

A Brief Introduction to HPCToolkit

Keren Zhou
Department of Computer Science
Rice University

<http://hpctoolkit.org>



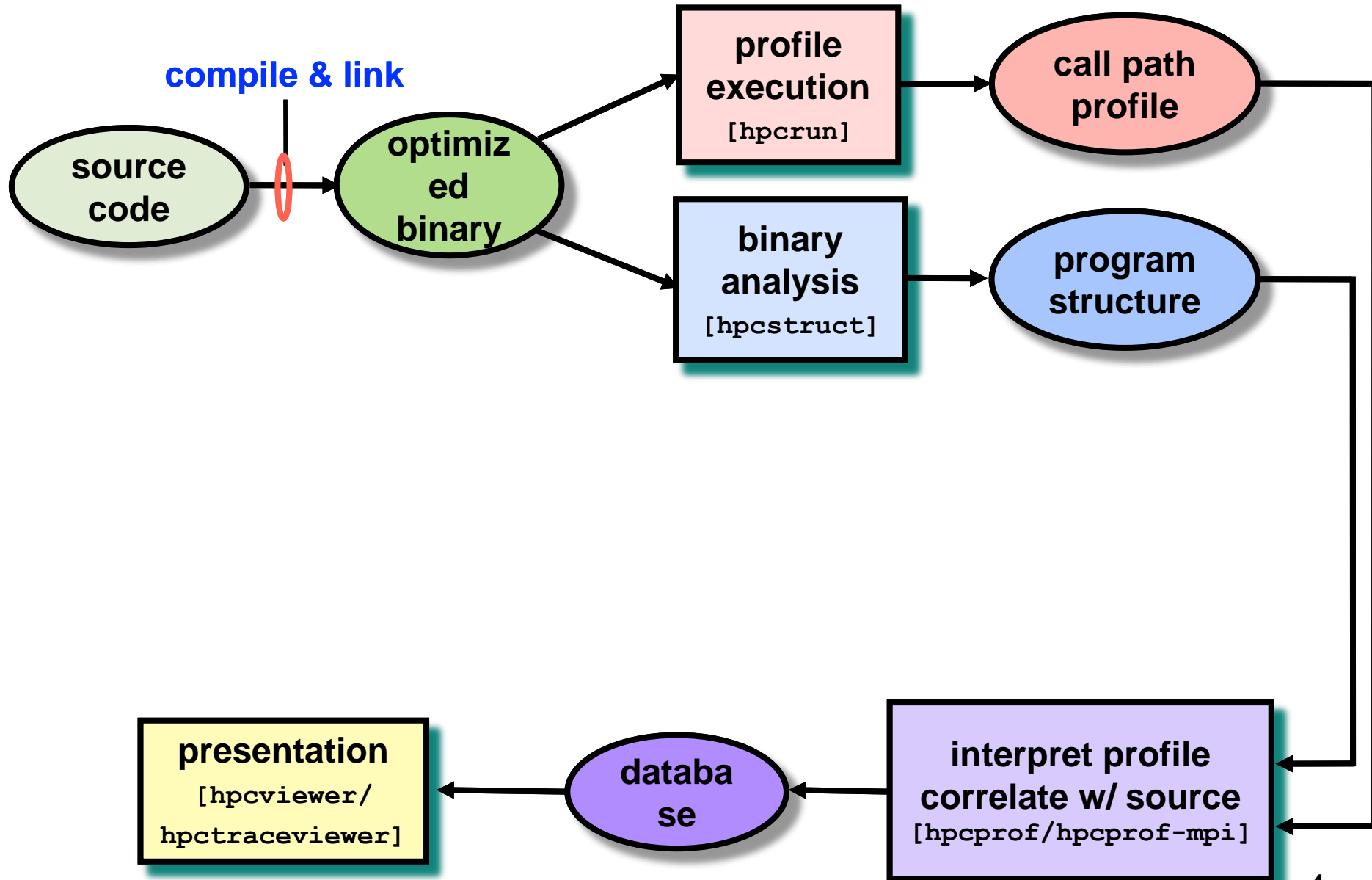
Outline

- **Overview of Rice's HPCToolkit**
- **OpenMP issues**
- **Using HPCToolkit's GUIs to analyze program performance**
- **Other capabilities**

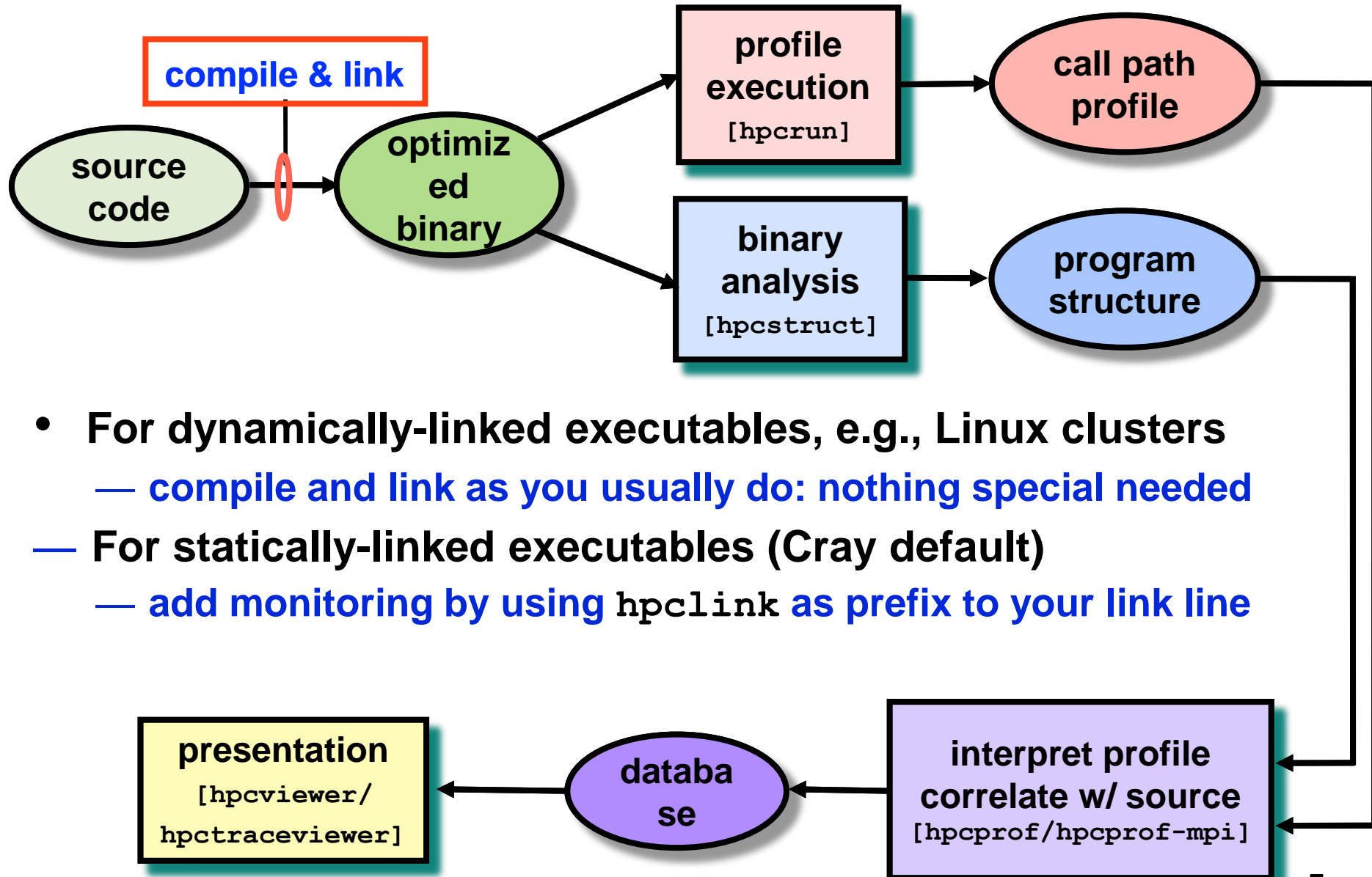
Rice University's HPCToolkit

- Employs binary-level measurement and analysis
 - observe **fully optimized**, **dynamically linked** executions
 - support **multi-lingual codes** with external binary-only libraries
- Uses sampling-based measurement (avoid instrumentation)
 - **controllable overhead**
 - **minimize** systematic error and avoid **blind spots**
 - enable data collection for **large-scale parallelism**
- Collects and correlates multiple derived performance metrics
 - diagnosis often requires more than one species of metric
- Associates metrics with both static and dynamic context
 - **loop nests**, **procedures**, **inlined code**, **calling context**
- Supports top-down performance analysis
 - identify costs of interest and drill down to causes
 - **up and down call chains**
 - **over time**

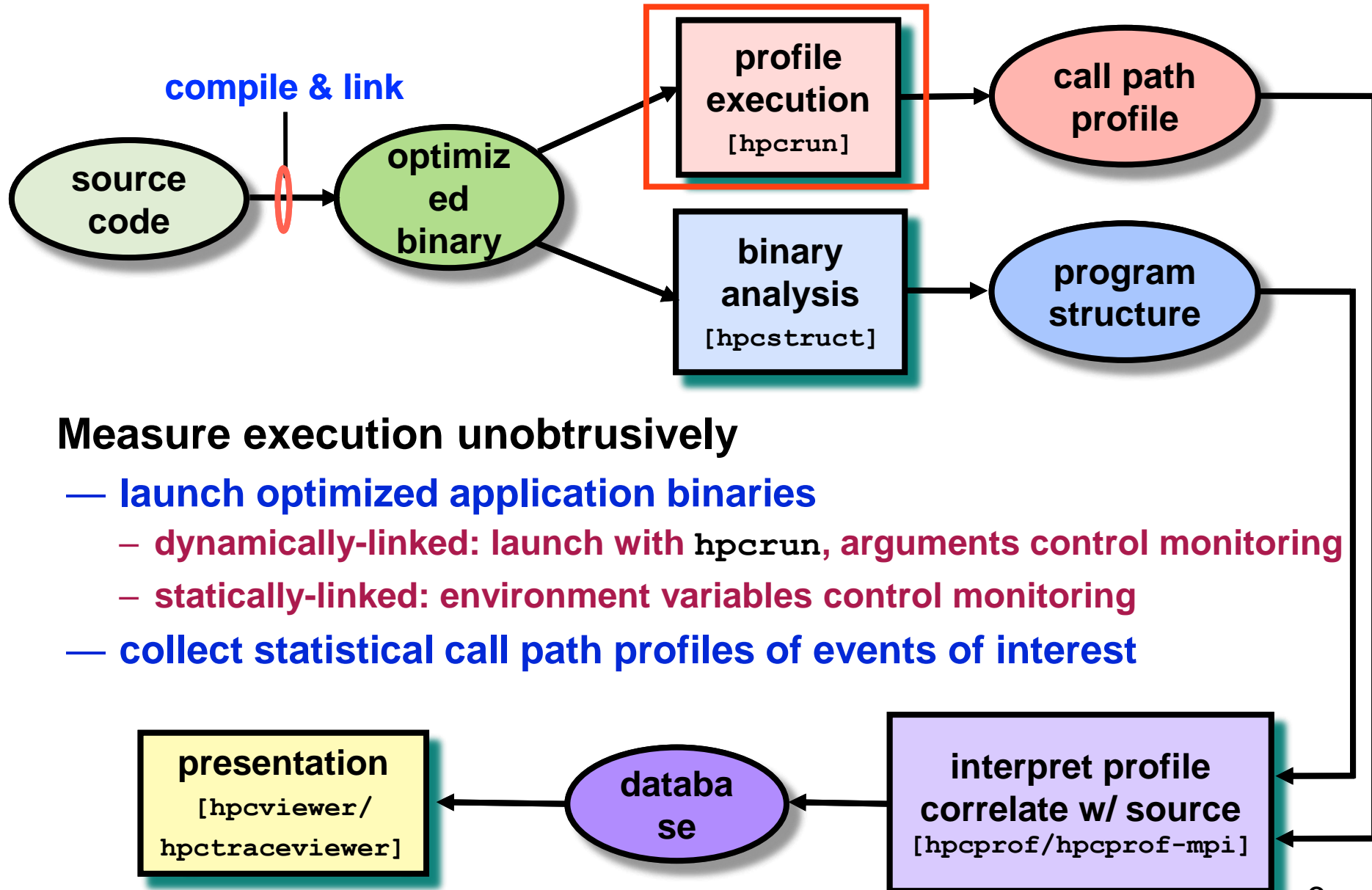
HPCToolkit Workflow



HPCToolkit Workflow



HPCToolkit Workflow



Measure execution unobtrusively

- **launch optimized application binaries**
 - **dynamically-linked: launch with `hpcrun`, arguments control monitoring**
 - **statically-linked: environment variables control monitoring**
- **collect statistical call path profiles of events of interest**

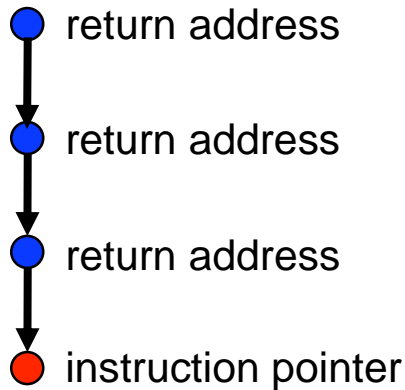
Call Path Profiling

Measure and attribute costs in context

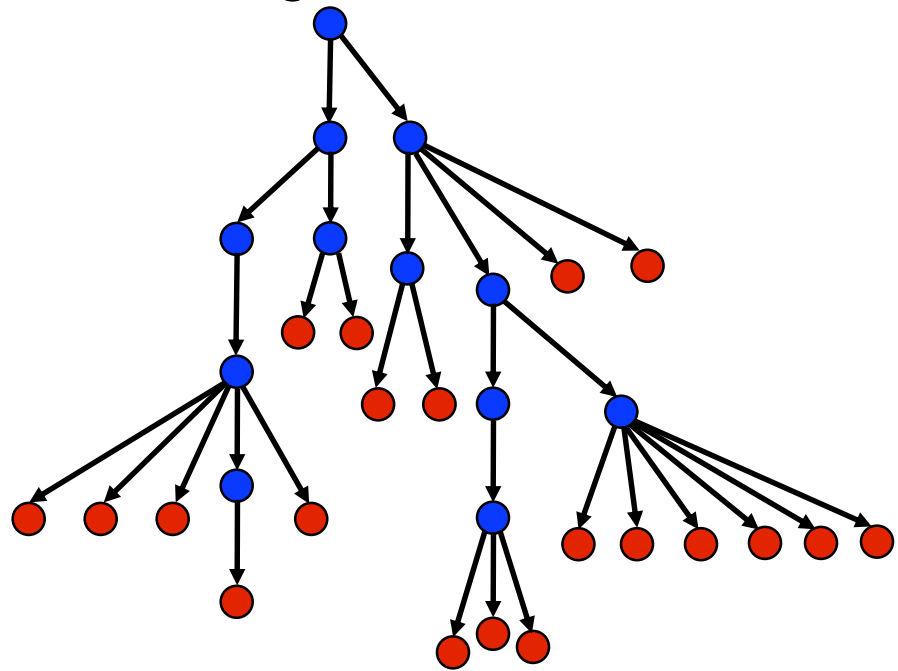
sample timer or hardware counter overflows

gather calling context using stack unwinding

Call path sample

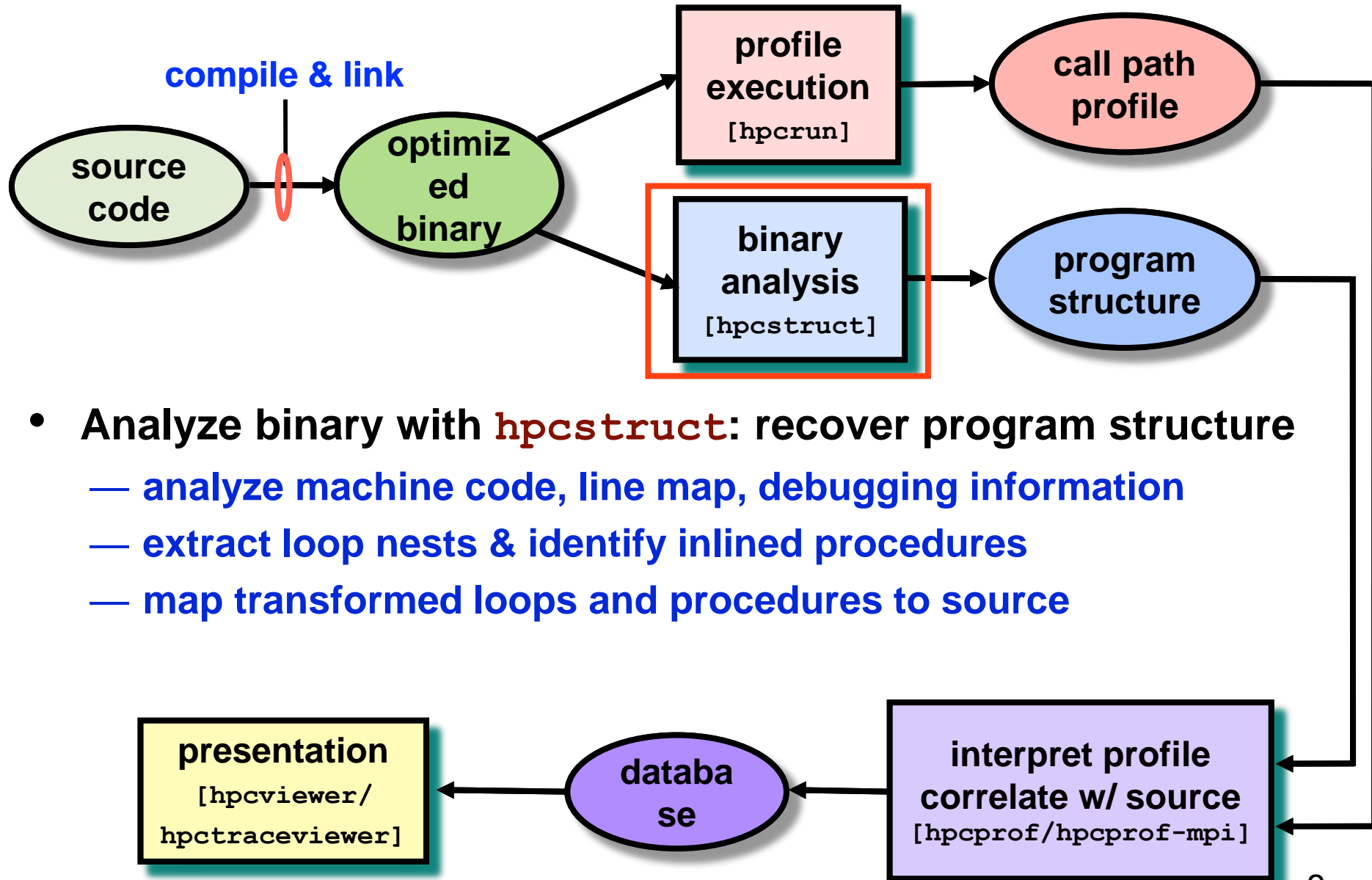


Calling context tree

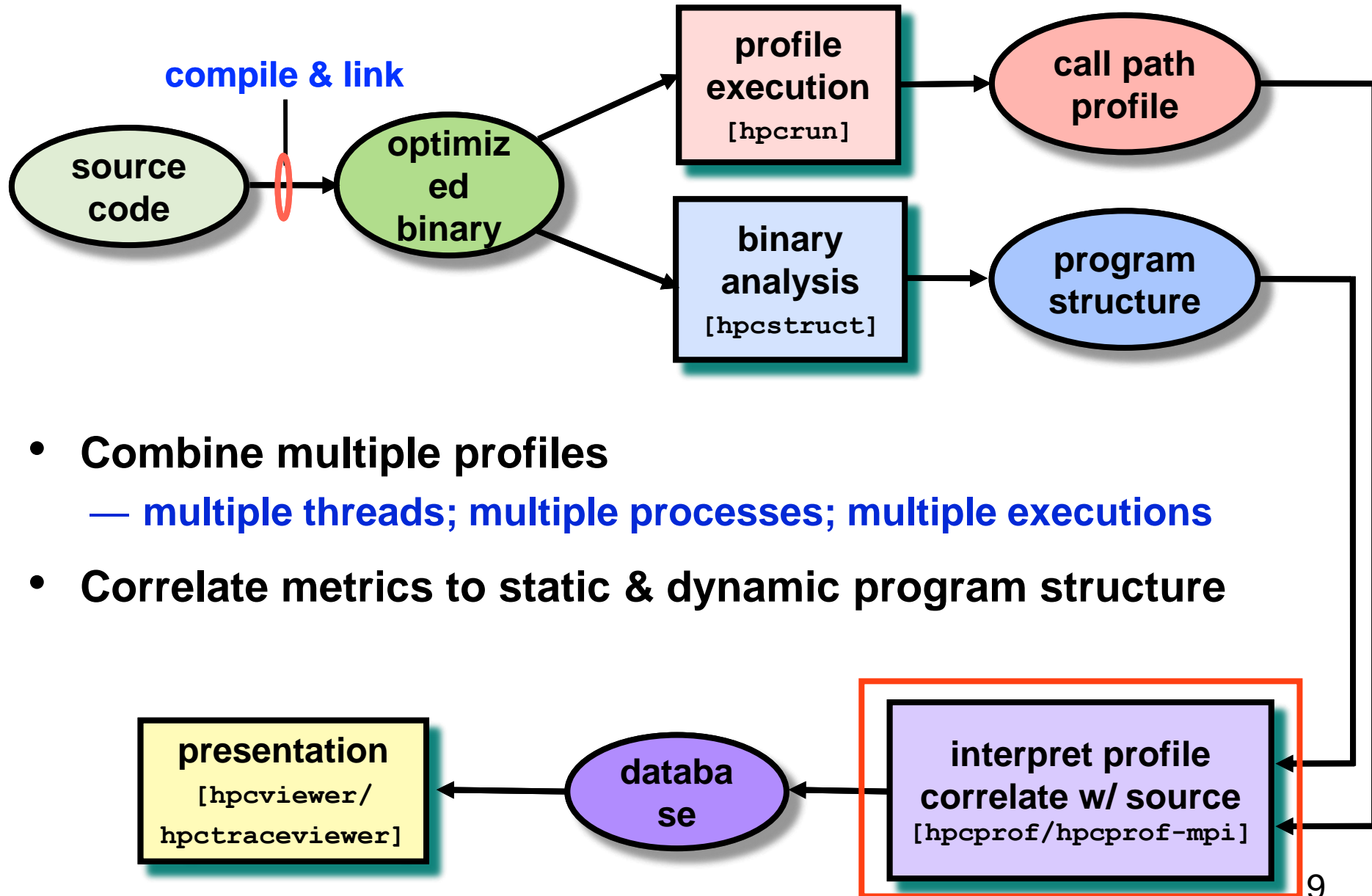


**Overhead proportional to sampling frequency...
...not call frequency**

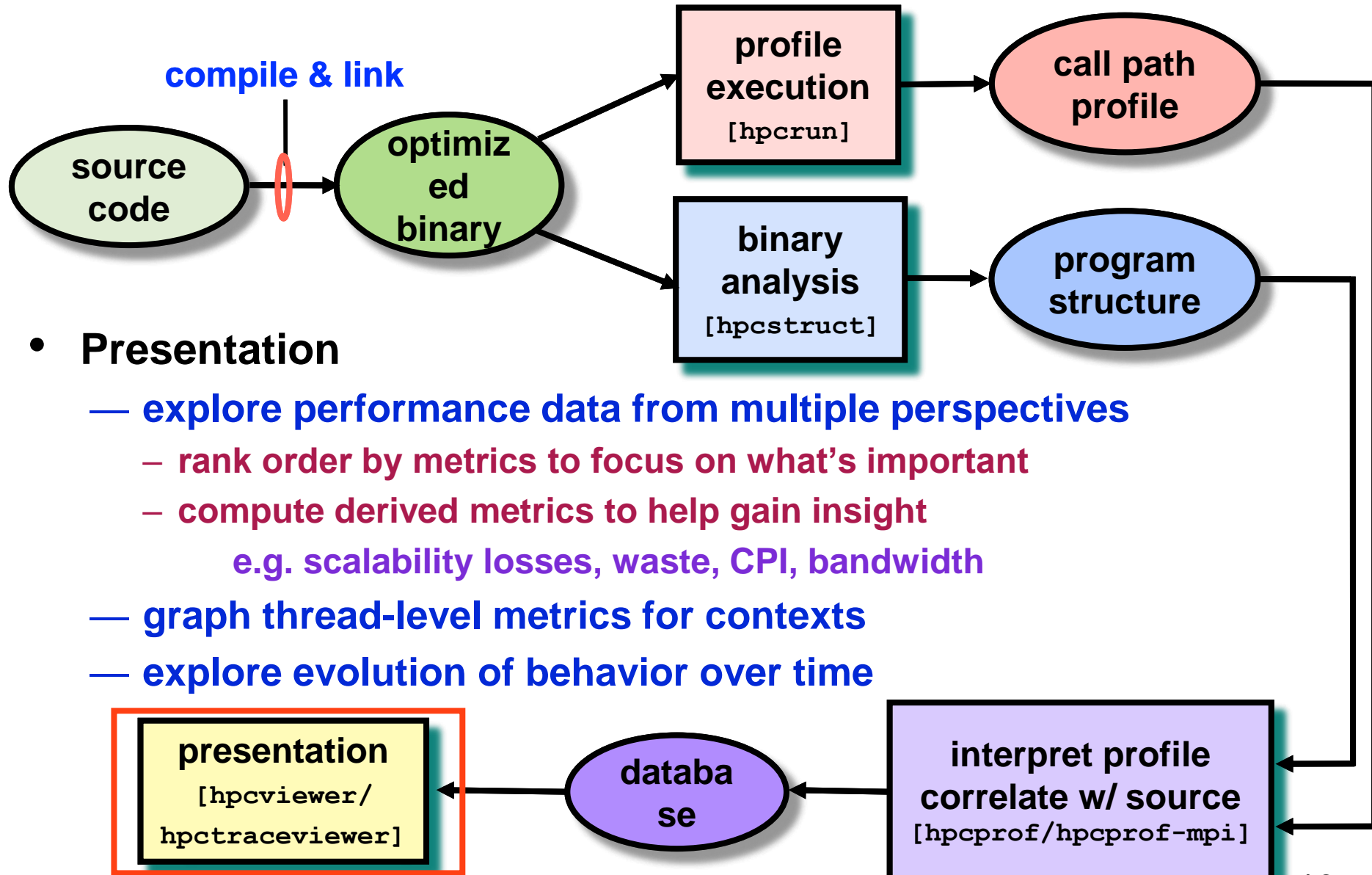
HPCToolkit Workflow



HPCToolkit Workflow



HPCToolkit Workflow



Code-centric Analysis with hpcviewer

The screenshot shows the hpcviewer application window titled "hpcviewer: lulesh-RAJA-parallel.exe". The interface includes a menu bar (File, Filter, View, Window, Help) and a toolbar. The main area is divided into several panes:

- source pane:** Displays the source code of the file "luleshRAJA-parallel.cxx". The code shows a function "forall" with a loop body.
- view control:** A set of buttons (Calling Context View, Callers View, Flat View) used to switch between different views of the code.
- metric display:** A set of icons (up/down arrows, a flame, a magnifying glass, and a list) used to filter and display metrics.
- navigation pane:** A tree view showing the execution flow of the program, starting from the "program root" and branching into various functions and loops.

The navigation pane is currently expanded to show the "forall" function, which is further expanded to show its internal structure, including a loop and several function calls.

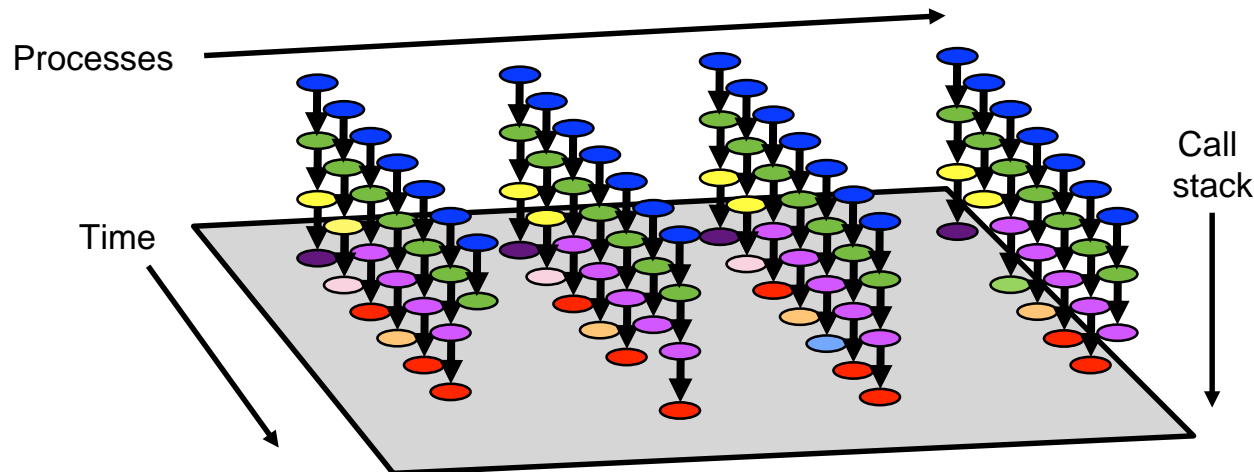
- function calls in full context
- inlined procedures
- inlined templates
- outlined OpenMP loops
- loops

The metric pane displays performance data for various functions and loops. The data is organized into three columns: the function name, the REALTIME (usec):Sum (I), and the REALTIME (usec):Sum (E). The data is sorted by the REALTIME (usec):Sum (I) column.

	REALTIME (usec):Sum (I)	REALTIME (usec):Sum (E)
Experiment Aggregate Metrics	2.26e+08 100 %	2.26e+08 100 %
<program root>	1.45e+08 63.9%	
497: main	1.45e+08 63.9%	6.01e+03 0.0%
loop at luleshRAJA-parallel.cxx: 3526	1.41e+08 63.8%	
3528: [I] LagrangeLeapFrog(Domain*)	1.44e+08 63.8%	
2715: [I] LagrangeNodal(Domain*)	8.70e+07 36.5%	
1554: [I] CalcForceForNodes(Domain*)	8.30e+07 36.7%	
1469: CalcVolumeForceForElems(Domain*)	8.25e+07 36.5%	
1454: [I] CalcHourglassControlForElems(Domain*, double*, double)	5.15e+07 22.8%	
1399: [I] CalcFBHourglassForceForElems(int*, double*, double*, double*, double*, double)	3.10e+07 13.7%	
1187: [I] void RAJA::forall<RAJA::IndexSet::ExecPolicy<RAJA::seq_segit, RAJA::omp_parallel_for_exec>>(const IndexSet& iset, const LoopBody& loop_body)	2.43e+07 10.8%	
405: [I] void RAJA::forall<RAJA::omp_parallel_for_exec, CalcFBHourglassForceForElems(int*, double*, double*, double*, double*, double)>(const IndexSet& iset, const LoopBody& loop_body)	2.43e+07 10.8%	
loop at forall_seq_any.hxx: 498	2.43e+07 10.8%	
505: [I] void RAJA::forall<CalcFBHourglassForceForElems(int*, double*, double*, double*, double*, double)>(const IndexSet& iset, const LoopBody& loop_body)	2.43e+07 10.8%	1.00e+03 0.0%
89: outline forall_omp_any.hxx:89 (0x423620)	2.42e+07 10.7%	3.91e+04 0.0%
loop at forall_omp_any.hxx: 90	2.42e+07 10.7%	3.41e+04 0.0%
91: [I] CalcFBHourglassForceForElems(int*, double*, double*, double*, double*, double)	2.42e+07 10.7%	9.84e+06 4.3%
1300: [I] CalcElemFBHourglassForce(double*, double*, double*, double*, double*, double)	1.11e+07 4.9%	1.11e+07 4.9%
1260: [I] CBRT(double)	3.27e+06 1.4%	2.00e+05 0.1%

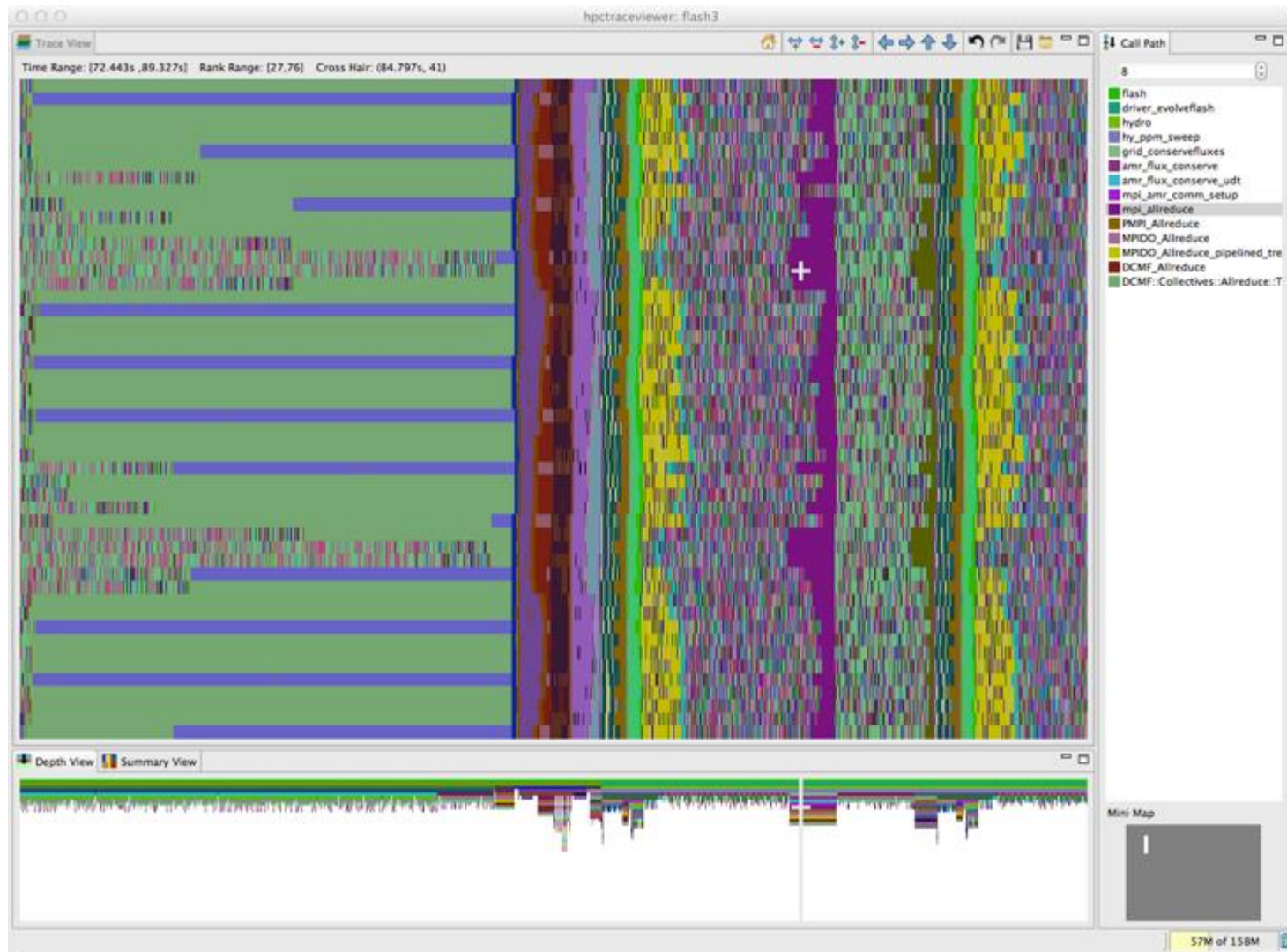
Understanding Temporal Behavior

- Profiling compresses out the temporal dimension
 - temporal patterns, e.g. serialization, are invisible in profiles
- What can we do? Trace call path samples
 - sketch:
 - 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



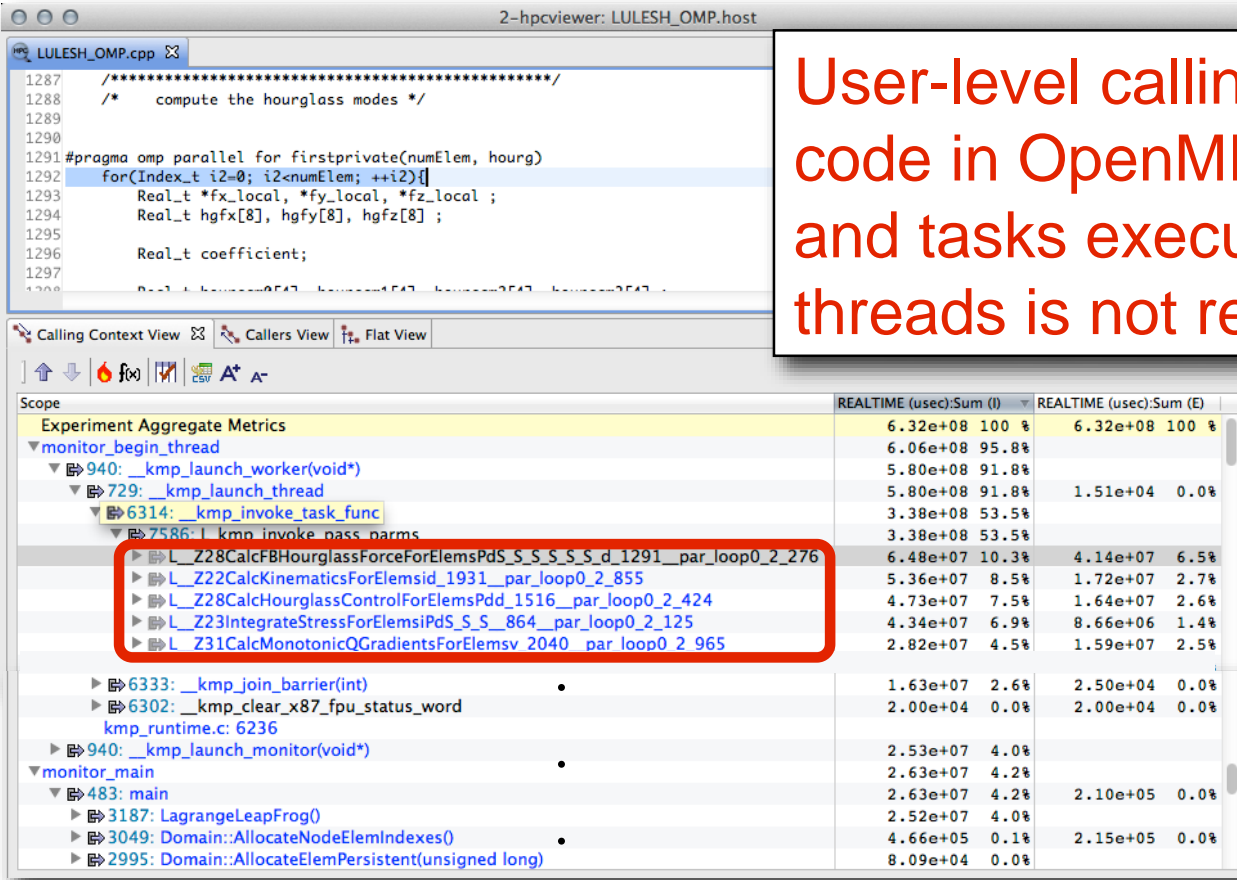
hpctraceviewer: detail of FLASH@256PE

Time-centric analysis: load imbalance among threads appears as different lengths of colored bands along the x axis



OpenMP: A Challenge for Tools

- Large gap between between threaded programming models and their implementations



User-level calling context for code in OpenMP parallel regions and tasks executed by worker threads is not readily available

- Runtime support is necessary for tools to bridge the gap

Challenges for OpenMP Node Programs

- **Tools provide implementation-level view of OpenMP threads**
 - **asymmetric threads**
 - **master thread**
 - **worker thread**
 - **run-time frames are interspersed with user code**
- **Hard to understand causes of idleness**
 - **long serial sections**
 - **load imbalance in parallel regions**
 - **waiting for critical sections or locks**

OMPT: An OpenMP Tools API

- **Goal: a standardized tool interface for OpenMP**
 - prerequisite for portable tools
 - missing piece of the OpenMP language standard
- **Design objectives**
 - enable tools to measure and attribute costs to application source and runtime system
 - support low-overhead tools based on asynchronous sampling
 - attribute to user-level calling contexts
 - associate a thread's activity at any point with a descriptive state
 - minimize overhead if OMPT interface is not in use
 - features that may increase overhead are optional
 - define interface for trace-based performance tools
 - don't impose an unreasonable development burden
 - runtime implementers
 - tool developers

LLNL's IuleshMPI_OMP (8 MPI x 3 OMP), 30, REALTIME@1000

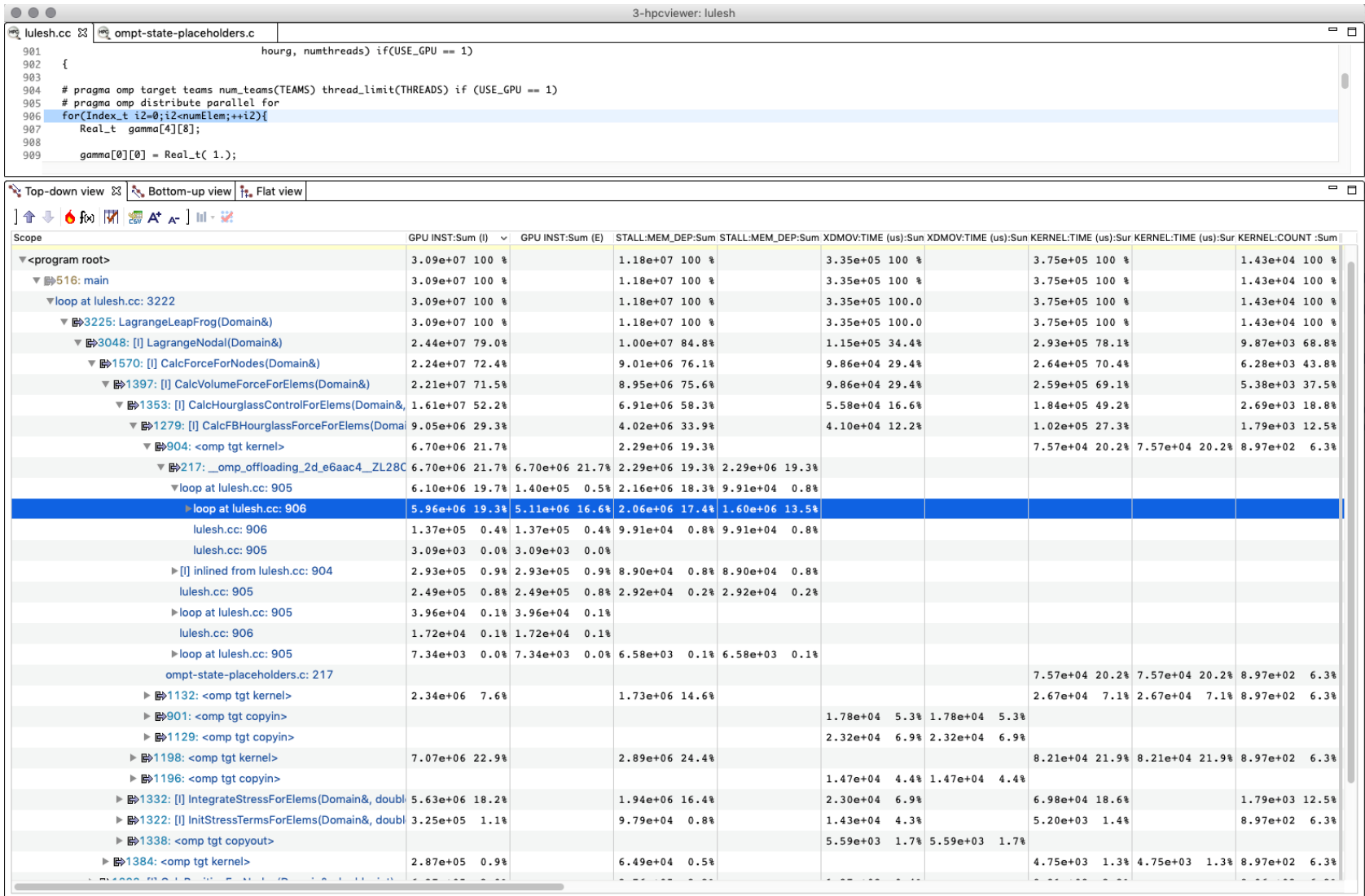


OpenMP Tool API Status

- HPCToolkit supports OpenMP 5.0 OMPT
- OMPT prototype implementations
 - LLVM (emerging: OpenMP 5.0)
 - interoperable with GNU, Intel compilers
 - IBM LOMP (currently targets OpenMP 4.5)
- Ongoing work
 - refining OpenMP 5.0 OMPT support in LLVM OpenMP
 - refining OpenMP 5.0 OMPT support in HPCToolkit
 - asynchronous call stack assembly for lightweight monitoring

HPCToolkit Capabilities for GPU Code

MPI + OpenMP 4.5 or CUDA GPU accelerated applications



Other Capabilities

- **Measure hardware counters using Linux perf_events**
 - available events can be listed with
 - `hpcrun -L`
 - launching a binary created by `hpcrun` with environment setting `HPCRUN_EVENT_LIST=LIST`
 - frequency based sampling: 300/s per thread or machine max
 - no need to set periods or frequencies unless you want precise control
 - hardware event multiplexing
 - measure more events than hardware counters
- **Kernel sampling**
 - measure activity in the Linux kernel in addition to your program
 - e.g., allocating and clearing memory pages
 - not available on BG/Q
 - measurement and attribution subject to system permissions
 - detailed attribution not available on NERSC or ANL systems

HPCToolkit at NERSC

- **NERSC cori**
 - a setup script or a set of module loads
 - `source /global/cscratch1/sd/kz21/env-static.sh`
 - `source /global/cscratch1/sd/kz21/env-shared.sh`
- **Man pages**
 - automatically added to MANPATH by the aforementioned command

HPCToolkit at ORNL

- **On Summit**
 - **module use /gpfs/alpine/csc322/world-shared/modulefiles**
 - **module load hpctoolkit**
- **Man pages**
 - **automatically added to MANPATH by the aforementioned command**

GUIs for your Laptop

- Download binary packages for HPCToolkit's user interfaces on your laptop
 - <http://hpctoolkit.org/download/hpcviewer>

Detailed HPCToolkit Documentation

<http://hpctoolkit.org/documentation.html>

- **Comprehensive user manual:**

- <http://hpctoolkit.org/manual/HPCToolkit-users-manual.pdf>

- **Quick start guide**

- essential overview that almost fits on one page

- **Using HPCToolkit with statically linked programs**

- a guide for using hpctoolkit on BG/Q and Cray platforms

- **The hpcviewer and hpctraceviewer user interfaces**

- **Effective strategies for analyzing program performance with HPCToolkit**

- analyzing scalability, waste, multicore performance ...

- **HPCToolkit and MPI**

- **HPCToolkit Troubleshooting**

- why don't I have any source code in the viewer?

- hpcviewer isn't working well over the network ... what can I do?

- **Installation guide**

Advice for Using HPCToolkit

Using HPCToolkit

- Add hpctoolkit's bin directory to your path using softenv
- Adjust your compiler flags (if you want full attribution to src)
 - add -g flag after any optimization flags
- Add hpclink as a prefix to your Makefile's link line
 - e.g. `hpclink CC -o myapp foo.o ... lib.a -lm ...`
- See what sampling triggers are available on Cray
 - use hpclink to link your executable
 - launch executable with environment variable `HPCRUN_EVENT_LIST=LIST`
 - you can launch this on 1 core of 1 node
 - no need to provide arguments or input files for your program
they will be ignored

Monitoring Large Executions

- Collecting performance data on every node is typically not necessary
- Can improve scalability of data collection by recording data for only a fraction of processes
 - set environment variable `HPCRUN_PROCESS_FRACTION`
 - e.g. collect data for 10% of your processes
 - set environment variable `HPCRUN_PROCESS_FRACTION=0.10`

Digesting your Performance Data

- Use hpcstruct to reconstruct program structure
 - e.g. `hpcstruct your_app`
 - creates `your_app.hpcstruct`
- Correlate measurements to source code with hpcprof and hpcprof-mpi
 - run hpcprof on the front-end to analyze data from small runs
 - run hpcprof-mpi on the compute nodes to analyze data from lots of nodes/threads in parallel
 - notes
 - much faster to do this on an x86_64 vis cluster (cooley) than on BG/Q
 - avoid expensive per-thread profiles with `--metric-db no`
- Digesting performance data in parallel with hpcprof-mpi
 - `qsub -A ... -t 20 -n 32 --mode c1 --proccount 32 --cwd `pwd` \`
`/projects/Tools/hpctoolkit/pkgs-vesta/hpctoolkit/bin/hpcprof-mpi \`
`-S your_app.hpcstruct \`
`-I /path/to/your_app/src/+ \`
`hpctoolkit-your_app-measurements.jobid`

Analysis and Visualization

- Use hpcviewer to open resulting database
 - warning: first time you graph any data, it will pause to combine info from all threads into one file
- Use hpctraceviewer to explore traces
 - warning: first time you open a trace database, the viewer will pause to combine info from all threads into one file
- Try out our user interfaces before collecting your own data
 - example performance data
<http://hpctoolkit.org/examples.html>

Installing HPCToolkit GUIs on your Laptop

- See <http://hpctoolkit.org/download/hpcviewer>
- Download the latest for your laptop (Linux, Mac, Windows)
 - hpctraceviewer
 - hpcviewer

A Note for Mac Users

When installing HPCToolkit GUIs on your Mac laptop, don't simply download and double click on the zip file and have Finder unpack them. Follow the Terminal-based installation directions on the website to avoid interference by Mac Security.