Using HPCToolkit to Measure and Analyze the Performance of GPU-accelerated Applications



John Mellor-Crummey Rice University

Tutorial
Mar-Apr 2021
NERSC and OLCF (Virtual)











Acknowledgments

Current funding

- DOE Exascale Computing Project (Subcontract 4000151982)
- DOE Labs: ANL (Subcontract 9F-60073), Tri-labs (LLNL Subcontract B645220)
- Industry: AMD, Total E&P Research & Technology

Team

- Rice University
 - HPCToolkit PI: Prof. John Mellor-Crummey
 - Research staff: Laksono Adhianto, Mark Krentel, Xiaozhu Meng, Scott Warren
 - Contractor: Marty Itzkowitz
 - Students: Jonathon Anderson, Aaron Cherian, Dejan Grubisic, Yumeng Liu, Keren Zhou
 - Recent summer interns: Vladimir Indjic, Tijana Jovanovic, Aleksa Simovic
- University of Wisconsin Madison
 - Dyninst PI: Prof. Barton Miller



HPCToolkit Team at the Tutorial



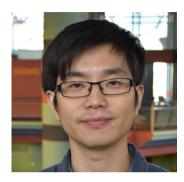
Laksono Adhianto



Mark Krentel



John Mellor-Crummey



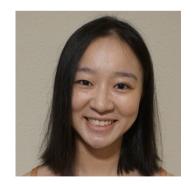
Xiaozhu Meng



Aaron Cherian



Dejan Grubisic



Yumeng Liu



Keren Zhou



DOE's Forthcoming Heterogeneous Exascale Platforms







Aurora compute nodes (ALCF)

- 2 Intel Xeon "Sapphire Rapids" processors
- 6 Intel Xe "Ponte Vecchio" GPUs
- 8 Slingshot endpoints
- Unified memory architecture

Frontier compute nodes (OLCF)

- 1 AMD EPYC CPU
- 4 purpose-built AMD Radeon Instinct GPUs
- Multiple Slingshot endpoints
- Unified memory architecture

El Capitan compute nodes (LLNL)

- Next-generation AMD EPYC "Genoa" CPU (5nm)
- Next-generation AMD Radeon Instinct GPUs

DOE's NVIDIA-based Heterogeneous Supercomputers



- Summit compute nodes (OLCF)
 - 2 IBM Power9 processors
 - 6 NVIDIA V100 GPUs
 - Dual-rail Mellanox EDR Infiniband Fat Tree



- Sierra compute nodes (LLNL)
 - 2 IBM Power9 processors
 - 4 NVIDIA V100 GPUs
 - Mellanox EDR Infiniband Fat Tree



- Perlmutter compute nodes (LBNL)
 - AMD Milan CPU
 - 4 NVIDIA A100 GPUs
 - Cray Slingshot Interconnect



Node-level Programming Models for Heterogeneous Supercomputers

- Native programming models from platform vendors
 - Intel DPC++

C++ + SYCL specification 1.2.1 + extensions: device queues, buffers, accessors, parallel_for, single_task, parallel_for_work_group, events

- CUDA
- HIP: AMD's CUDA-like model
- Directive-based models
 - OpenACC

• OpenMP

offload structured blocks and device functions, work sharing loops, data environment, data mappings

C++ template-based models

- RAJA
- Kokkos

parallelism loops, iteration spaces, execution policies, traversal templates, lambda functions, n-dimensional array abstractions, and lambda functions



Global Programming Models

- Message passing
 - MPI
- Partitioned global address space programming models
 - languages
 - Coarray Fortran, Coarray C++, Chapel, UPC
 - libraries
 - UPC++, GASNet, OpenSHMEM, Global Arrays
- Object-based
 - Charm++

Sources: Frontier Spec Sheet https://www.olcf.ornl.gov/wp-content/uploads/2019/05/frontier_specsheet.pdf https://docs.nersc.gov/development/programming-models



Performance Analysis Challenges for GPU-accelerated Supercomputers

Myriad performance concerns

- Computation: need extreme-scale data parallelism to keep GPUs busy
- Data movement costs within and between memory spaces
- Internode communication
- I/O

Many ways to hurt performance

 insufficient parallelism, load imbalance, serialization, replicated work, data copies, synchronization, lack of locality, ...

Hardware and execution model complexity

- Multiple compute engines with vastly different characteristics, capabilities, and concerns
- Multiple memory spaces with different performance characteristics
- Asynchronous execution



Measurement Challenges for GPU-accelerated Supercomputers

Extreme-scale parallelism

Serialization within tools will disrupt parallel performance

Dependent on third-party measurement interfaces

- Hardware
 - CPU hardware performance monitoring unit
 - GPU hardware counters and PC sampling
- Software
 - Glibc LD_AUDIT for tracking dynamic loading of shared libraries
 - Linux perf_events for kernel measurement
 - GPU monitoring and instrumentation libraries from vendors

Multiple measurement modalities and interfaces

- Sampling on the CPU
- Callbacks when GPU operations are launched and (sometimes) completed
- GPU event stream, including PC sampling measurements
- Frequent GPU kernel launches require a low-overhead measurement substrate



Engineering Challenges for Performance Tools

Complex applications

- Compositions of programming models
- > 100 dynamic libraries
- Application binaries exceeding 5GB
- HPC libraries that intercept system calls (mmap, munmap, open, close)
- Quirky application characteristics
 - NAMD: exit initiated by a non-initial thread
 - Kull: forking non-readable helper application

Dynamic library loading

- Implicit system locks on dynamic library state
- RUNPATH: library-specific library load path
- Early threads in library init constructors
- Nested dynamic library loading

Provisioning thread local state

Implicit lock when creating or destroying thread local storage

Process fork

atfork handlers trigger thread destructors

Interactions with vendor tool substrates

- Libraries lack documentation of their actions, e.g. creating threads
- Callbacks for submission and completion on unspecified (and sometimes different) threads

Interaction between tools and software stack

- Interaction of signals with everything
- Managing monitoring when forking

Lack of vendor tooling and documentation

Non-standard GPU binary formats that lack public documentation



Other GPU Performance Tools

Features

- Trace view
 - A series of events that happen over time on each process, thread, and GPU stream
- Profile view
 - A correlation of performance metrics with program contexts

Tools

- GPU vendors
 - Nsight Systems, Nsight Compute, nvprof, ROCProfiler, Intel VTune
- Third party
 - TAU, VampirTrace, ARM Map



Shortcomings of Other Tools for Complex GPU-accelerated Programs

- They lack a comprehensive profile view to analyze
 - CPU calling contexts (including inlined frames) where GPU operations are invoked
 - Understanding where, how and why GPU kernels arose from instantiation of nested templates
 - Understanding costs of GPU APIs (e.g., cudaMemcpy) invoked from many different contexts
 - Sophisticated GPU calling contexts
 - OpenMP Target, Kokkos, and RAJA generate GPU code with many small procedures
 - Loop-level performance information on CPUs and GPUs
- At best, existing tools only attribute runtime cost to a flat profile view of functions executed on GPUs



Outline

- Performance measurement and analysis challenges for GPU-accelerated supercomputers
- Introduction to HPCToolkit performance tools
 - Overview of HPCToolkit components and their workflow
 - HPCToolkit's graphical user interfaces
- Analyzing the performance of GPU-accelerated supercomputers with HPCToolkit
 - Overview of HPCToolkit's GPU performance measurement capabilities
 - Collecting measurements
 - Analysis and attribution
 - Scalable analysis of performance data
- Status, ongoing work, final remarks



Rice University's HPCToolkit Performance Tools

• Employs binary-level measurement and analysis

- Observes executions of fully optimized, dynamically-linked parallel applications
- Supports multi-lingual codes with external binary-only libraries

Collects sampling-based measurements of CPU

- Controllable overhead
- Minimize systematic error and avoid blind spots
- Enable data collection for large-scale parallelism

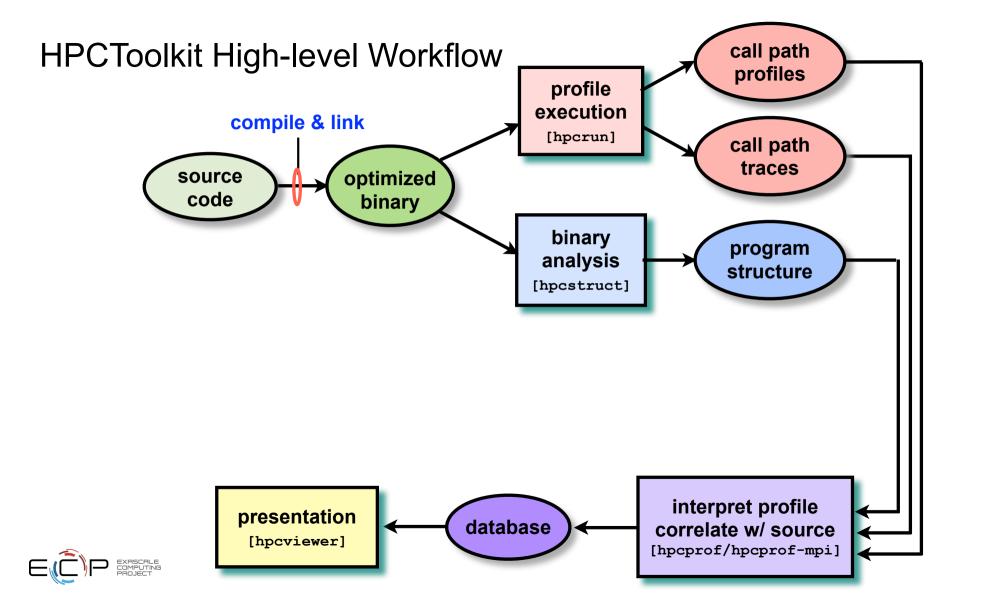
Measures GPU performance using vendor APIs

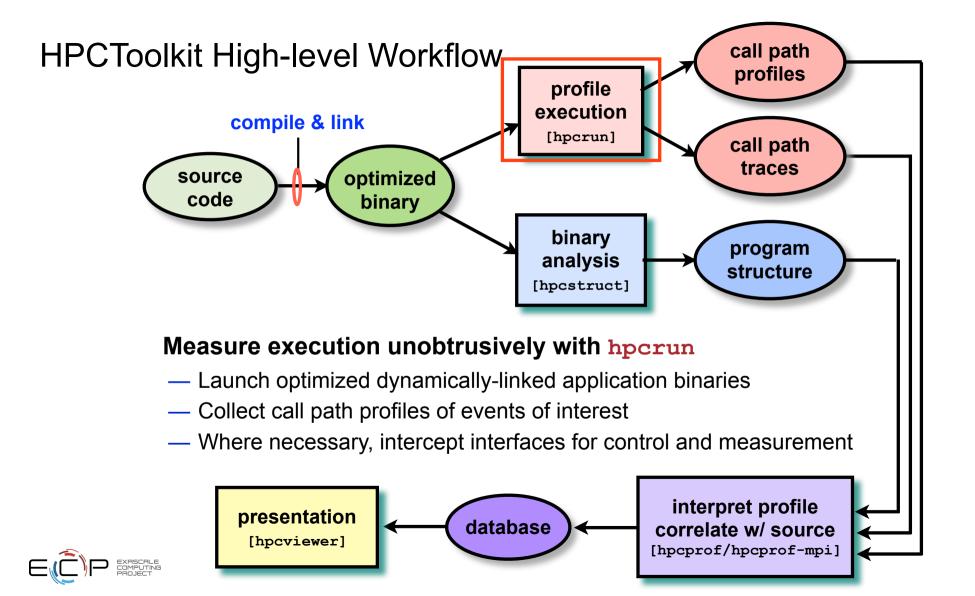
- Register callbacks to monitor launch/completion of GPU operations
- Monitor asynchronous GPU operations using activity APIs from NVIDIA and AMD
- Collect fine-grain measurements of GPU code using PC sampling (NVIDIA) and instrumentation (Intel GTPin)

Associates metrics with both static and dynamic context

- Loop nests, procedures, inlined code, calling contexts on both CPU and GPU
- Enables one to specify and compute derived CPU and GPU performance metrics of your choosing
 - Diagnosis often requires more than one species of metric
- Supports top-down performance analysis
 - Identify costs of interest and drill down to causes: up and down call chains, over time

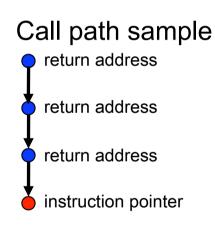


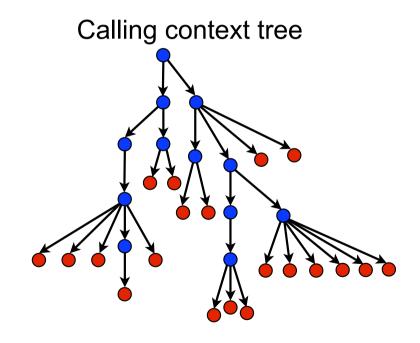




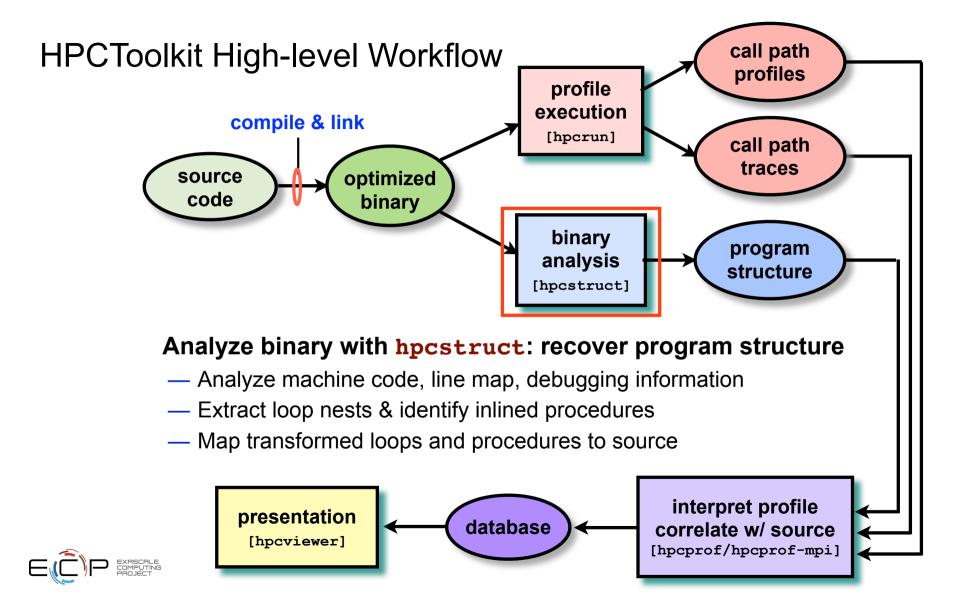
Call Path Profiling

- Measure and attribute costs in context
 - Sample timer or hardware counter overflows
 - Gather CPU calling context using stack unwinding







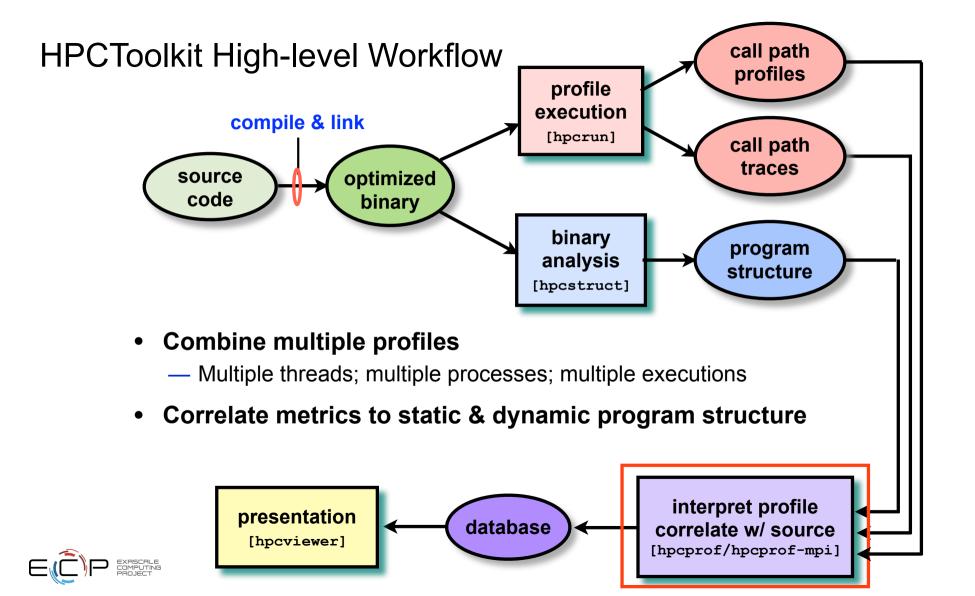


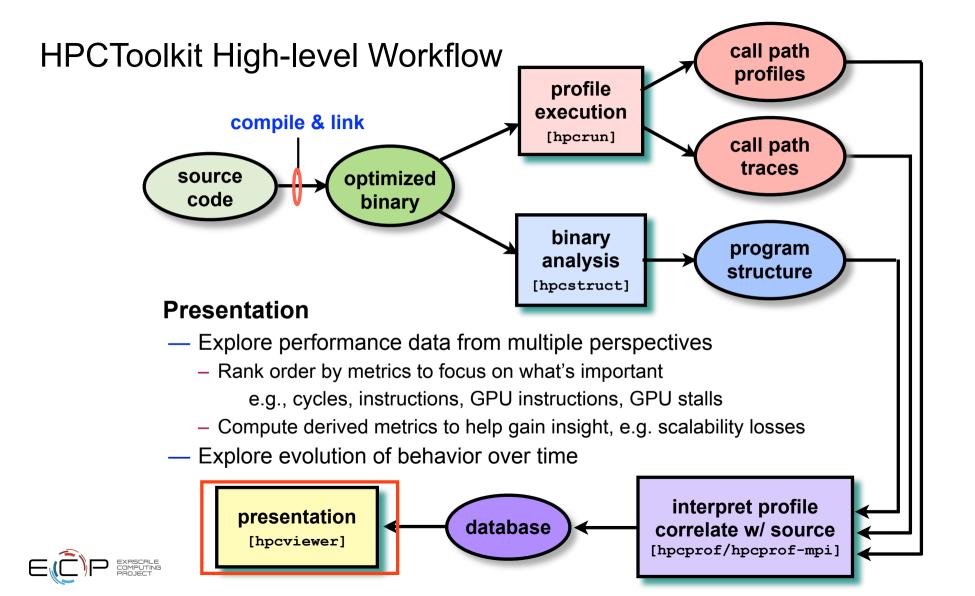
Dyninst: A Toolkit for Binary Analysis and Instrumentation

Architectures Code Gen API Patch API X86 64 Symtab API Parse API Power/BE Power/LE **DataFlow API** Instruction API **ARM** AMD Vega StackWalker API ProcControl API **CUDA** Intel GPU

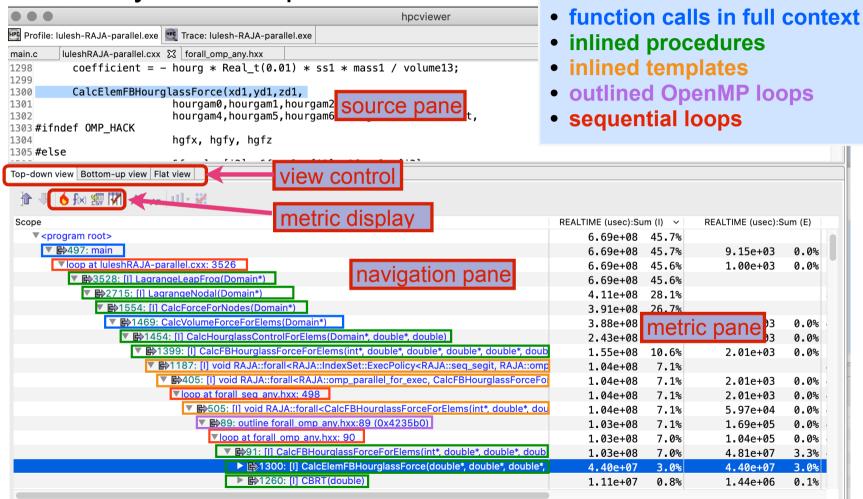


Lead Institution: University of Wisconsin – Madison





Code-centric Analysis with hpcviewer





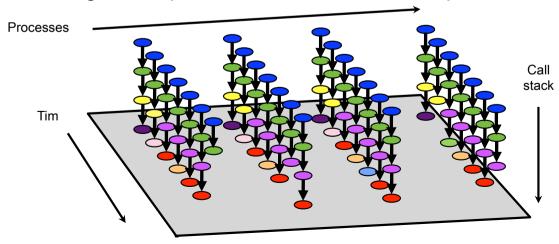
Understanding Temporal Behavior

Profiling compresses out the temporal dimension

- Temporal patterns, e.g. serial sections and dynamic load imbalance are invisible in profiles

What can we do? Trace call path samples

- N times per second, take a call path sample of each thread
- Organize the samples for each thread along a time line
- View how the execution evolves left to right
- What do we view? assign each procedure a color; view a depth slice of an execution





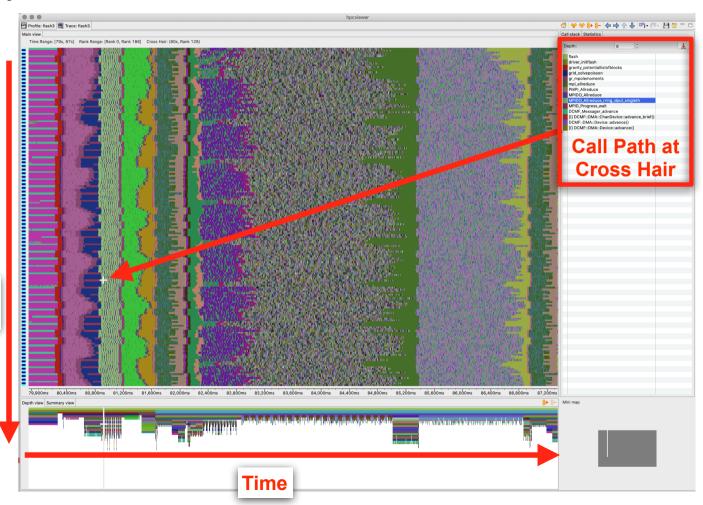
Time-centric Analysis of Call Path Traces

Detail of a trace of Flash3 - block structured AMR code written in Fortran

• 256 ranks

Depth, trace, and call path views

Ranks/ Threads



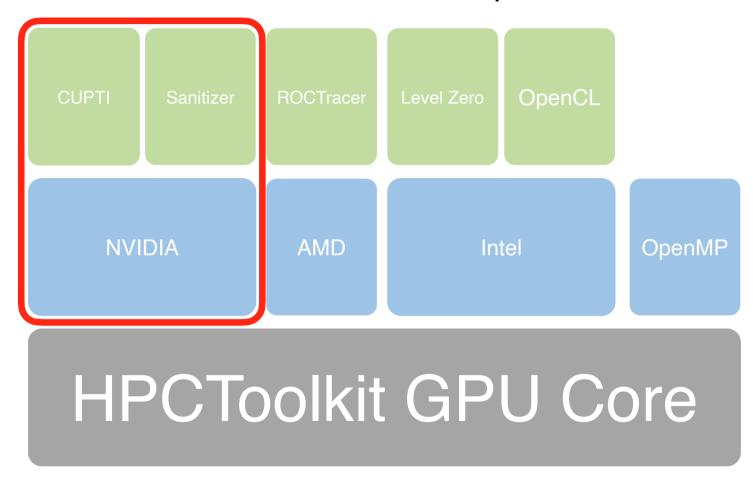


Outline

- Performance measurement and analysis challenges for GPU-accelerated supercomputers
- Introduction to HPCToolkit performance tools
 - Overview of HPCToolkit components and their workflow
 - HPCToolkit's graphical user interfaces
- Analyzing the performance of GPU-accelerated supercomputers with HPCToolkit
 - Overview of HPCToolkit's GPU performance measurement capabilities
 - Collecting measurements
 - Analysis and attribution
 - Scalable analysis of performance data
- Status, ongoing work, final remarks



HPCToolkit for GPU-accelerated Computations



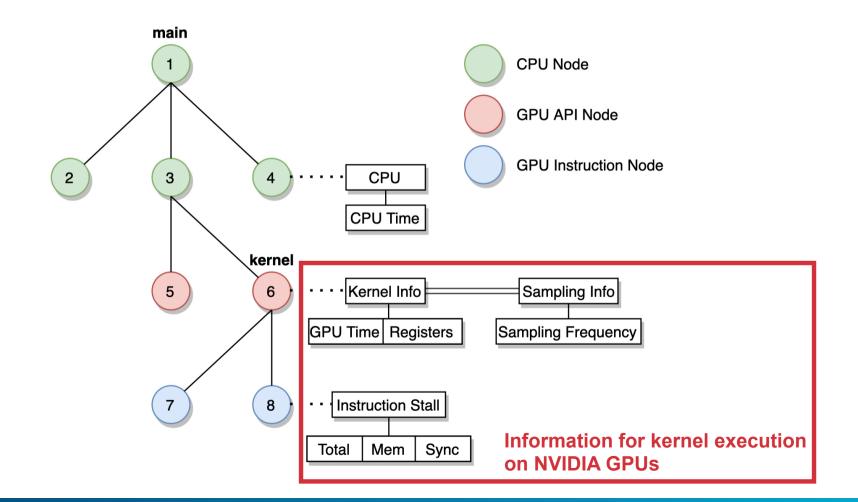


Highlights of HPCToolkit's Support for GPU-accelerated Codes

- It unwinds the CPU call stack to identify the CPU calling context for each GPU API invocation
- It employs novel data structures for fast and non-blocking inter-thread communication
- It employs binary analysis of GPU code to attribute fine-grain performance measurements to functions, inlined functions and templates, loops, and statements
 - NVIDIA, Intel, and AMD GPU binaries
- It uses a novel technique to reconstruct an approximate GPU calling context tree for computations from instruction-level measurements
- On NVIDIA GPUs: derives a rich set of metrics from PC samples from a single execution
- It performs scalable analysis of sparse representations of performance measurements and produces sparse representations tailored for graphical user interfaces

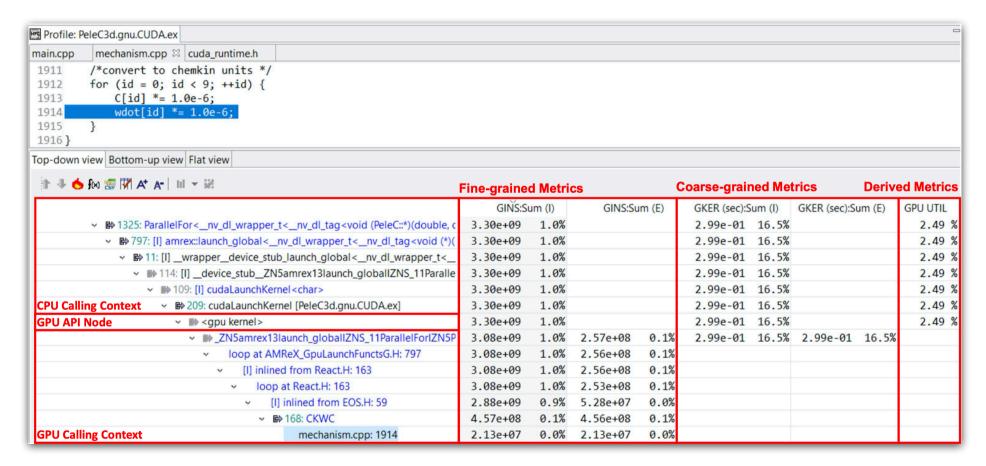


HPCToolkit's Sparse Representation of Measurements at Run-time





HPCToolkit's Code-Centric Profiles of GPU-accelerated Code



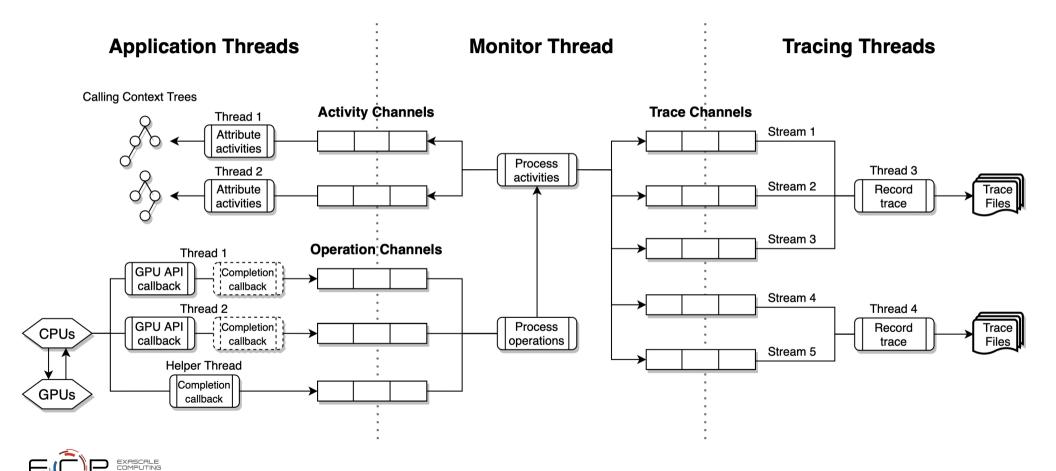


GPU Performance Measurement

- Three categories of threads
 - Application Threads (*N* per process)
 - Launch kernels, move data, and synchronize GPU calls
 - Monitor Thread (1 per process)
 - Monitor GPU events and collect GPU measurements
 - Tracing Threads (1 for every K GPU streams)
- Interactions
 - Create correlation: An application thread T creates a correlation record when it launches a
 kernel and tags the kernel with a correlation ID C, notifying the monitor thread that C belongs to T
 - Attribute measurements: The monitor thread collects measurements associated with C and communicates measurement records back to thread T
 - Record traces: The monitor thread sends activity traces to tracing threads to record in a separate trace file per GPU stream (NVIDIA, AMD) or device queue (Intel, AMD)

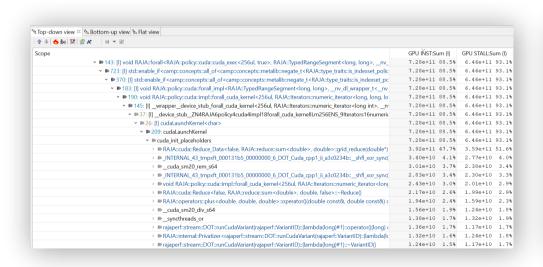


HPCToolkit's Runtime Monitoring Infrastructure for OpenCL

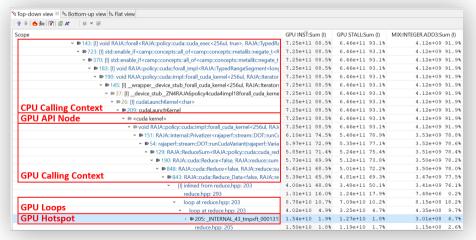


Approximation of GPU Calling Contexts to Understand Performance

- GPU code from C++ template-based programming models is complex
- NVIDIA GPUs collect flat PC samples
- Flat profiles for instantiations of complex C++ templates are inscrutable

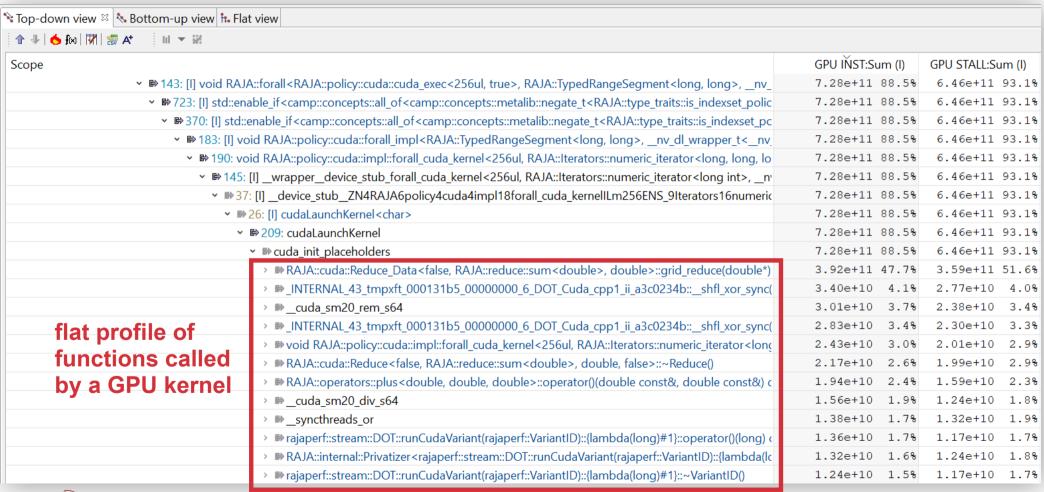


- HPCToolkit reconstructs approximate GPU calling contexts
 - Reconstruct call graph from machine code
 - Infer calls at call sites
 - PC samples of call instructions indicate calls
 - Use call counts to apportion costs to call sites
 - PC samples in a routine





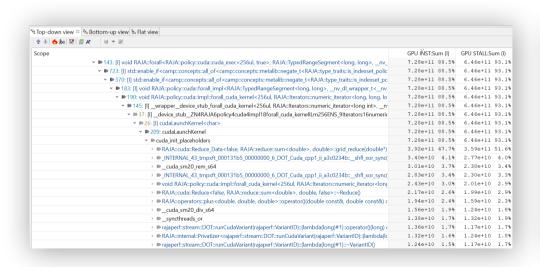
Approximation of GPU Calling Contexts to Understand Performance



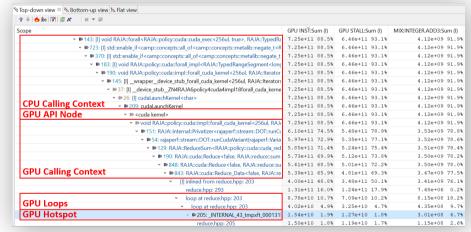


Approximation of GPU Calling Contexts to Understand Performance

- GPU code from C++ template-based programming models is complex
- NVIDIA GPUs collect flat PC samples
- Flat profiles for instantiations of complex C++ templates are inscrutable

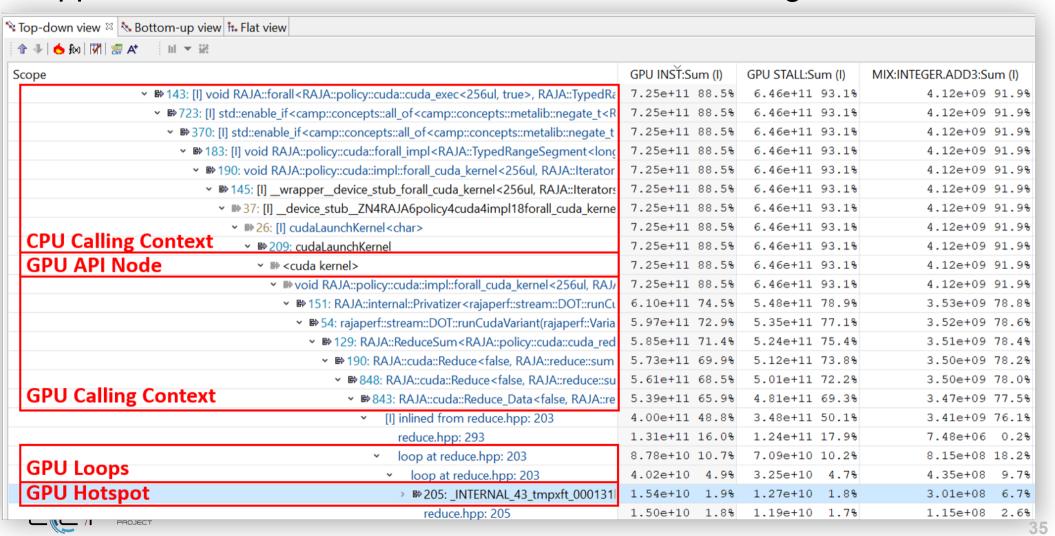


- HPCToolkit reconstructs approximate GPU calling contexts
 - Reconstruct call graph from machine code
 - Infer calls at call sites
 - PC samples of call instructions indicate calls
 - Use call counts to apportion costs to call sites
 - PC samples in a routine





Approximate Performance Attribution to GPU Calling Contexts



Reconstruction of GPU Calling Context Trees

Problem

Vendor GPU monitoring APIs don't collect call paths inside GPU kernels

Challenges

- GPU functions may be invoked from different call sites
- Need to decide how to attribute costs to each call site

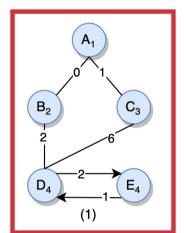
Solution

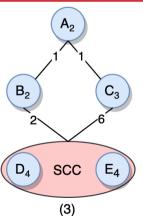
- Reconstruct GPU calling context tree from flat instruction samples and static GPU call graph

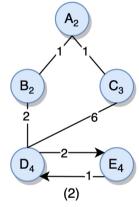


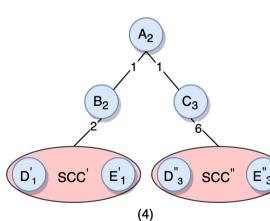
6/11/2020

- Construct a GPU static call graph based on functions and call instructions. Initialize call edge counts using counts or samples of call instructions.
- 2. For call graphs based on samples: if a function has samples and no incoming call edge has a non-zero weight, assign each of its incoming call edges a weight of 1; repeat for call edges of callers until at least one incoming call edge has samples.
- 3. Identify strongly connected components (SCCs) using Tarjan's algorithm. Rewire call graph, removing SCC internal structure and linking external calls to SCC.
- 4. Build CCT by splitting call graph. Like gprof, assume that every call to a function has equal cost. Apportion costs of each function among its call sites according to ratios of calls from each call site.



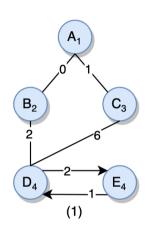


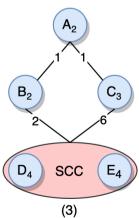


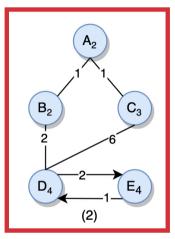


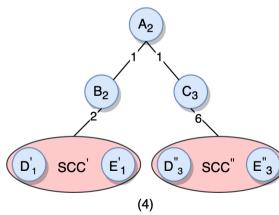


- 1. Construct a GPU static call graph based on functions and call instructions. Initialize call edge counts using counts or samples of call instructions.
- 2. For call graphs based on samples: if a function has samples and no incoming call edge has a non-zero weight, assign each of its incoming call edges a weight of 1; repeat for call edges of callers until at least one incoming call edge has samples.
- 3. Identify strongly connected components (SCCs) using Tarjan's algorithm. Rewire call graph, removing SCC internal structure and linking external calls to SCC.
- 4. Build CCT by splitting call graph. Like gprof, assume that every call to a function has equal cost. Apportion costs of each function among its call sites according to ratios of calls from each call site.



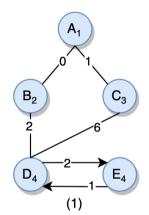


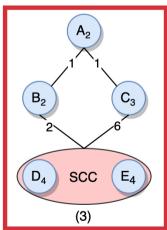


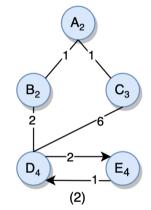


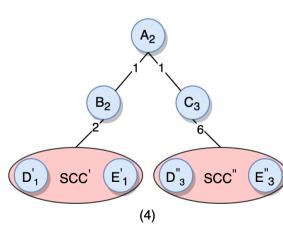


- 1. Construct a GPU static call graph based on functions and call instructions. Initialize call edge counts using counts or samples of call instructions.
- 2. For call graphs based on samples: if a function has samples and no incoming call edge has a non-zero weight, assign each of its incoming call edges a weight of 1; repeat for call edges of callers until at least one incoming call edge has samples.
- 3. Identify strongly connected components (SCCs) using Tarjan's algorithm. Rewire call graph, removing SCC internal structure and linking external calls to SCC.
- 4. Build CCT by splitting call graph. Like gprof, assume that every call to a function has equal cost. Apportion costs of each function among its call sites according to ratios of calls from each call site.



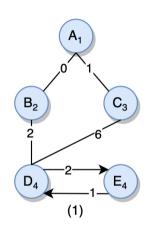


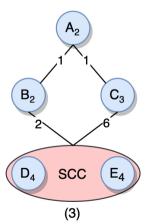


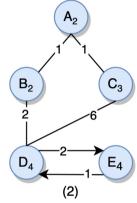


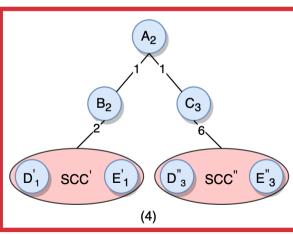


- 1. Construct a GPU static call graph based on functions and call instructions. Initialize call edge counts using counts or samples of call instructions.
- 2. For call graphs based on samples: if a function has samples and no incoming call edge has a non-zero weight, assign each of its incoming call edges a weight of 1; repeat for call edges of callers until at least one incoming call edge has samples.
- 3. Identify strongly connected components (SCCs) using Tarjan's algorithm. Rewire call graph, removing SCC internal structure and linking external calls to SCC.
- 4. Build CCT by splitting call graph. Like gprof, assume that every call to a function has equal cost. Apportion costs of each function among its call sites according to ratios of calls from each call site.







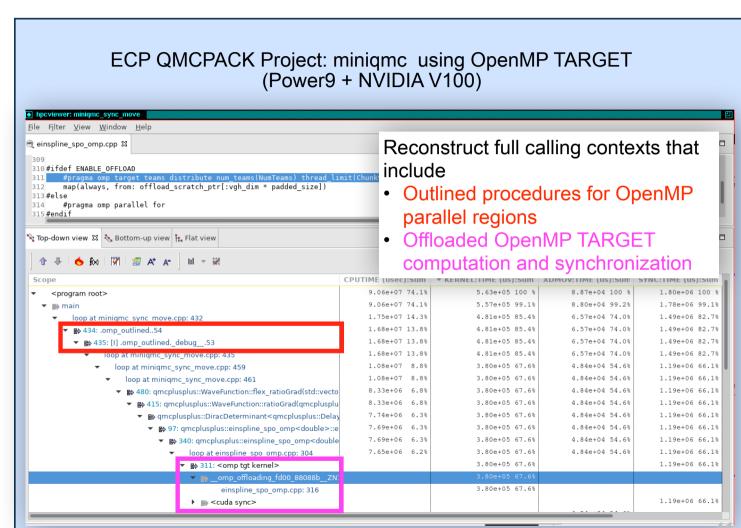




Support for OpenMP TARGET

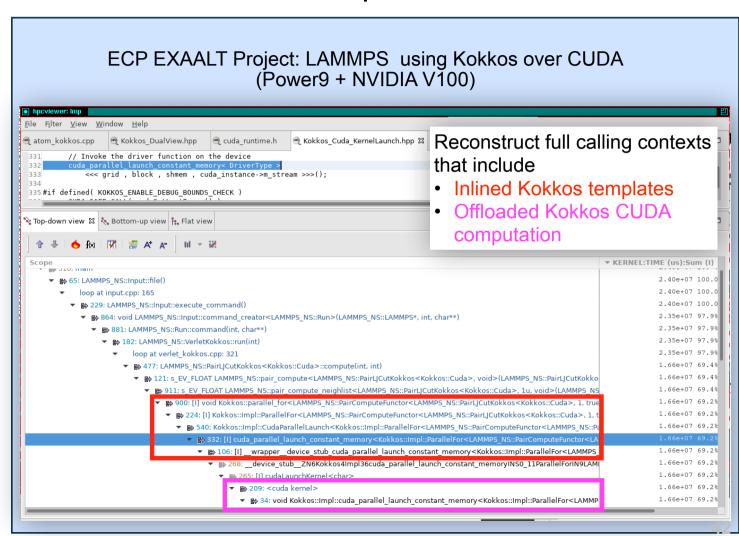
- HPCToolkit implementation of OMPT OpenMP API
 - host monitoring
 - leverages callbacks for regions, threads, tasks
 - employs OMPT API for call stack introspection
 - GPU monitoring
 - leverages callbacks for device initialization, kernel launch, data operations
 - reconstruction of userlevel calling contexts
- Leverages implementation of OMPT in LLVM OpenMP and libomptarget





Support for RAJA and and Kokkos C++ Template-based Models

- RAJA and Kokkos provide portability layers atop C++ template-based programming abstractions
- HPCToolkit employs binary analysis to recover information about procedures, inlined functions and templates, and loops
 - Enables both developers and users to understand complex template instantiation present with these models





Deriving GPU Metrics

Problem

- GPU PC sampling cannot be used in the same pass with metric collection
- Nsight-compute runs nine passes to collect multiple metrics for kernels

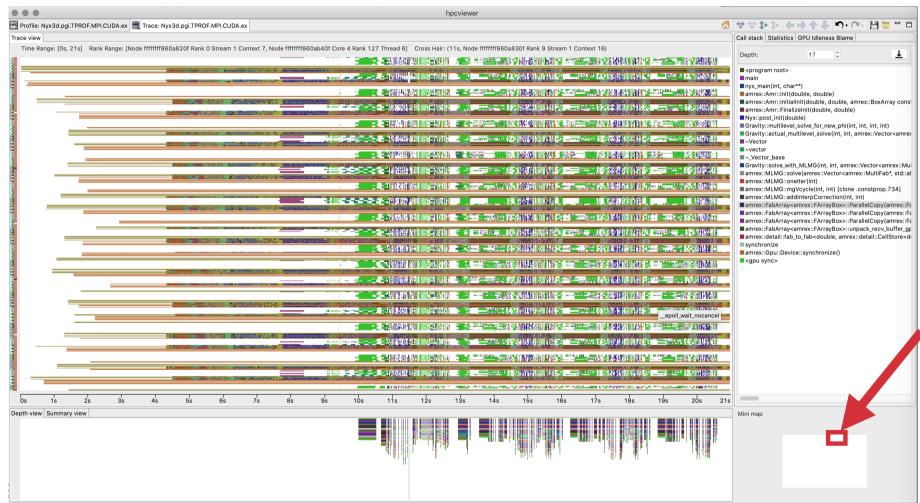
Our approach

- Measure a single pass of an execution and collect PC samples
- Derive multiple metrics using PC samples and other activity records
 - e.g., GPU SM utilization, GPU occupancy, ...



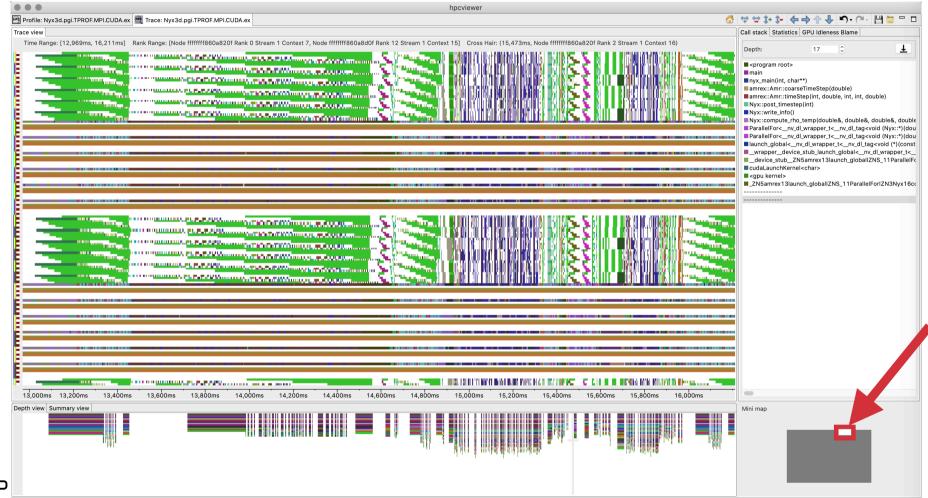
6/11/2020

Nyx with CUDA: Trace of Multi-rank Multi-GPU Executions





Nyx with CUDA: Trace of Multi-rank Multi-GPU Executions

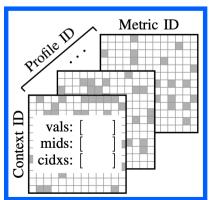


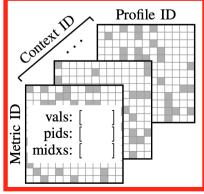


Scalable Analysis of Performance Data

- When to reduce profile data?
 - After termination: Linux perf, NVIDIA nvvp, and Paraver record detailed traces
 - At termination
 - Scalasca, Tau, Vampir use MPI to unify profile data into CUBE format
 - HPCToolkit saves separate profiles and traces per thread
- Scalable analysis of performance data using out-of-core algorithms
 - Inspect profiles and balance across ranks by aggregate size
 - Unify call stacks from all threads
 - Overlay static information on calling context trees: procedures, inline functions, loops, stmts
 - Generate computed statistics: aggregate and per profile
 - Write out two sparse outputs
 - profile-major-sparse database
 - calling-context-major-sparse database
 - Implementation: MPI + OpenMP







Is Using Sparse Formats Important?

Assess the space savings of sparse profiles

- AMD2006 CPU
 - 1 metric
 - 9 metrics, including some rare metrics
- Nyx GPU
- LAMMPS GPU
- Findings
 - as much as 21x space reduction for measurements
 - as much 337x reduction for output data

		Size (
Dataset		Dense	Sparse	Ratio
AMG2006 (1)	M D	659.0 7370.0	911.0 836.0	0.723> 8.819>
AMG2006 (9)	M D	$21.7 \\ 2290.0$	$11.1 \\ 33.5$	1.956 68.34
Nyx	M D	$5890.0 \\ 130\mathrm{GiB}$	$278.0 \\ 601.0$	$21.14 \Rightarrow 221.5 \Rightarrow$
LAMMPS	M D	$85.5 \\ 8250.0$	$5.23 \\ 24.5$	16.35 × 336.9 ×



Scalable Analysis of Performance Data: 64K profiles of AMG2006

Input

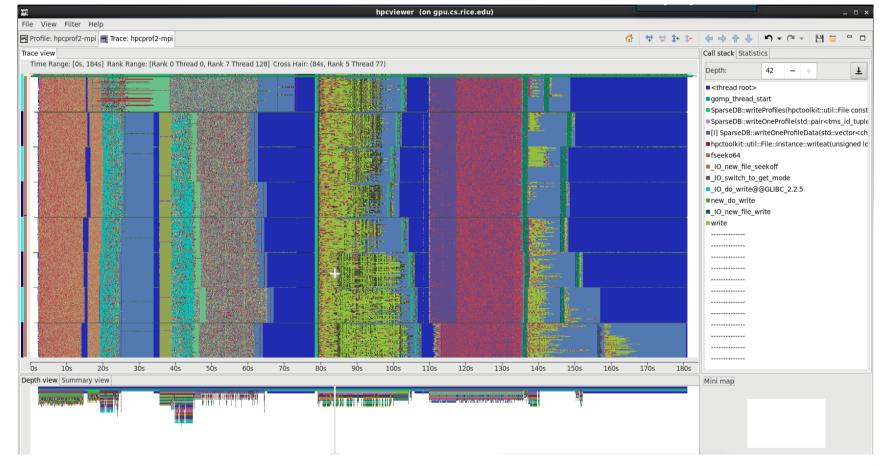
- 5GB profiles
- 225GB traces

Analysis

- 8 KNL nodes
- 1 rank / node
- 128T / rank

Execution time

• 184s





Outline

- Performance measurement and analysis challenges for GPU-accelerated supercomputers
- Introduction to HPCToolkit performance tools
 - Overview of HPCToolkit components and their workflow
 - HPCToolkit's graphical user interfaces
- Analyzing the performance of GPU-accelerated supercomputers with HPCToolkit
 - Overview of HPCToolkit's GPU performance measurement capabilities
 - Collecting measurements
 - Analysis and attribution
 - Scalable analysis of performance data
- Status, ongoing work, final remarks



Status for Various GPUs

Vendor	Coarse-grain measurement	Fine-grain measurement	Tracing	Binary analysis: loops, inlined code
NVIDIA	CUPTI	PC sampling	CUPTI	nvdisasm + Dyninst
Intel	OpenCL and Level 0	GTPin instrumentation	OpenCL callbacks	IGA + Dyninst
AMD	Roctracer	emerging Dyninst instrumenter	Roctracer	emerging Dyninst decoder



Detailed Performance Analysis Requires Support at Many Levels

		•
	Hardware and Software Stack Components	Partners
•	Hardware must include support for fine-grain measurement and attribution	GPU vendors
	performance counters are not enough; NVIDIA's PC sampling approximates our needs	
•	System software must provide appropriate interfaces for introspection and analysis	
	 e.g. Linux perf_events supports sample-based performance monitoring even in the kernel 	Red Hat
	e.g. dynamic loader (ld.so) provides LD_AUDIT interface for monitoring and control of dynamic library operations	
	 elfutils must support NVIDIA's extended line maps in CUDA 11.2+ GPU binaries 	
	GPU vendor software stacks (kernel driver, runtime, tools API)	GPU vendors
•	Compiler must compute high-quality DWARF information	Vendors and LLVM
	 associate each machine instruction with full call chains involving inlined templates and functions 	community
•	Runtime must maintain information needed to map computations to a source-level view	OpenMP Language
	OpenMP's OMPT helps bridge the vast gap between the implementation and user-level view	Committee and LLVM Community
•	Performance tools must gather measurements using multiple modalities and map them to source	
	precise attribution when possible	
	 reconstruct approximate attribution when precise attribution is unavailable 	Wisconsin's Dyninst Project
	GPU calling context	
	loops in CPU and GPU code	
	attribute inefficiencies from where they are observed back to their causes	
	/I	

Ongoing Work

Interface

- Emerging GPU Performance Advisor tool for NVIDIA GPUs
 - attributes instruction stalls with backward slicing, analyzes code, offers advice about most promising improvements
- Integrated user interface that supports both profiles and traces
 - Automated serialization analysis of CPU and GPU traces

Internals

- Collecting GPU hardware counters, which will support Roofline analysis
- Updating measurement and analysis support for NVIDIA GPUs (emerging CUPTI, more info about inlining)
- Extending HPCToolkit to support analysis of machine learning frameworks: Pytorch, Tensorflow
- · Improving scalability of measurement and analysis
- Developing instrumentation to assess performance on Intel GPUs
- Refining implementation of monitoring for Intel's Level 0
- Improving binary analysis of AMD GPU binaries



Final Remarks

- Nice to work with national labs and have early involvement in big procurements
 - Amplifies our ability to affect vendor hardware and software in the near term
- Software development challenges are myriad
 - Developing tools for three GPU software stacks at the same time is ridiculous
 - Building capabilities ahead of current vendor hardware and software
 - AMD and Intel software is a work in progress
 - instability and API-breaking changes are common
 - Relying on vendor closed-source components is a challenge
 - standards specify only an API, but internals matter for tools that see all
 - undocumented behaviors about things that matter
 - missing capabilities, e.g. need excellent DWARF mappings for optimized GPU code
 - NVIDIA serializes kernels to facilitate measurement with PC sampling

