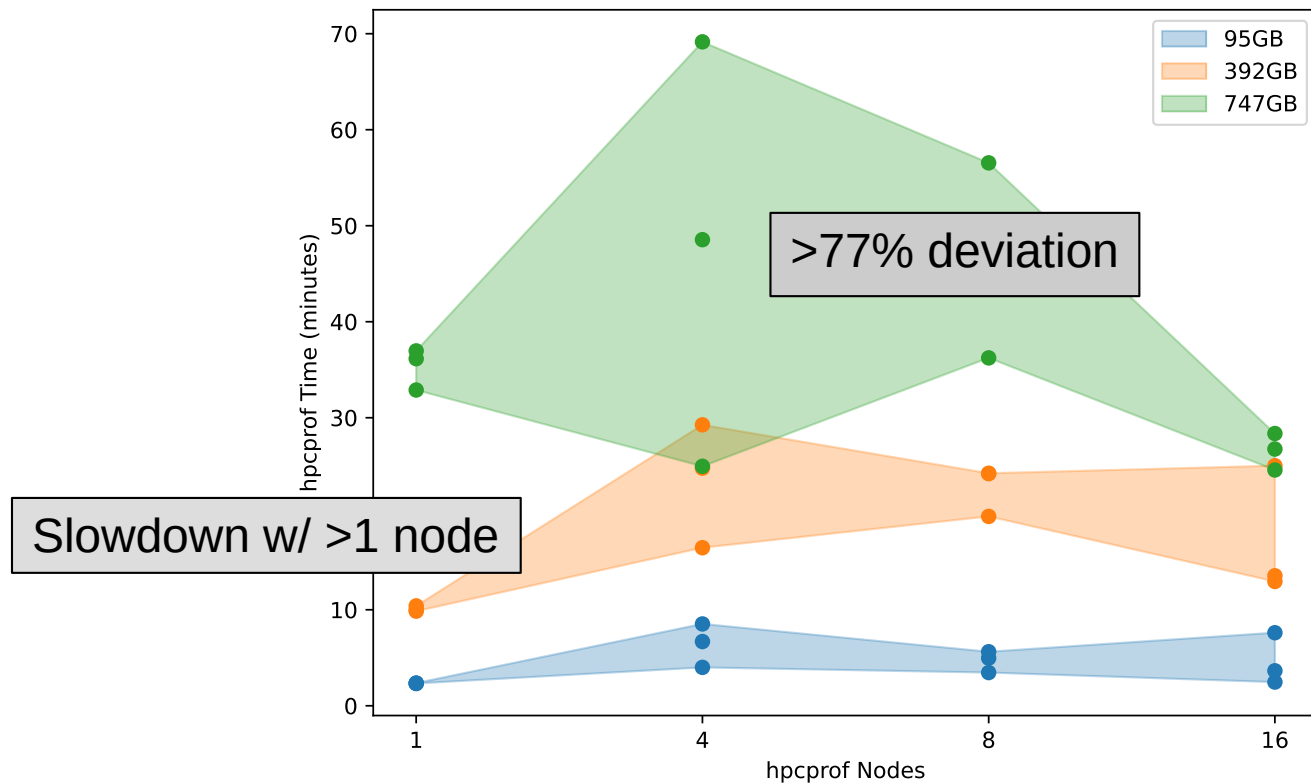
HPCToolkit Improvements for Frontier I

- HPCToolkit **built and (mostly) ran with no change**, but we made improvements for Frontier
- **hpcrun**
 - Bug fix: corrected support for Cray OpenMP
 - Improved how CPU threads are mapped to MPI ranks
 - now using job launcher variables: independent of any MPI version
 - Added support to boost resolution of CPU traces
 - collect CPU callstack when offloading GPU operations
- **hpcstruct**
 - Use OpenMP to inspect profiles to identify CPU and GPU binaries involved in the execution
 - Reduces preparation for binary analysis of an hpctoolkit database to seconds

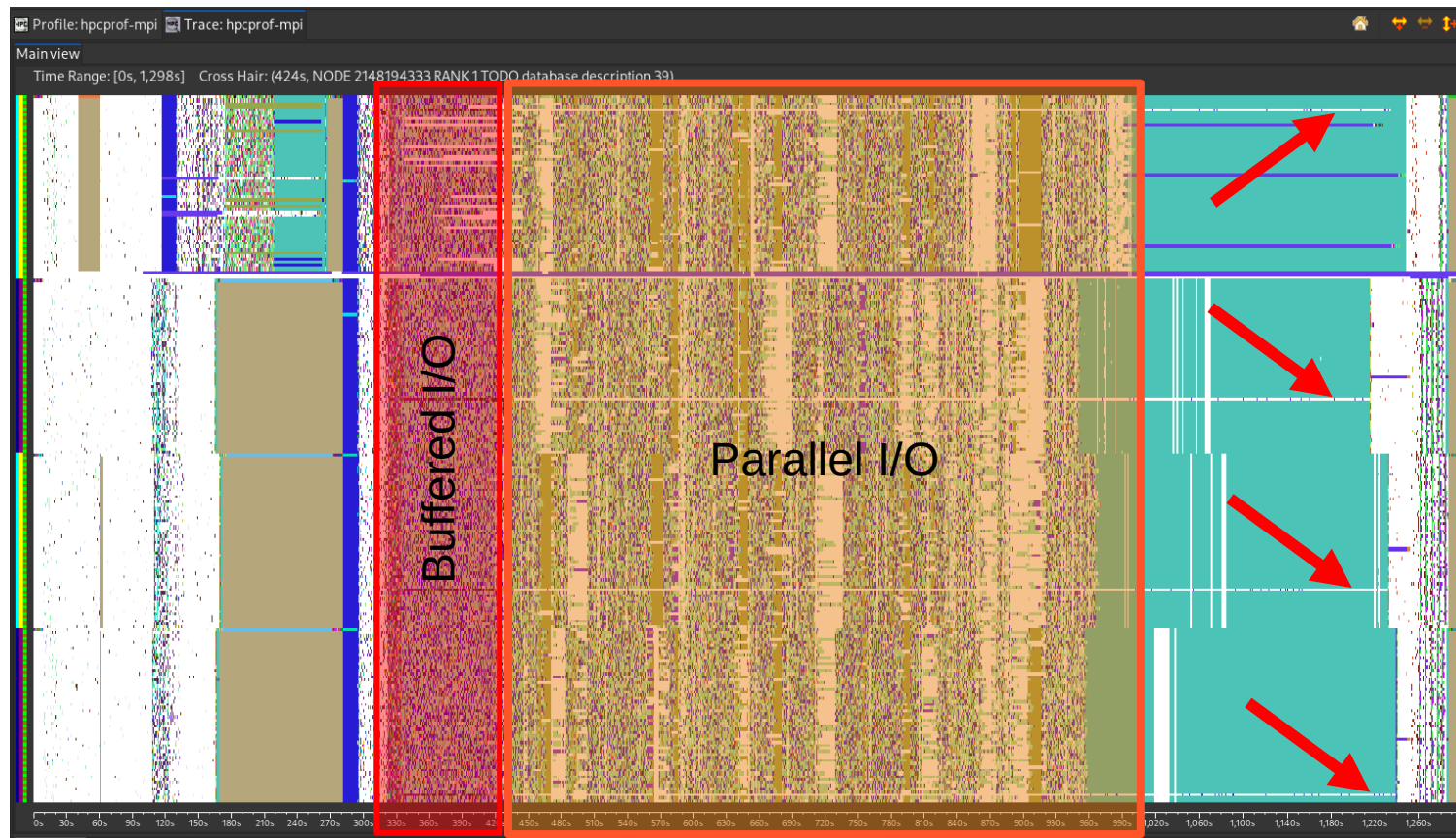
HPCToolkit Improvements for Frontier II

- **hpcprof**
 - Bug fix: handle truncated reads on network file systems
- **hpcviewer**
 - Reduced memory footprint when viewing large profiles and traces
 - Improved performance and usability of graphing metrics for many execution contexts
 - Corrected handling of corner cases at trace start and end

Remaining Issue: Poor Multi-node Performance of hpcprof



Remaining Issue: Serialized I/O



First Taste of Exascale –

First Taste of Exascale – False Start

- Large-scale jobs on Frontier **regularly fail**, and **fail immediately**
 - MPI bootstrap barrier times out after 2 minutes and fails the job
 - Unrelated to node crashes (which happen less immediately)
- OLCF Help: “Orion can’t serve libraries + executable for many nodes under 2 minutes”
- Recommended solution: copy all libraries to node-local NVMe first
 - User Guide “example” batch script is 60 lines long
 - Uses LD_LIBRARY_PATH, conflicts with Spack’s RPATH
- Spark discussion for alternative solutions. Spindle?
- Data on following slides collected in April, failure occurs **much more frequently** since

First Taste of Exascale – Leviathan

- Strong-scaling LAMMPS, **large** problem, **short** simulation
 - Up to 8192 nodes, ~1.6 EFLOPs (theoretical peak)
 - Lennard-Jones, up to 1.77 trillion atoms (216M atoms/node)
 - Runs for 7900 timesteps, 18.4 minutes (under measurement)
- Analysis
 - Final size: **0.793 TB**
 - Analysis time, single node: **36 minutes**



First Taste of Exascale – Smaug

- Strong-scaling LAMMPS, **small** problem, **long** simulation
 - Up to 8192 nodes, ~1.6 EFLOPs (theoretical peak)
 - Lennard-Jones, up to 34 billion atoms (4.1M atoms/node)
 - Runs for 110 thousand timesteps, 19.7 minutes (under measurement)
 - Boosted CPU trace resolution, approximates Leviathan for 4.3 hours
- Analysis
 - Final size: **5.68 TB**
 - Analysis time, single node: **1 hour, 40 minutes**



- Smaug configuration, 4096 nodes
 - Later run, after hpcrun improvements
 - Total size: 2.85 TB
- X11 over SSH tunnel from Frontier
 - `ssh -Y frontier.olcf.ornl.gov`
- Rendering on a virtual desktop at Rice
- Hpcviewer 2023.05



hpcviewer (on login03)

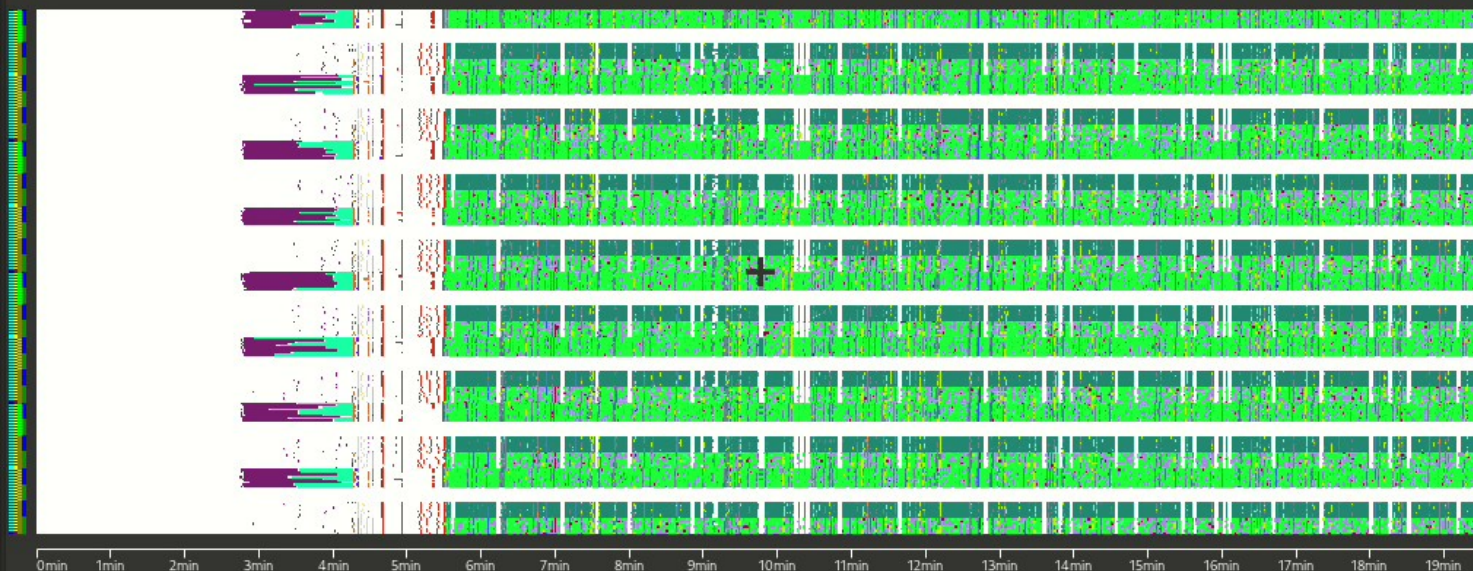
File View Filter Help

hpc Profile: Imp hpc Trace: Imp



Main view

Time Range: [0min, 19.52min] Cross Hair: (9.75min, NODE 2148202570 RANK 16292 GPUCONTEXT 0 GPUSTREAM 1)



Depth: 6 - + 🔽

Call stack Statistics "1

- <program root>
- main
- LAMMPS_NS::Input::file() [liblamr
- LAMMPS_NS::Input::execute_com
- LAMMPS_NS::Run::command(int, i
- LAMMPS_NS::VerletKokkos::run(ir
- LAMMPS_NS::CommKokkos::forw
- MPIABI_Send [libmpiwrapper.so]

Depth view Summary view



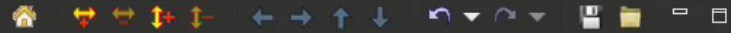
Mini map



hpcviewer (on login03)

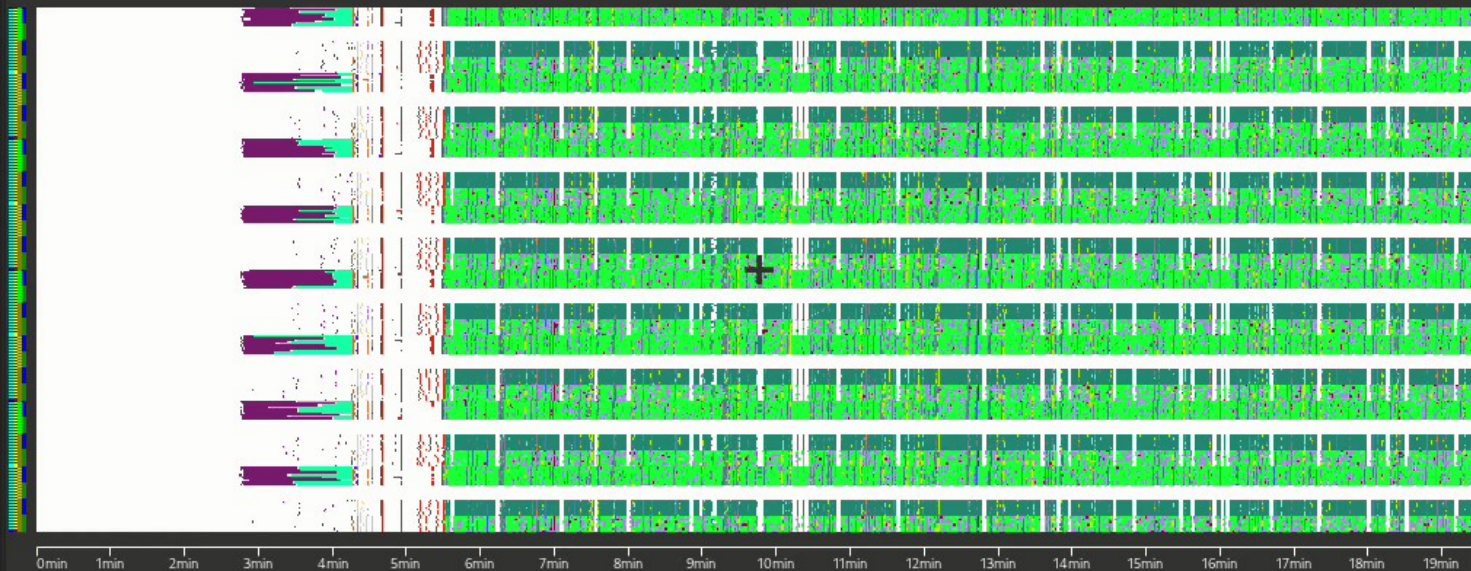
File View Filter Help

Profile: Imp Trace: Imp



Main view

Time Range: [0min, 19.52min] Cross Hair: (9.75min, NODE 2148202570 RANK 16292 GPUCONTEXT 0 GPUSTREAM 1)



Depth: 6

Call stack Statistics "1

- <program root>
- main
- LAMMPS_NS::Input::file() [liblamr
- LAMMPS_NS::Input::execute_com
- LAMMPS_NS::Run::command(int, i
- LAMMPS_NS::VerletKokkos::run(ir
- LAMMPS_NS::CommKokkos::forw
- MPIABI_Send [libmpiwrapper.so]

Depth view Summary view



Mini map



hpcviewer (on login03)

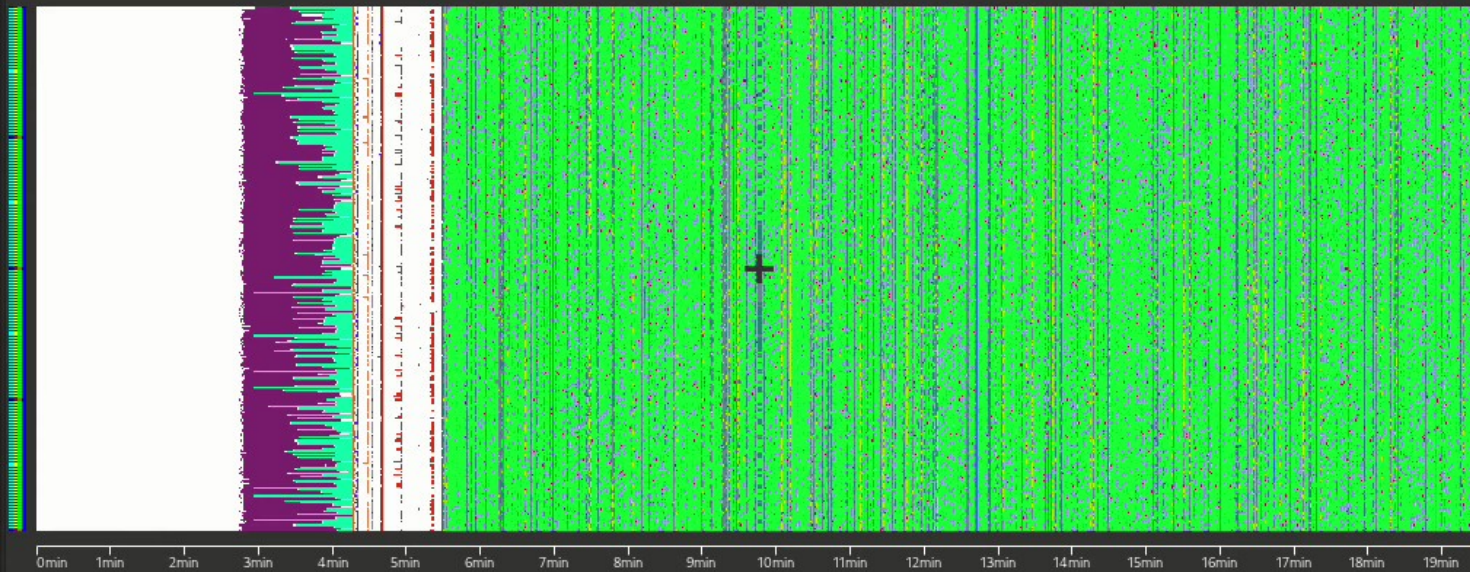
File View Filter Help

Profile: Imp Trace: Imp



Main view

Time Range: [0min, 19.52min] Cross Hair: (9.75min, NODE 2148202570 CORE 33 RANK 16292 THREAD 0)



Depth: 6

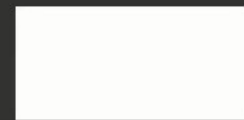
Call stack Statistics

- <program root>
- main
- LAMMPS_NS::Input::file() [liblamr
- LAMMPS_NS::Input::execute_com
- LAMMPS_NS::Run::command(int, i
- LAMMPS_NS::VerletKokkos::run(ir
- LAMMPS_NS::CommKokkos::forw
- MPIABI_Send [libmpiwrapper.so]

Depth view Summary view



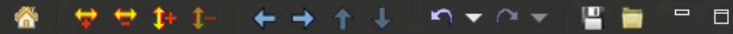
Mini map



hpcviewer (on login03)

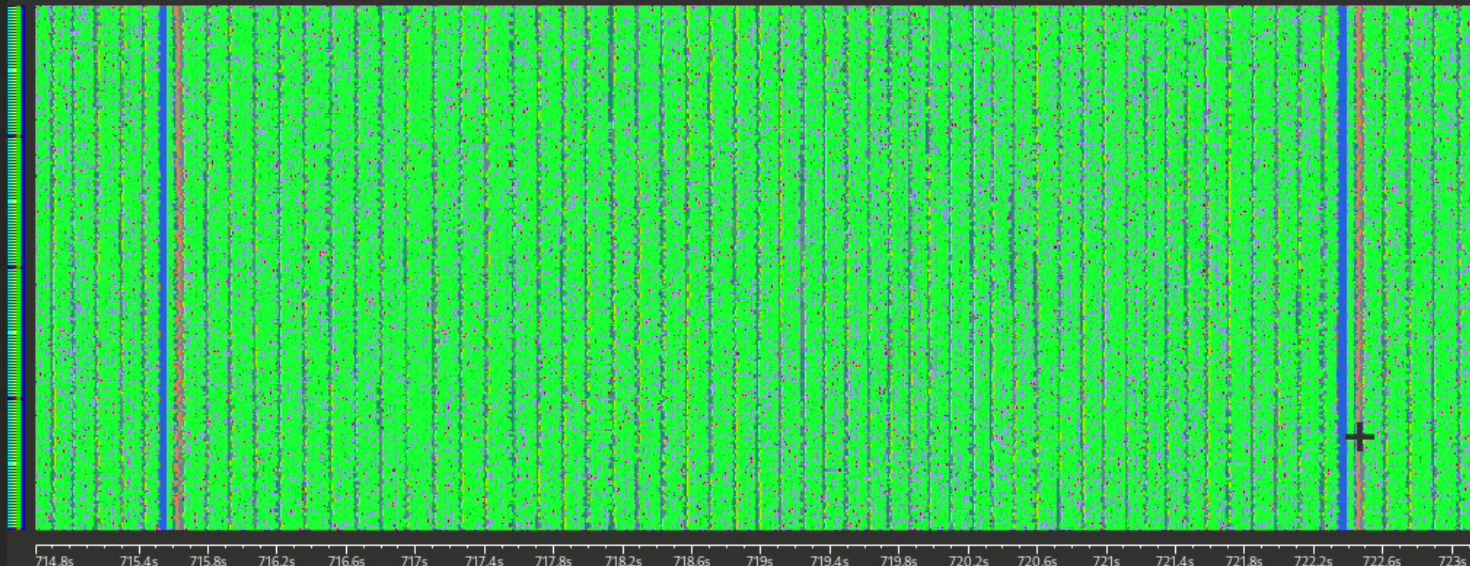
File View Filter Help

Profile: Imp Trace: Imp



Main view

Time Range: [714.8s, 723.16s] Cross Hair: (722.47s, NODE 2148146297 CORE 25 RANK 26851 THREAD 0)



Depth: 6

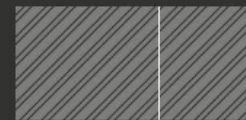
Call stack Statistics

- <program root>
- main
- LAMMPS_NS::Input::file() [liblamr
- LAMMPS_NS::Input::execute_com
- LAMMPS_NS::Run::command(int, i
- LAMMPS_NS::VerletKokkos::run(ir
- LAMMPS_NS::AtomKokkos::sort()
- LAMMPS_NS::AtomVec::copy(int, i

Depth view Summary view



Mini map



hpcviewer (on login03)

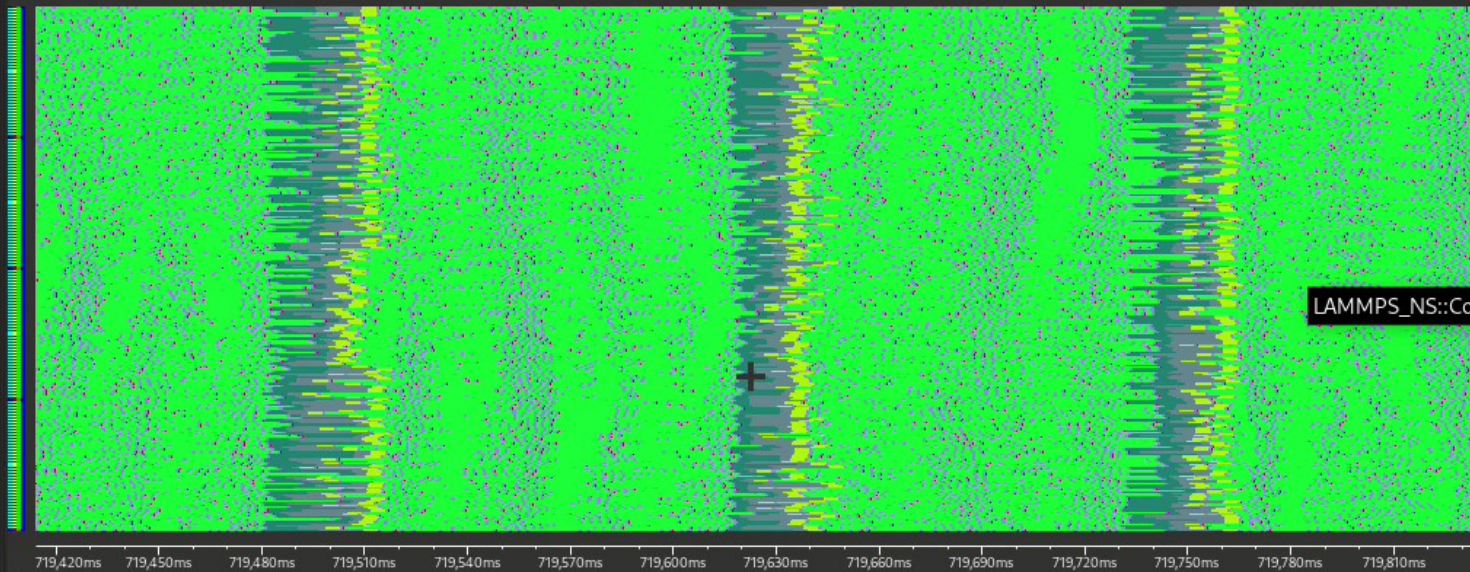
File View Filter Help

Profile: Imp Trace: Imp



Main view

Time Range: [719,414ms, 719,836ms] Cross Hair: (719,623ms, NODE 2148145257 CORE 33 RANK 23028 THREAD 0)



Depth: 6 - +

Call stack Statistics "1

- <program root>
- main
- LAMMPS_NS::Input::file() [liblamr
- LAMMPS_NS::Input::execute_com
- LAMMPS_NS::Run::command(int, i
- LAMMPS_NS::VerletKokkos::run(ir
- LAMMPS_NS::CommKokkos::exchange() [liblammps.so.0]
- LAMMPS_NS::CommKokkos::exch
- Kokkos::DualView<int*, Kokkos::L
- Kokkos::deep_copy<int*, Kokkos::l
- Kokkos::Impl::DeepCopyHIP(void*,
- hipMemcpyAsync [libamdhip64.sc

Depth view Summary view



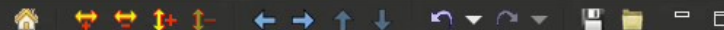
Mini map



hpcviewer (on login03)

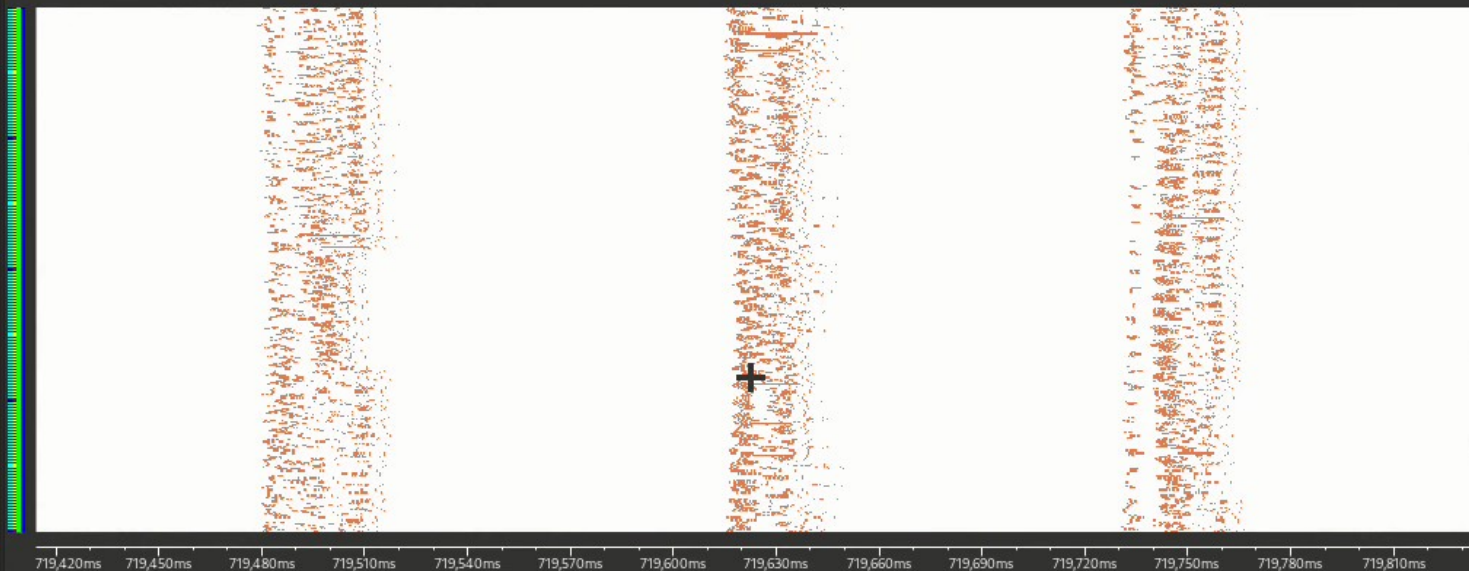
File View Filter Help

hpc Profile: Imp hpc Trace: Imp



Main view

Time Range: [719,414ms, 719,836ms] Cross Hair: (719,623ms, NODE 2148145257 RANK 23028 GPUCONTEXT 0 GPUSTREAM 0)



Depth view Summary view

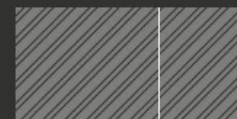


Depth: 79 - +

Call stack Statistics "1

- <program root>
- main
- LAMMPS_NS::Input::file() [liblan
- LAMMPS_NS::Input::execute_co
- LAMMPS_NS::Run::command(in
- LAMMPS_NS::VerletKokkos::rur
- LAMMPS_NS::CommKokkos::ex
- LAMMPS_NS::CommKokkos::ex
- Kokkos::DualView<int*, Kokkos:
- Kokkos::deep_copy<int*, Kokkos:
- Kokkos::Impl::DeepCopyHIP(voi
- hipMemcpyAsync [libamdhip64
- <gpu copy>

Mini map



hpcviewer (on login03)

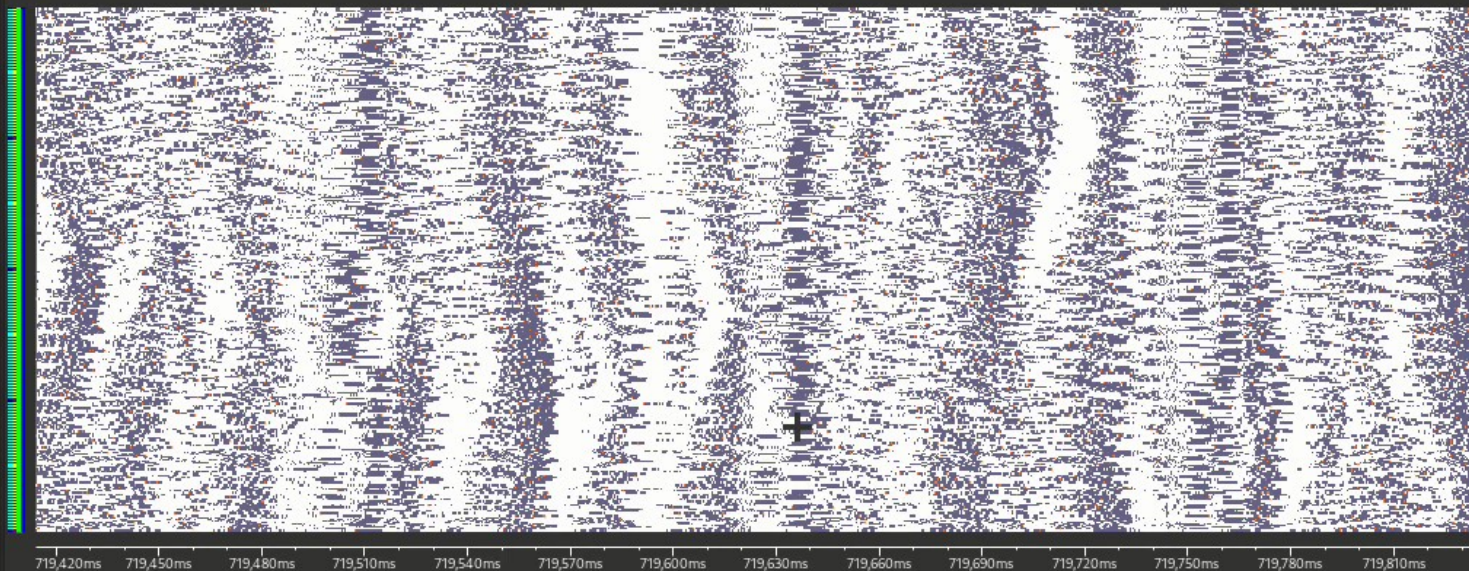
File View Filter Help

HPC Profile: Imp HPC Trace: Imp



Main view

Time Range: [719,414ms, 719,836ms] Cross Hair: (719,637ms, NODE 2148180086 RANK 26123 GPUCONTEXT 0 GPUSTREAM 1)



Depth view Summary view



Depth: 79 - +

Call stack Statistics "1"

Procedure	%
sched_yield, <no activit	63.9%
<gpu kernel>	34.9%
<gpu copy>	1.3%
<gpu memset>	0.0%

Mini map



Work in Progress

- Eliminating the I/O bottleneck in hpcprof
- Building a new remote hpcserver, companion to hpcviewer
 - Will load up the entire database in the memory of several compute nodes
 - Serves from memory, directly to an hpcviewer client on your laptop
- Collaborating with AMD on a new ROCm tools interface
 - Awaiting support for PC samples for non-root users
- Developing a Python API and library for performance data analysis
 - Supports automated and exploratory analysis of large-scale performance data
 - To be used for regression testing and validation
- Testing and tuning HPCToolkit on Sunspot (TDS for Aurora)