

Kokkos Tools:

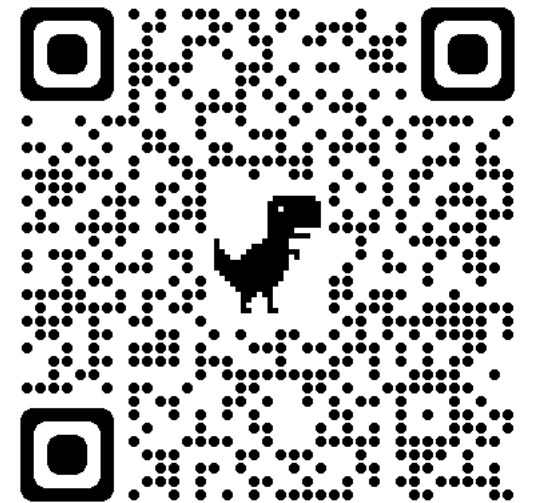
Kokkos support in the TAU and APEX portable performance measurement tools

Kevin A. Huck

Oregon Advanced Computing Institute for Science and Society (OACISS)



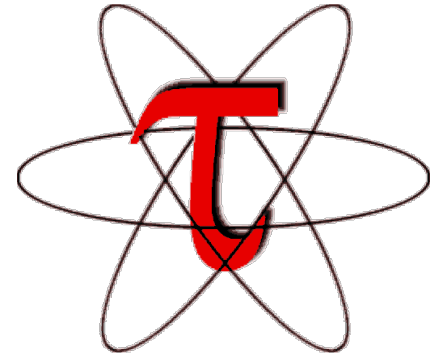
<http://www.nic.uoregon.edu/~khuck/kokkos/KUG2023>



TAU Performance System

TAU Performance System

- **T**uning and **A**nalysis **U**tilities (29+ year project)
- Integrated performance toolkit:
 - Multi-level performance instrumentation
 - Highly configurable
 - Widely ported performance profiling / tracing system
 - Portable (java, python) visualization / exploration / analysis tools
- Supports all major HPC programming models
- MPI/SHMEM, OpenMP, OpenACC, CUDA, HIP, SYCL/OneAPI, **Kokkos...**
- Support for ML/AI frameworks: TensorFlow, pyTorch, Horovod
- Integrated with PAPI, LIKWID for hardware counter support
- <https://tau.uoregon.edu> or <https://github.com/UO-OACISS/tau2> (public mirror)



Performance Measurement

■ Timers

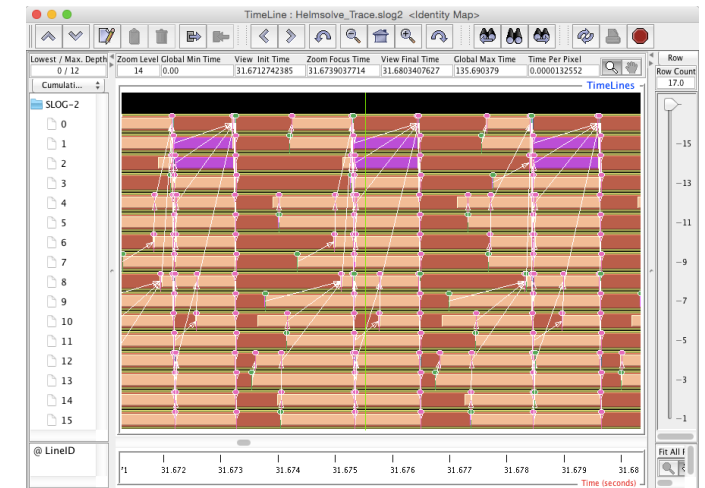
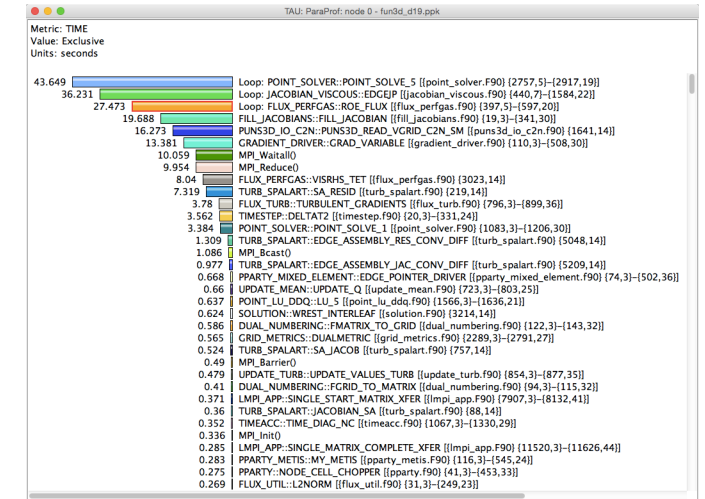
- Requires instrumentation of some kind
 - Manual, automated
 - Source, compiler provided, binary
 - **Library callbacks**, API wrappers, weak symbol replacement
- Simple to implement

■ Sampling

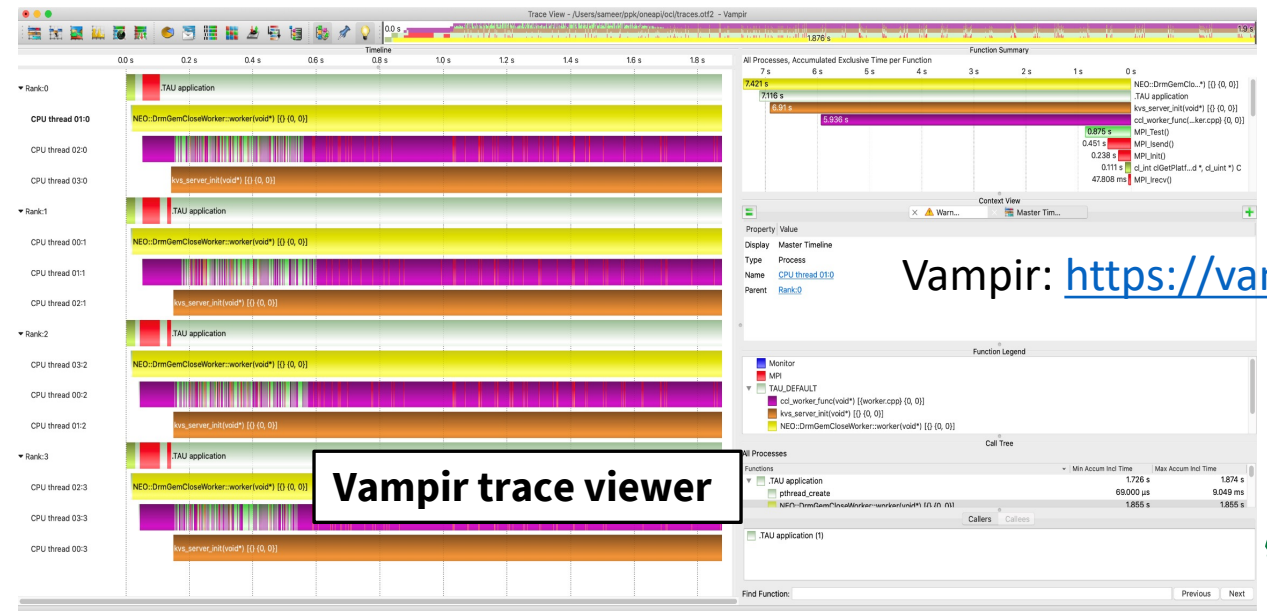
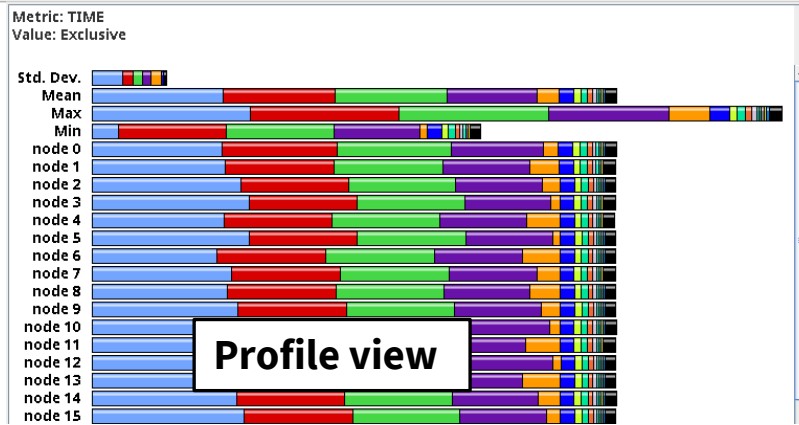
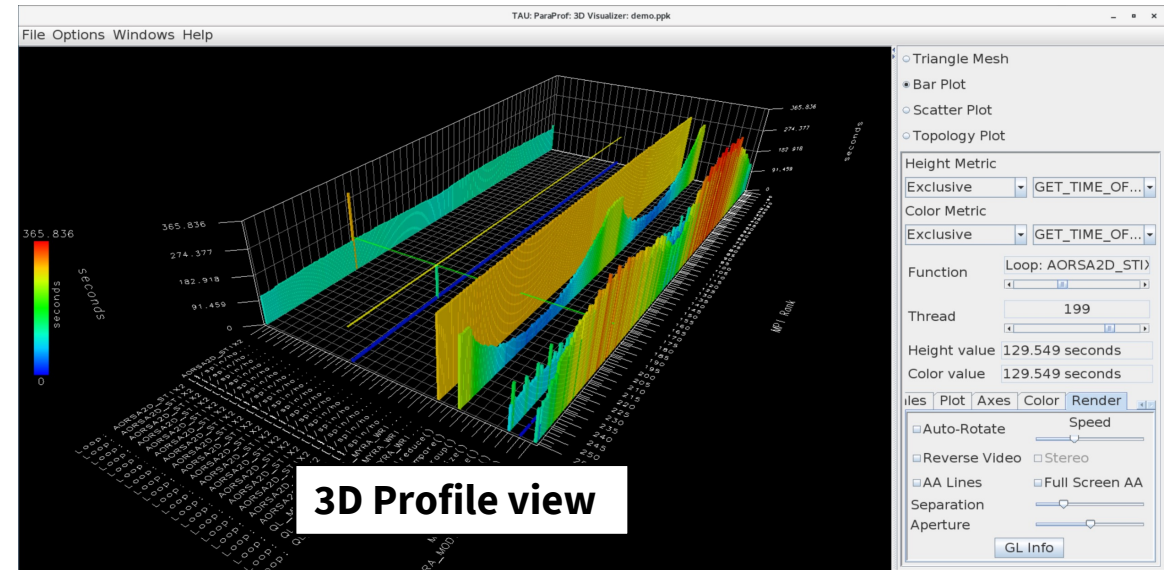
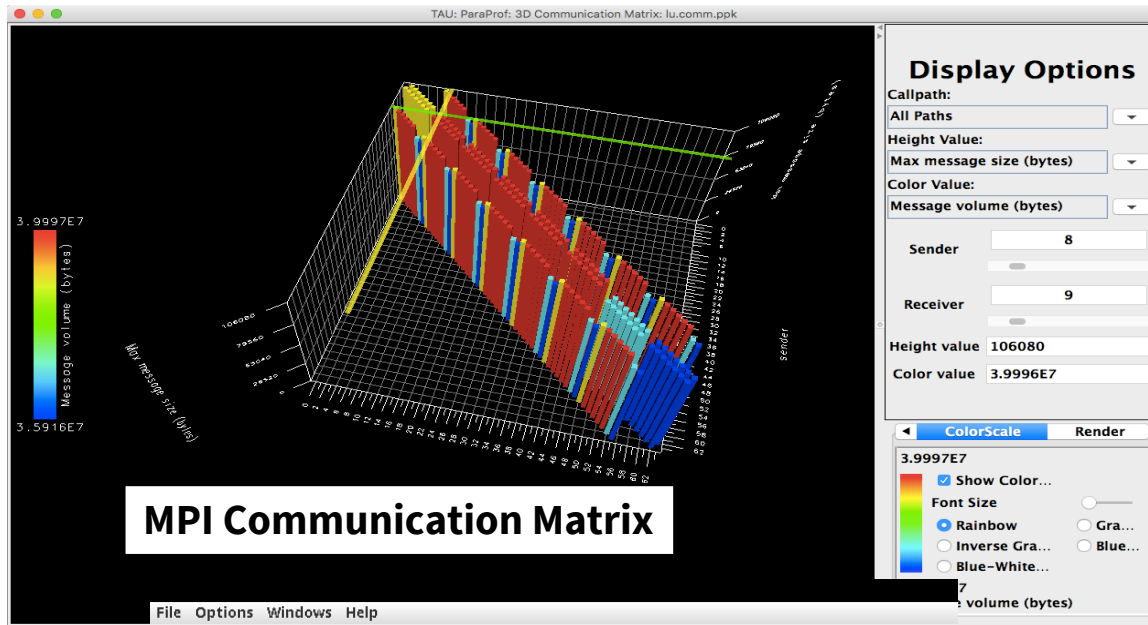
- Requires specialized system libraries / support
 - Periodic signals, signal handler
 - Call stack unwinding
- No modification to executable/library needed
- Potential to interfere with system support (signal handlers)
- Can mix with timers to generate a hybrid profile

Profiling and Tracing

- **Profiling:** how much time was spent in each measured function on each thread in each process?
 - Collapses the time axis
 - No ordering or causal event information
 - Small summary per thread/process, regardless of execution time – only grows with number of timers & threads/processes
- **Tracing:** record all function entry & exit events on a timeline
 - Detailed view of what happened
 - The longer the program runs, the bigger the trace



TAU Analysis Tools: ParaProf, Vampir



Vampir: <https://vampir.eu>

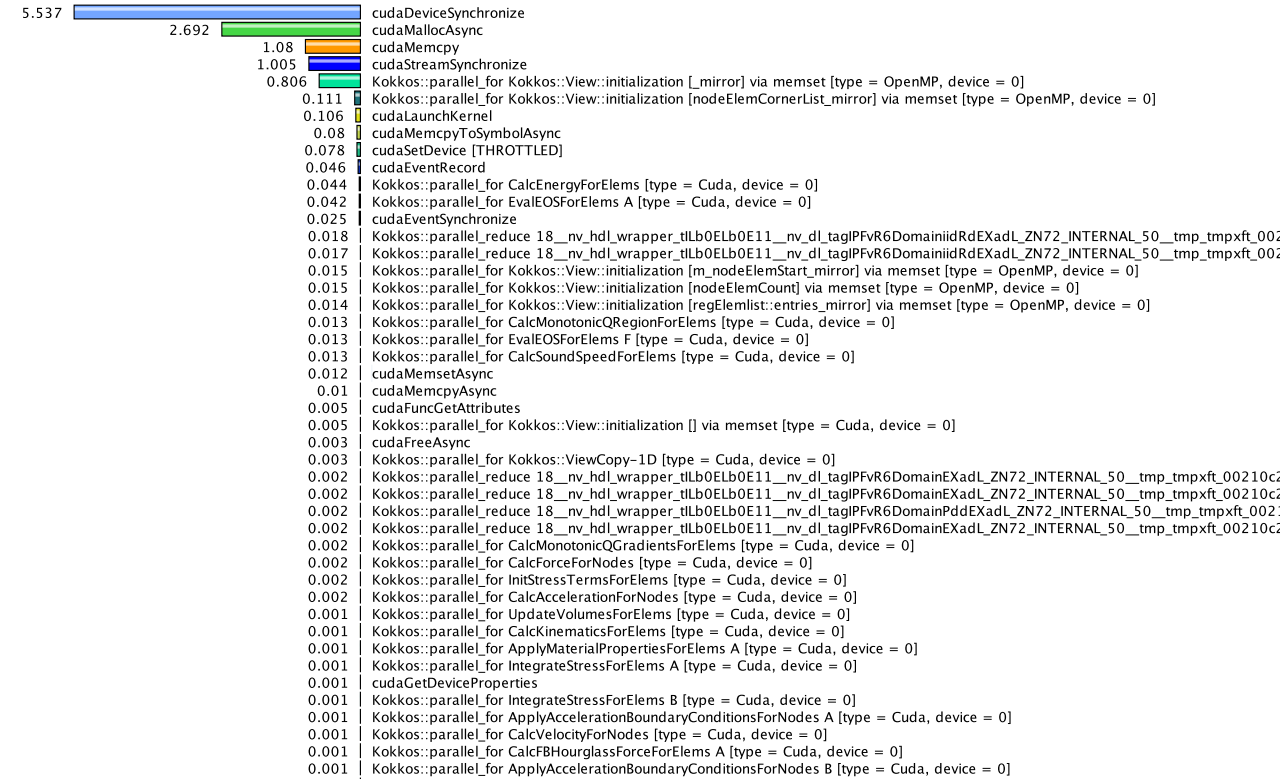
Vampir trace viewer

Kokkos support in TAU – since February, 2017

- TAU implements the Kokkos Profiling API (`Kokkos_Profiling_C_Interface.h`)
- TAU sets an environment variable `KOKKOS_PROFILE_LIBRARY` to tell Kokkos that it should enable profiling and enable function callbacks to the TAU implementations
- TAU implements
 - `kokkosp_[init|finalize]_library`
 - `kokkosp_[begin|end]_parallel_[for|scan|reduce]`
 - `kokkosp_[push|pop]_profile_region`
- Names for regions are passed to the tools to provide intelligent labels
- In addition, TAU also implements support for native Pthreads, OpenMP, OpenACC, CUDA, HIP, SYCL back-end measurement – no code changes necessary
- Fun fact: if you have a Raja application, and Raja is configured with `-DRAJA_ENABLE_RUNTIME_PLUGINS`, Raja implements the same callback API!

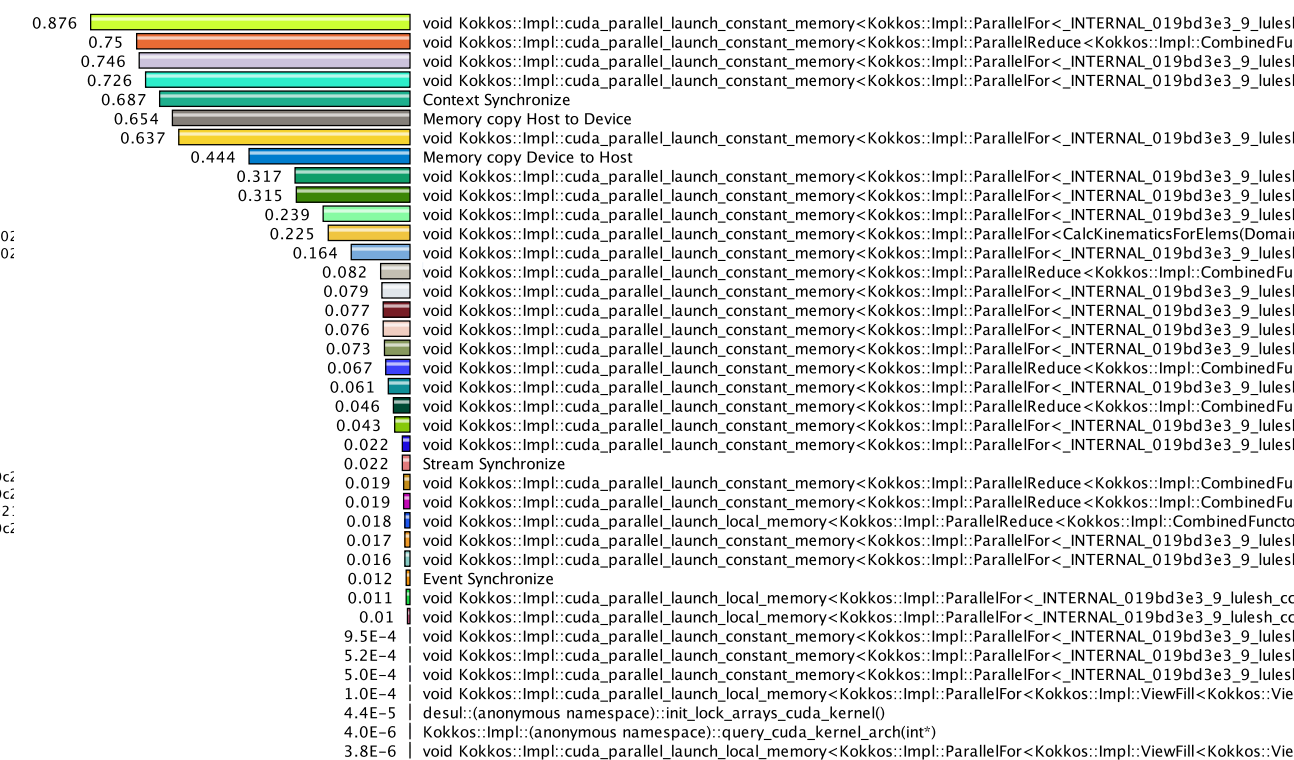
TAU Example – Kokkos Lulesh (from kokkos-miniapps)

Metric: TAUGPU_TIME
Value: Exclusive
Units: seconds



Main thread launching kernels

Metric: TAUGPU_TIME
Value: Exclusive
Units: seconds



Virtual thread with CUDA activity

PerfStubs side note...

- PerfStubs is a “frictionless” instrumentation library
 - <https://github.com/UO-OACISS/perfstubs>
 - One source file, three headers
 - Provides a plugin interface for performance tools
 - Can be compiled away if desired
- Integrated into several libraries (so far) as a git submodule
 - CAMTIMERS
 - PETSc
 - Ginkgo
 - ADIOS2
 - Others?
- Provides runtime integration with TAU & APEX

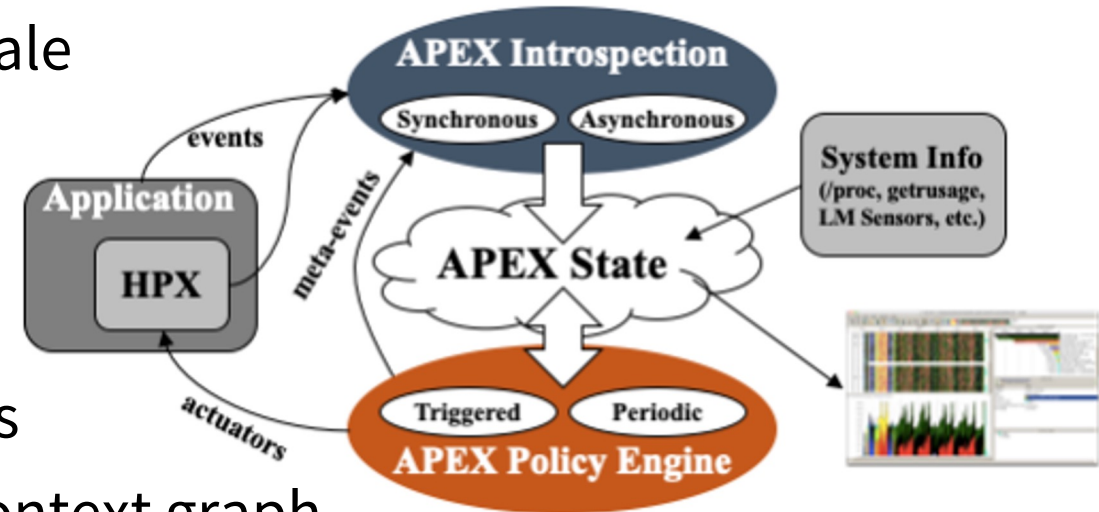
Boehme, Huck, Madsen, Weidendorfer,
“The Case for a Common Instrumentation Interface for HPC Codes”
<https://doi.org/10.1109/ProTools49597.2019.00010>, 2019

APEX

Autonomic Performance Environment for Exascale (APEX)



- Autonomic Performance Environment for eXascale
- Performance Measurement
- **Runtime Adaptation**
- Designed for AMT runtimes (HPX)
 - but works with conventional parallel models

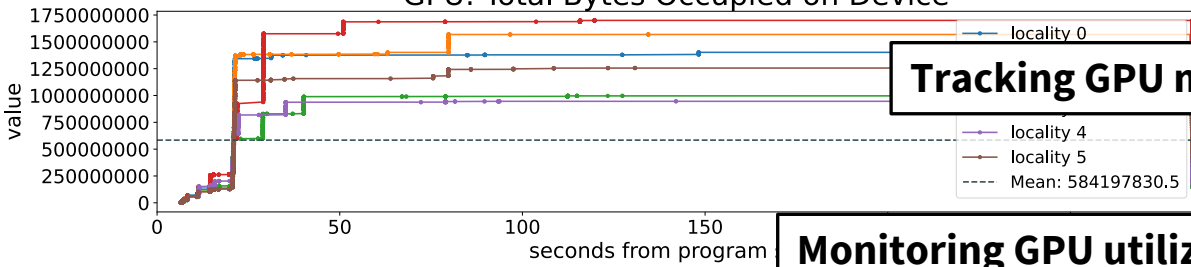


- Focus on **task dependency** graph, not calling context graph
- Supports HPX, C/C++ threads, OpenMP, OpenACC, Kokkos, Raja, CUDA, HIP, SYCL, StarPU... Working on YAKL, Iris
- <https://github.com/UO-OACISS/apex> and <https://github.com/khuck/apex-tutorial>
- Active Harmony* (Nelder Mead), Simulated Annealing, hill climbing for parametric search methods

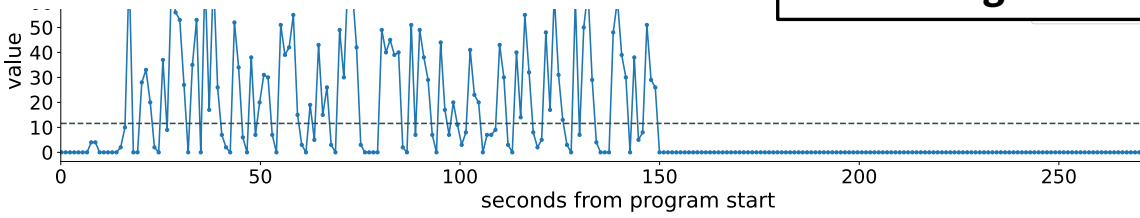
<https://doi.org/10.1109/ESPM256814.2022.00008> : “Broad Performance Measurement Support for Asynchronous Multi-Tasking with APEX”, Huck, ESPM, 2022

APEX example – Octo-Tiger (Octree astrophysics in HPX, Kokkos)

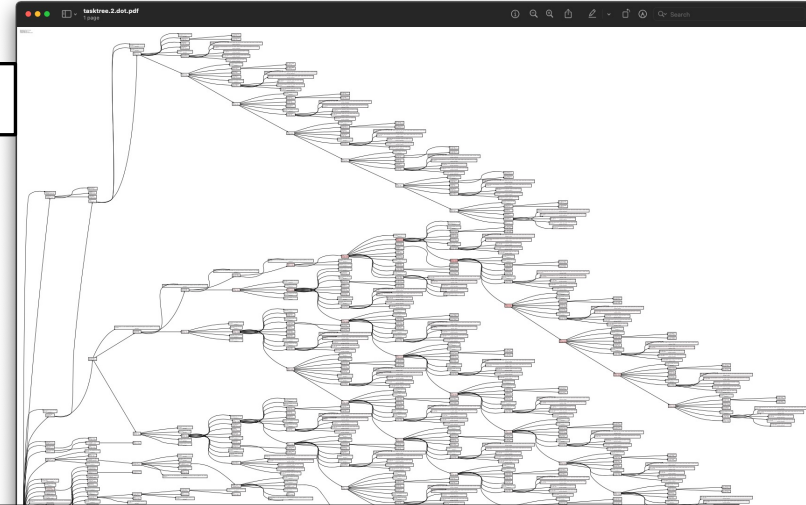
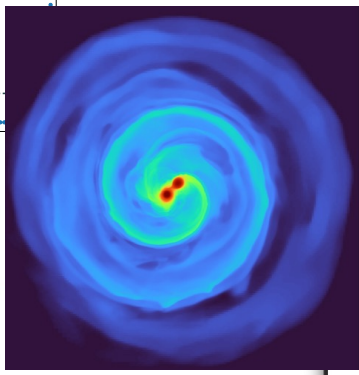
GPU: Total Bytes Occupied on Device



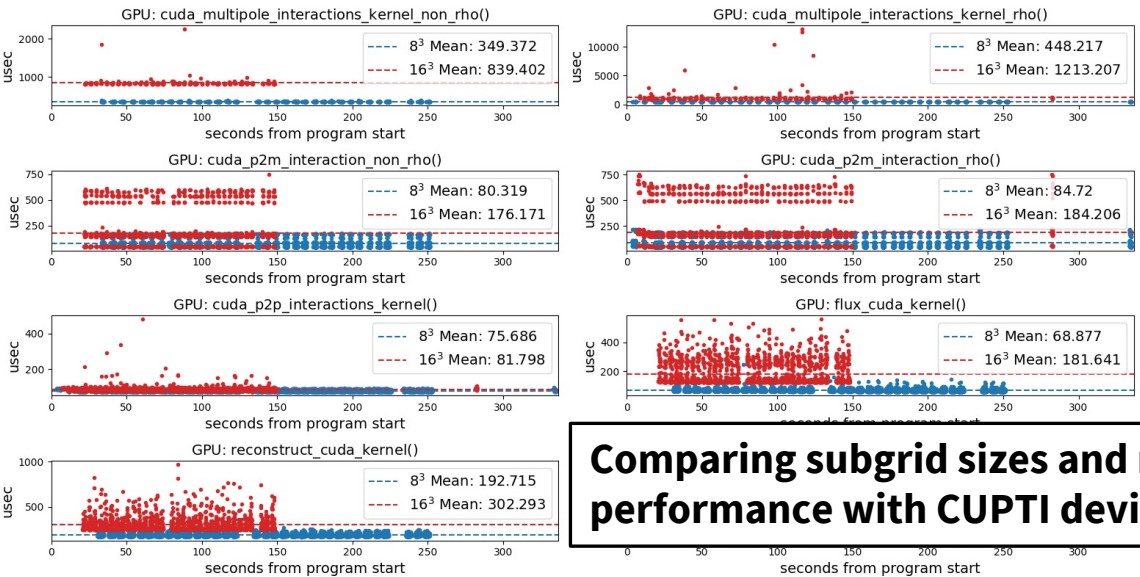
Tracking GPU memory usage with CUPTI



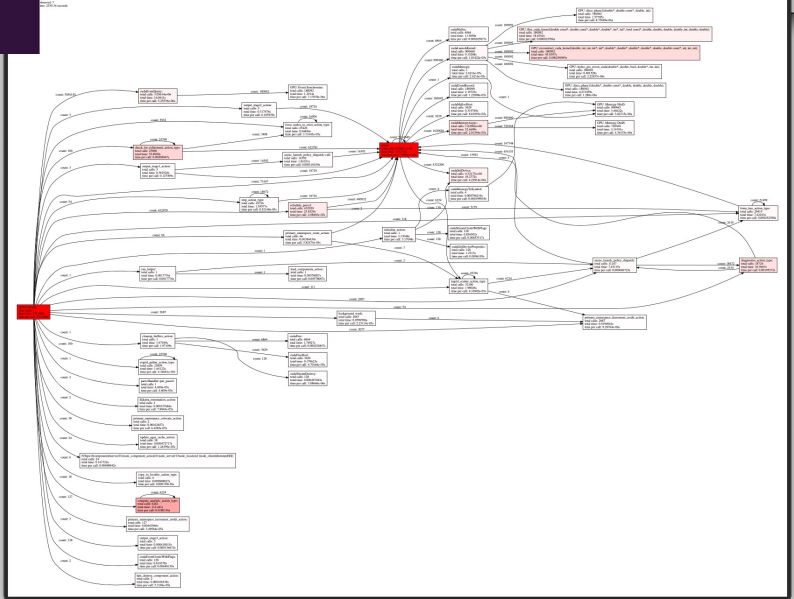
Monitoring GPU utilization with NVML library



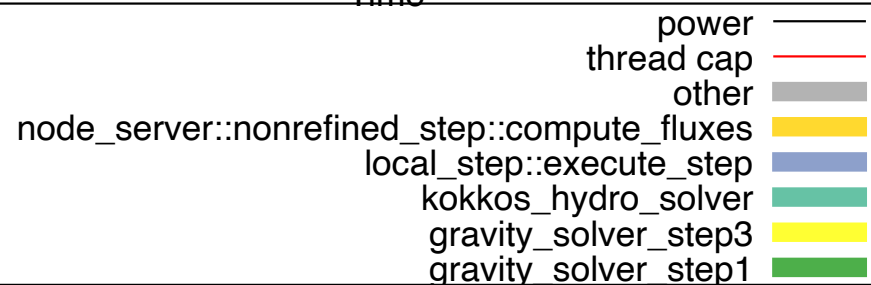
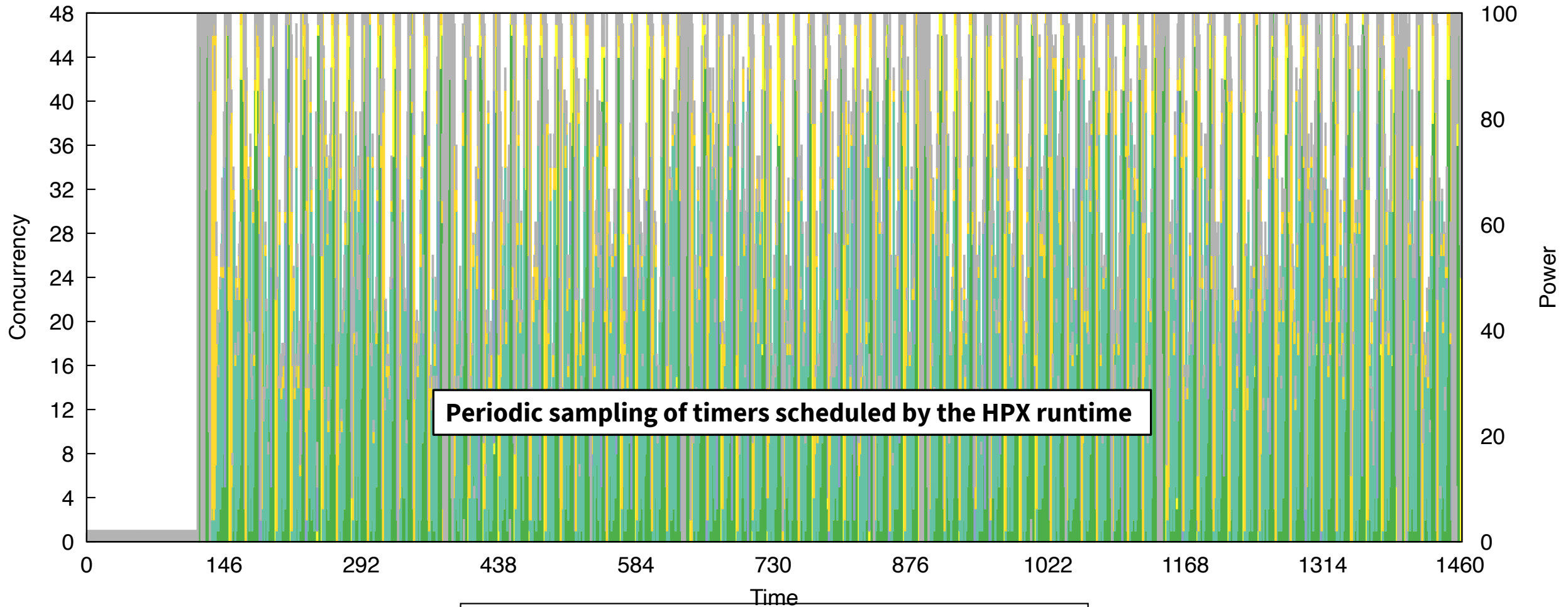
Full task tree (above) and task graph (below) showing task dependencies



Comparing subgrid sizes and relative kernel performance with CUPTI device activity



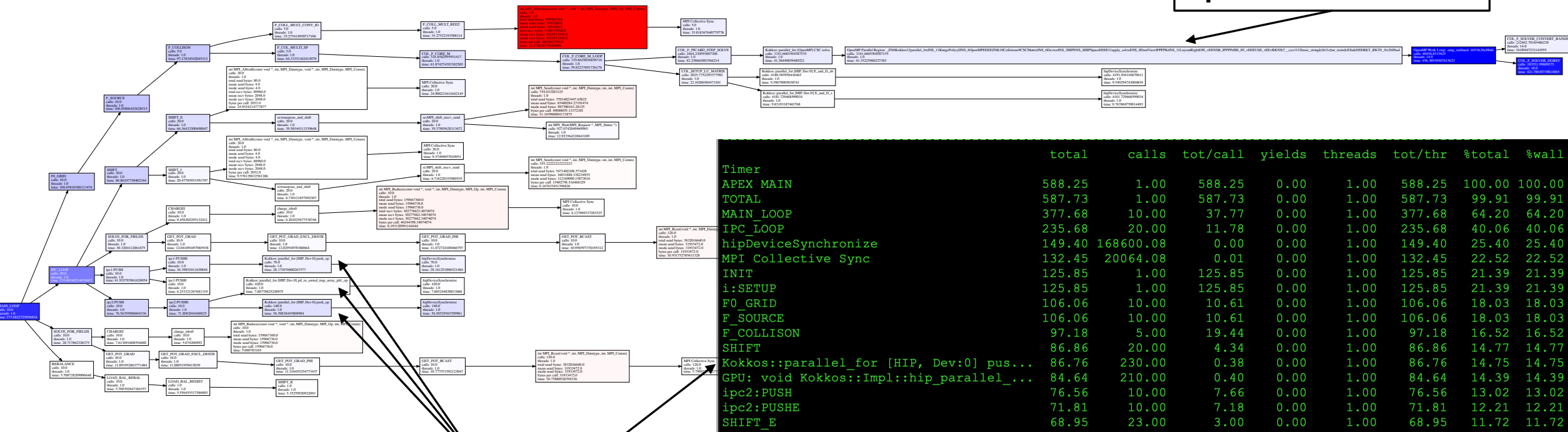
Octotiger on Fugaku – recent result



Example: XGC (tokamak plasma fusion PIC) on Frontier, 512 ranks

- Uses support for MPI, OpenMP-Tools, PerfStubs, Kokkos, Hip
- Post-processing view of MAIN_LOOP subtree, only with accumulated times > 5.0 seconds (only 72 nodes of 6298 of full tree)
- Red: MPI, blue: other, intensity = % of total subtree

OpenMP Kokkos Kernels

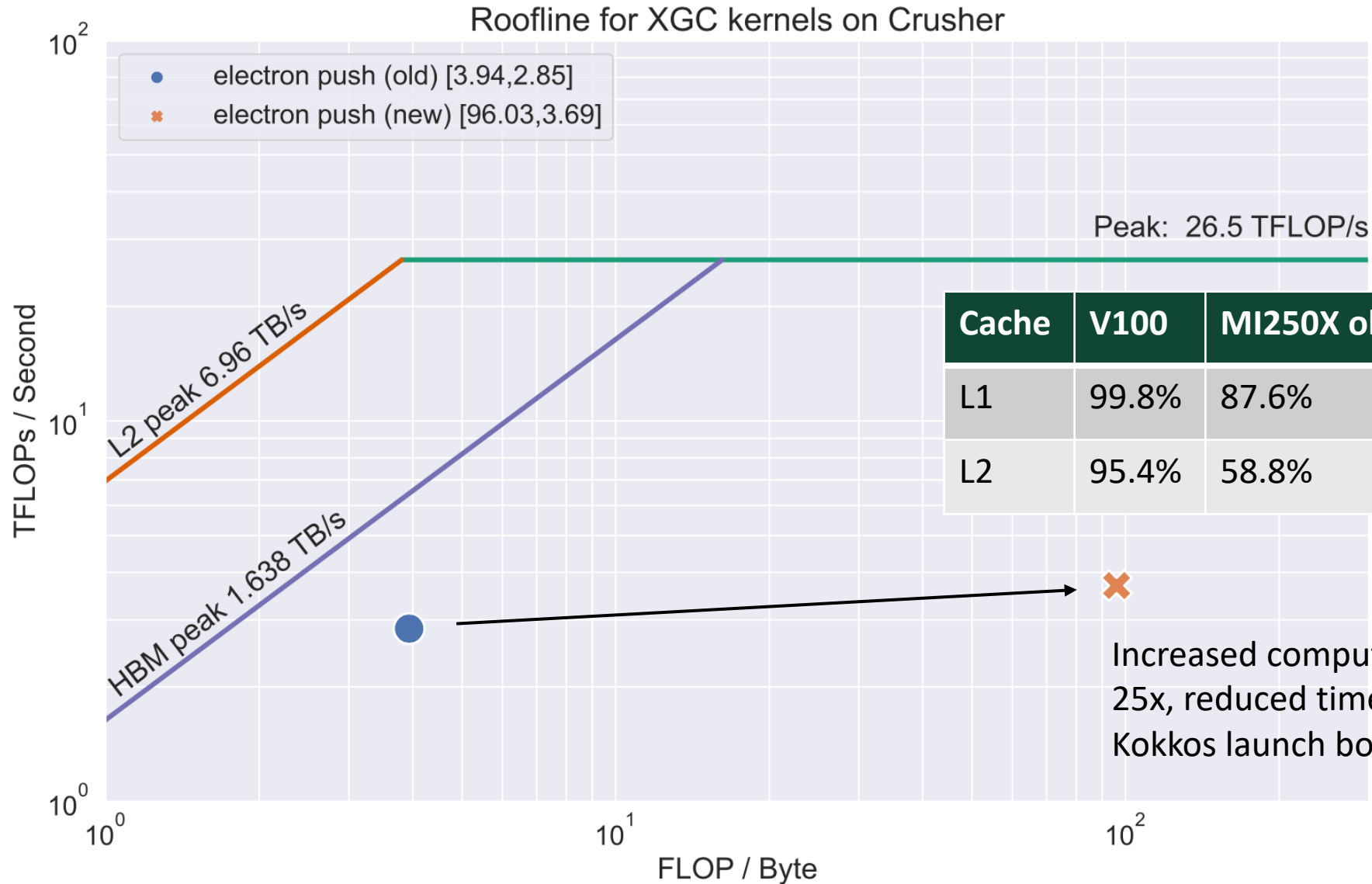


	total	calls	tot/call	yields	threads	tot/thr	%total	%wall
Timer								
APEX MAIN	588.25	1.00	588.25	0.00	1.00	588.25	100.00	100.00
TOTAL	587.73	1.00	587.73	0.00	1.00	587.73	99.91	99.91
MAIN_LOOP	377.68	10.00	37.77	0.00	1.00	377.68	64.20	64.20
IPC_LOOP	235.68	20.00	11.78	0.00	1.00	235.68	40.06	40.06
hipDeviceSynchronize	149.40	168600.08	0.00	0.00	1.00	149.40	25.40	25.40
MPI Collective Sync	132.45	20064.08	0.01	0.00	1.00	132.45	22.52	22.52
INIT	125.85	1.00	125.85	0.00	1.00	125.85	21.39	21.39
i;SETUP	125.85	1.00	125.85	0.00	1.00	125.85	21.39	21.39
F0 GRID	106.06	10.00	10.61	0.00	1.00	106.06	18.03	18.03
F_SOURCE	106.06	10.00	10.61	0.00	1.00	106.06	18.03	18.03
F_COLLISION	97.18	5.00	19.44	0.00	1.00	97.18	16.52	16.52
SHIFT	86.86	20.00	4.34	0.00	1.00	86.86	14.77	14.77
Kokkos::parallel_for [HIP, Dev:0] pus...	86.76	230.00	0.38	0.00	1.00	86.76	14.75	14.75
GPU: void Kokkos::Impl::hip_parallel_...	84.64	210.00	0.40	0.00	1.00	84.64	14.39	14.39
ipc2:PUSH	76.56	10.00	7.66	0.00	1.00	76.56	13.02	13.02
ipc2:PUSHE	71.81	10.00	7.18	0.00	1.00	71.81	12.21	12.21
SHIFT E	68.95	23.00	3.00	0.00	1.00	68.95	11.72	11.72

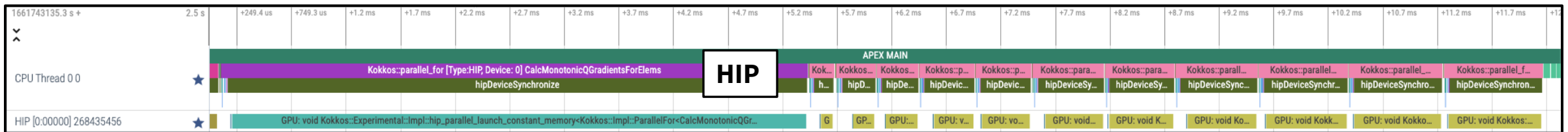
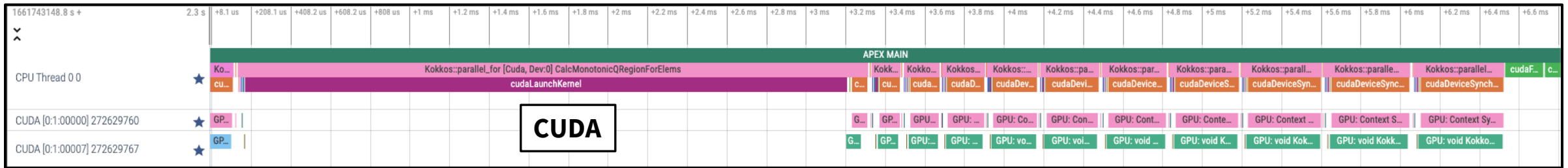
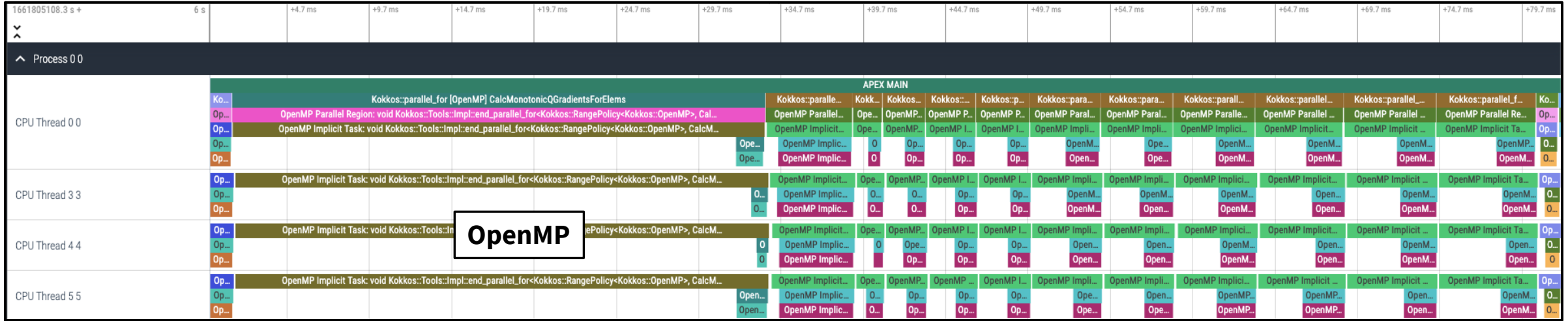
HIP Kokkos Kernels

XGC: Push Kernel on Crusher/Frontier, Kokkos helped generate roofline

Kokkos + APEX + PAPI + rocprofiler = 😊



Kokkos Lulesh and APEX Tracing – OpenMP, CUDA, HIP back ends



Kokkos Support in APEX

- APEX implements the same profiling API that TAU does, *and...*
- APEX provides autotuning (search) support
- Kokkos provides the ability to autotune with:
 - **-DKokkos_ENABLE_TUNING=ON**
- Automatically provides input and context variables for parallel_for, parallel_reduce, parallel_scan, parallel_copy.
 - TeamPolicy: team size and vector length
 - MDRangePolicy: tile sizes
 - RangePolicy*: block size

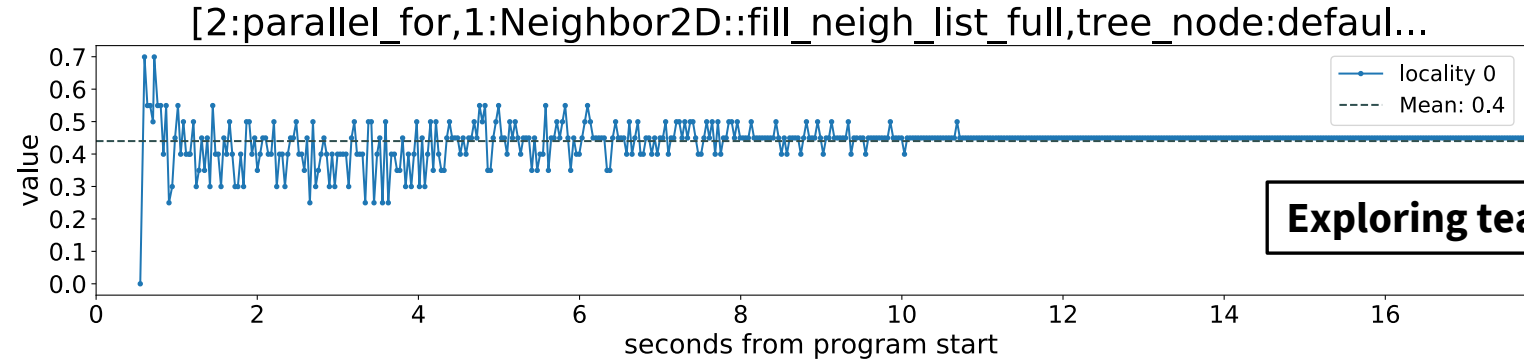
*(in a long-dormant development fork/branch...would be nice to have because many kernels use RangePolicy)

APEX Autotuning of ExaMiniMD Neighbor2D::fill_neigh_list_full kernel

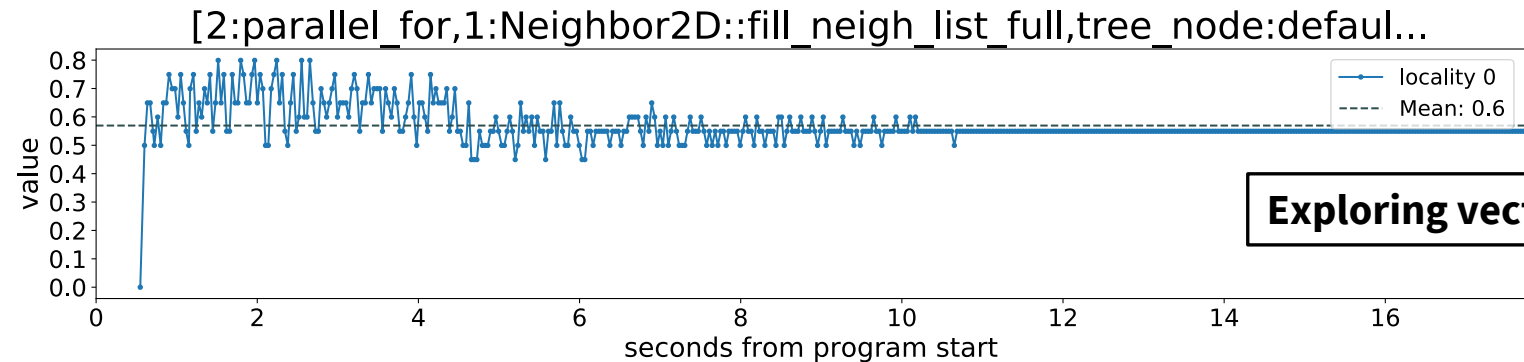
Is it worth it? ...maybe?

Baseline	17.4972s
Tuning	17.7294s
Tuned	17.3790s

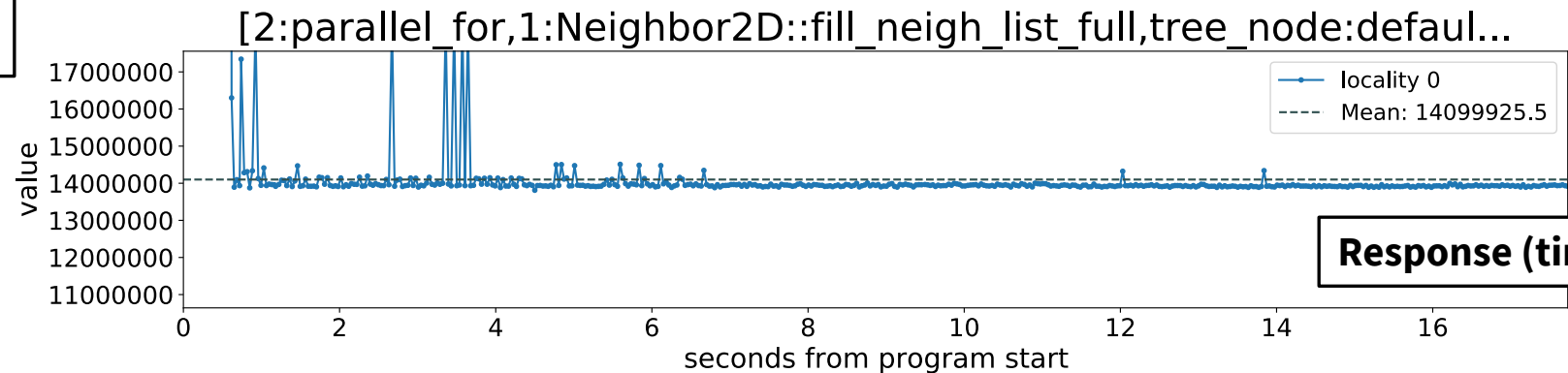
Only one kernel is TeamPolicy in ExaMiniMD - all of the rest of the kernels are Range policies...



Exploring team sizes



Exploring vector lengths



Response (time)

TAU or APEX?

- Use TAU when:
 - Advanced MPI or SHMEM measurements
 - Sampling support
 - HW/OS context (per-OS thread measurements)
 - Broader HW support
 - Python/ML/AI support
 - TAU plugin support
- Use APEX when:
 - Support for asynchronous tasking
 - Focus on algorithmic task dependency, not HW/OS
 - Runtime autotuning / feedback & control support

Kokkos Wishlist

- Autotuning:
 - Support for Range policy
 - Access to non-normalized input variables (all are mapped to [0.0 ... 1.0])
 - I think I know why this was done, but it's confusing/misleading without a map
- Default labels with source info 😊
 - Or at least a function pointer



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.





Thanks! Questions?