# Measurement and Analysis of GPU-accelerated Applications with HPCToolkit

Jonathon Anderson, Yuning Xia, John Mellor-Crummey
Rice University

Consortium for the Advancement of Scientific Software
Birds of a Feather Days
February 11, 2026

# DOE's GPU-Accelerated Exascale Platforms



- Frontier compute nodes (OLCF)
  - 1 AMD EPYC "Trento" CPU
  - 4 MI250X AMD Radeon Instinct GPUs
  - 4 Slingshot 11 endpoints
  - Unified memory architecture

- Aurora compute nodes (ALCF)
  - 2 Intel Xeon "Sapphire Rapids" processors
  - 6 Intel Data Center GPU Max 1500
  - 8 Slingshot 11 endpoints
  - Unified memory architecture

- El Capitan compute nodes (LLNL)
  - 4 AMD MI300A APU
  - 4 Slingshot 11 endpoints
  - Unified memory architecture

# Tuning HPC Applications



**Profile the Application**

Collect detailed performance data on an application execution

**Analyze Performance Data**

Aggregate performance data from all application threads and generate a database of "analysis results"
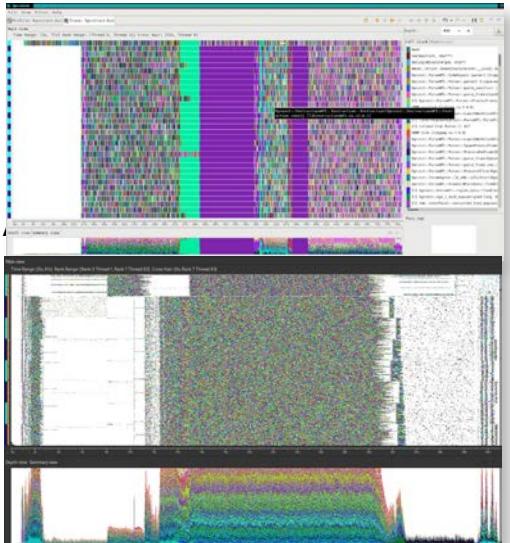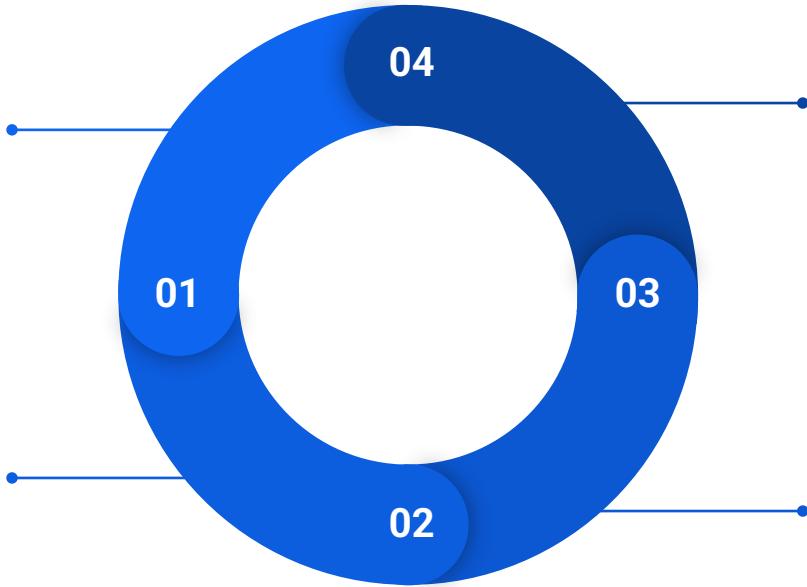
**04**

**01**

**03**

**02**

**Optimize and Improve**

Remove performance bottlenecks and further improve the application

**Inspect and Identify**

Inspect locations with poor performance and identify root causes

# Tuning HPC Applications

# Why HPCToolkit?

- **Widely applicable**: many parallel programming models within & across nodes

- **Easy**: profiles unmodified application binaries

- **Fast:** low-overhead measurement; parallel analysis of large performance data

- **Scalable**: measure and analyze GPU-accelerated executions at large scale

- **Informative**: learn where an application spends its time and why
  - call path profiles associate metrics with application source code contexts
  - hierarchical traces reveal execution dynamics
  - LEO identifies and quantifies root causes of GPU performance losses

- **Broad audience**: developers of applications, frameworks, runtimes & tools

- **Multiplatform**: unlike vendor tools, works with a wide range of CPUs and GPUs
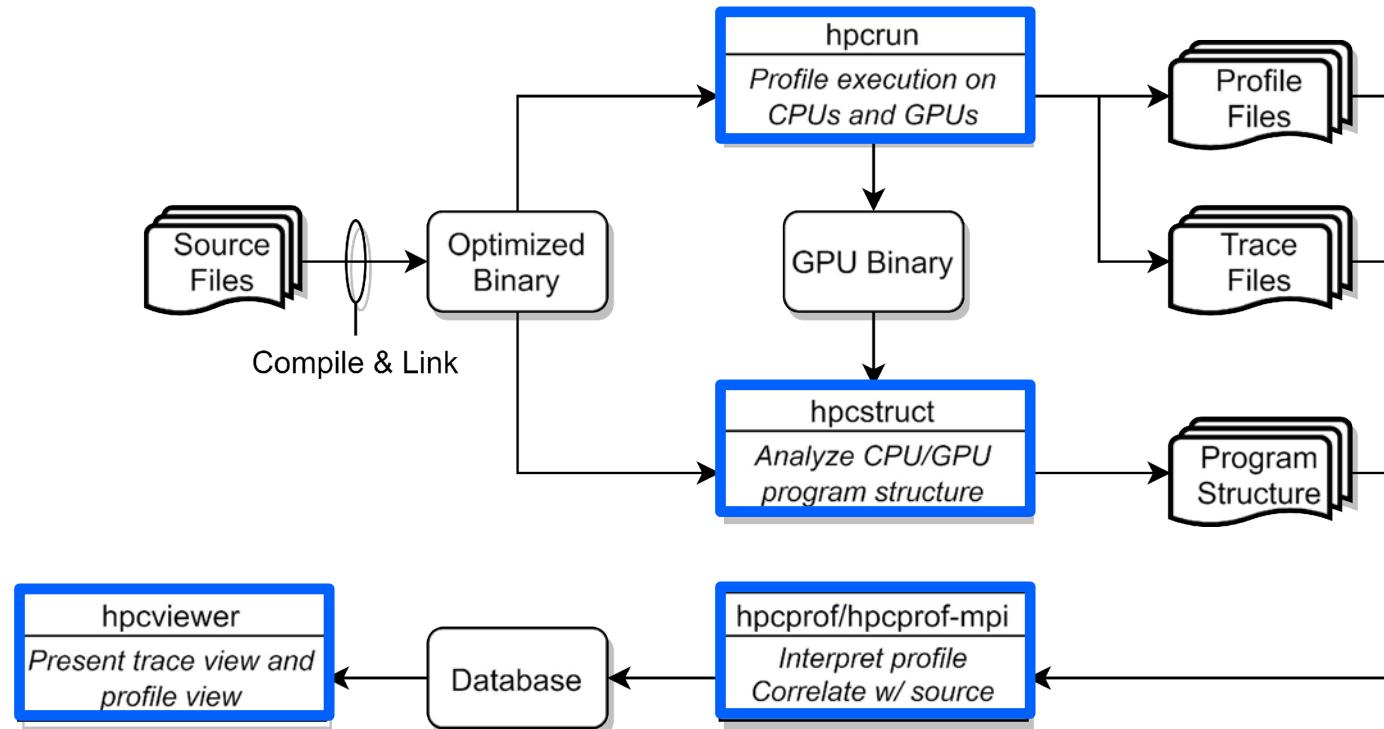
# How is HPCToolkit Different from Vendor Tools?

- More scalable tracing than vendor tools
  - — measure exascale executions across many nodes and GPUs
  - — GUI supports interactive exploration of TB of performance data
- Scalable, parallel post-mortem analysis vs. non-scalable in-GUI analysis
- Detailed reconstruction of calling context profiles within GPU kernels
- Identifies and quantifies root causes of GPU performance losses

# Today's Agenda

- Introduce HPCToolkit tools and workflow

- Illustrate HPCToolkit's use with some case studies

- Live demos

  - profiling and tracing

  - instruction-based performance metrics

- Instruction-level performance monitoring

  - explaining instruction-level performance

  - improving analysis and attribution of PC samples

  - automating analysis of GPU bottlenecks

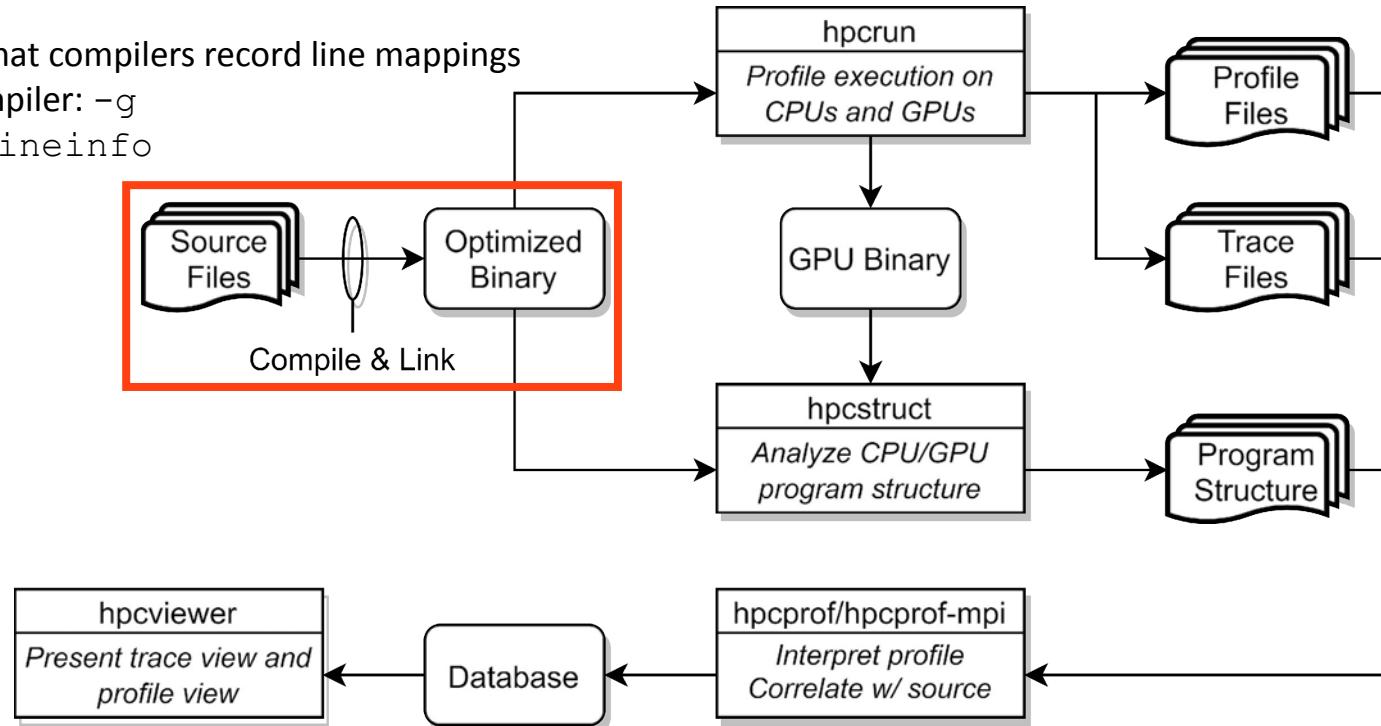- Discussion about needs, problems, and suggestions

# HPCToolkit's Workflow for GPU-accelerated Applications

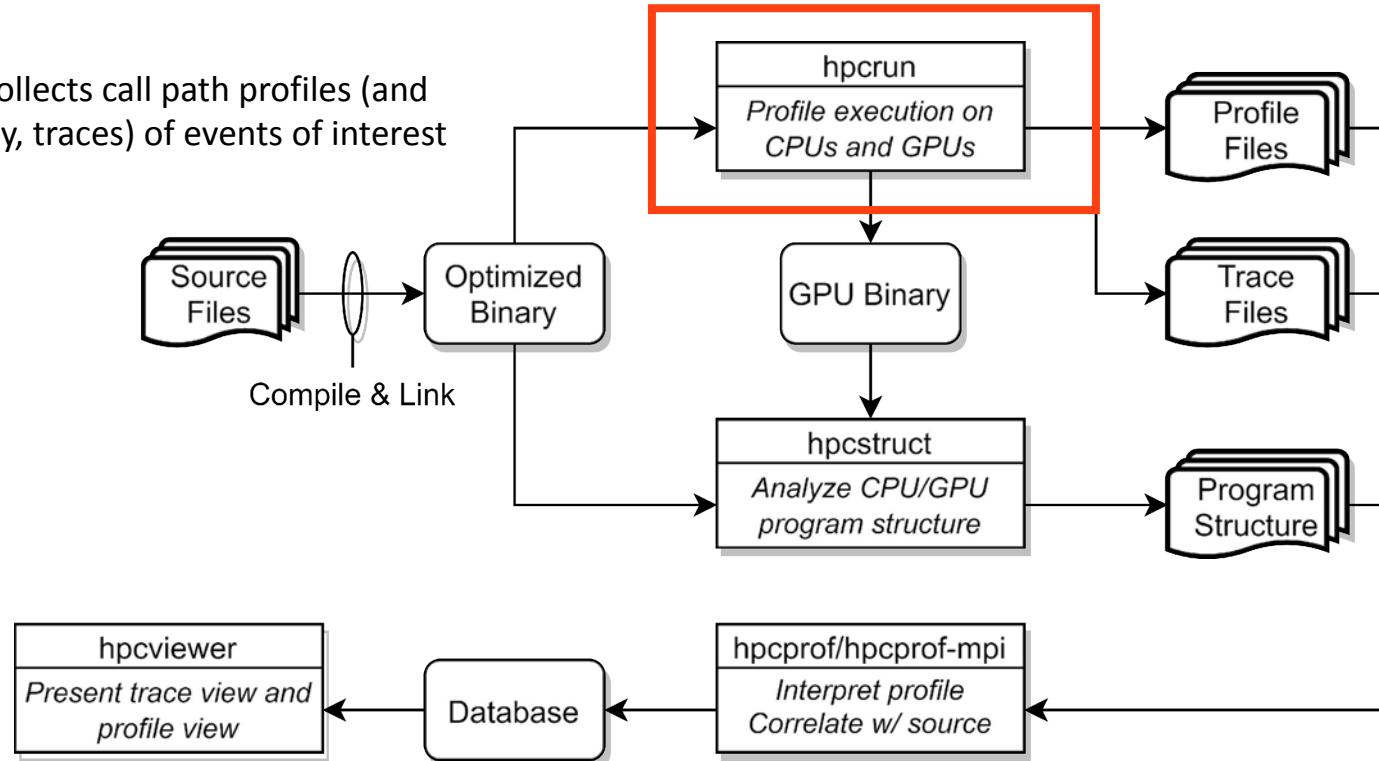# HPCToolkit's Workflow for GPU-accelerated Applications

Step 1:
- Ensure that compilers record line mappings
- host compiler: $-g$
- nvcc: $-lineinfo$

# HPCToolkit's Workflow for GPU-accelerated Applications

Step 2:
- *hpcrun* collects call path profiles (and optionally, traces) of events of interest
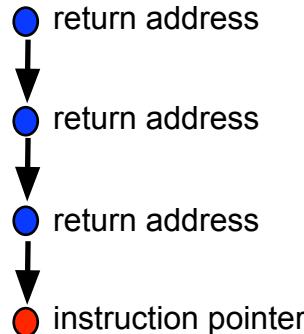
# Measurement of CPU and GPU-accelerated Applications

- Sampling using Linux timers and hardware counters on the CPU
- Callbacks when GPU operations are launched
- Event stream or callbacks for GPU operation completion
- PC Samples: AMD, NVIDIA, Intel
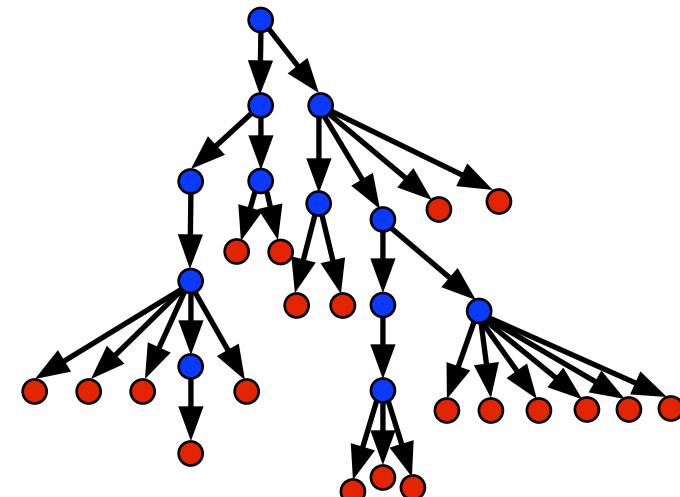- Binary instrumentation of GPU kernels: Intel

# Call Stack Unwinding to Attribute Costs in Context

- Unwind when timer or hardware counter overflows

  —measurement overhead $\propto$ sampling frequency rather than call frequency

- Unwind to capture context for events such as GPU kernel launches

Call path sample

Calling context tree

- return address

- return address

- return address

- instruction pointer

# hpcrun: Measure CPU and/or GPU activity
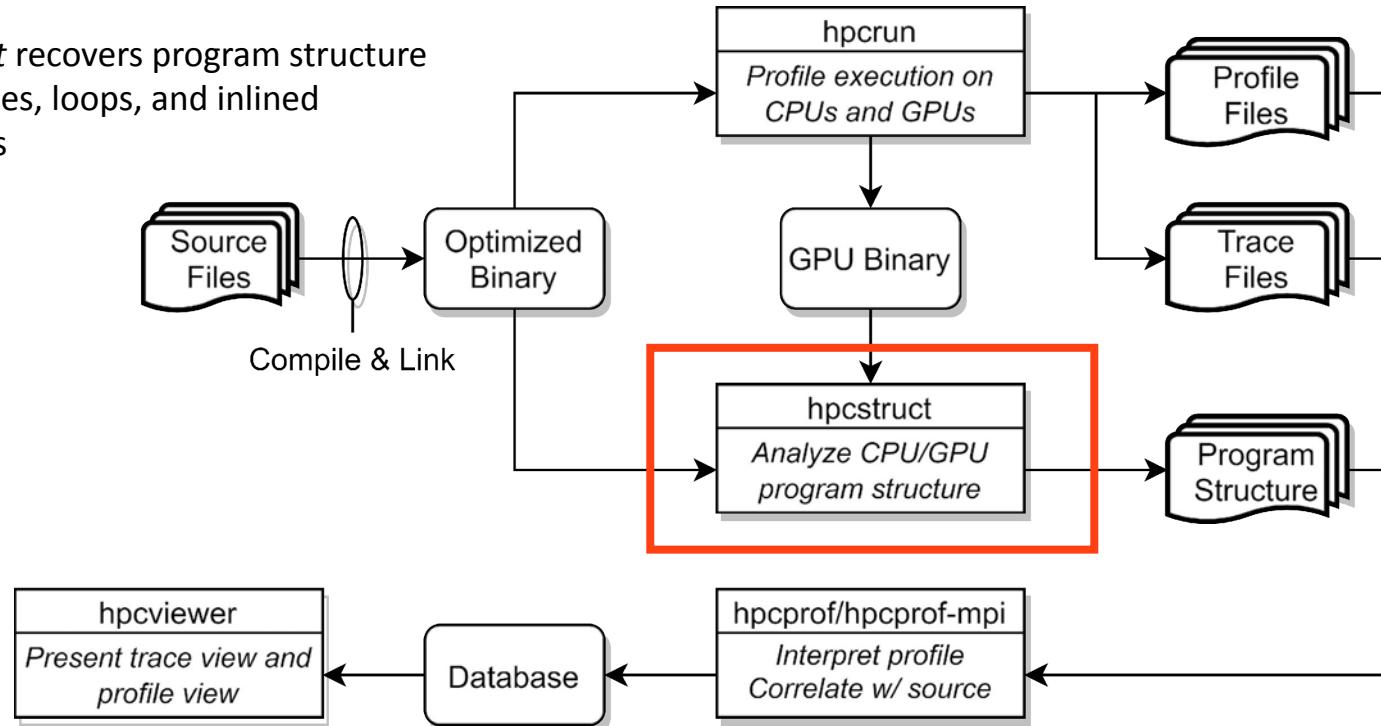
- CPU and GPU profiling
  - `hpcrun -e CPUTIME -e gpu=XXX <app> …`

- CPU and GPU profiling and tracing
  - `hpcrun -e CPUTIME -e gpu=xxx -tt <app>`

- GPU PC sampling
  - `hpcrun -e gpu=yyy,pc <app>`

- Measuring MPI programs
  - `srun -n <ranks> … hpcrun -e CPUTIME -e gpu=XXX <app>`

$xxx \in \{cuda, rocm, opencl, level0\}$
$yyy \in \{cuda, rocm, level0\}$

# HPCToolkit's Workflow for GPU-accelerated Applications

Step 3:
- *hpcstruct* recovers program structure about lines, loops, and inlined functions

# hpcstruct: Analyze CPU and GPU Binaries Using Multiple Threads

- Usage

  ```
  hpcstruct [--gpucfg yes] <measurement-directory>
  ```
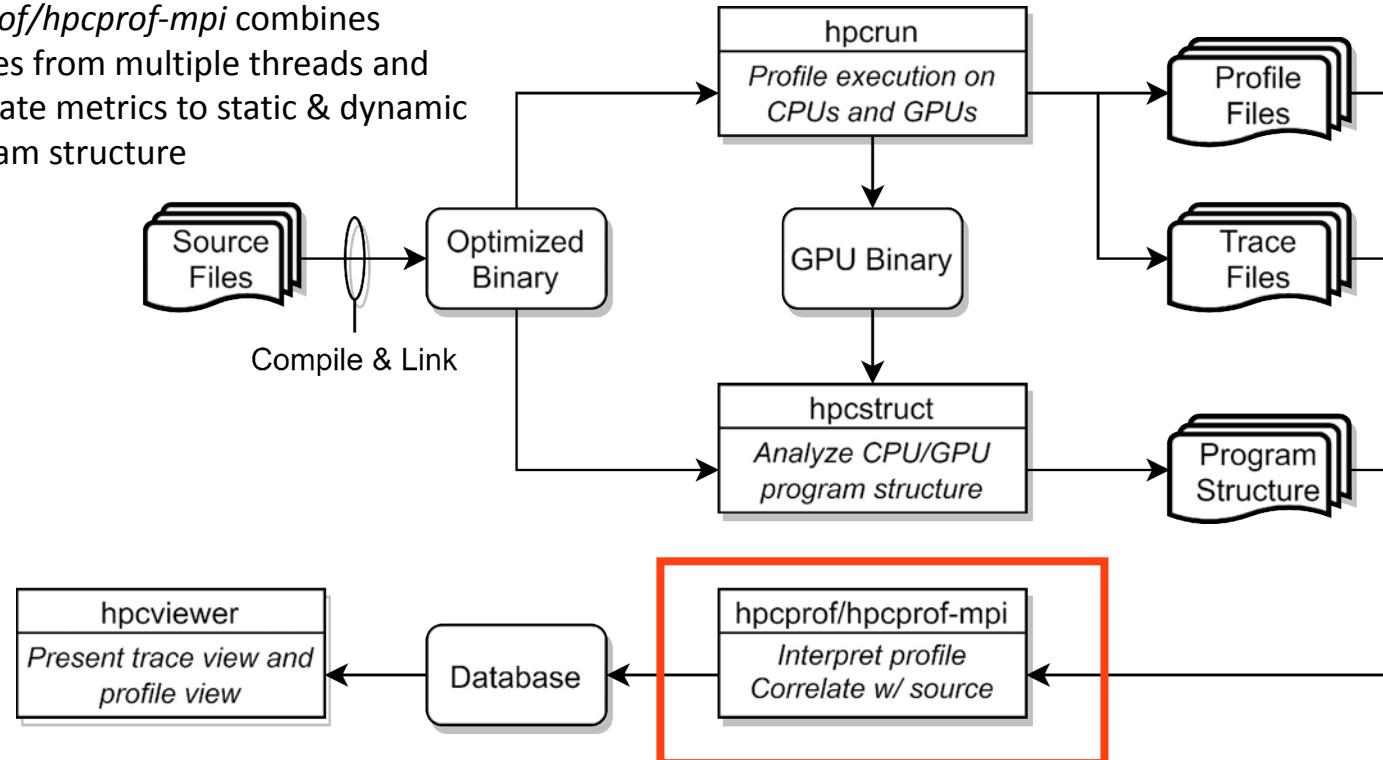
- What it does
    - Recover program structure information
        - Files, functions, inlined templates or functions, loops, source lines
    - In parallel, analyze all CPU and GPU binaries that were measured by HPCToolkit
        — typically analyze large application binaries with 16 threads
        — typically analyze multiple small application binaries concurrently with 2 threads each
    - Cache binary analysis results for reuse when analyzing other executions

  NOTE: **--gpucfg yes** needed only for detailed analysis of PC samples

# HPCToolkit's Workflow for GPU-accelerated Applications

Step 4:

- *hpcprof/hpcprof-mpi* combines profiles from multiple threads and correlate metrics to static & dynamic program structure

# hpcprof/hpcprof-mpi: Associate Measurements with Program Structure

- Analyze data from modest-scale executions with multithreading
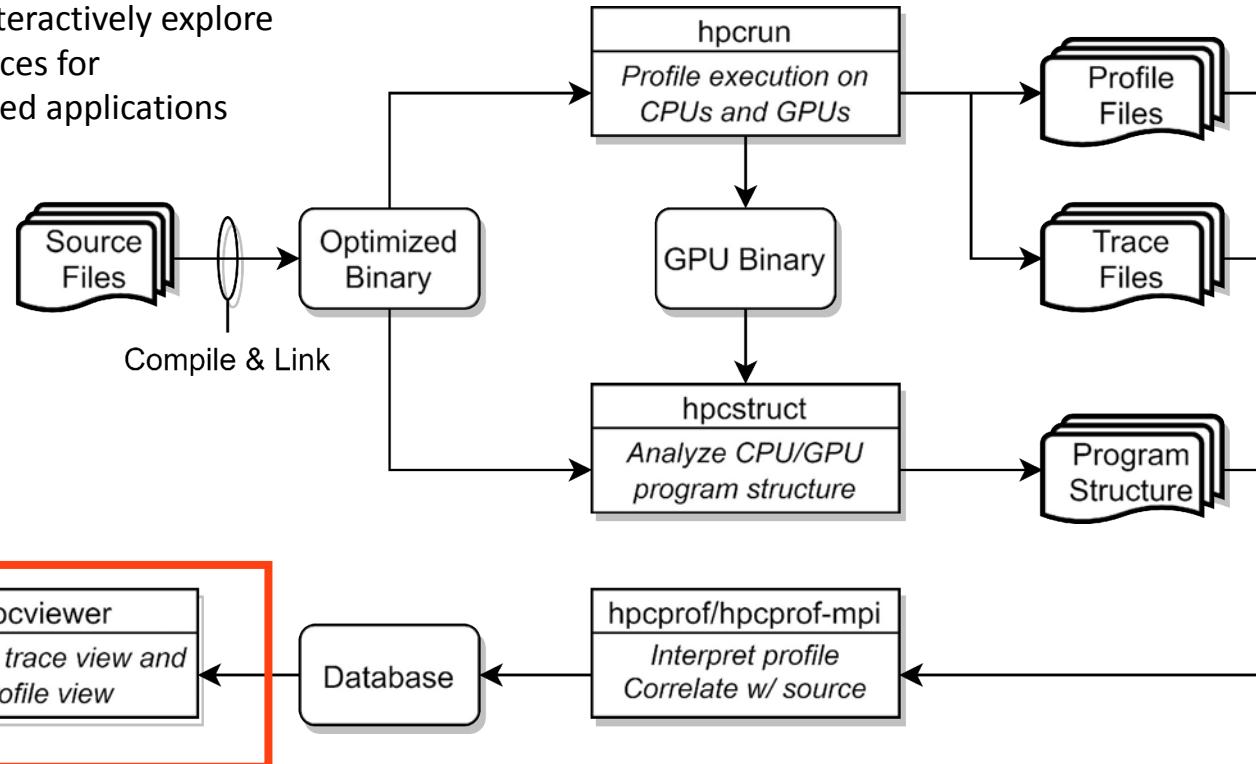
```
hpcprof <measurement-directory>
```

- Analyze data from large-scale executions with distributed-memory parallelism + multithreading

```
srun -n ${NODES} --ppn 1 —depth=128 \
  hpcprof-mpi <measurement-directory>
```

# HPCToolkit's Workflow for GPU-accelerated Applications

Step 4:

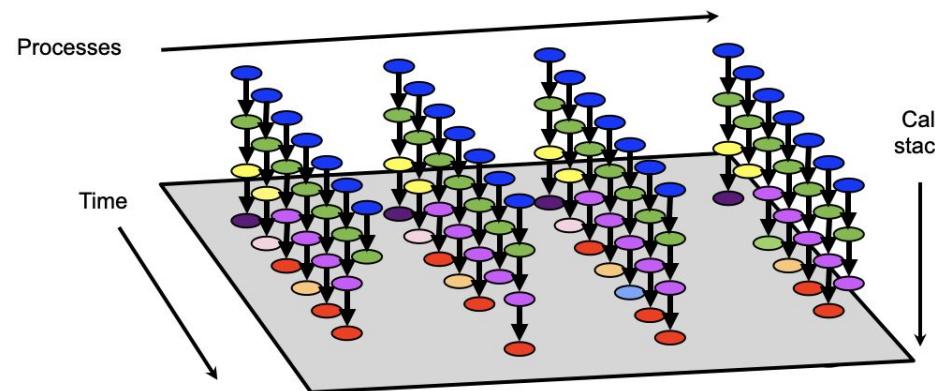- *hpcviewer* - interactively explore profile and traces for GPU-accelerated applications

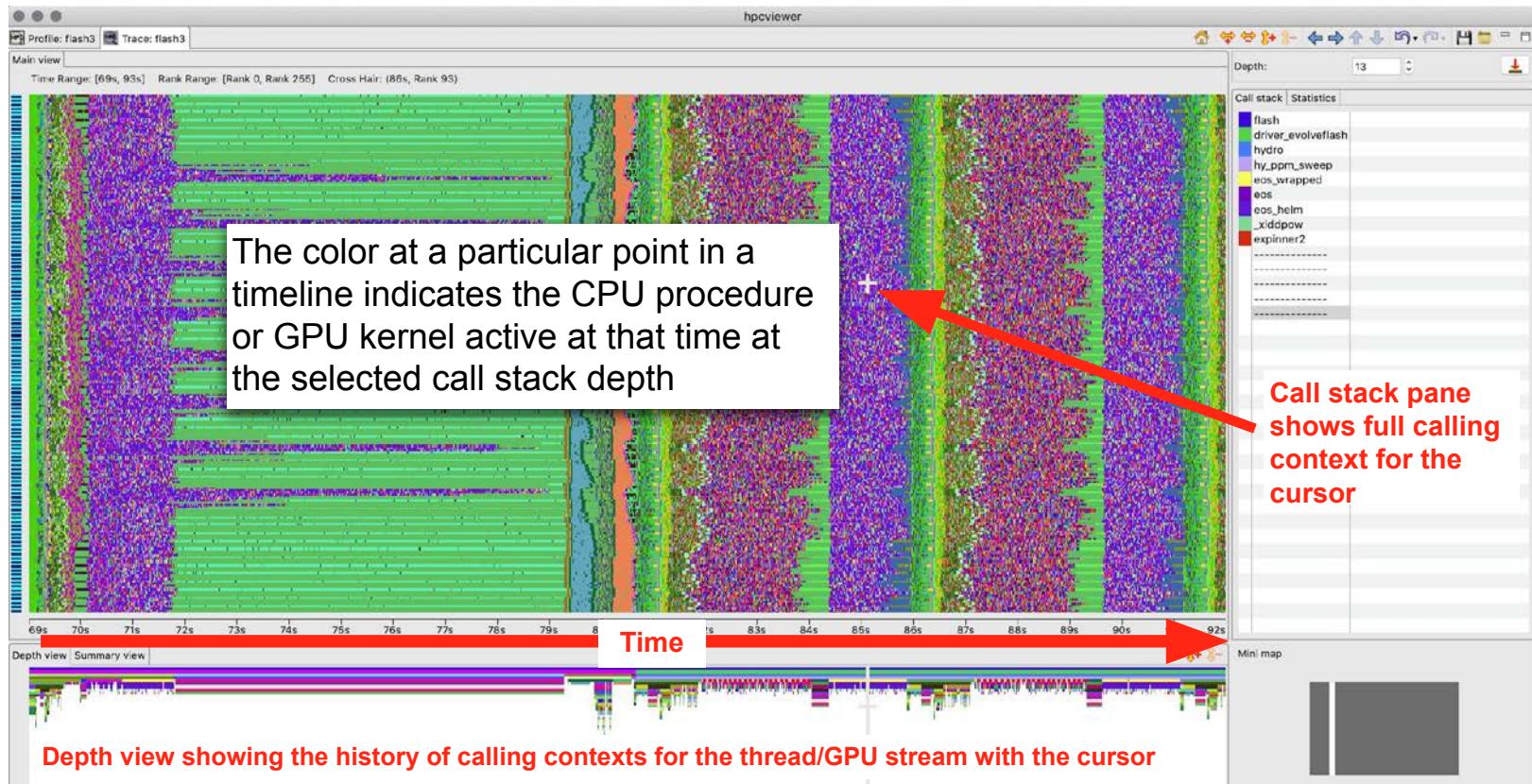# Code-centric Analysis with hpcviewer

# 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

# Understanding hpcviewer's Trace View



**MPI ranks, OpenMP Threads, GPU streams**

The color at a particular point in a timeline indicates the CPU procedure or GPU kernel active at that time at the selected call stack depth

**Call stack pane shows full calling context for the cursor**

**Time**

**Depth view showing the history of calling contexts for the thread/GPU stream with the cursor**

A multi-level call stack based view of execution over time

**Minimap indicates part of execution trace shown**

# Today's Agenda

- Introduce HPCToolkit tools and workflow
- Illustrate HPCToolkit's use with some case studies
- Live demos
  - profiling and tracing
  - instruction-based performance metrics
- Instruction-level performance monitoring
  - explaining instruction-level performance
  - improving analysis and attribution of PC samples
  - automating analysis of GPU bottlenecks
- Discussion about needs, problems, and suggestions

# Case Studies

- ExaWind (Nalu-Wind + AMRWind)  - Wind turbine and wind farm simulation

- PeleLMeX - Adaptive mesh hydrodynamics code for low mach number reacting flows

- GAMESS  (OpenMP) - general ab initio quantum chemistry package

- LAMMPS (Kokkos) - classical molecular dynamics code with a focus on materials modeling

At Exascale!

# ExaWind: Modeling Turbine Wake Formation



Figure credit: Jon Rood, NREL

# ExaWind: Wakes from Three Turbines over Time



Figure credit: Jon Rood, NREL

# ExaWind: Visualization of a Wind Farm Simulation



Figure credit: Jon Rood, NREL

# ExaWind: Execution Traces on Frontier Collected with HPCToolkit

Traces on roughly 64K MPI ranks + 8K GPUs for ~17minutes

Before: MPI waiting (bad), shown in red          After: MPI overhead negligible*



Figure credits: Jon Rood, NREL          *replaced non-blocking send/recv with ialltoallv

# ExaWind Testimonials for HPCToolkit

*I just wanted to mention we've been using HPCToolkit a lot for our ExaWind application on Frontier, which is a hugely complicated code, and your profiler is one of the only ones we've found that really lets us easily instrument and then browse what our application is doing at runtime including GPUs. As an example, during a recent hackathon we had, we improved our large scale performance by 24x by understanding our code better with HPCToolkit and running it on 1000s of nodes while profiling. We also recently improved upon this by 10% for our total runtime.*

*- Jon Rood NREL (5/31/2024)*

*One big thing for us is that we can't overstate how complicated ExaWind is in general, and how complicated it is to build, so finding out that HPCToolkit could easily profile our entire application without a ton of instrumentation during the build process, and be able to profile it on a huge amount of Frontier with line numbers and visualizing the trace was really amazing to us.*

*- Jon Rood NREL (6/3/2024)*

# Full SAF Case: Load Balancing Issues

Confirmed on Frontier production runs: 60% speedup

## Load balancing ON

Grid from checkpoint file before regridding, Avg. Time/dt = 92.3 s

After regridding, Avg. Time/dt = 62.5s
- Down to 40s when setting amr.max_grid_size=32

Load balancing OFF
(pink-> Cvode calls)

Figure credit: PeleLMEX Team, NREL Hackathon, February 2025

# LAMMPS on Frontier: Executions with Kernel Duration of Milliseconds



8K nodes, 64K MPI ranks + 64K GPU tiles

# LAMMPS on Frontier: Executions with Kernel Duration of Milliseconds

# LAMMPS on Frontier: Executions with Kernel Duration of Milliseconds

# LAMMPS on Frontier: 8K nodes, 64K MPI ranks + 64K GPU tiles

## Kernel duration of microseconds

# Case Study: GAMESS

- General Atomic and Molecular Electronic Structure System (GAMESS)

  — general *ab initio* quantum chemistry package

- Calculates the energies, structures, and properties of a wide range of chemical systems


- Experiments

  - GPU-accelerated nodes at a prior Perlmutter hackathon

    - Single node with 4 GPUs

    - Five nodes with 20 GPUs

**Perlmutter node at a glance**
AMD Milan CPU
4 NVIDIA A100 GPUs
256 GB memory

# Time-centric Analysis: GAMESS 4 ranks, 4 GPUs on Perlmutter



GAMESS original        All CPU threads and GPU streams

# Time-centric Analysis: GAMESS 4 ranks, 4 GPUs on Perlmutter



GAMESS original

All CPU threads and GPU streams

# Time-centric Analysis: GAMESS 4 ranks, 4 GPUs on Perlmutter



GAMESS original        All GPU streams; whole execution

# Time-centric Analysis: GAMESS 4 ranks, 4 GPUs on Perlmutter



GPU load imbalance due to triangular iteration spaces

GAMESS original          GPU streams: 1 iteration

GAMESS original

CPU Threads and GPU Streams

# Time-centric Analysis: GAMESS 5 nodes, 40 ranks, 20 GPUs on Perlmutter

GAMESS improved with better manual distribution of work in input

# Time-centric Analysis: GAMESS 5 nodes, 40 ranks, 20 GPUs on Perlmutter



GAMESS improved adding Rank 0 Thread 0 to GPU streams

# Time-centric Analysis: GAMESS 5 nodes, 40 ranks, 20 GPUs on Perlmutter



1 CPU Stream, 2 GPU Streams: 6 Iterations

# Time-centric Analysis: GAMESS 5 nodes, 40 ranks, 20 GPUs on Perlmutter

# Time-centric Analysis: GAMESS 5 nodes, 40 ranks, 20 GPUs on Perlmutter

# Time-centric Analysis: GAMESS 5 nodes, 40 ranks, 20 GPUs on Perlmutter

```
1096 C
1097         IJ=1-INC
1098         DO 150 I=2,NA
1099             IJ=IJ+INC
1100             IM1=I-1
1101             DO 140 J=1,IM1
1102                 IJ=IJ+INC
1103                 AIJ=A(IJ)
1104                 IF(AIJ.EQ.ZERO) GO TO 140
1105                     CALL DAXPY(MB,AIJ,B(I,1),NA,AB(J,1),NAB)
1106                     CALL DAXPY(MB,AIJ,B(J,1),NA,AB(I,1),NAB)
1107     140     CONTINUE
1108     150 CONTINUE
1109         RETURN
1110         END
```

Top-down view | Bottom-up view | Flat view
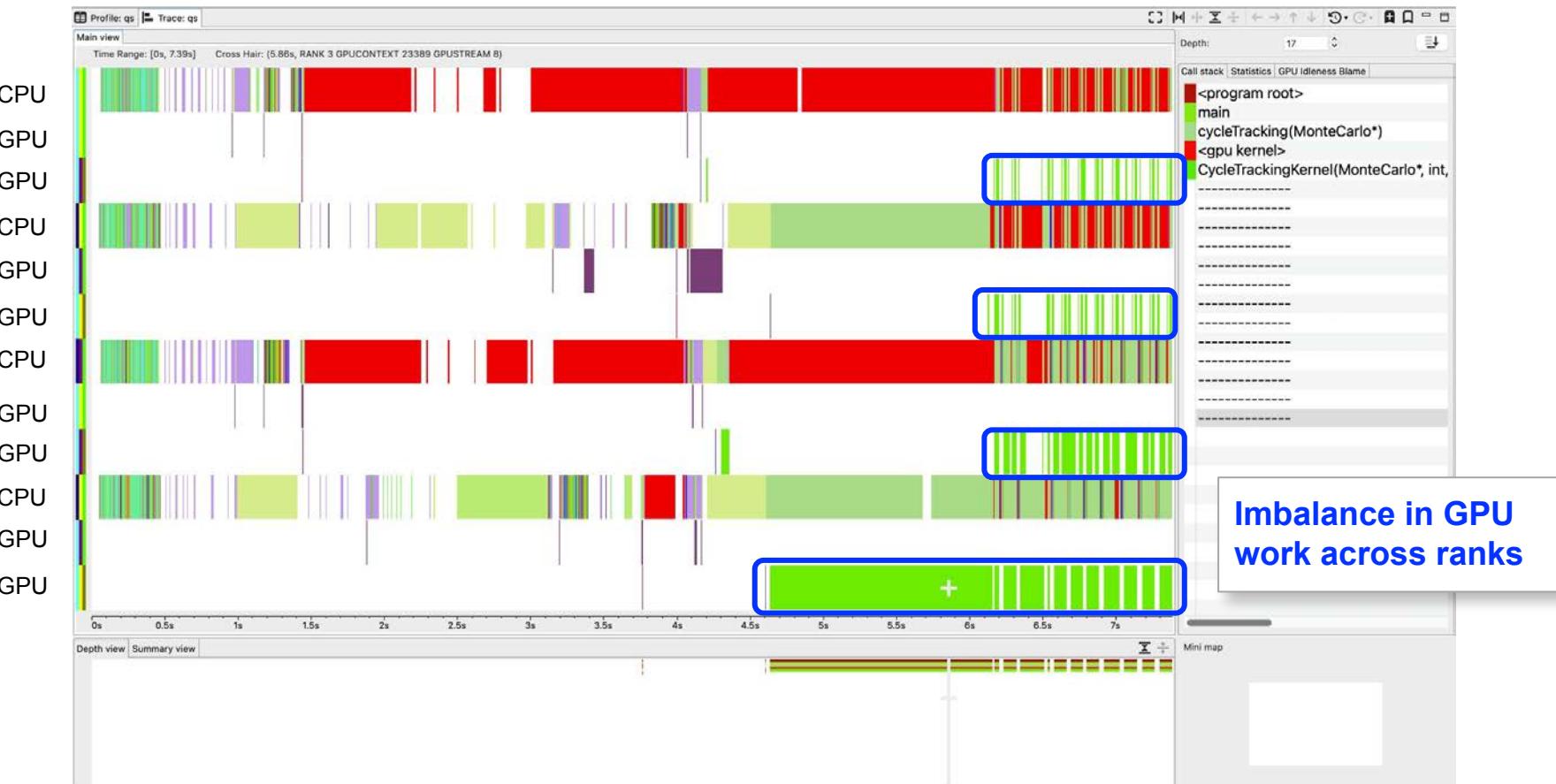
RICE

nne Leadership
outing Facility

# Today's Agenda

- Introduce HPCToolkit tools and workflow
- Illustrate HPCToolkit's use with some case studies
- Live demos
  - profiling and tracing
  - instruction-based performance metrics
- Instruction-level performance monitoring
  - explaining instruction-level performance
  - improving analysis and attribution of PC samples
  - automating analysis of GPU bottlenecks
- Discussion about needs, problems, and suggestions

# Case Study: Quicksilver

- Proxy application that represents some elements of LLNL's Mercury code

- Solves a simplified dynamic Monte Carlo particle transport problem

  - Attempts to replicate memory access patterns, communication patterns, and branching or divergence of Mercury for problems using multigroup cross sections

- Parallelization: MPI, OpenMP, CUDA, and HIP

- Performance Issues

  - load imbalance (for canned example)

  - latency bound table look-ups

  - a highly branchy/divergent code path

  - poor vectorization potential

# Quicksilver: Trace view

# Quicksilver: Detailed analysis within a Kernel using PC Sampling

# Quicksilver: Attribution to Code within a Kernel

# Today's Agenda

- Introduce HPCToolkit tools and workflow
- Illustrate HPCToolkit's use with some case studies
- Live demos
  - profiling and tracing
  - instruction-based performance metrics
- Instruction-level performance monitoring
  - explaining instruction-level performance
  - improving analysis and attribution of PC samples
  - automating analysis of GPU bottlenecks
- Discussion about needs, problems, and suggestions

# Work in Progress: Present GPU Metrics using a Donut Graph

- Figure shows forthcoming top-down display of CPU metrics
- Plan similar top-down display of GPU metrics
  - Issues
  - Exposed stalls
    - Memory
    - Pipeline
    - Ifetch
    - …
  - Hidden stalls

# WIP: Enabling Instruction-level Metrics for Complex Kernels

- Calling context is great for visualization, but slow to analyze for complex GPU kernels!

  - Every plausible calling context must be expanded to attribute performance

  - Example: Quicksilver has one kernel with 145 separate GPU functions
    - Call graph below, **red** nodes have multiple plausible calling contexts

  - Sometimes >100,000 calling contexts for a single function, e.g. cuda_div

# WIP: Enabling Instruction-level Metrics for Complex Kernels

- Next generation: Graph-based calling contexts
  - Performance is attributed to nodes
  - Distribution factors recorded on edges

- Efficient representation for GPU code
  - Avoid reconstruction of duplicate contexts
  - Match GPU measurement capabilities

- Clean model for both CPU and GPU performance
  - Not a hybrid model, one graph handles all

# WIP: Enabling Instruction-level Metrics for Complex Kernels

- Case study: Quicksilver

  - 409 unique calling contexts

  - 145 graph nodes (1 per GPU function)

- Preliminary results

  - 48% estimated size reduction

  - 1.3x potential speed up

# WIP: Leo - Platform-independent analysis with PC Sampling

# WIP: Leo - Platform-independent analysis with PC Sampling

- GPU PC sampling collects a wealth of information
  - Stall reasons (AMD, Intel, NVIDIA)
  - Hidden stalls (NVIDIA, AMD)
  - Register utilization (NVIDIA, AMD)
  - Compute unit utilization (AMD)
    - Pipeline utilization
    - Wavefront occupancy
    - Thread utilization (divergence)
- Leo pinpoints and quantifies root causes of GPU stalls
  - Ranks code locations suffering from GPU instruction stalls
  - Dumps machine code into assembly code
  - Performs dataflow analysis to trace stalls back to potential root causes

```
Database: /data/per-kernel/Algorithm_REDUCE_SUM/amd/hpctoolkit-raja-perf.exe-database
Measurements: /data/per-kernel/Algorithm_REDUCE_SUM/amd/hpctoolkit-raja-perf.exe-measurements

PROGRAM TOTALS
--------------------------------------------------------------------------------------
  Total Execution Time: 0.0037s
  Total Stall Cycles:  1,873,805,312
  Total GPU Cycles:    1,946,157,056
  Overall Stall Ratio: 96.3%
  Kernels Analyzed:    1
  Kernels Skipped:     0

TOP 1 KERNELS BY STALL CYCLES
--------------------------------------------------------------------------------------
#  Kernel                     Time (s)      Stall Cycles    Stall %     Occupancy     Sample Rate
--------------------------------------------------------------------------------------
1  .text                        0.0020       343,932,928     97.0%     30% (vgpr)    181.3G/s / 28%
--------------------------------------------------------------------------------------

KERNEL #1: .text
  GPU Binary: c59f0354ed847c14b96a717c2709e313.gpubin
  Occupancy: 30% (12/40 waves/CU, limited by vgpr) [VGPRs=40, SGPRs=50, LDS=64B]
======================================================================================
                               Leo GPU Performance Analysis
======================================================================================
WARNING: Very low profile coverage: 66/15359 (0.4%)
Kernel: .text                               Architecture: AMD MI300
Total Stall Cycles: 423,624,704
======================================================================================

STALL ANALYSIS (PC Sampling → Back-slicing → Root Cause)
--------------------------------------------------------------------------------------
Stall Location         Stall Opcode       Root Cause Location         Root Opcode            Cycles % Total Speedup
--------------------------------------------------------------------------------------
REDUCE_SUM-Hip.cpp:104:25   s_waitcnt vmcnt(0)  <-- REDUCE_SUM-Hip.cpp:104:28   global_load_dwordx2    277,872,640   65.6%   2.71x
amd_device_functions.h:778:9 s_barrier          <-- amd_device_functions.h:778:9 s_barrier             44,040,192   10.4%   1.11x
functional.hpp:242          s_waitcnt vmcnt(9)  <-- block_load_func.hpp:258     global_load_dwordx2     16,777,216    4.0%   1.04x
block_load_func.hpp:258     global_load_dwordx2 <-- block_load_func.hpp:258     v_addc_co_u32_e32       16,055,075    3.8%   1.04x
(base) JMC22-5:LEO johnmc$
```

# Discussion

- What performance issues do you need help with?
- What would help you investigate these issues?
- Are there any features that you would like added?

# HPCToolkit Resources

- Documentation
  - User manual for HPCToolkit: https://hpctoolkit.gitlab.io/hpctoolkit
  - Cheat sheet: https://gitlab.com/hpctoolkit/hpctoolkit/-/wikis/HPCToolkit-cheat-sheet
  - User manual for hpcviewer: https://hpctoolkit.gitlab.io/hpctoolkit/users/hpcviewer/hpcviewer.html
  - Tutorial videos
    - http://hpctoolkit.org/training.html
      - recorded demo of GPU analysis of Quicksilver: https://youtu.be/vixa3hGDuGg
      - recorded tutorial presentation including demo with GPU analysis of GAMESS: https://vimeo.com/781264043
- Software
  - Download hpcviewer GUI binaries for your laptop, desktop, cluster, or supercomputer
    - OS: Linux, Windows, MacOS
    - Processors: x86_64, aarch64, ppc64le
    - http://hpctoolkit.org/download.html
  - Install HPCToolkit on your Linux desktop, cluster, or supercomputer using Spack
    - http://hpctoolkit.org/software.html

RICE