

Performance Measurement and Analysis using the TAU Performance System

Allen D. Malony, Sameer Shende, Kevin Huck, Camille Coti, Wyatt Spear, Jeffrey S. Vetter, Mary Hall, Kei Davis, Ian Foster, Scott Klasky, Kerstin Kleese van Dam, Todd Munson

Funded through ECP 2.3.2.10 (PROTEAS-TUNE) and ECP 2.2.6.08 (CODAR)

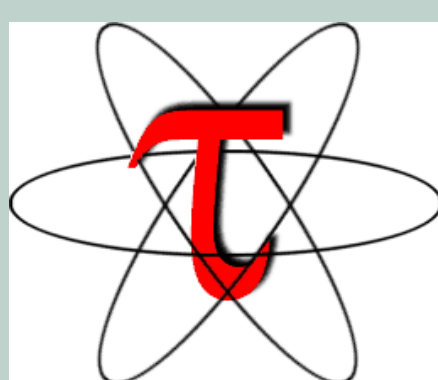
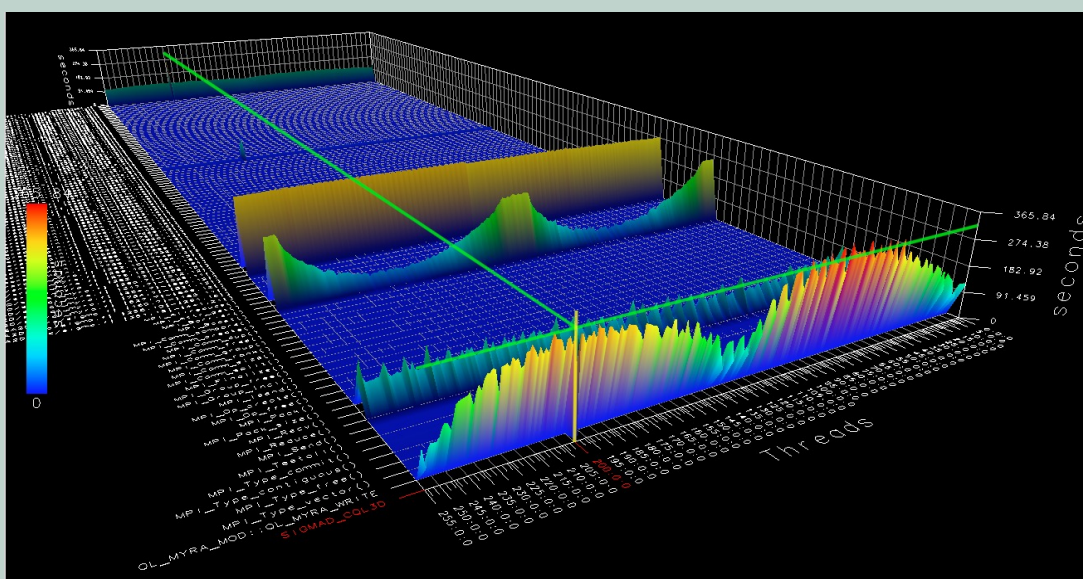


BERKELEY LAB
Bringing Science Solutions to the World



What is TAU?

- Tuning and Analysis Utilities (25+ year project)
- <http://tau.uoregon.edu>
- Comprehensive performance profiling and tracing
 - ✓ Scalable, flexible, portable
 - ✓ Targets all parallel programming/execution paradigms
- Integrated performance toolkit
 - ✓ Instrumentation, measurement, analysis, visualization
 - ✓ Timers, samples, counters, integrated tool callback support, hardware counter support
 - ✓ Widely-ported performance profiling and tracing system
 - ✓ Performance data management and data mining
 - ✓ Open source (BSD-style license)



PROTEAS-TUNE Project Goals

(PROgramming Toolchain for Emerging Architectures and Systems)

Programmer productivity and performance portability are two of the most important challenges facing applications targeting future exascale computing platforms. The PROTEAS-TUNE project is a strategic response to the continuous changes in architectures and hardware (e.g., heterogeneous computing, deep memory hierarchies, nonvolatile memory) that are defining the landscape for emerging ECP systems. PROTEAS-TUNE is a flexible programming framework and integrated toolchain that will provide ECP applications the opportunity to work with programming abstractions and to evaluate solutions that address the exascale programming challenges they face.

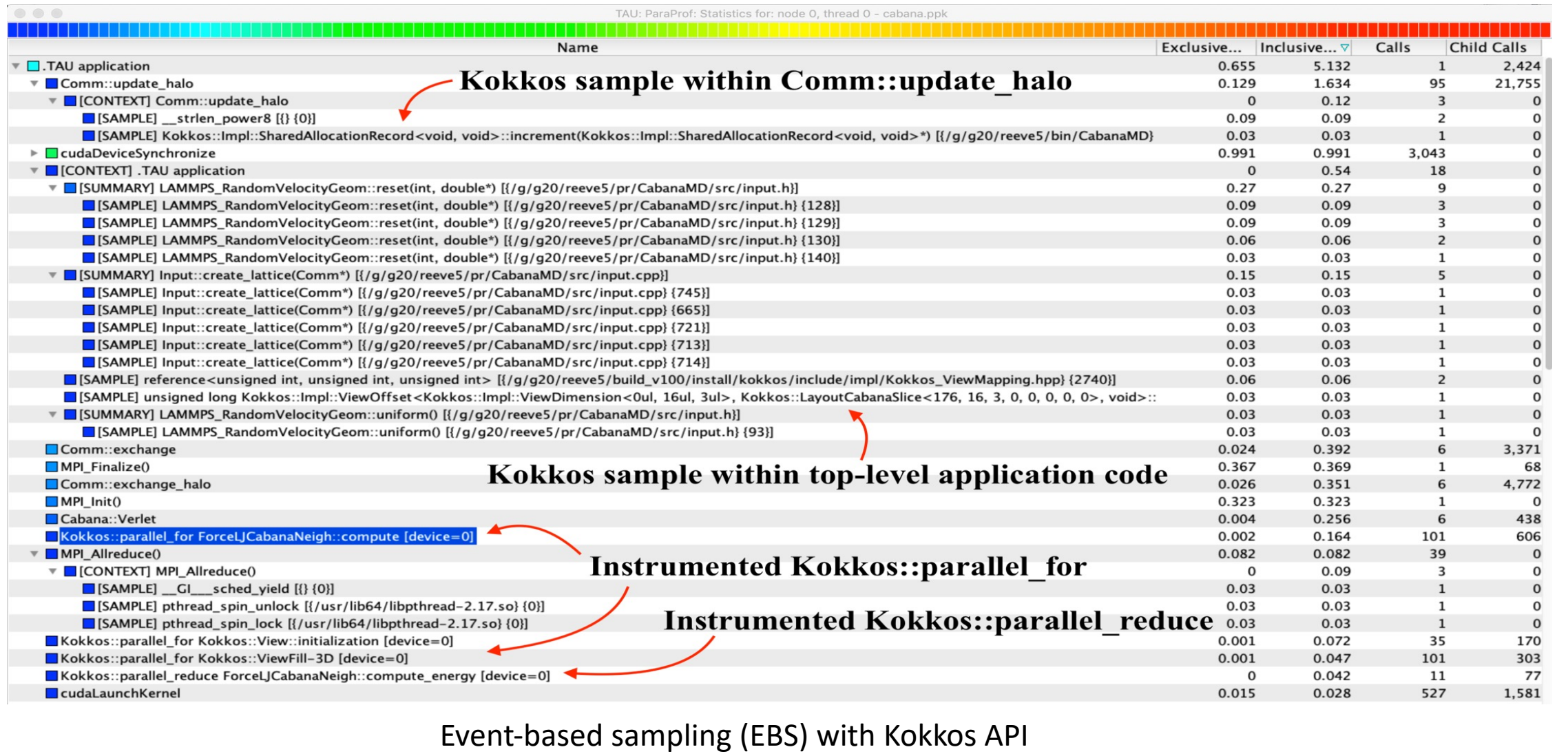
Key Capabilities: LLVM; OpenACC, CUDA, HIP, OpenCL, OneAPI; Performance tools with TAU; Expertise and software systems for heterogeneous computing (GPUs, FPGAs, Manycore) and deep memory hierarchies including nonvolatile memory; Performance portability metrics, tools, and strategies.

New Capabilities in TAU from PROTEAS-TUNE

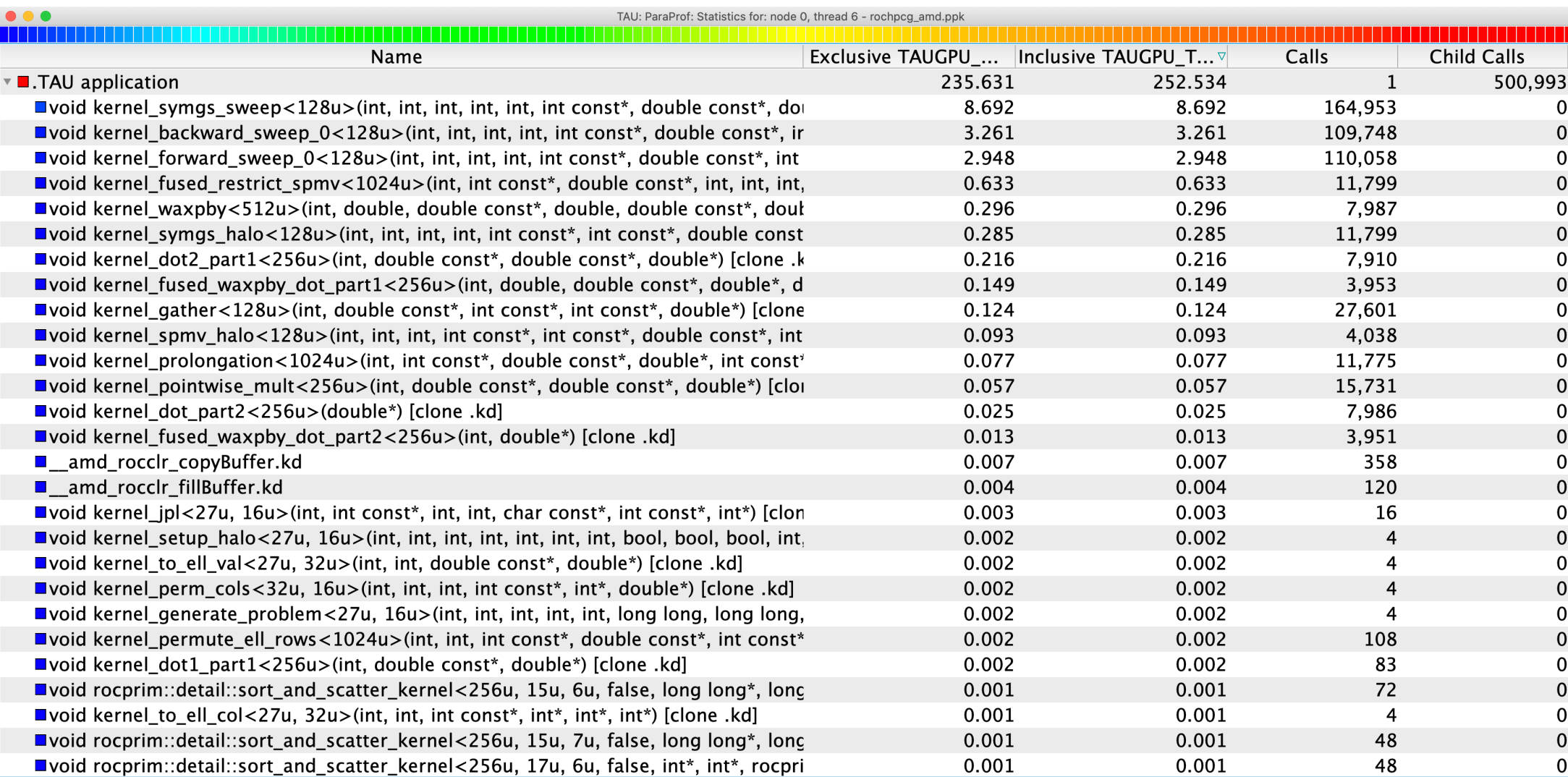
- **Updated support for latest CUDA, AMD GPUs and software**
- **clang/clang++ TAU plugin** : Selective instrumentation of C and C++ code, choice of functions to instrument based on name and / or source file – now with support for hipcc.
- **Monitoring support** : updated monitoring plugin provides support for capturing hardware and OS state, NVIDIA Monitoring Library (NVML) support, and broad PAPI counter support (capture all available metrics globally, in one run, using multiplexing).
- **Python3+CUDA** : Enhanced dynamic measurement support for complex deep/machine learning platforms such as TensorFlow and PyTorch.
- **Intel OneAPI** : Implemented support for Intel oneAPI, including Level Zero, time spent in kernels. on GPU and time spent in OneAPI calls.
- **OpenMP** : Updating OMPT offload event support provided by AMD 5.1.0 clang / hipcc compilers.
- **F18/Flang Instrumentation** : TAU support for the Flang Fortran compiler was added. PDT- and Compiler-based instrumentation support for Flang was implemented and tested on x86_64 Linux platforms.
- **OpenACC** : Updated profiling support for Clacc, developing OpenACC support for f18 Fortran.

Event-based Sampling (EBS):

CabanaMD on an IBM AC922 with NVIDIA V100 GPUs



AMD HIP: Kernel execution on GPUs: rochpcg



TAU / PROTEAS-TUNE: Next Steps

- **CUDA 11** : Finish new Perfworks APIs for CUDA/CUPTI 10+ to replace deprecated CUPTI metric support.
- **OpenMP / OpenACC** : Continue to explore and implement prototype measurement for OpenMP and OpenACC regions executed on target devices.
- **ROCm/HIP** : Continued support for AMD GPUs
- **oneAPI** : Continue to update and develop support for Intel GPUs
- **TAU compiler wrappers** : new instrumentation support based on libclang and libtooling libraries from LLVM

Related Project: APEX

Autonomic Performance Environment for Exascale

<https://github.com/UO-OACISS/apex>

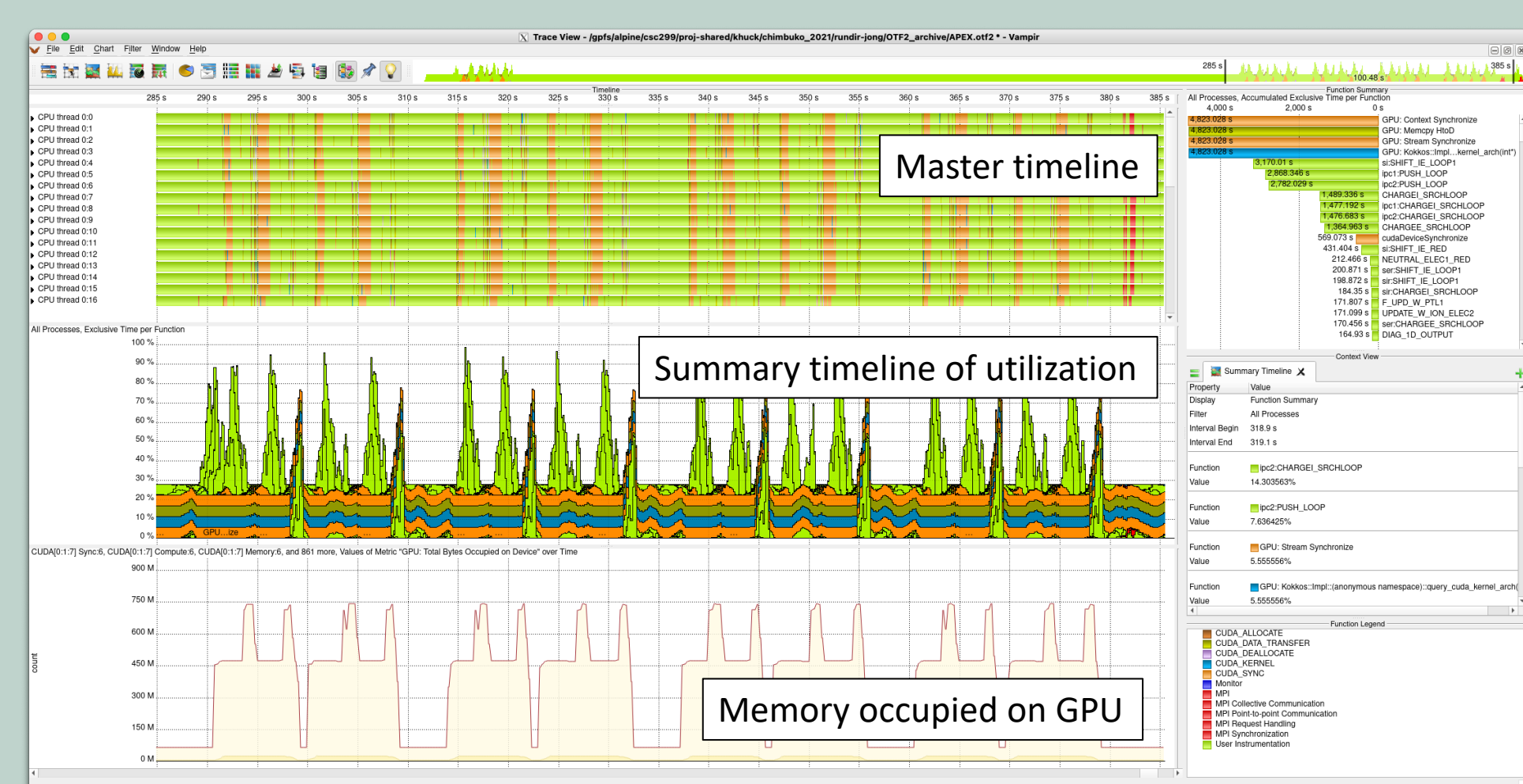
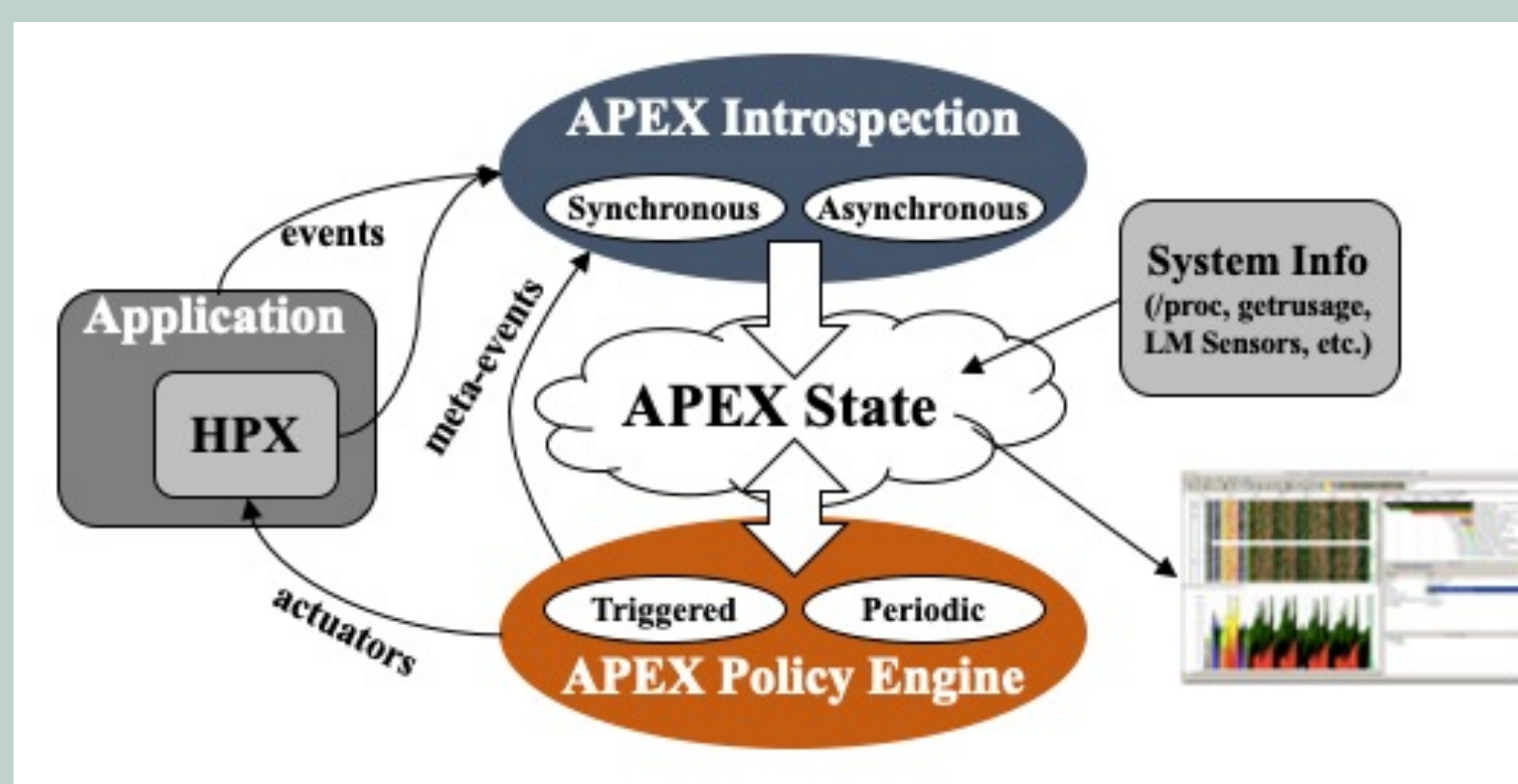


Figure: XGC executed on Summit, measured with APEX, visualized in Vampir, including CUDA and Kokkos activity. APEX supports several asynchronous threading and tasking programming models, including C++ std::async, HPX, OpenMP, OpenACC, Kokkos, Raja, CUDA, HIP. Monitoring support similar to the TAU plugin, providing utilization statistics for the filesystem, network, devices, CPU, GPU, and memory.

This research was supported by the Exascale Computing Project (17-SC-20-SC), a collaborative effort of two U.S. Department of Energy organizations (Office of Science and the National Nuclear Security Administration) responsible for the planning and preparation of a capable exascale ecosystem, including software, applications, hardware, advanced system engineering, and early testbed platforms, in support of the nation's exascale computing imperative.