

Gaining Insight into Parallel Program Performance using HPCToolkit

John Mellor-Crummey
Department of Computer Science
Rice University
johnmc@rice.edu



<http://hpctoolkit.org>



Challenges for Computational Scientists

- **Execution environments and applications are rapidly evolving**
 - **architecture**
 - rapidly changing multicore microprocessor designs
 - increasing scale of parallel systems
 - growing use of accelerators
 - **applications**
 - transition from MPI everywhere to threaded implementations
 - add additional scientific capabilities
 - maintain multiple variants or configurations
- **Steep increase in development effort to deliver performance, evolvability, and portability**
- **Computational scientists need to**
 - **assess weaknesses in algorithms and their implementations**
 - **improve scalability of executions within and across nodes**
 - **adapt to changes in emerging architectures**

Performance tools can play an important role as a guide

Performance Analysis Challenges

- **Complex architectures are hard to use efficiently**
 - multi-level parallelism: multi-core, ILP, SIMD instructions
 - multi-level memory hierarchy
 - result: gap between typical and peak performance is huge
- **Complex applications present challenges**
 - measurement and analysis
 - understanding behaviors and tuning performance
- **Supercomputer platforms compound the complexity**
 - unique hardware
 - unique microkernel-based operating systems
 - multifaceted performance concerns
 - computation
 - communication
 - I/O

Performance Analysis Principles

- **Without accurate measurement, analysis is irrelevant**
 - avoid systematic measurement error
 - measure actual executions of interest, not an approximation
 - fully optimized production code on the target platform
- **Without effective analysis, measurement is irrelevant**
 - quantify and attribute problems to source code
 - compute insightful metrics
 - e.g., “scalability loss” or “waste” rather than just “cycles”
- **Without scalability, a tool is irrelevant for supercomputing**
 - large codes
 - large-scale threaded parallelism within and across nodes

Performance Analysis Goals

- Programming model independent tools
- Accurate measurement of complex parallel codes
 - large, multi-lingual programs
 - fully optimized code: loop optimization, templates, inlining
 - binary-only libraries, sometimes partially stripped
 - complex execution environments
 - dynamic loading (Linux clusters) vs. static linking (Cray, Blue Gene)
 - SPMD parallel codes with threaded node programs
 - batch jobs
- Effective performance analysis
 - insightful analysis that pinpoints and explains problems
 - correlate measurements with code for actionable results
 - support analysis at the desired level
 - intuitive enough for application scientists and engineers
 - detailed enough for library developers and compiler writers
- Scalable to petascale and beyond

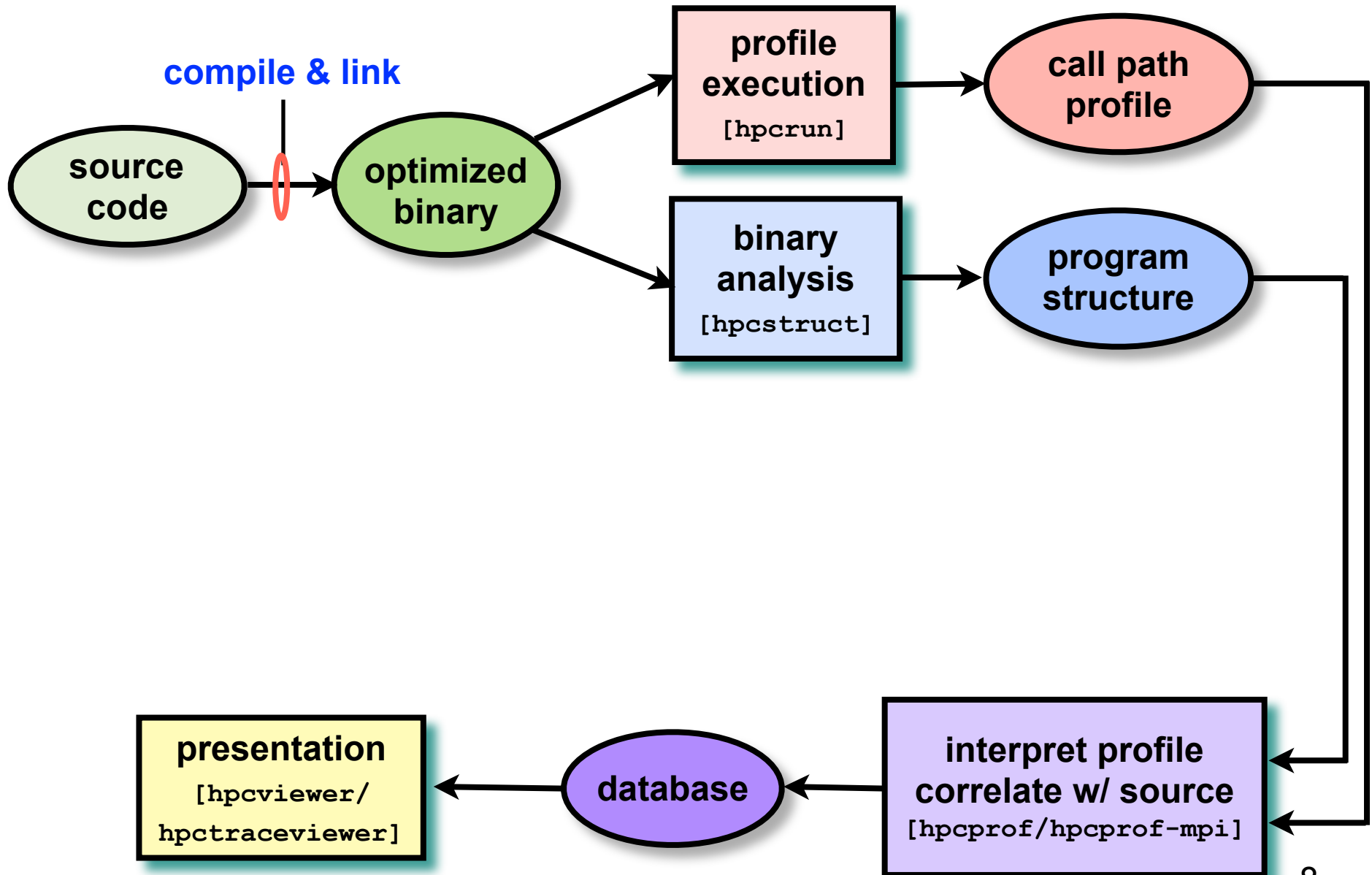
HPCToolkit Design Principles

- **Employ binary-level measurement and analysis**
 - observe **fully optimized**, **dynamically linked** executions
 - support **multi-lingual codes** with external binary-only libraries
- **Use sampling-based measurement (avoid instrumentation)**
 - **controllable overhead**
 - **minimize** systematic error and avoid blind spots
 - enable data collection for **large-scale parallelism**
- **Collect and correlate multiple derived performance metrics**
 - diagnosis typically requires more than one species of metric
- **Associate metrics with both static and dynamic context**
 - **loop nests**, **procedures**, **inlined code**, **calling context**
- **Support top-down performance analysis**
 - natural approach that minimizes burden on developers

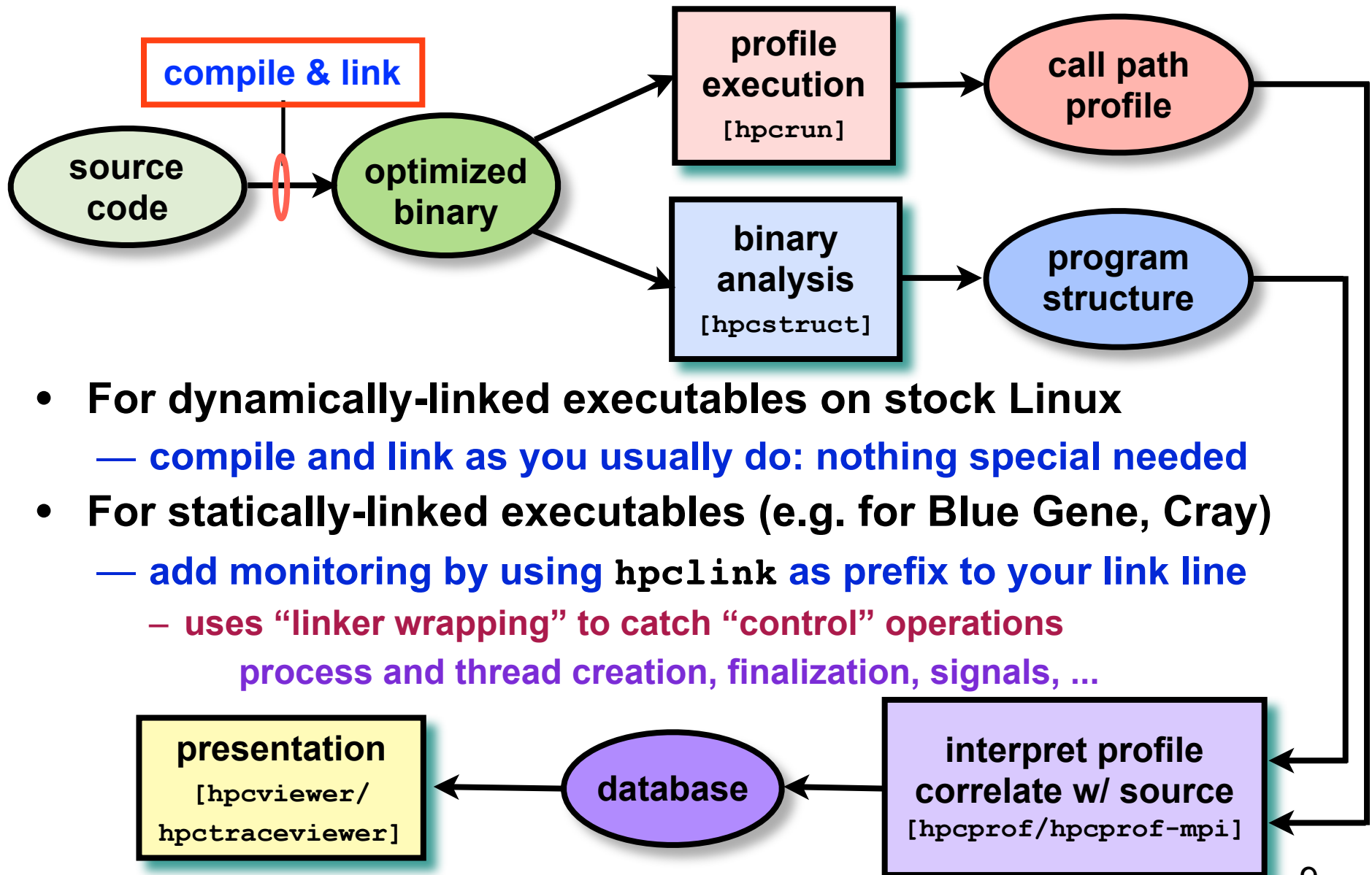
Outline

- **Overview of Rice's HPCToolkit**
- **Accurate measurement**
- **Effective performance analysis**
- **Pinpointing scalability bottlenecks**
 - scalability bottlenecks on large-scale parallel systems
 - scaling on multicore processors
- **Understanding temporal behavior**
- **Assessing process variability**
- **Understanding threading, GPU, locks, and memory hierarchy**
 - blame shifting
 - attributing memory hierarchy costs to data
- **Summary and conclusions**

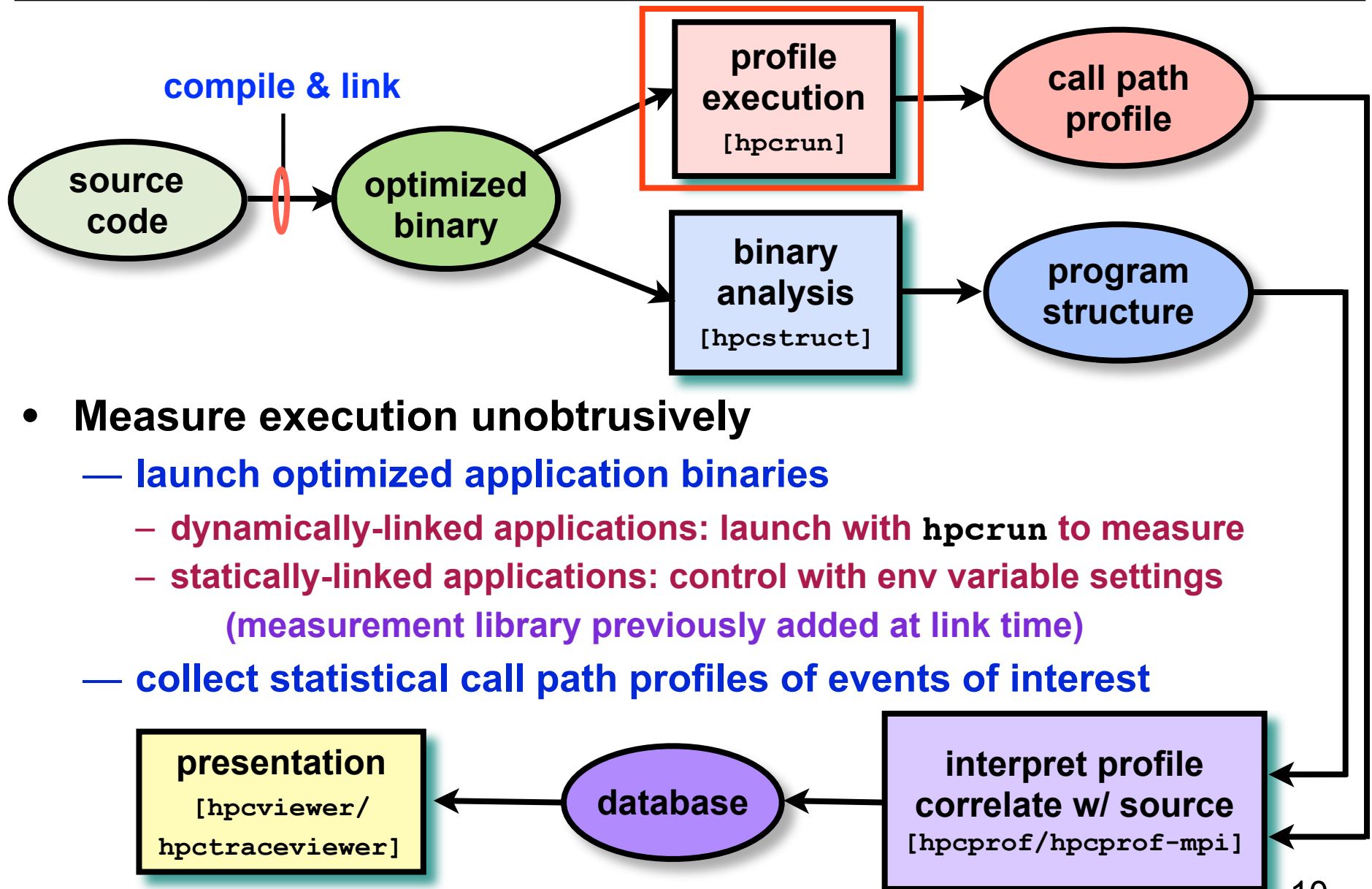
HPCToolkit Workflow



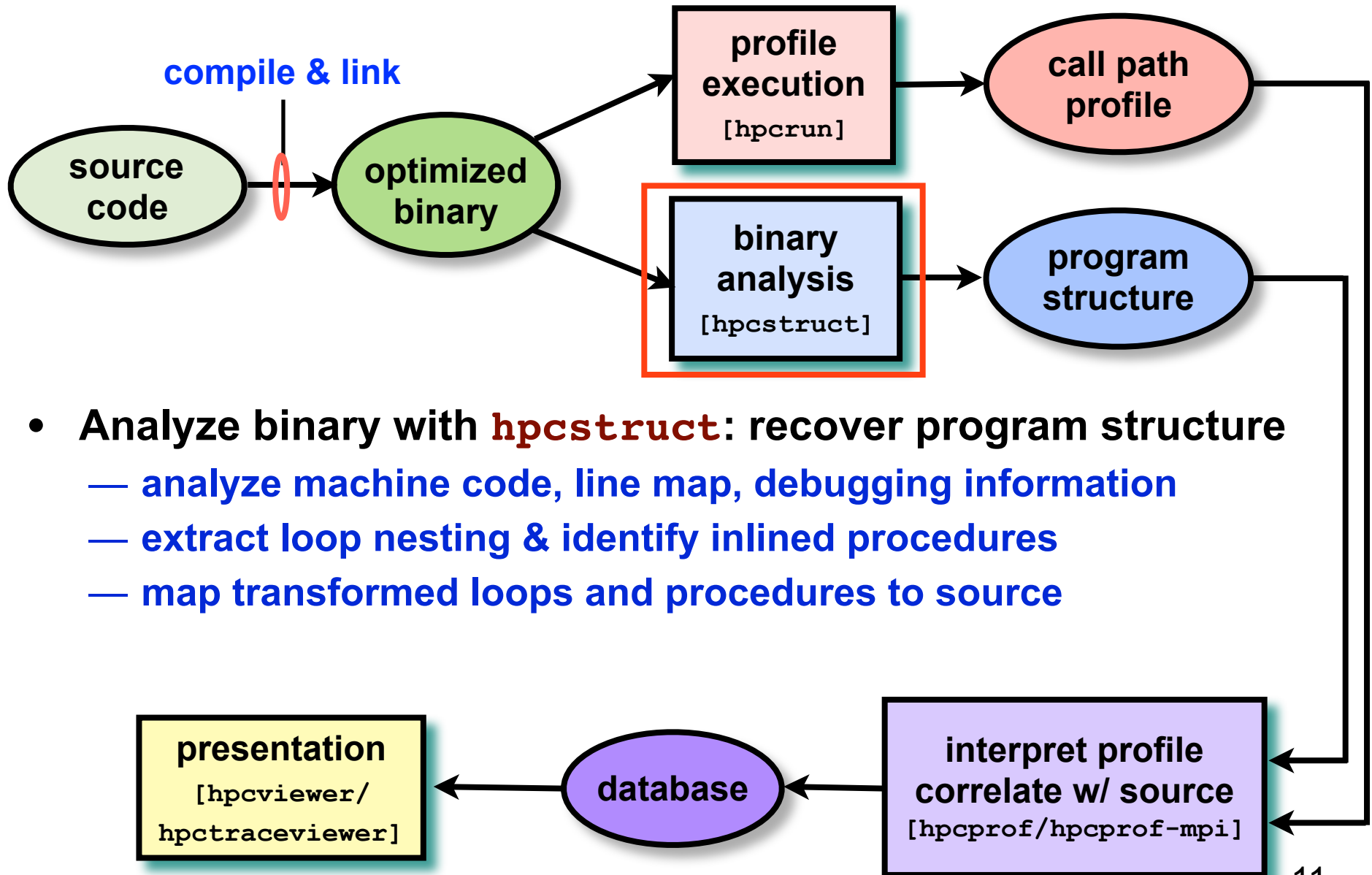
HPCToolkit Workflow



HPCToolkit Workflow

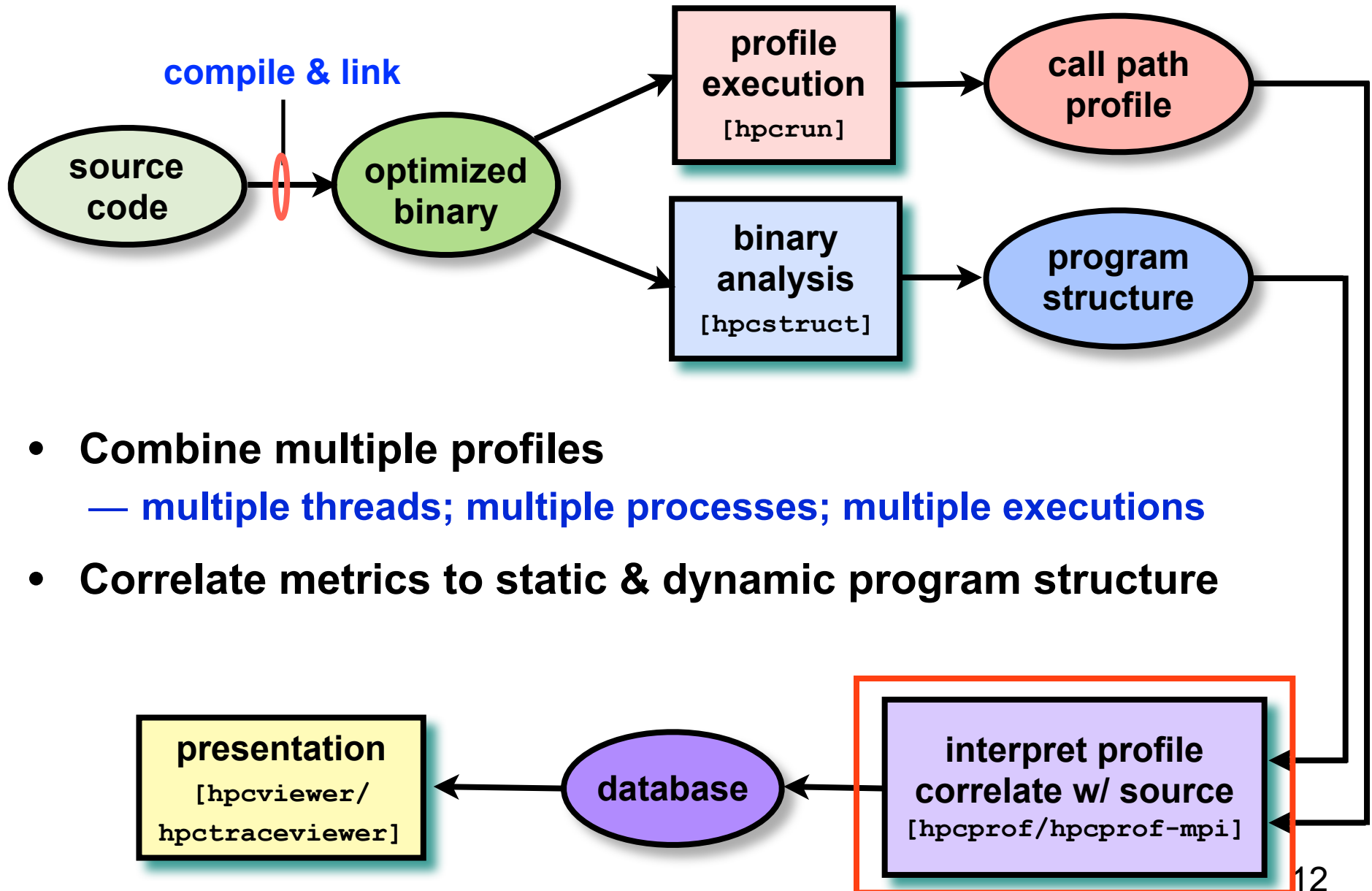


HPCToolkit Workflow

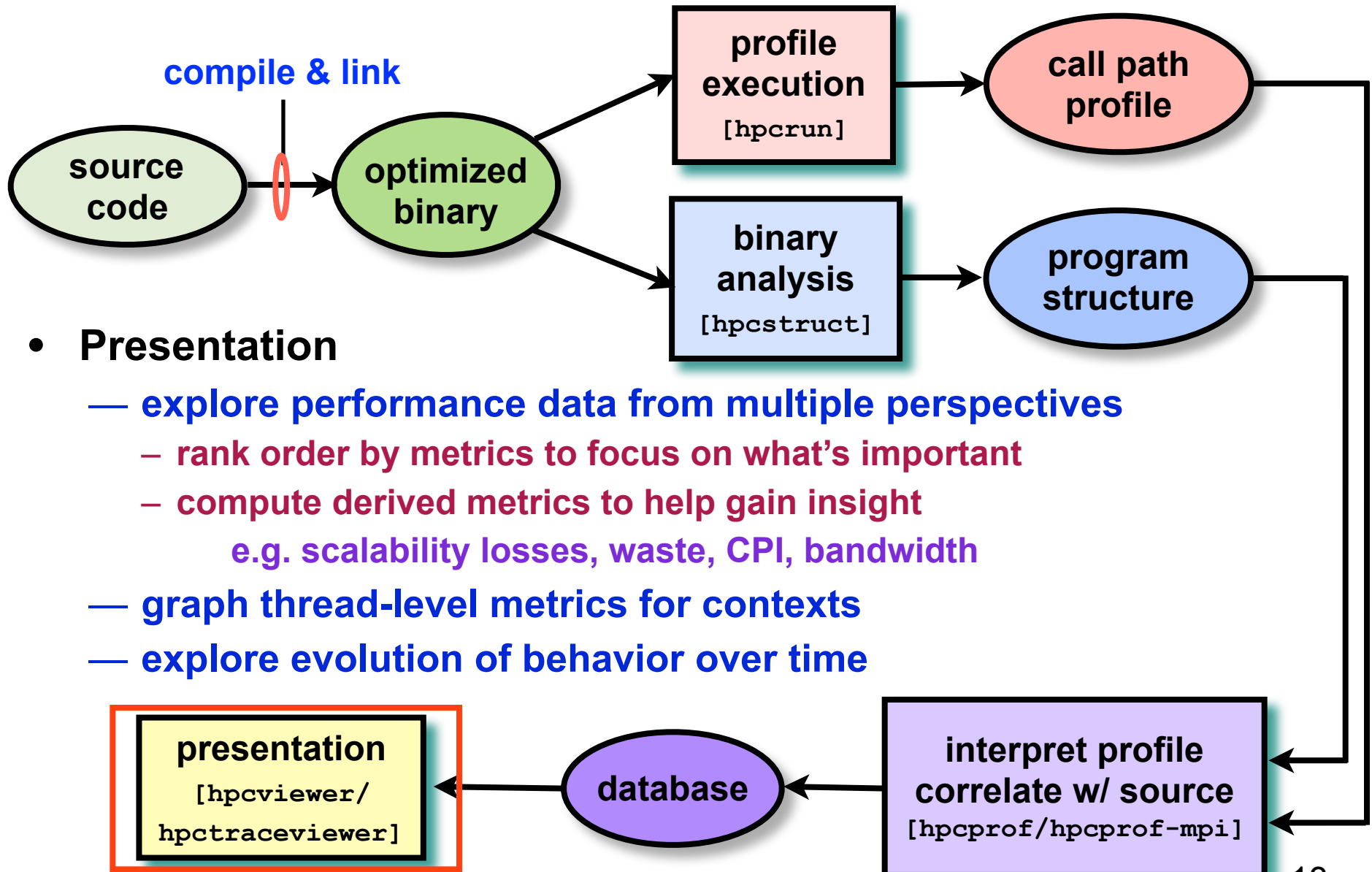


- Analyze binary with **hpcstruct**: recover program structure
 - analyze machine code, line map, debugging information
 - extract loop nesting & identify inlined procedures
 - map transformed loops and procedures to source

HPCToolkit Workflow



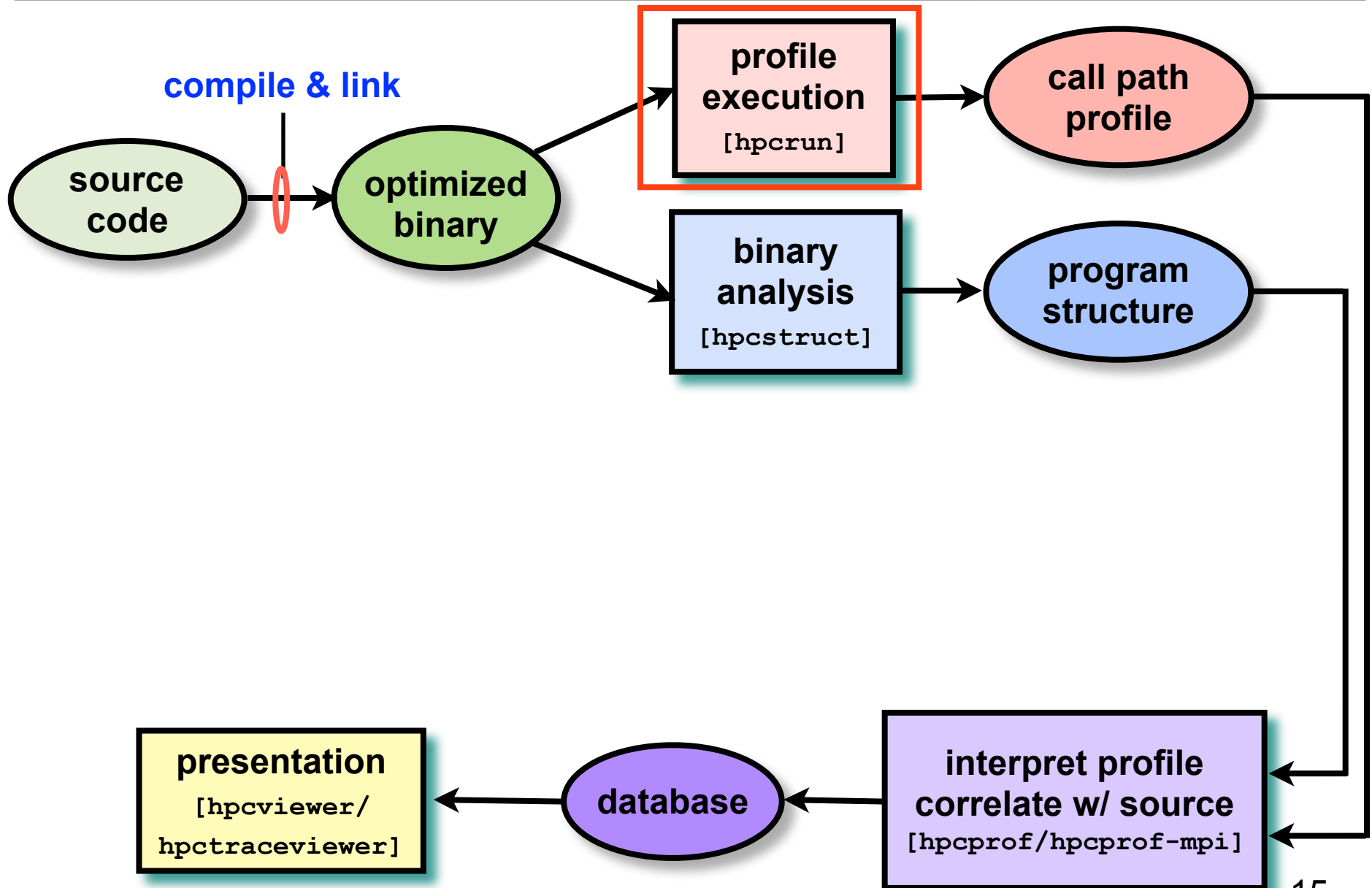
HPCToolkit Workflow



Outline

- Overview of Rice's HPCToolkit
- Accurate measurement
- Effective performance analysis
- Pinpointing scalability bottlenecks
 - scalability bottlenecks on large-scale parallel systems
 - scaling on multicore processors
- Understanding temporal behavior
- Assessing process variability
- Understanding threading, GPU, locks, and memory hierarchy
 - blame shifting
 - attributing memory hierarchy costs to data
- Summary and conclusions

Measurement



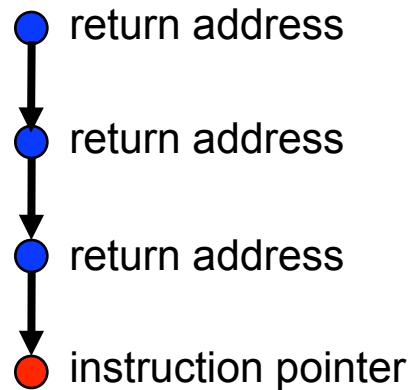
Call Path Profiling

Measure and attribute costs in context

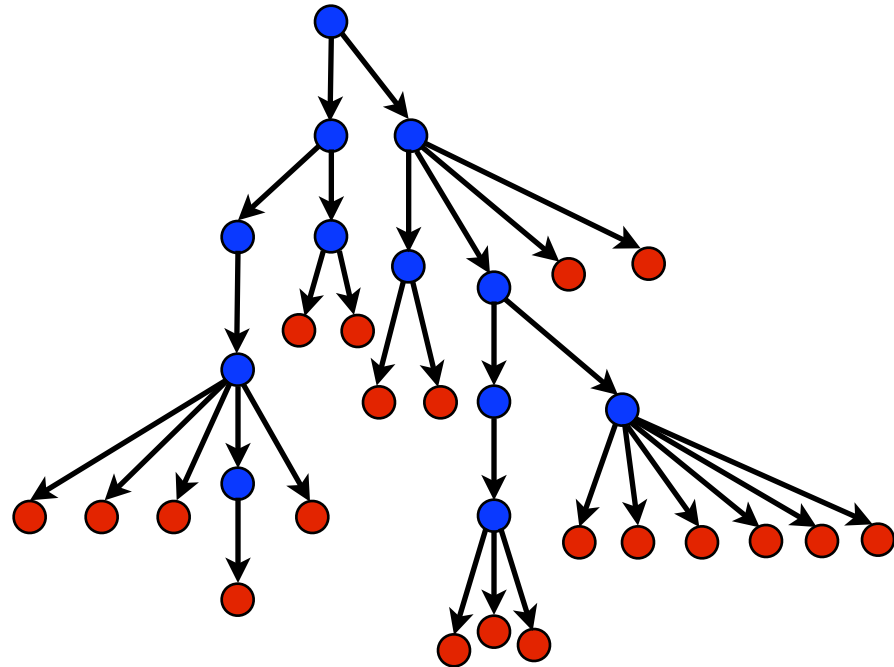
sample timer or hardware counter overflows

gather calling context using stack unwinding

Call path sample



Calling context tree



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

Why Sampling?

The performance uncertainty principle implies that the accuracy of performance data is inversely correlated with the degree of performance instrumentation – AI Malony, PhD Thesis 1991

Instrumentation of MADNESS with TAU

Method	Number of Profiled Events	Runtime (seconds)	Overhead (%)
Uninstrumented		654s	
Compiler-based Instrumentation	1321	19625s	2901%
Regular Source Instrumentation	183	748s	14.4%
Source Instrumentation with headers (-optHeaderInst)	806	1628s	150%
-optHeaderInst and selective instrumentation (auto)	539	685s	4.7%
callpath depth 2, -optHeaderInst and selective instrumentation (auto)	1773	693s	6%
callpath depth 100, -optHeaderInst and selective instrumentation (auto)	8535	893s	36.5%

Figure source: <http://www.nic.uoregon.edu/tau-wiki/MADNESS>

Why Sampling?

The performance uncertainty principle implies that the accuracy of performance data is inversely correlated with the degree of performance instrumentation – AI Malony, PhD Thesis 1991

Instrumentation of MADNESS with TAU

Method	Number of Profiled Events	Runtime (seconds)	Overhead (%)
Uninstrumented		654s	
Compiler-based Instrumentation	1321	19625s	2901%
Dynamic instrumentation	100	710s	11.4%
callpath depth 100, -optHeaderInst and selective instrumentation (auto)	8535	893s	36.5%

Each of these instrumentation approaches ignores any functions in libraries available only in binary form

Figure source: <http://www.nic.uoregon.edu/tau-wiki/MADNESS>

full instrumentation slows execution by 30x!

Novel Aspects of Our Approach

- **Unwind fully-optimized and even stripped code**
 - use on-the-fly binary analysis to support unwinding
- **Cope with dynamically-loaded shared libraries on Linux**
 - note as new code becomes available in address space
 - problematic for instrumentation-based tools, unless using Dyninst or Pin
- **Integrate static & dynamic context information in presentation**
 - dynamic call chains including procedures, inlined functions, loops, and statements

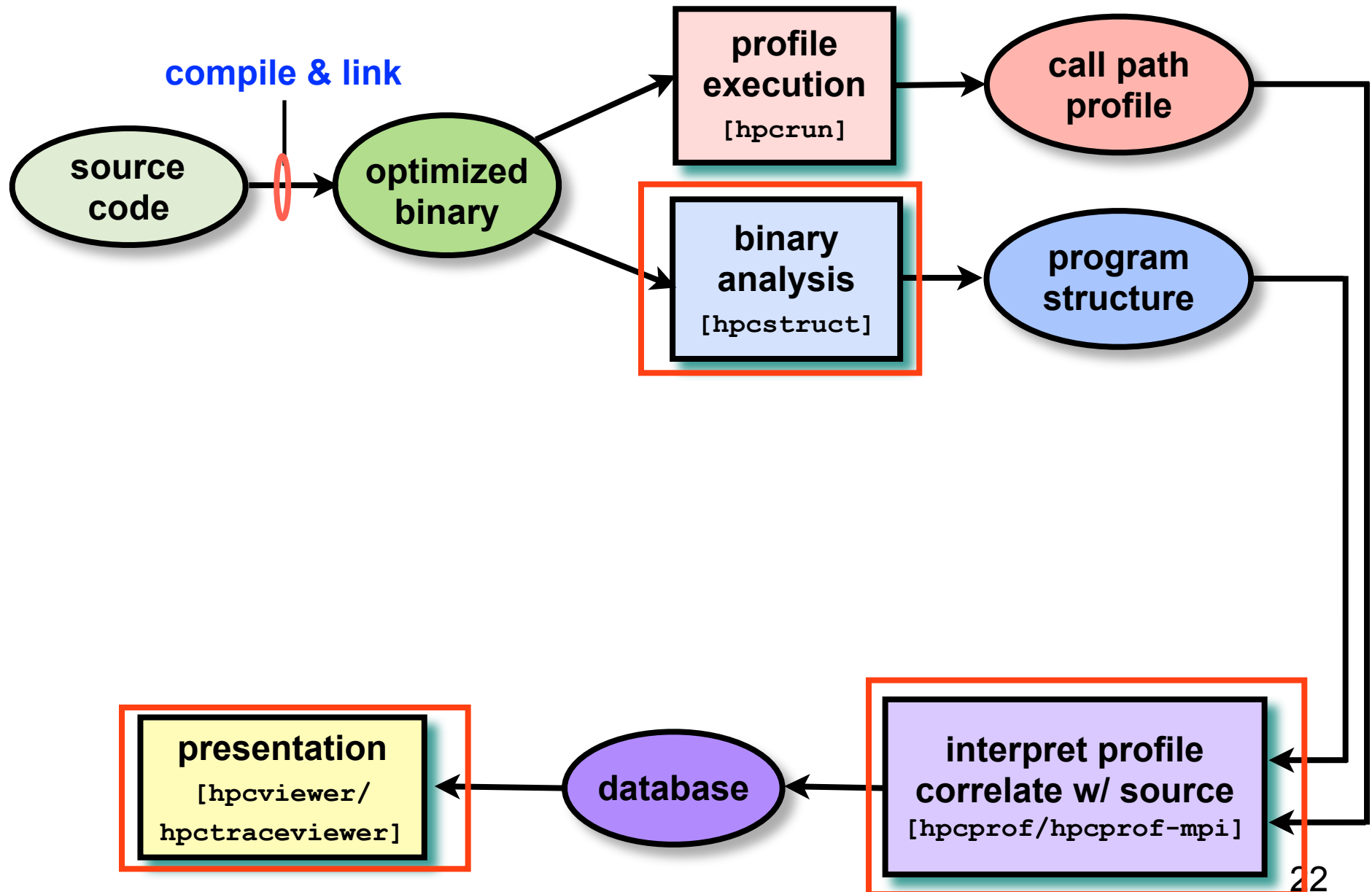
Measurement Effectiveness

- **Accurate**
 - **PFLOTRAN on Cray XT @ 8192 cores**
 - 148 unwind failures out of 289M unwinds
 - 5e-5% errors
 - **Flash on Blue Gene/P @ 8192 cores**
 - 212K unwind failures out of 1.1B unwinds
 - 2e-2% errors
 - **SPEC2006 benchmark test suite (sequential codes)**
 - fully-optimized executables: Intel, PGI, and Pathscale compilers
 - 292 unwind failures out of 18M unwinds (Intel Harpertown)
 - 1e-3% error
- **Low overhead**
 - **e.g. PFLOTRAN scaling study on Cray XT @ 512 cores**
 - measured cycles, L2 miss, FLOPs, & TLB @ 1.5% overhead
 - **suitable for use on production runs**

Outline

- Overview of Rice's HPCToolkit
- Accurate measurement
- Effective performance analysis
- Pinpointing scalability bottlenecks
 - scalability bottlenecks on large-scale parallel systems
 - scaling on multicore processors
- Understanding temporal behavior
- Assessing process variability
- Understanding threading, GPU, locks, and memory hierarchy
 - blame shifting
 - attributing memory hierarchy costs to data
- Summary and conclusions

Effective Analysis



Recovering Program Structure

- **Analyze an application binary**
 - **identify object code procedures and loops**
 - decode machine instructions
 - construct control flow graph from branches
 - identify natural loop nests using interval analysis
 - **map object code procedures/loops to source code**
 - leverage line map + debugging information
 - discover inlined code
 - account for many loop and procedure transformations

Unique benefit of our binary analysis

- **Bridges the gap between**
 - **lightweight measurement of fully optimized binaries**
 - **desire to correlate low-level metrics to source level abstractions**

Analyzing Results with hpcviewer

The screenshot displays the hpcviewer application window titled "hpcviewer: MOAB: mbperf_iMesh 200 B (Barcelona 2360 SE)". The interface is divided into several panes:

- source pane:** Shows the source code for `mbperf_iMesh.cpp`. The code includes comments and a `SequenceCompare` class definition. A red box highlights the `source pane` label.
- view control:** A toolbar with buttons for "Calling Context View", "Callers View", and "Flat View". A red box highlights the `view control` label.
- metric display:** A toolbar with icons for navigation and metrics. A red box highlights the `metric display` label.
- navigation pane:** A tree view showing the execution scope. It includes nodes for `main`, `testB`, and various inlined functions and loops. A red box highlights the `navigation pane` label.
- metric pane:** A table displaying performance metrics for different scopes. A red box highlights the `metric pane` label.

costs for

- inlined procedures
- loops
- function calls in full context

Scope	PAPI_L1_DCM (I)	PAPI_TOT_CYC (I)
main	8.63e+08 100 %	1.13e+11 100 %
testB(void*, int, double const*, int const*)	8.35e+08 96.7%	1.10e+11 97.6%
inlined from mbperf_iMesh.cpp: 261	6.81e+08 78.9%	0.98e+11 86.5%
loop at mbperf_iMesh.cpp: 280-313	3.43e+08	0.9%
imesh_getvtxarrcoords_	3.20e+08 37.1%	2.18e+10 19.3%
MBCore::get_coords(unsigned long const*, int, double*)	3.20e+08 37.1%	2.16e+10 19.1%
loop at MBCore.cpp: 681-693	3.20e+08 37.1%	2.16e+10 19.1%
inlined from stl_tree.h: 472	2.04e+08 23.7%	9.38e+09 8.3%
loop at stl_tree.h: 1388	2.04e+08 23.6%	9.37e+09 8.3%
inlined from TypeSequenceManager.hpp: 27	1.78e+08 20.6%	8.56e+09 7.6%
TypeSequenceManager.hpp: 27	1.78e+08 20.6%	8.56e+09 7.6%

Principal Views

- **Calling context tree view - “top-down” (down the call chain)**
 - associate metrics with each dynamic calling context
 - high-level, hierarchical view of distribution of costs
 - example: quantify initialization, solve, post-processing
- **Caller’s view - “bottom-up” (up the call chain)**
 - apportion a procedure’s metrics to its dynamic calling contexts
 - understand costs of a procedure called in many places
 - example: see where PGAS put traffic is originating
- **Flat view - ignores the calling context of each sample point**
 - aggregate all metrics for a procedure, from any context
 - attribute costs to loop nests and lines within a procedure
 - example: assess the overall memory hierarchy performance within a critical procedure

Toolchain Demo: Lulesh Serial Code

Source Code Attribution: LULESH A++

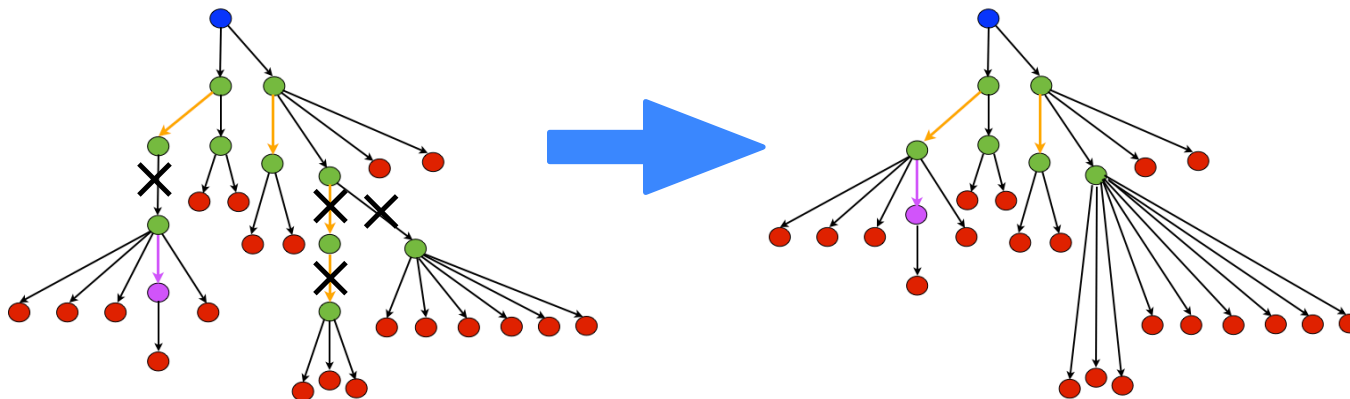
The screenshot displays the hpcviewer interface for the lulesh-app. The top pane shows the source code of LULESH-App.C. Two sections are highlighted: a blue box around the `AreaFace` function (lines 919-929) and a red box around a loop in the `CharacteristicLength` function (lines 938-940). A purple arrow points to the subtraction operation in line 922 of the `AreaFace` function. The bottom pane shows the 'Calling Context View' with a tree structure of the call stack. The path `main > loop at LULESH-App.C: 1802 > loop at LULESH-App.C: 1817 > inlined from LULESH-App.C: 1619 > 1626: LagrangeElements > inlined from LULESH-App.C: 1040 > 1046: CalcKinematicsForElems > loop at LULESH-App.C: 1003 > inlined from LULESH-App.C: 936 > loop at LULESH-App.C: 938 > inlined from LULESH-App.C: 922` is highlighted. The right pane shows a table of performance metrics for various scopes.

Performance attribution to inlined code and loops

Scope	PAPL_TOT_CYC.[0,0] (I)	PAPL_TOT_CYC.[0,0] (E)	PAPL_RES_STL.[0,0] (I)	PAPL_RES_STL.[0,0] (E)	P
Experiment Aggregate Metrics	7.87e+12 100 %	7.87e+12 100 %	5.64e+11 100 %	5.64e+11 100 %	
main	7.87e+12 100 %	1.10e+08 0.0%	5.64e+11 100 %	1.20e+07 0.0%	
loop at LULESH-App.C: 1802	7.86e+12 99.8%		5.63e+11 99.8%		
loop at LULESH-App.C: 1817	7.86e+12 99.8%	5.20e+07 0.0%	5.63e+11 99.8%	4.00e+06 0.0%	
inlined from LULESH-App.C: 1619	7.86e+12 99.8%		5.63e+11 99.8%		
1626: LagrangeElements	4.20e+12 53.3%	3.18e+09 0.0%	2.83e+11 50.2%	8.16e+08 0.1%	
inlined from LULESH-App.C: 1040	2.40e+12 30.4%		1.75e+11 31.1%		
1046: CalcKinematicsForElems	2.40e+12 30.4%	1.07e+10 0.1%	1.75e+11 31.1%	5.48e+08 0.1%	
loop at LULESH-App.C: 1003	2.40e+12 30.4%	1.06e+10 0.1%	1.75e+11 31.1%	5.46e+08 0.1%	
inlined from LULESH-App.C: 936	9.25e+11 11.8%	4.91e+09 0.1%	6.60e+10 11.7%	6.40e+07 0.0%	
loop at LULESH-App.C: 938	9.18e+11 11.7%	4.29e+09 0.1%	6.53e+10 11.6%	5.60e+07 0.0%	
inlined from LULESH-App.C: 922	6.25e+11 7.9%	2.84e+09 0.0%	5.06e+10 9.0%	5.60e+07 0.0%	
922: operator-(doubleArray const&, doubleArray const&)	5.04e+10 0.6%	8.22e+08 0.0%	4.36e+09 0.8%		
925: operator+(doubleArray const&, doubleArray const&)	4.80e+10 0.6%	8.32e+08 0.0%	4.14e+09 0.7%		
927: operator*(doubleArray const&, doubleArray const&)	4.40e+10 0.6%	7.64e+08 0.0%	4.70e+09 0.8%		
923: operator-(doubleArray const&, doubleArray const&)	4.07e+10 0.5%	4.56e+08 0.0%	4.93e+09 0.9%		
924: operator-(doubleArray const&, doubleArray const&)	3.95e+10 0.5%	4.62e+08 0.0%	5.26e+09 0.9%	2.00e+07 0.0%	
922: doubleArray::doubleArray(doubleArray const&, int)	3.95e+10 0.5%	2.63e+09 0.0%	2.78e+09 0.5%		
927: operator*(doubleArray const&, doubleArray const&)	3.92e+10 0.5%	2.72e+08 0.0%	4.44e+09 0.8%	1.20e+07 0.0%	
927: operator*(doubleArray const&, doubleArray const&)	3.87e+10 0.5%	2.72e+08 0.0%	4.44e+09 0.8%	1.20e+07 0.0%	

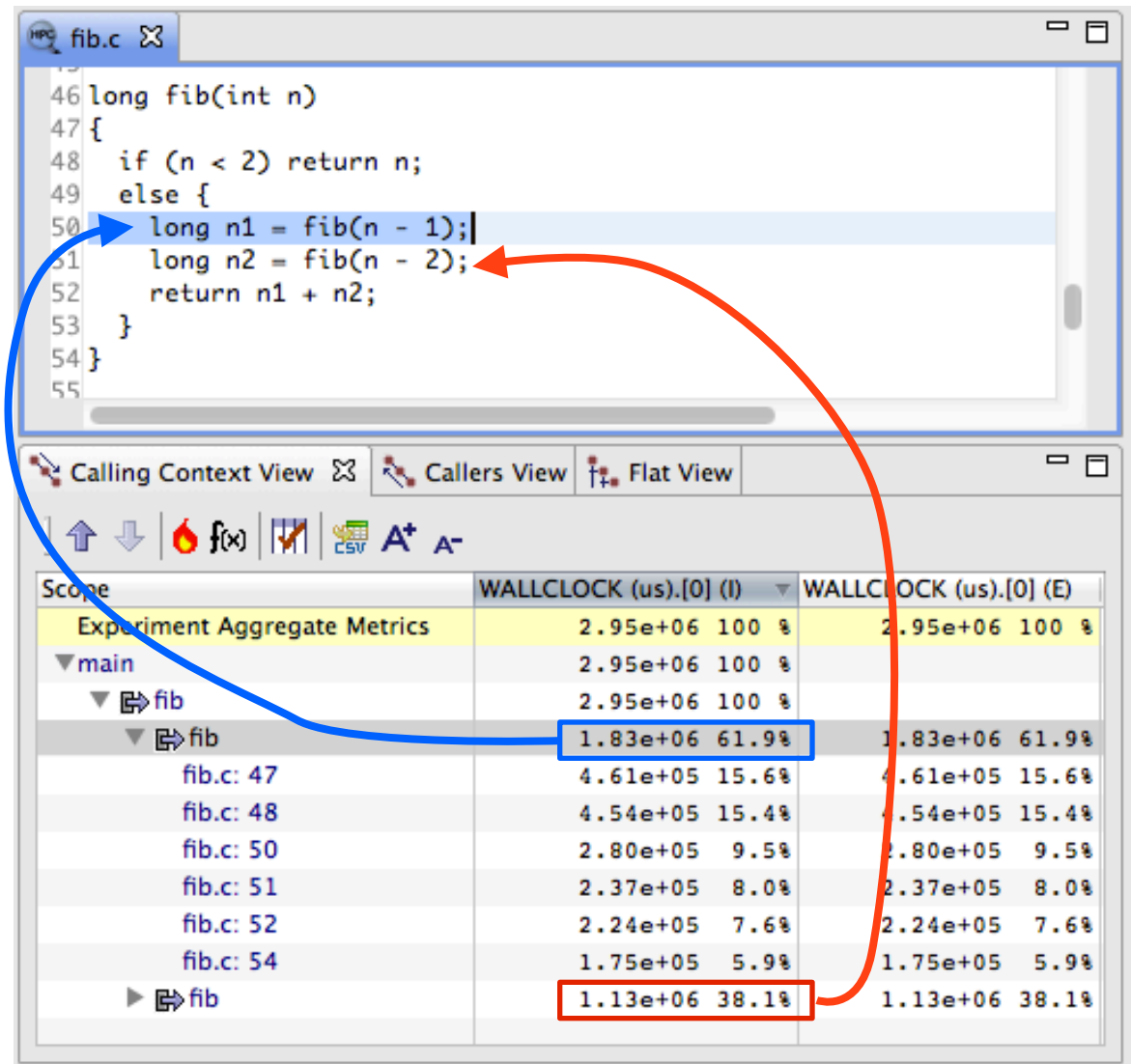
Handling Call Chains with Recursion

- **Problem:** some recursive algorithms, e.g., quicksort have many long and unique call chains
 - each sample can expose a unique call chain
 - space overhead can be significant for recursive computations that have many unique call chains, e.g. broad and deep trees
 - for parallel programs, the total space overhead can be especially problematic when thread-level views are merged
- **Approach**
 - collapse recursive chains to save space
 - preserve one level of recursion so high-level properties of the recursive solution remain available



Example: Recursive Fibonacci

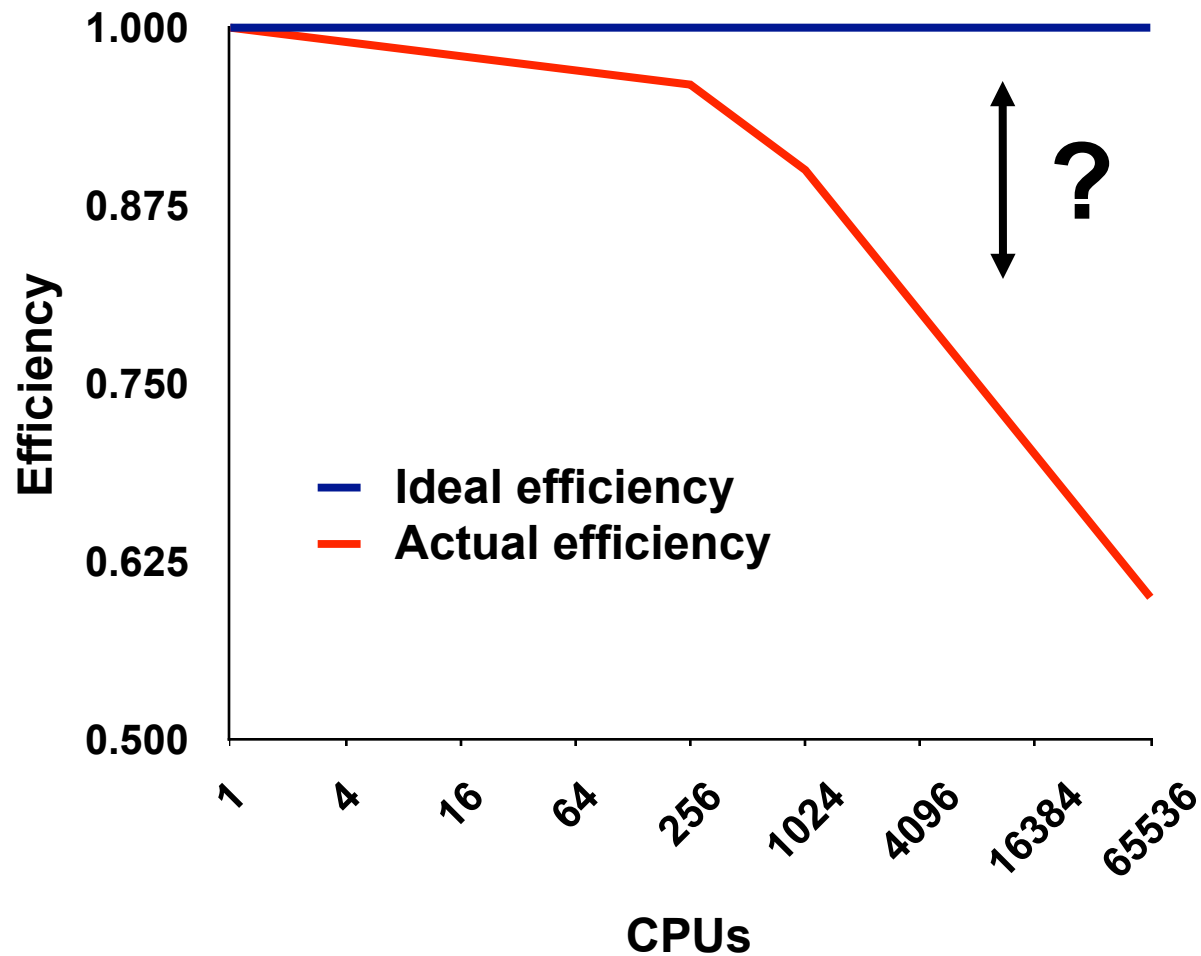
- Compact representation
- Summarizes costs for each subtree in the recursion
- $T_{\text{fib}(n-1)} / T_{\text{fib}(n-2)} = 1.619$
(within .1% of the golden ratio)



Outline

- Overview of Rice's HPCToolkit
- Accurate measurement
- Effective performance analysis
- Pinpointing scalability bottlenecks
 - scalability bottlenecks on large-scale parallel systems
 - scaling on multicore processors
- Understanding temporal behavior
- Assessing process variability
- Understanding threading, GPU, locks, and memory hierarchy
 - blame shifting
 - attributing memory hierarchy costs to data
- Summary and conclusions

The Problem of Scaling



Note: higher is better

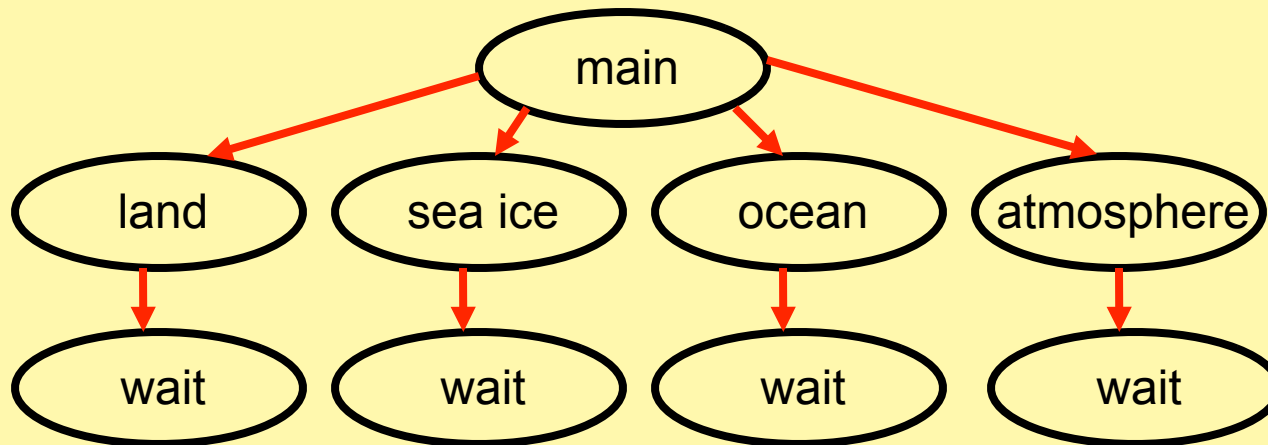
Goal: Automatic Scaling Analysis

- Pinpoint scalability bottlenecks
- Guide user to problems
- Quantify the magnitude of each problem
- Diagnose the nature of the problem

Challenges for Pinpointing Scalability Bottlenecks

- **Parallel applications**
 - modern software uses layers of libraries
 - performance is often context dependent

Example climate code skeleton

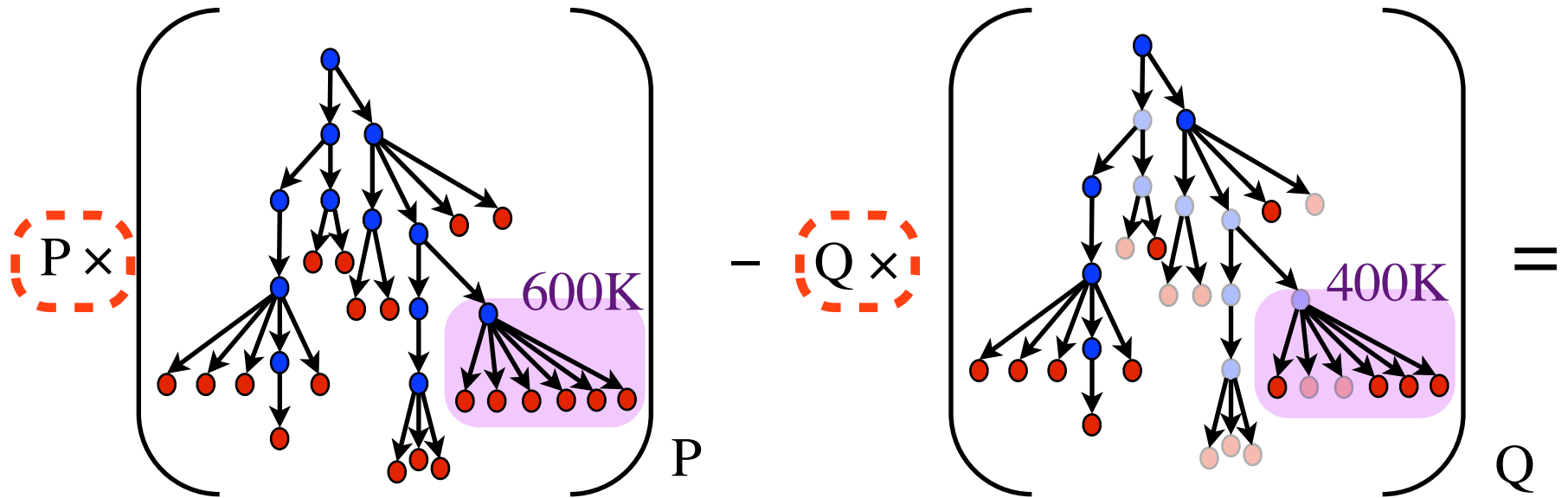


- **Monitoring**
 - bottleneck nature: computation, data movement, synchronization?
 - **2 pragmatic constraints**
 - acceptable data volume
 - low perturbation for use in production runs

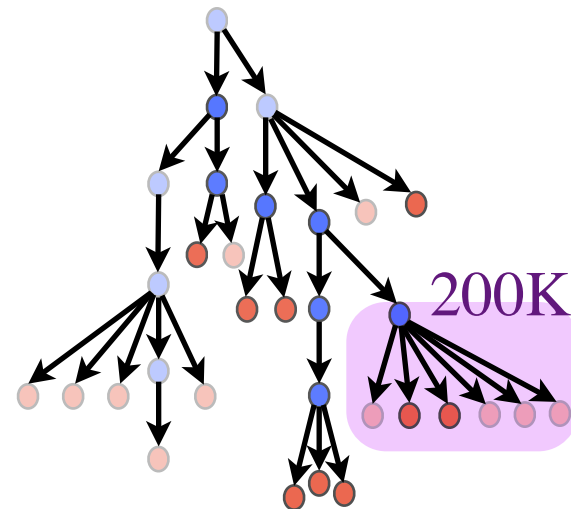
Performance Analysis with Expectations

- You have performance expectations for your parallel code
 - strong scaling: linear speedup
 - weak scaling: constant execution time
- Put your expectations to work
 - measure performance under different conditions
 - e.g. different levels of parallelism or different inputs
 - express your expectations as an equation
 - compute the deviation from expectations for each calling context
 - for both inclusive and exclusive costs
 - correlate the metrics with the source code
 - explore the annotated call tree interactively

Pinpointing and Quantifying Scalability Bottlenecks

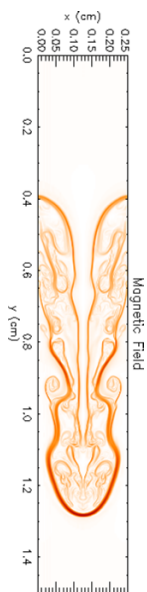


coefficients for analysis
of strong scaling

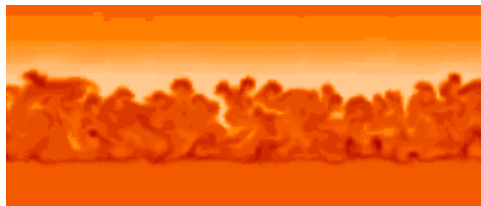


Scalability Analysis Demo

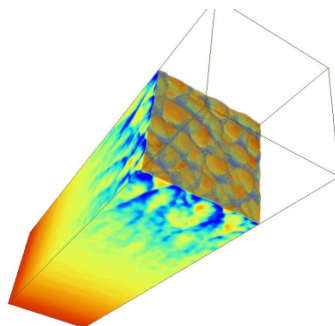
Code: University of Chicago FLASH
Simulation: white dwarf detonation
Platform: Blue Gene/P
Experiment: 8192 vs. 256 processors
Scaling type: weak



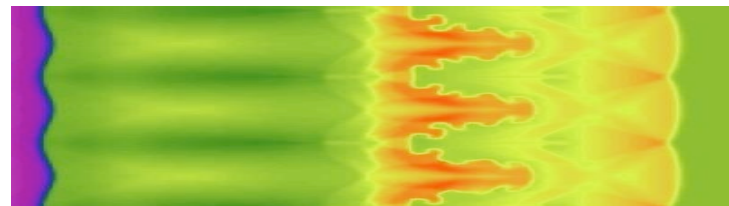
*Magnetic
Rayleigh-Taylor*



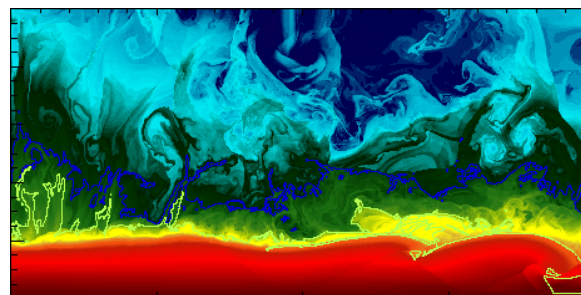
Nova outbursts on white dwarfs



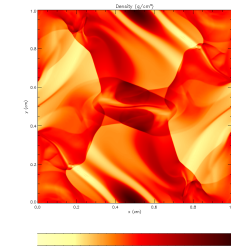
Cellular detonation



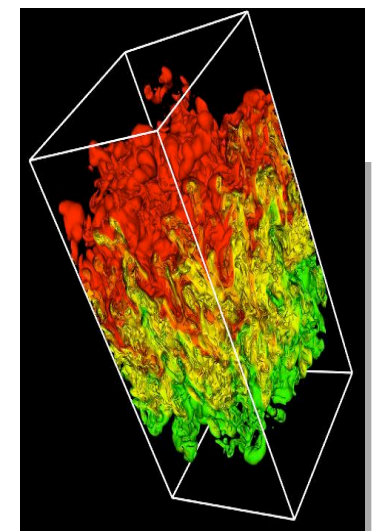
Laser-driven shock instabilities



Helium burning on neutron stars



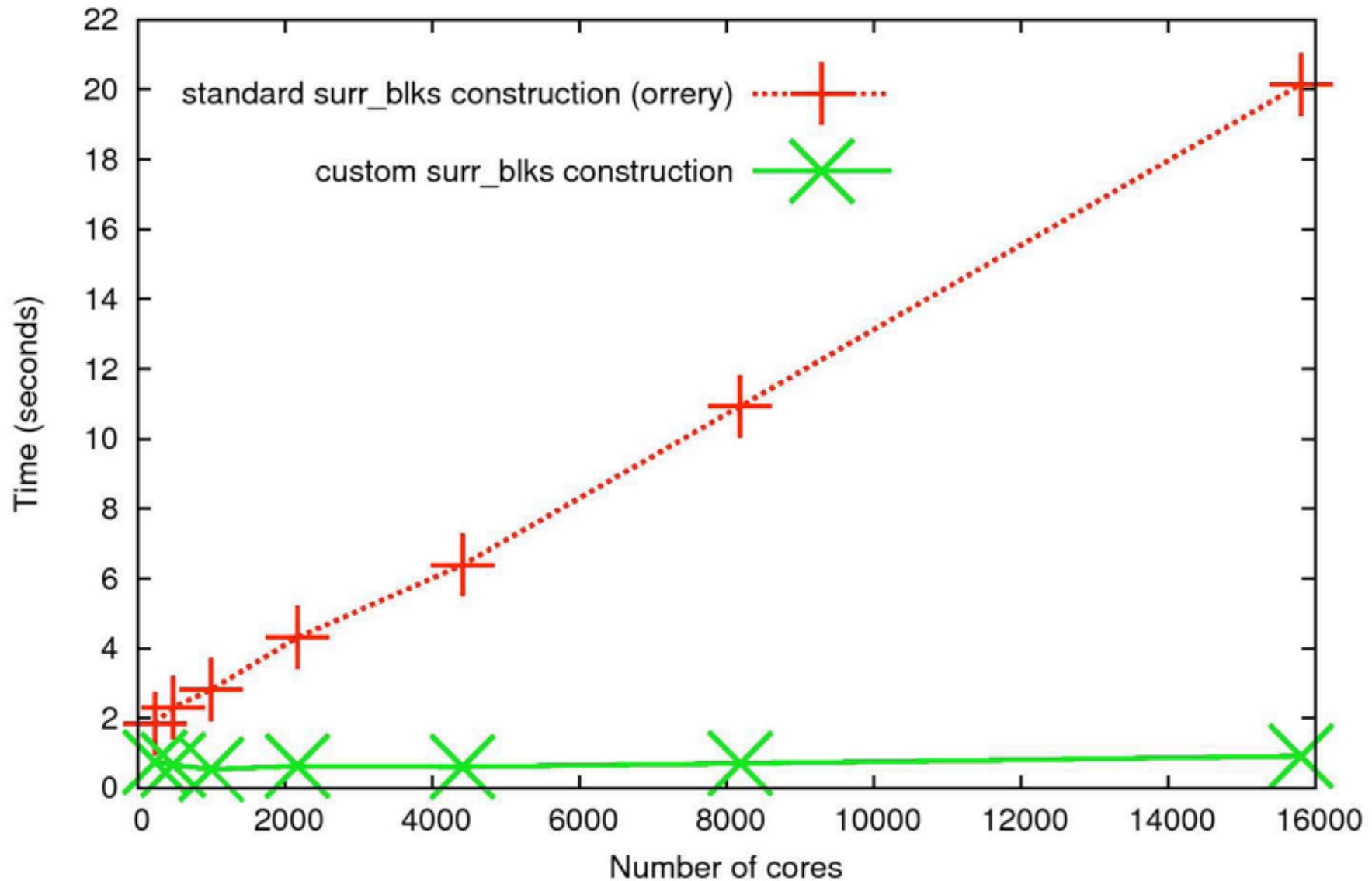
*Orzag/Tang MHD
vortex*



Rayleigh-Taylor instability

Figures courtesy of FLASH Team, University of Chicago

Improved Flash Scaling of AMR Setup



Graph courtesy of Anshu Dubey, U Chicago

Scaling on Multicore Processors

- **Compare performance**
 - single vs. multiple processes on a multicore system
- **Strategy**
 - differential performance analysis
 - subtract the calling context trees as before, unit coefficient for each

S3D: Multicore Losses at the Procedure Level

The screenshot shows the hpcviewer interface. The top pane displays the source code for the subroutine `rhsf`. The bottom pane shows a performance table with columns for Scope, 1-core (ms) (I), 1-core (ms) (E), 8-core(1) (ms) (I), 8-core(1) (ms) (E), and Multicore Loss. The `rhsf` row is highlighted, and its Multicore Loss of 13.0% is circled in blue. A red arrow points from the 1-core (ms) (I) value of 1.11e08 to the 8-core(1) (ms) (I) value of 1.77e08.

```
1 subroutine rhsf( q, rhs )
2 !-----
3 ! Changes
4 ! Ramanan Sankaran - 01/04/05
5 ! 1. Diffusive fluxes are computed without having to convert units.
6 ! Ignore older comments about conversion to CGS units.
7 ! This saves a lot of flops.
8 ! 2. Mixavg and Lewis transport modules have been made interchangeable
9 ! by adding dummy arguments in both.
10 !-----
11 !           Author: James Sutherland
12 !           Date:   April, 2002
13 !-----
14 ! This routine calculates the time rate of change for the
15 ! momentum, continuity, energy, and species equations.
16 !
```

Scope	1-core (ms) (I)	1-core (ms) (E)	8-core(1) (ms) (I)	8-core(1) (ms) (E)	Multicore Loss
Experiment Aggregate Metrics	1.11e08 100 %	1.11e08 100 %	1.88e08 100 %	1.88e08 100 %	7.64e07 100 %
rhsf	1.07e08 96.5%	6.60e06 5.9%	1.77e08 94.1%	1.65e07 8.8%	9.92e06 13.0%
diffflux_proc_looptool	2.86e06 2.6%	2.86e06 2.6%	8.12e06 4.3%	8.12e06 4.3%	5.27e06 6.9%
integrate_erk_jstage_lt	1.09e08 98.1%	1.25e06 1.1%	1.84e08 97.9%	5.94e06 3.2%	4.70e06 6.1%
GET_MASS_FRAC.in.VARIABLES_M	1.49e06 1.3%	1.49e06 1.3%	6.08e06 3.2%	6.08e06 3.2%	4.59e06 6.0%
ratx	1.01e07 9.1%	1.00e07 9.0%	4.41e07 23.5%	1.40e07 7.4%	3.95e06 5.2%
qssa	3.52e06 3.2%	3.52e06 3.2%	5.71e06 3.0%	5.71e06 3.0%	2.18e06 2.9%
ratt	3.26e07 29.2%	1.48e07 13.3%	4.38e07 23.3%	1.66e07 8.8%	1.76e06 2.3%
CALC_INV_AVG_MOL_WT.in.THER	9.70e05 0.9%	9.70e05 0.9%	2.68e06 1.4%	2.68e06 1.4%	1.70e06 2.2%
computeheatflux_looptool	1.46e06 1.3%	1.46e06 1.3%	2.88e06 1.5%	2.88e06 1.5%	1.41e06 1.8%
rdwdot	3.09e06 2.8%	3.09e06 2.8%	4.33e06 2.3%	4.33e06 2.3%	1.24e06 1.6%

Execution time
increases 1.65x in
subroutine rhsf

subroutine rhsf
accounts for 13.0% of
the multicore scaling
loss in the execution

S3D: Multicore Losses at the Loop Level

The screenshot displays the hpcviewer interface. The top pane shows Fortran code with a loop at lines 197-223 highlighted. The bottom pane shows a performance table with columns for Scope, 1-core (ms) (I), 1-core (ms) (E), 8-core(1) (ms) (I), 8-core(1) (ms) (E), and Multicore Loss. A red arrow points from the highlighted code loop to the corresponding row in the table.

Execution time increases 2.8x in the loop that scales worst

loop contributes 6.9% of the scaling loss for the whole execution

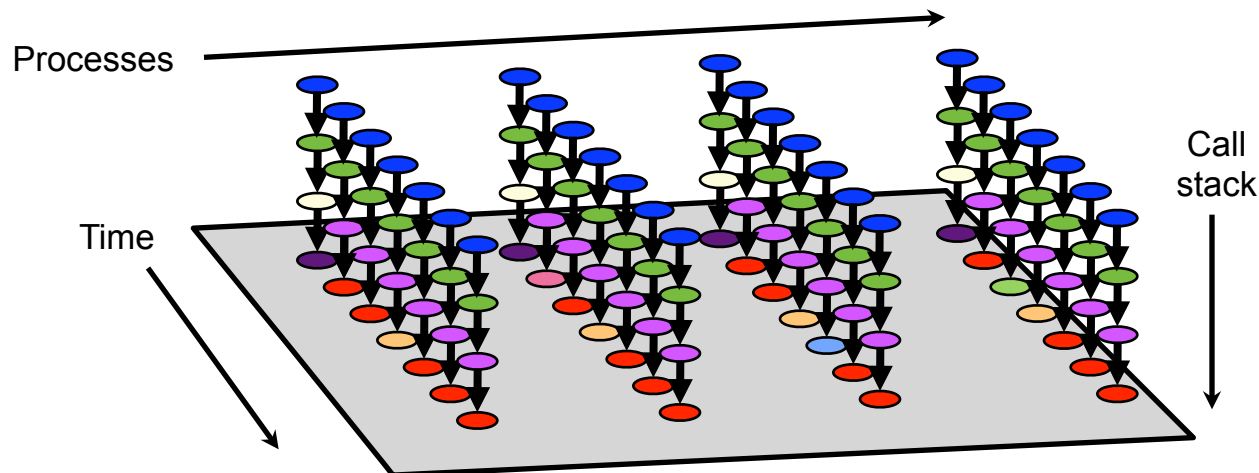
Scope	1-core (ms) (I)	1-core (ms) (E)	8-core(1) (ms) (I)	8-core(1) (ms) (E)	Multicore Loss
loop at diffflux_gen_uj.f: 197-223	2.86e06 2.6%	2.86e06 2.6%	8.12e06 4.3%	8.12e06 4.3%	5.27e06 6.9%
loop at integrate_erk_jstage_lt_ge	1.09e08 98.1%	1.25e06 1.1%	1.84e08 97.9%	5.94e06 3.2%	4.70e06 6.1%
loop at variables_m.f90: 88-99	1.49e06 1.3%	1.49e06 1.3%	6.08e06 3.2%	6.08e06 3.2%	4.60e06 6.0%
loop at rhsf.f90: 516-536	2.70e06 2.4%	1.31e06 1.2%	6.49e06 3.5%	3.72e06 2.0%	2.41e06 3.1%
loop at rhsf.f90: 538-544	3.35e06 3.0%	1.45e06 1.3%	7.06e06 3.8%	3.82e06 2.0%	2.36e06 3.1%
loop at rhsf.f90: 546-552	2.56e06 2.3%	1.47e06 1.3%	5.86e06 3.1%	3.42e06 1.8%	1.96e06 2.6%
loop at thermchem_m.f90: 127-1	8.00e05 0.7%	8.00e05 0.7%	2.28e06 1.2%	2.28e06 1.2%	1.48e06 1.9%
loop at heatflux_lt_gen.f: 5-132	1.46e06 1.3%	1.46e06 1.3%	2.88e06 1.5%	2.88e06 1.5%	1.41e06 1.8%
loop at rhsf.f90: 576	6.65e05 0.6%	6.65e05 0.6%	1.87e06 1.0%	1.87e06 1.0%	1.20e06 1.6%
loop at getrates.f: 504-505	8.00e06 7.2%	8.00e06 7.2%	8.74e06 4.7%	8.74e06 4.7%	7.35e05 1.0%
loop at derivative_x.f90: 213-690	1.78e06 1.6%	1.78e06 1.6%	2.47e06 1.3%	2.47e06 1.3%	6.95e05 0.9%

Outline

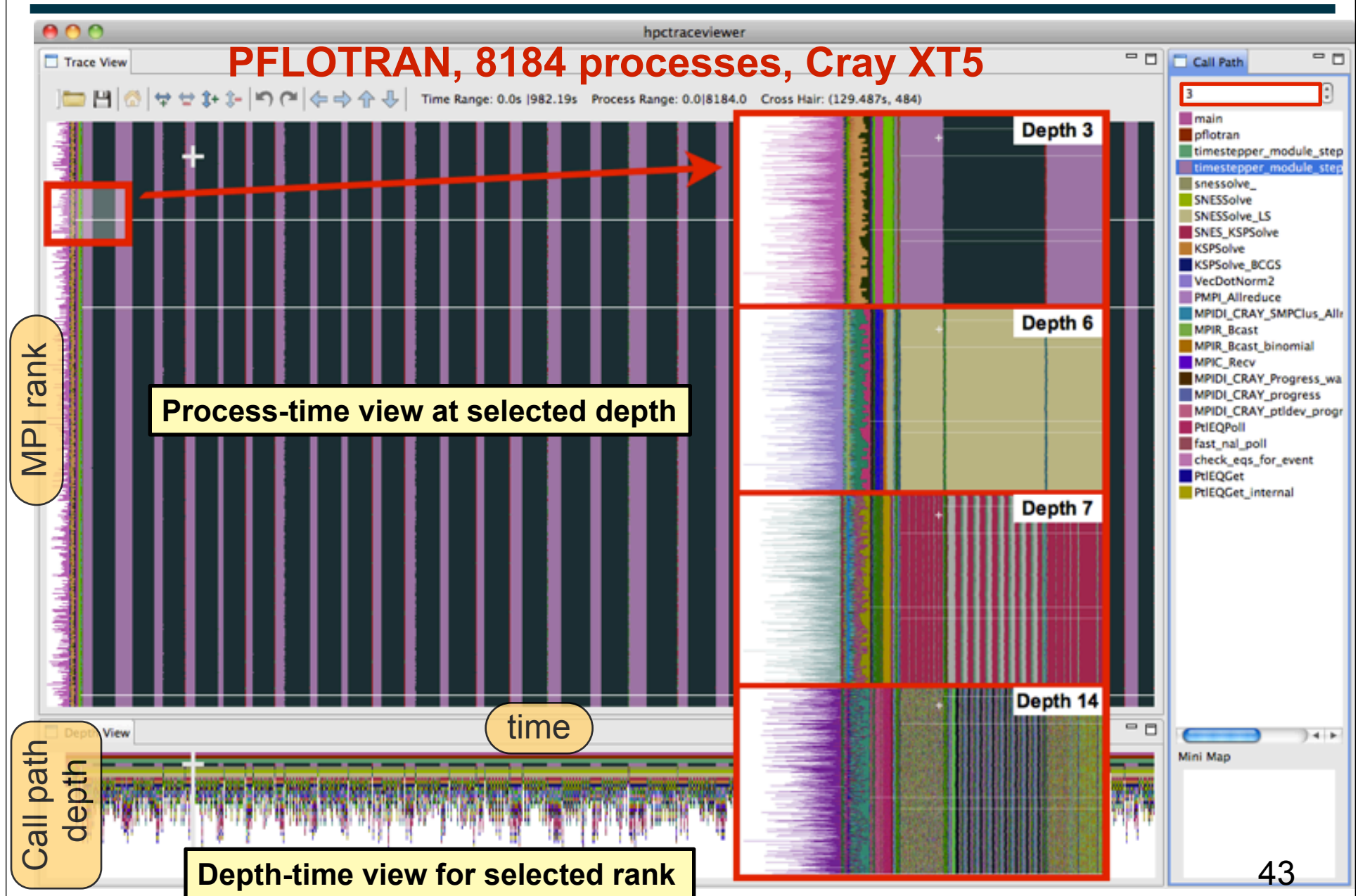
- Overview of Rice's HPCToolkit
- Accurate measurement
- Effective performance analysis
- Pinpointing scalability bottlenecks
 - scalability bottlenecks on large-scale parallel systems
 - scaling on multicore processors
- Understanding temporal behavior
- Assessing process variability
- Understanding threading, GPU, locks, and memory hierarchy
 - blame shifting
 - attributing memory hierarchy costs to data
- Summary and conclusions

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

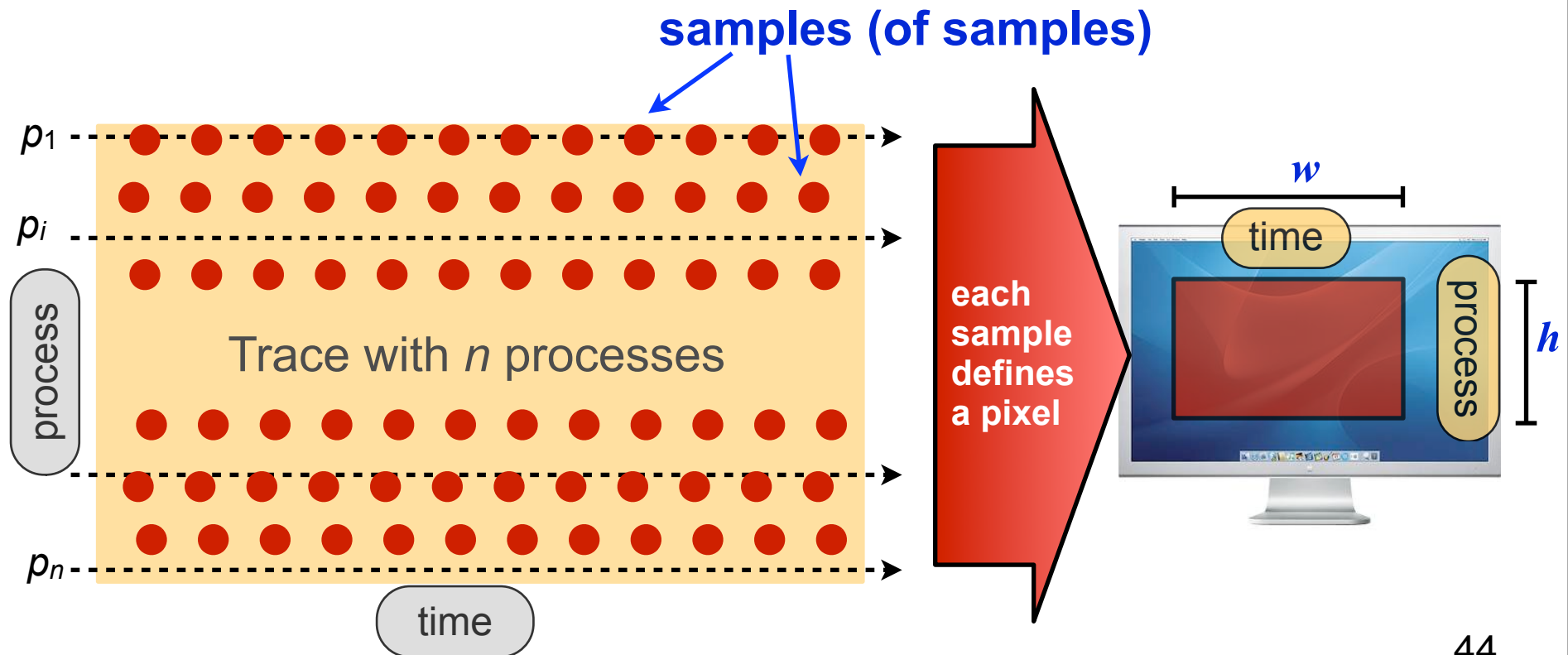


Exposes Temporal Call Path Patterns

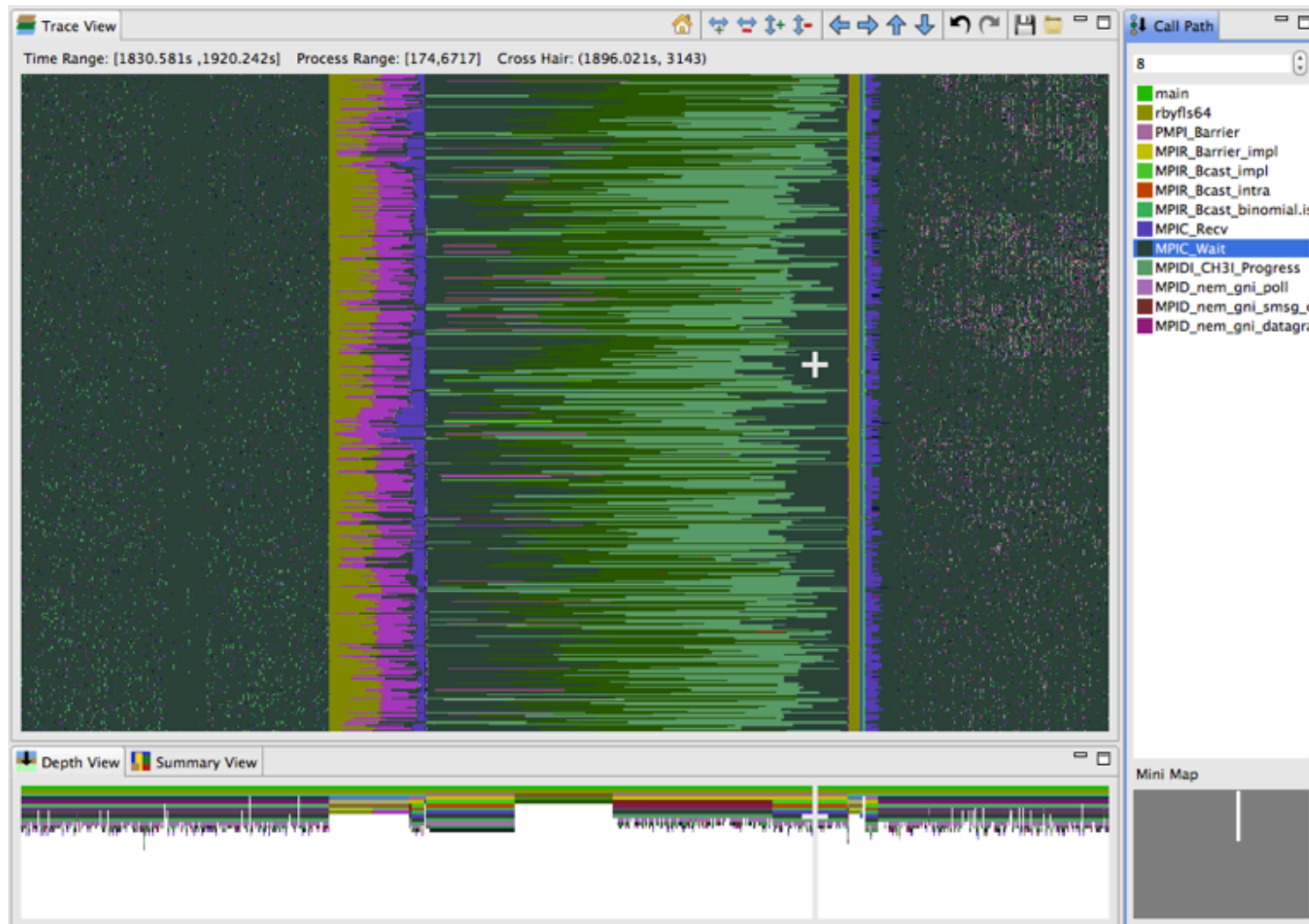


Presenting Large Traces on Small Displays

- How to render an arbitrary portion of an arbitrarily large trace?
 - we have a display window of dimensions $h \times w$
 - typically many more processes (or threads) than h
 - typically many more samples (trace records) than w
- Solution: sample the samples!



MPBS: 16K cores @ 50 min



NOTES:

- (1) I/O in this execution to /dev/null (to show we can scale without burning hours writing application data files)
- (2) panel above shows a zoomed view of an execution detail.

Outline

- Overview of Rice's HPCToolkit
- Accurate measurement
- Effective performance analysis
- Pinpointing scalability bottlenecks
 - scalability bottlenecks on large-scale parallel systems
 - scaling on multicore processors
- Understanding temporal behavior
- Assessing process variability
- Understanding threading, GPU, locks, and memory hierarchy
 - blame shifting
 - attributing memory hierarchy costs to data
- Summary and conclusions

Example: Massive Parallel Bucket Sort (MPBS)

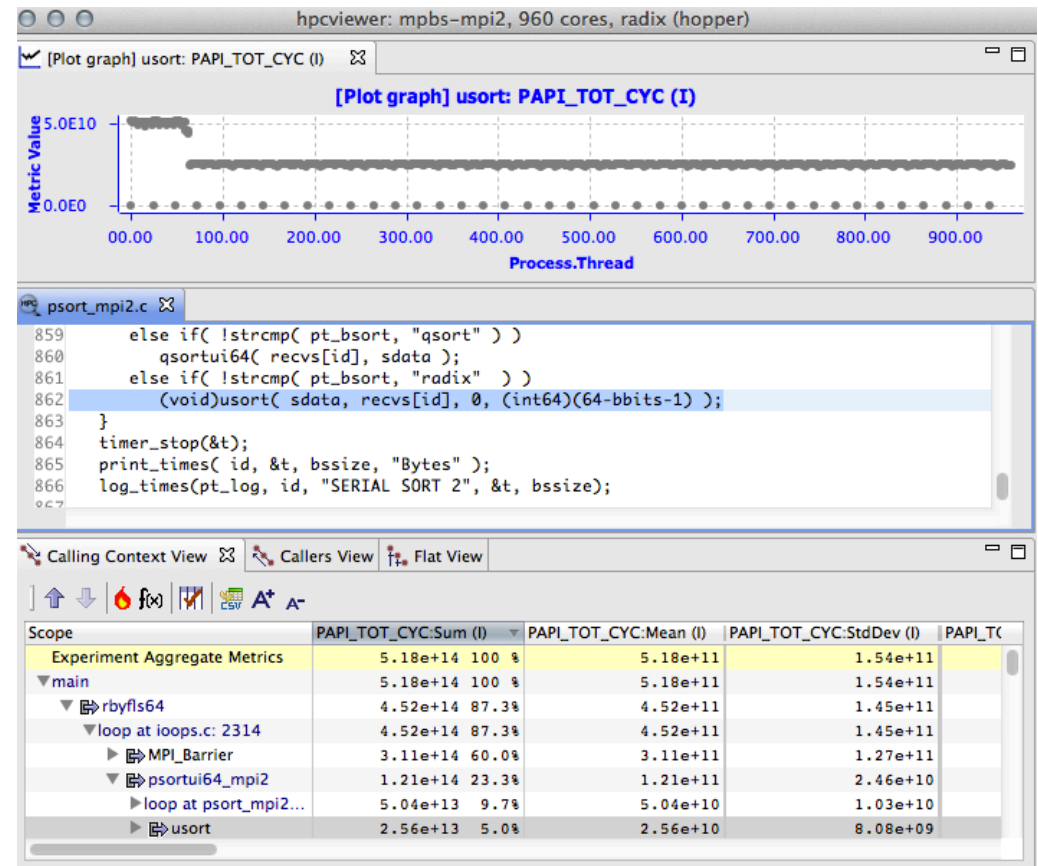
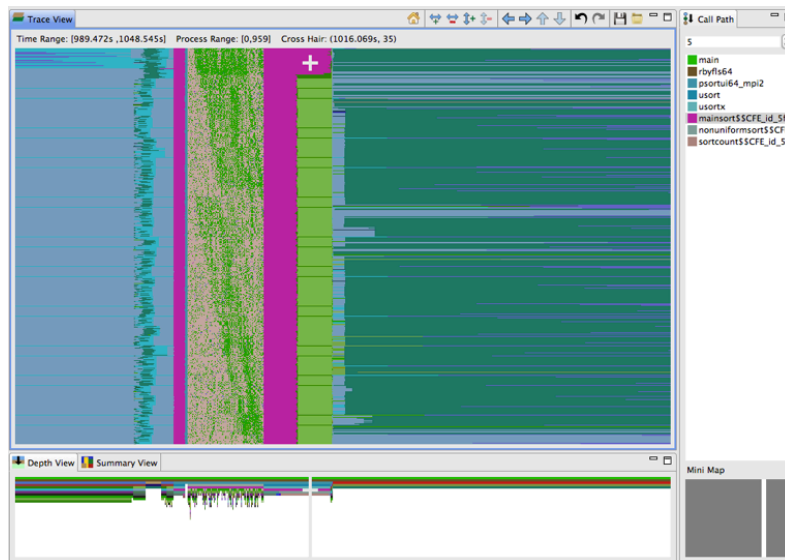
Program execution consists of two phases

- **Produces a large number of files**
 - each file has a fixed numbered sequence of buckets
 - each bucket has a fixed number of records
 - each record is a 4, 8, or 16-byte integer
 - each file produced by sequentially filling each bucket with integer records
 - most significant bits set to bucket number
 - file complete when all buckets filled and file written to disk
- **Performs a two-stage sort on the contents of all files**
 - records are sorted for a given bucket number across all of the generated files
 - then written to a single file
 - this is repeated for each bucket
 - this yields a single sorted file as a result

Sample execution: radix sort, 960 cores, 512MB/core

MPBS @ 960 cores, radix sort

Two views of load imbalance since not on a 2^k cores



Outline

- Overview of Rice's HPCToolkit
- Accurate measurement
- Effective performance analysis
- Pinpointing scalability bottlenecks
 - scalability bottlenecks on large-scale parallel systems
 - scaling on multicore processors
- Understanding temporal behavior
- Assessing process variability
- Understanding threading, GPU, locks, and memory hierarchy
 - blame shifting
 - attributing memory hierarchy costs to data
- Summary and conclusions

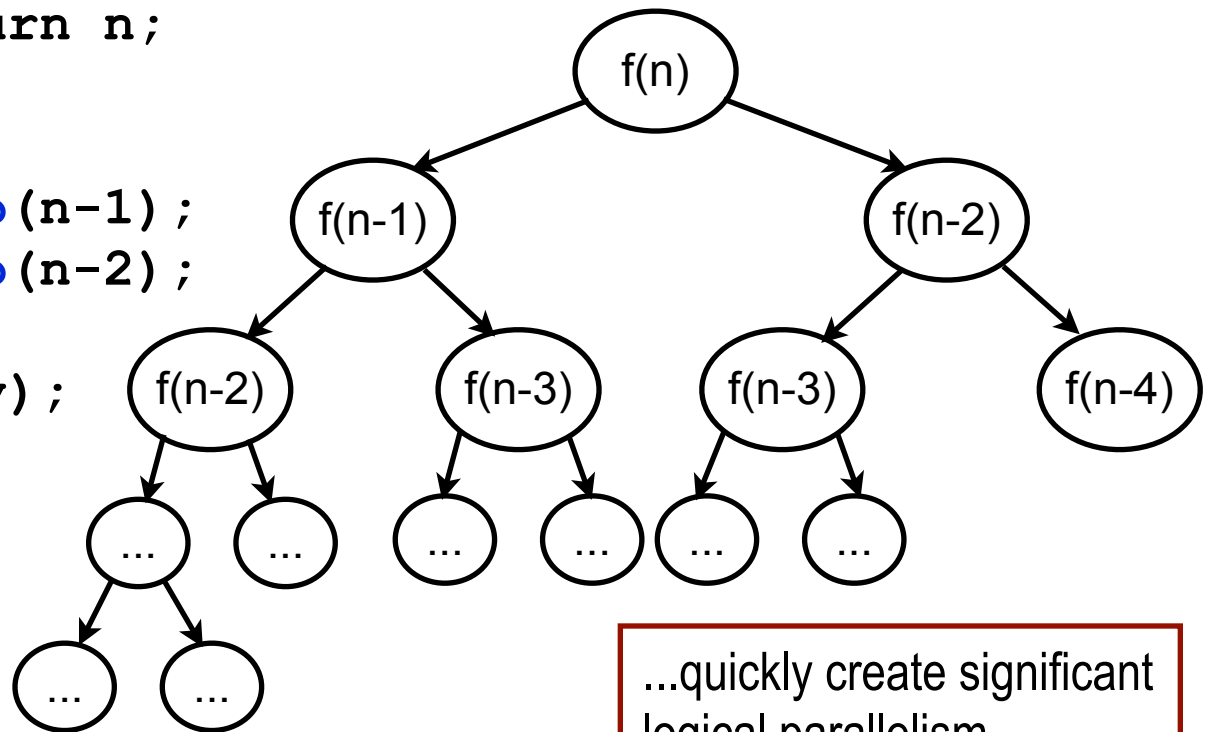
Blame Shifting

- **Problem:** in many circumstances sampling measures symptoms of performance losses rather than causes
 - worker threads waiting for work
 - threads waiting for a lock
 - MPI process waiting for peers in a collective communication
- **Approach:** shift blame for losses from victims to perpetrators
- **Flavors**
 - active measurement
 - analysis only

Cilk: A Multithreaded Language

```
cilk int fib(n) {  
  if (n < 2) return n;  
  else {  
    int x, y;  
    x = spawn fib(n-1);  
    y = spawn fib(n-2);  
    sync;  
    return (x + y);  
  }  
}
```

asynchronous calls
create logical tasks that
only block at a **sync**...



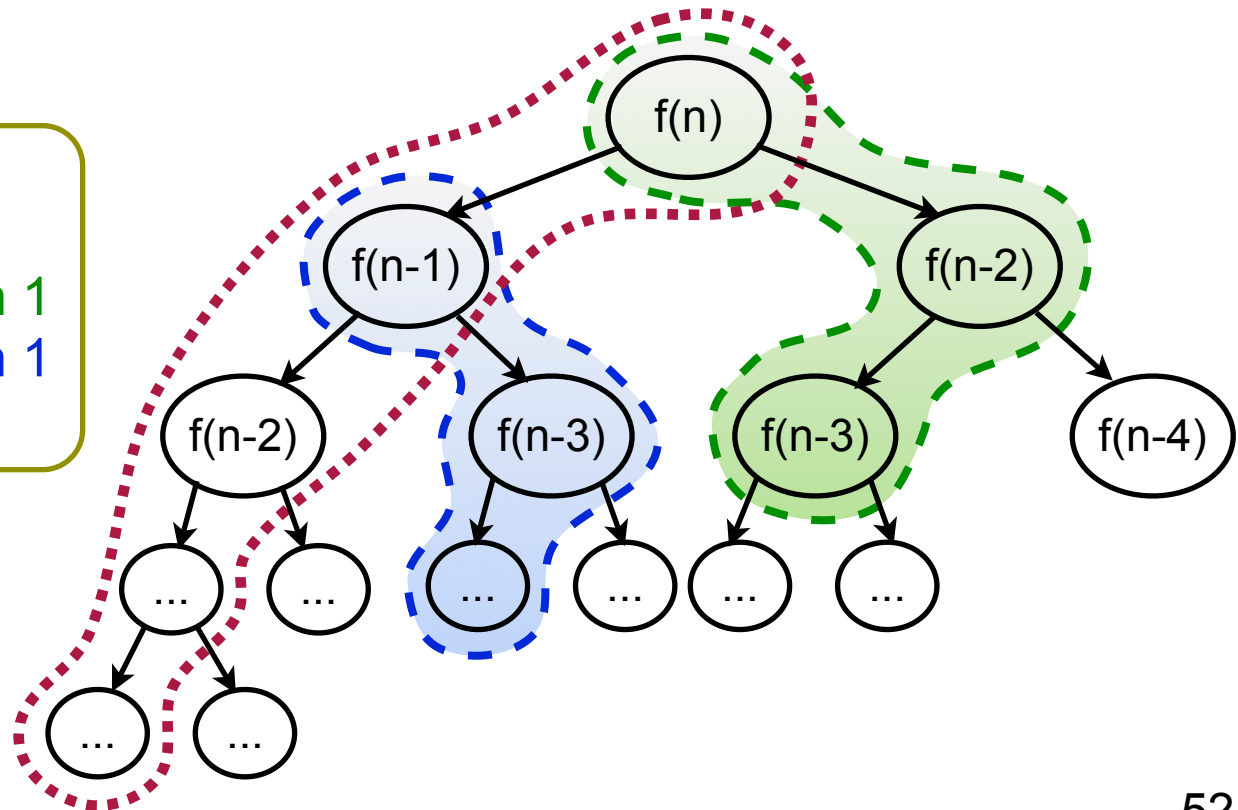
...quickly create significant
logical parallelism.

Cilk Program Execution using Work Stealing

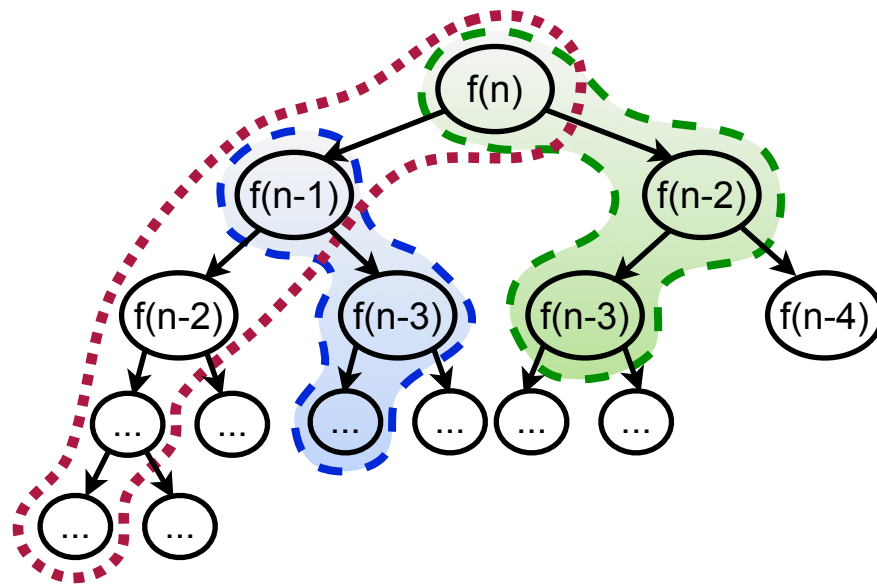
- Challenge: Mapping logical tasks to compute cores
- Cilk approach:
 - lazy thread creation plus work-stealing scheduler
 - **spawn**: a potentially parallel task is available
 - an idle thread steals tasks from a random working thread

Possible Execution:

thread 1 begins
thread 2 steals from 1
thread 3 steals from 1
etc...



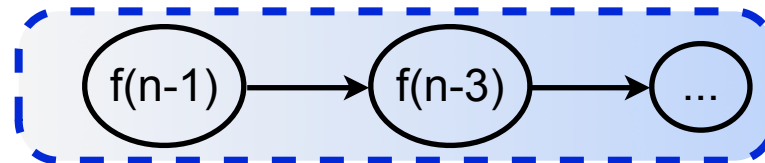
Wanted: Call Path Profiles of Cilk



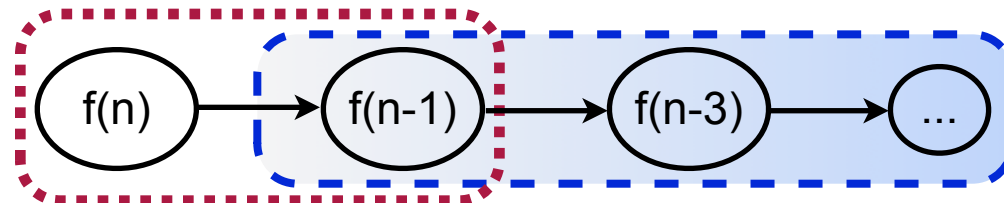
thread 1
thread 2
thread 3

Work stealing *separates*
user-level calling contexts in
space and time

- Consider **thread 3**:
 - physical call path:



- logical call path:



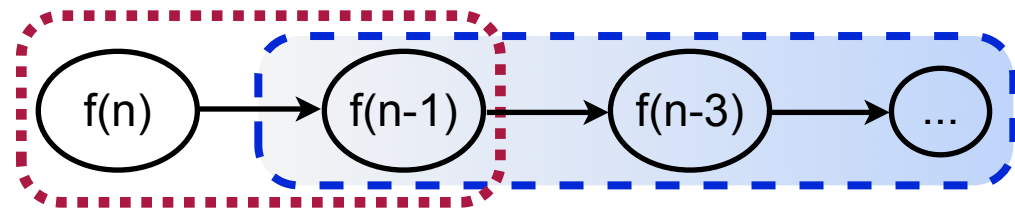
Logical call path profiling: Recover *full* relationship
between *physical* and *user-level* execution

Effective Performance Analysis

Three Complementary Techniques:

- Recover *logical calling contexts* in presence of work-stealing

```
cilk int fib(n) {  
  if (n < 2) {...}  
  else {  
    int x, y;  
    x = spawn fib(n-1);  
    y = spawn fib(n-2);  
    sync;  
    return (x + y);  
  }  
}
```



high parallel overhead from
creating many small tasks

- Quantify *parallel idleness* (insufficient parallelism)
- Quantify *parallel overhead*
- Attribute *idleness* and *overhead* to *logical contexts*
— at the source level

Measuring & Attributing Parallel Idleness

- **Metrics: Effort = “work” + “idleness”**
 - associate metrics with user-level calling contexts
 - **insight:** attribute idleness to its cause: context of *working* thread
 - a thread looks past itself when ‘bad things’ happen to others
- **Work stealing-scheduler: one thread per core**
 - maintain W (# working threads) and I (# idling threads)
 - slight modifications to work-stealing run time
 - atomically incr/decr W when thread exits/enters scheduler
 - when a sample event interrupts a working thread
 - $I = \text{\#cores} - W$
 - apportion others' idleness to me: I / W
- **Example: Dual quad-cores; on a sample, 5 are *working*:**



for each $\mathcal{W} += 1$ $\sum \mathcal{W} = 5$
worker: $\mathcal{I} += 3/5$ $\sum \mathcal{I} = 3$



idle: drop sample
(it's in the scheduler!)

Parallel Overhead

- **Parallel overhead**
 - **when a thread works on something other than user code**
 - **(we classify waiting for work as idleness)**
- **Pinpointing overhead with call path profiling**
 - **impossible, without prior arrangement**
 - **work and overhead are both machine instructions**
 - **insight: have compiler tag instructions as overhead**
 - **quantify samples attributed to instructions that represent ovhd**
 - **use post-mortem analysis**

Blame Shifting: Lulesh OpenMP Code

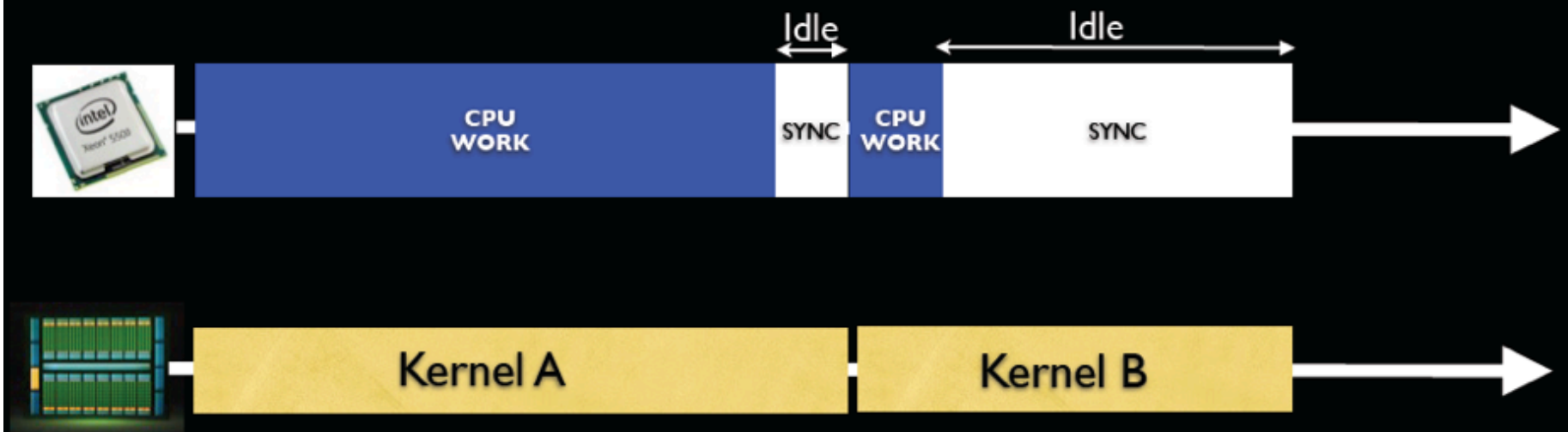
HPCToolkit OpenMP Metrics Explained

- **p_req_core_idleness**
 - idleness for each parallel region is measured with respect to the maximum number of threads ever requested for a parallel region. The number of threads for a parallel region is specified by `omp_set_num_threads`, `OMP_NUM_THREADS`, or (by default) the number of cores on the node.
- **p_all_core_idleness**
 - idleness for each parallel region is measured with respect to the total number of cores on the node.
- **p_all_thread_idleness**
 - idleness for each parallel region is measured with respect to the number of threads employed for that parallel region.
- **p_work**
 - useful work performed by the thread
- **p_overhead**
 - work performed by the thread on behalf of the OpenMP runtime system shared library

Blame Shifting for Hybrid Codes

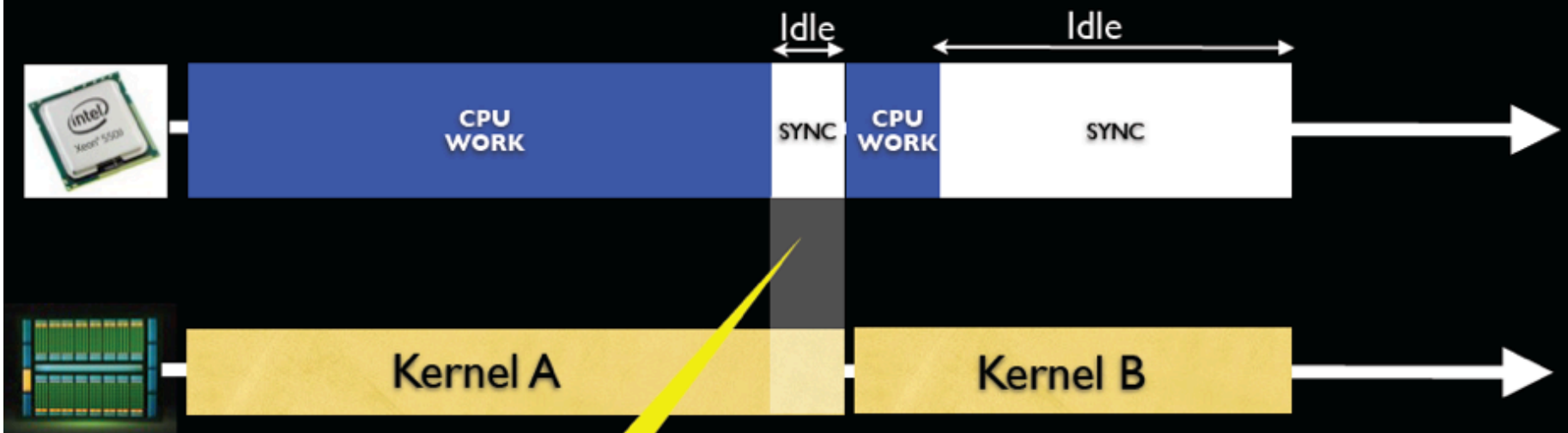
- If GPU is idle, code executing on CPU is responsible for not offloading (enough) work to GPU
 - ✦ Attribute blame to CPU code executing while GPU is idle
- If CPU is idle waiting for GPU kernel(s) to finish, executing GPU kernel(s) are responsible for CPU idleness
 - ✦ Attribute proportional blame to each such kernels
- Credit codes that are well overlapped

Performance Expectations for Hybrid Code with Blame Shifting



Top GPU-kernel may not be the best candidate for tuning

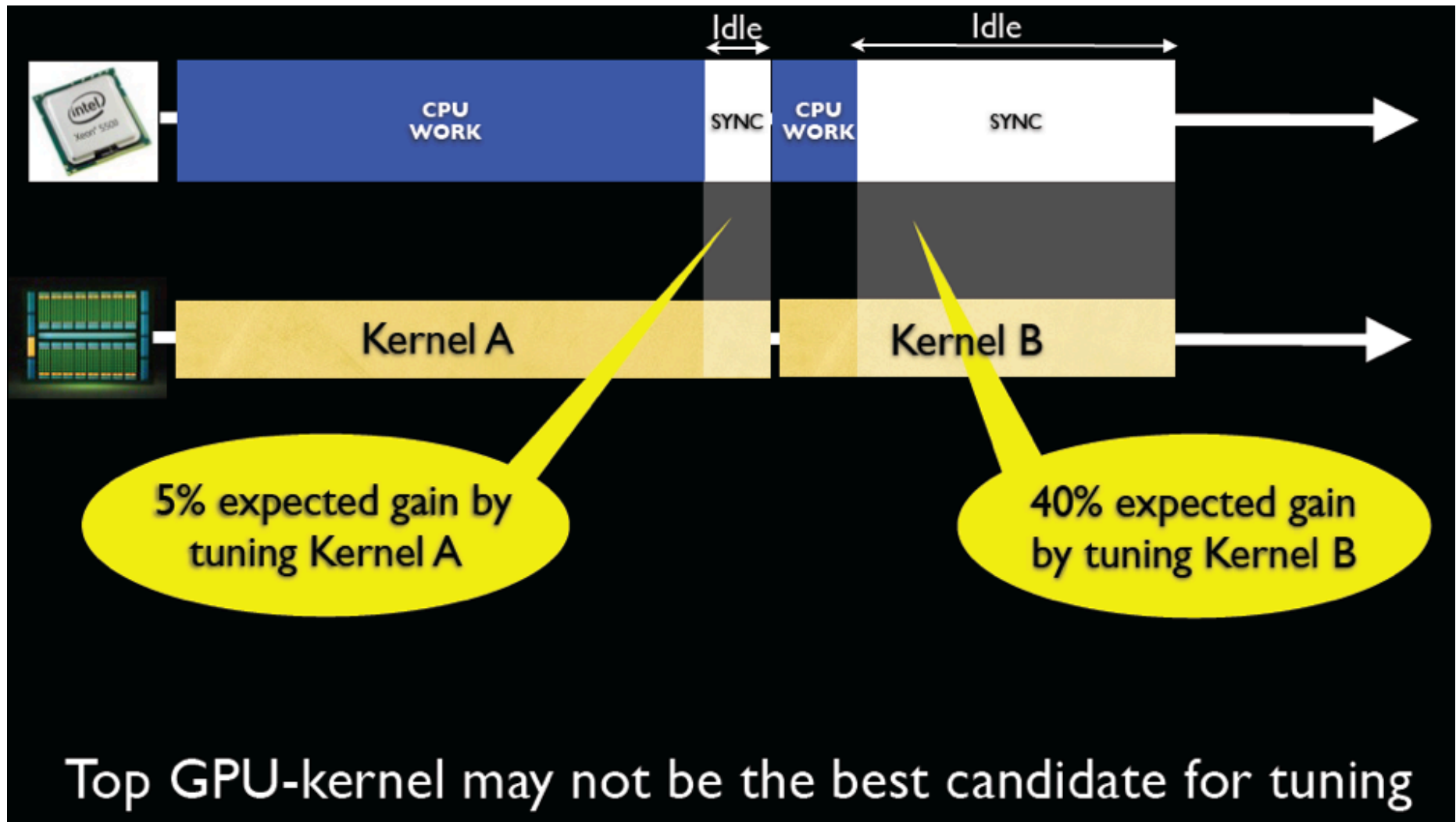
Performance Expectations for Hybrid Code with Blame Shifting



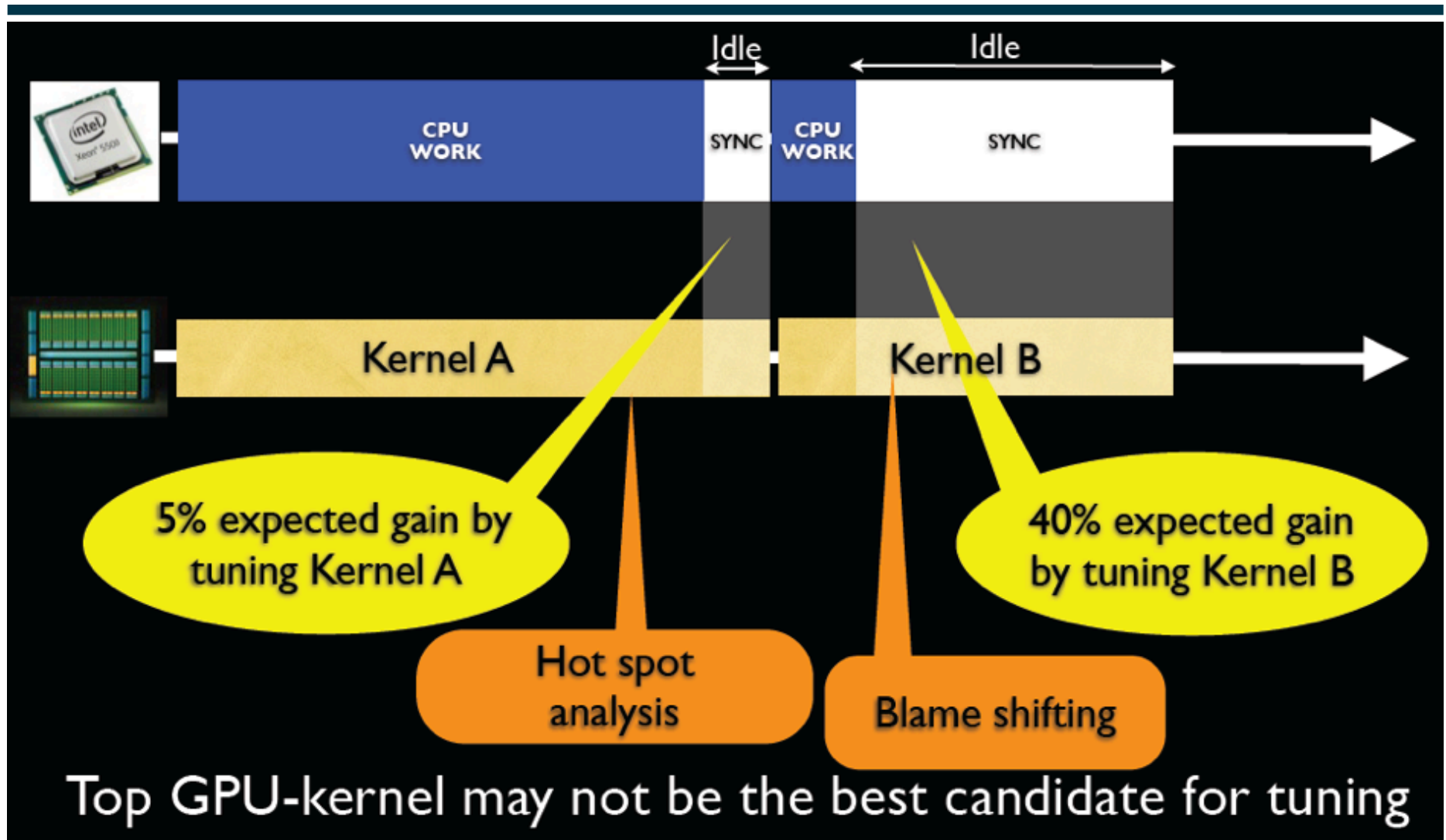
5% expected gain by
tuning Kernel A

Top GPU-kernel may not be the best candidate for tuning

Performance Expectations for Hybrid Code with Blame Shifting



Performance Expectations for Hybrid Code with Blame Shifting



HPCToolkit GPU Metrics Explained

- **CPU_IDLE (CI)**
 - When a sample event occurs in a CPU context C, this metric is incremented for C if the CPU thread is waiting for some GPU activity to finish.
- **CPU_IDLE_CAUSE (CIC)**
 - When a sample event occurs while the CPU is waiting in a context C, this metric is incremented for each context G that launched a kernel active on a GPU.
- **GPU_IDLE_CAUSE (GIC)**
 - When a sample event occurs in a CPU context C, this metric is incremented for C when there are no active GPU kernels.
- **OVERLAPPED_CPU (OC)**
 - When a sample event occurs in a CPU context C, this metric is incremented for C when CPU thread is not waiting for a GPU that has some unfinished activity.
- **OVERLAPPED_GPU (OG)**
 - When a sample event occurs in a CPU context C, this metric is incremented this metric is incremented for each context G that launched a kernel active on the GPU if the CPU thread is not waiting for GPU.
- **GPU_ACTIVITY_TIME (GAT)**
 - This metric is increased by T for the GPU context that launched a kernel K, where T is the time K spent executing.
- **H_TO_D_BYTES (H2D)**
 - This metric is incremented by bytes transferred from CPU to GPU, and attributed to the calling context where the host to device memory copy was invoked.
- **D_TO_H_BYTES (D2H)**
 - This metric is incremented by bytes transferred from GPU to CPU and attributed to the calling context where device to host memory copies were invoked.

Note, that we don't have a GPU_IDLE metric (unlike CPU_IDLE), because when the GPU is idle, there is clearly no code executing on it, contrary to that when CPU is idle, it makes sense to show where the CPU was idling.

Hybrid Code Demo:
Lulesh: CPU/GPU blame shifting
LAMMPS: Tracing

LAMMPS Slow GPU Copies on Keeneland

- **From Keeneland support staff:**
“My first guess is that those nodes had GPUs that weren't seated correctly -- instead of PCIe x16, they only had PCIe x8 or less”
- **Sample related error log messages:**
 - PCIE (needs GPU reseal)
 - kid036 : GPU 0 has incorrect PCIe width
 - kid036 : GPU 0 has low bandwidth (< 5.0GB/s) : 3.08614
 - kid058 : GPU 0 has incorrect PCIe width
 - kid058 : GPU 0 has low bandwidth (< 5.0GB/s) : 0.405049
 - kid105 : GPU 0 has incorrect PCIe width
 - kid105 : GPU 0 has low bandwidth (< 5.0GB/s) : 0.813523

64 MPI processes



Blame Shifting to Understand Lock Contention

- **Lock contention causes idleness**
 - explicitly threaded programs (Pthreads, etc)
 - implicitly threaded programs (critical sections in OpenMP, Cilk...)
- **Use “blame-shifting” to shift blame from victim to perpetrator**
 - use shared state (locks) to communicate blame
- **How it works**
 - consider spin-waiting
 - sample a working thread:
 - charge to ‘work’ metric
 - sample an idle thread
 - accumulate in idleness counter assoc. with lock (atomic add)
 - working thread releases a lock
 - atomically swap 0 with lock’s idleness counter
 - exactly represents contention while that thread held the lock
 - unwind the call stack to attribute lock contention to a calling context

Lock contention in MADNESS

```
578     add(MEMFUN_OBJT(memfunT)& obj,  
579         memfunT memfun,  
580         const arg1T& arg1, const arg2T& arg2, const arg3T& arg3, const TaskAttributes&  
581         Future<REMFUTURE(MEMFUN_RETURNT(memfunT))> result;  
582         add(new TaskMemfun<memfunT>(result,obj,memfun,arg1,arg2,arg3,attr));  
583         return result;  
584     }
```

quantum chemistry; MPI + pthreads

Calling Context View Callers View Flat View

16 cores; 1 thread/core (4 x Barcelona)

µs

Scope	...	% idleness (all/E).%	idleness (all/E)
Experiment Aggregate Metrics		2.35e+01 100 %	1.57e+09 100 %
▼ pthread_spin_unlock		2.35e+01 100.0	
▼ madness::Spinlock::unlock() const		2.35e+01 100.0	
▼ inlined from worldmutex.h: 142		1.78e+01 75.6%	
▼ madness::ThreadPool::add(madness::PoolTaskInterface*)		1.78e+01 75.6%	
▼ inlined from worldtask.h: 581		7.35e+00 31.2%	4.92e+08 31.2%
▶ madness::Future<> madness::WorldObject<>::task<>		7.35e+00 31.2%	4.92e+08 31.2%
▼ inlined from worldtask.h: 569		4.56e+00 19.4%	3.09e+07 19.4%
▶ madness::Future<> madness::WorldObject<>::task<>		4.56e+00 19.4%	3.09e+07 19.4%
▶ inlined from worlddep.h: 68		1.53e+00 6.5%	1.02e+07 6.5%
▼ inlined from worldtask.h: 570		1.49e+00 6.3%	9.97e+07 6.3%
▶ madness::Future<> madness::WorldObject<>::task<>		1.49e+00 6.3%	9.97e+07 6.3%
▶ inlined from worldtask.h: 558		1.38e+00 5.9%	9.26e+07 5.9%
▶ madness::Future<> madness::WorldTaskQueue::add<>(ma		6.72e-01 2.9%	4.49e+07 2.9%

lock contention accounts for **23.5%** of execution time.

Adding futures to shared global work queue.

Outline

- Overview of Rice's HPCToolkit
- Accurate measurement
- Effective performance analysis
- Pinpointing scalability bottlenecks
 - scalability bottlenecks on large-scale parallel systems
 - scaling on multicore processors
- Understanding temporal behavior
- Assessing process variability
- Understanding threading, GPU, locks, and memory hierarchy
 - blame shifting
 - attributing memory hierarchy costs to data
- Summary and conclusions

Data Centric Analysis

- **Goal: associate memory hierarchy performance losses with data**
- **Approach**
 - **intercept allocations to associate with their data ranges**
 - **associate latency with data using “instruction-based sampling” on AMD Opteron CPUs**
 - **identify instances of loads and store instructions**
 - **identify the data structure an access touches based on L/S address**
 - **measure the total latency associated with each L/S**
 - **present quantitative results using hpcviewer**

Data Centric Analysis of S3D

The screenshot displays the hpcviewer interface for the s3d_f90.x application. The top pane shows Fortran source code, with line 389, `yspec(:) = yspecies(i, j, k, :)`, highlighted. A blue arrow points from a text box to this line. The bottom pane shows the 'Calling Context View' with a tree of execution scopes. A red box highlights the entry for 'loop at chemkin_m.f90: 389', and a red arrow points from a text box to it. A table of performance metrics is visible on the right side of the bottom pane.

yspecies latency for this loop is 14.5% of total latency in program

41.2% of memory hierarchy latency related to yspecies array

Scope	LATENCY.[0,0] (v)	#(LD+ST).[0,0] (I)	#(LD+ST).[0,0] (E)	CACHE_MISS.[0,0] (I)
Experiment Aggregate Metrics	1.38e+05 100 %	5.02e+04 100 %	5.02e+04 100 %	9.92e+03 100 %
ALLOCATE_VARIABLES_ARRAYS.in.VARIABLES_M	5.68e+05 41.2%	9.40e+03 18.7%		3.14e+03 31.6%
solve_driver	5.68e+05 41.2%	9.40e+03 18.7%		3.14e+03 31.6%
loop at solve_driver.f90: 137	5.66e+05 41.0%	9.32e+03 18.6%		3.11e+03 31.4%
integrate	5.66e+05 41.0%	9.32e+03 18.6%		3.11e+03 31.4%
integrate_erk_jstage_it	5.36e+05 38.9%	8.51e+03 17.0%		2.92e+03 29.5%
loop at integrate_erk_jstage_it.gen.f: 47	5.36e+05 38.9%	8.51e+03 17.0%		2.92e+03 29.5%
rhsf	5.17e+05 37.4%	7.99e+03 15.9%		2.78e+03 28.0%
REACTION_RATE.in.CHEMKIN_M	2.57e+05 18.6%	1.24e+03 2.5%		1.20e+03 12.1%
REACTION_RATE_BOUNDS.in.CHEMKIN_M	2.57e+05 18.6%	1.24e+03 2.5%	1.10e+03 2.2%	1.20e+03 12.1%
loop at chemkin_m.f90: 385	2.57e+05 18.6%	1.24e+03 2.5%		1.20e+03 12.1%
loop at chemkin_m.f90: 386	2.57e+05 18.6%	1.24e+03 2.5%		1.20e+03 12.1%
loop at chemkin_m.f90: 387	2.57e+05 18.6%	1.24e+03 2.5%		1.20e+03 12.1%
loop at chemkin_m.f90: 389	2.00e+05 14.5%	1.10e+03 2.2%	1.10e+03 2.2%	1.06e+03 10.7%
chemkin_m.f90: 389	2.00e+05 14.5%	1.10e+03 2.2%	1.10e+03 2.2%	1.06e+03 10.7%

74M of 433M

Outline

- Overview of Rice's HPCToolkit
- Accurate measurement
- Effective performance analysis
- Pinpointing scalability bottlenecks
 - scalability bottlenecks on large-scale parallel systems
 - scaling on multicore processors
- Understanding temporal behavior
- Assessing process variability
- Understanding threading, GPU, locks and memory hierarchy
 - blame shifting
 - attributing memory hierarchy costs to data
- Summary and conclusions

Summary

- **Sampling provides low overhead measurement**
- **Call path profiling + binary analysis + blame shifting = insight**
 - scalability bottlenecks
 - where insufficient parallelism lurks
 - sources of lock contention
 - load imbalance
 - temporal dynamics
 - bottlenecks in hybrid code
 - problematic data structures
- **Other capabilities**
 - attribute memory leaks back to their full calling context

Status

- **Operational today on**
 - 64- and 32-bit x86 systems running Linux (including Cray XT/E/K)
 - IBM Blue Gene/P/Q
 - IBM Power7 systems running Linux
- **Emerging capabilities**
 - **NVIDIA GPU**
 - measurement and reporting using GPU hardware counters
 - data centric analysis
- **Available as open source software at hpctoolkit.org**

Ongoing Work

- **Standardize OpenMP tools API**
 - **enable first-class support for BG/Q OpenMP implementation**
- **Visualization of massive traces**
 - **parallel trace server**
- **Harden support for GPU and hybrid codes**

Some Challenges Ahead

- **Support characteristics of emerging hardware and software**
 - **heterogeneous hardware**
 - manycore, CPU+GPU
 - dynamic power and frequency scaling
 - **software**
 - one-sided communication
 - asynchronous operations
 - dynamic parallelism
 - adaptation
 - failure recovery
 - new programming models
- **Augment monitoring capabilities throughout the stack**
 - **hardware, OS, runtime, language-level API**
- **Transition from descriptive to prescriptive feedback**
- **Guide online adaptation and tuning**

Anecdotal Comparison with Tau and Vampir

Program	Original time	HPCToolkit		TAU		VampireTrace	CCP
		Profiling	Tracing	Profiling	Tracing	Tracing	Tracing
LAMMPS	21.0	23.2 (10.5%)	23.9 (13.8%)	31.3 (49.0%)	36.6 (74.3%)	7846.6 (37265%)	21.0 (0%)
LULESH	17.0	18.2 (7%)	18.3 (7.7%)	27.7 (63%)	27.6 (62.4%)	912.5(5268%)	17.6 (3.5%)

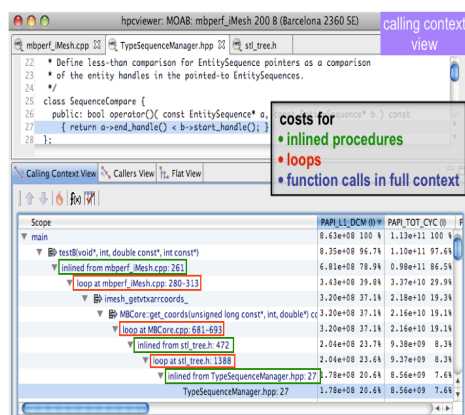
Table 1. Time Comparison.

Program	HPCToolkit		TAU		VampireTrace	CCP
	Profiling	Tracing	Profiling	Tracing	Tracing	Tracing
LAMMPS	17MB	67MB	98MB (5.8x)	5.2GB (79.5x)	85GB (1299x)	152MB (2.3x)
LULESH	268KB	4.0MB	1.2MB (4.6x)	175MB (43.8x)	559MB (139.8x)	11MB (2.8x)

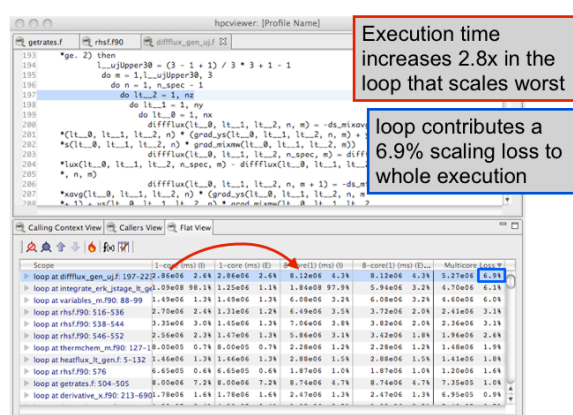
Table 2. Size Comparison.

NOTE: Despite HPCToolkit's need to wrap GPU interfaces for hybrid codes, which increases overhead, HPCToolkit's space and time overhead is still much lower than other tools

HPCToolkit Capabilities at a Glance



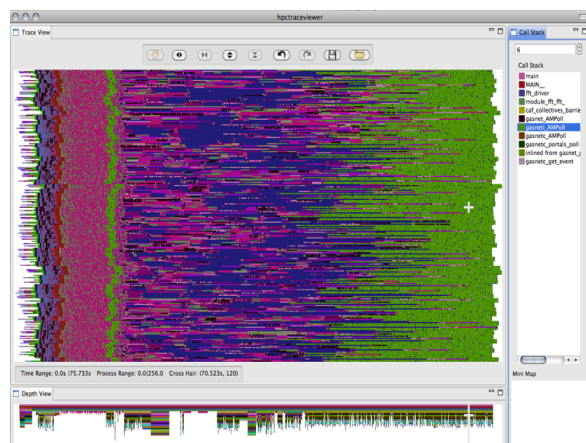
Attribute Costs to Code



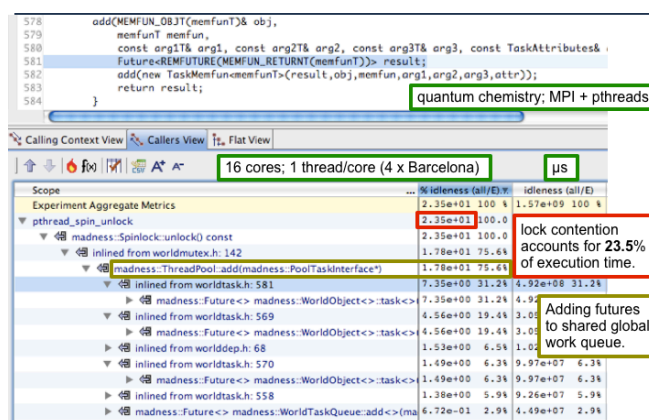
Pinpoint & Quantify Scaling Bottlenecks



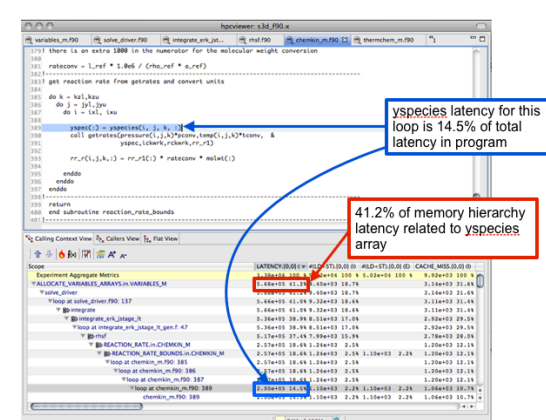
Assess Imbalance and Variability



Analyze Behavior over Time



Shift Blame from Symptoms to Causes



Associate Costs with Data

hpctoolkit.org

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/P and Cray XT
 - The hpcviewer user interface
 - 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**

Using HPCToolkit

- Add hpctoolkit's bin directory to your path
 - use `hpctoolkit`
- Perhaps adjust your compiler flags for your application
 - sadly, most compilers throw away the line map unless `-g` is on the command line. add `-g` flag after any optimization flags if using anything but the Cray compilers/ Cray compilers provide attribution to source without `-g`.
- Add `hpclink` as a prefix to your Makefile's link line
 - e.g. `hpclink mpixlf -o myapp foo.o ... lib.a -lm ...`
- Decide what hardware counters to monitor
 - statically-linked executables (e.g., Cray, Blue Gene)
 - use `hpclink` to link your executable
 - launch executable with environment var `HPCRUN_EVENT_LIST=LIST`
(BG/P hardware counters supported)
 - dynamically-linked executables (e.g., Linux)
 - use `hpcrun -L` to learn about counters available for profiling
 - use `papi_avail`
you can sample any event listed as “profilable”

Using Profiling and Tracing Together

- When tracing, good to have an event that represents a measure of time
 - e.g., **WALLCLOCK** or **PAPI_TOT_CYC**
- Turn on tracing while sampling using one of the above events
 - **Cray XT/E/K: set environment variable in your launch script**
`setenv HPCRUN_EVENT_LIST "PAPI_TOT_CYC@3000000"`
`setenv HPCRUN_TRACE 1`
`aprun your_app`
 - **Linux: use hpcrun**
`hpcrun -e PAPI_TOT_CYC@3000000 -t your_app`
 - **Blue Gene/P at ANL: pass environment settings to cqsub**
`cqsub -p YourAllocation -q prod-devel -t 30 -n 2048 -c 8192 \`
`--mode vn --env HPCRUN_EVENT_LIST=WALLCLOCK@1000 \`
`--env HPCRUN_TRACE=1 your_app`

Monitoring Using Hardware Counters

- **Cray XT/E/K:** set environment variable in your launch script

```
setenv HPCRUN_EVENT_LIST "PAPI_TOT_CYC@3000000  
PAPI_L2_MISS@400000 PAPI_TLB_MISS@400000  
PAPI_FP_OPS@400000"  
aprun your_app
```
- **Linux:** use hpcrun

```
hpcrun -e PAPI_TOT_CYC@3000000 -e PAPI_L2_MISS@400000 \  
-e PAPI_TLB_MISS@400000 -e PAPI_FP_OPS@400000 \  
your_app
```
- **Blue Gene/P at ANL:** pass environment settings to cqsub

```
cqsub -p YourAllocation -q prod-devel -t 30 -n 2048 -c 8192 \  
--mode vn --env HPCRUN_EVENT_LIST=WALLCLOCK@1000 \  
your_app
```

Analysis and Visualization

- Use hpcstruct to reconstruct program structure
 - e.g. `hpcstruct your_app`
 - creates `your_app.hpcstruct`
- Use hpcprof to correlate measurements to source code
 - run `hpcprof` on the front-end node
 - run `hpcprof-mpi` on the compute nodes to analyze data in parallel
- Use hpcviewer to open resulting database
- Use hpctraceviewer to explore traces (collected with `-t` option)

Memory Leak Detection with HPCToolkit

- **Statically linked code**
 - `hpclink --memleak -o your_app foo.o ... lib.a -lm ...`
 - at launch time
 - `setenv HPCTOOLKIT_EVENT_LIST=MEMLEAK`
 - `your_app`
- **Dynamically linked code**
 - `hpcrun -e MEMLEAK your_app`