



Center for Scalable Application Development Software

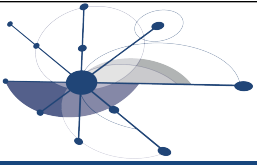
A Slice of CScADS: Performance Tools for Petascale Platforms

John Mellor-Crummey

Laksono Adhianto, Mike Fagan, Michael Franco,
Mark Krentel, Reed Landrum, Xu Liu, Nathan Tallent

Department of Computer Science
Rice University





CScADS Co-PIs and Senior Personnel

John Mellor-Crummey, Keith Cooper



Peter Beckman, Ewing Lusk



Jack Dongarra

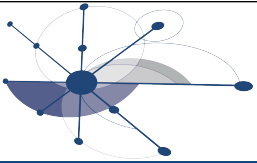


Bart Miller



Katherine Yelick



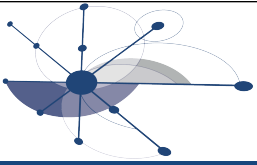


CScADS Mission

- Provide open source software systems, tools, and components that address a spectrum of needs
 - directly usable by application experts
 - support development of enabling technologies by the CS community
- Target architectures of critical interest to DOE
 - Cray XT
 - Blue Gene/P
 - multicore processors in general
- Engage DOE application teams and vendors
- Engage the research community in SciDAC challenges

SciDAC-2 Mission

- Develop comprehensive scientific computing software infrastructure to enable petascale science
- Develop new generation of data management and knowledge discovery tools for large data sets



CScADS Research and Development

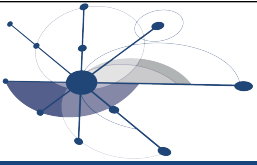
Vertical integration across the petascale software stack

- System software for leadership computing platforms
- Communication libraries
- Math libraries
- Open source compilers
- Performance tools and infrastructure
- Application engagement: analysis and tuning
 - e.g., Annual CScADS Workshop on Leadership Computing
 - experts worked with approximately 200 INCITE and SciDAC code developers to help them scale to DOE's largest systems



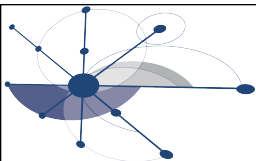
Key Performance Questions

- Why doesn't my application scale as well as I hoped?
- How can I identify bottlenecks in multithreaded node programs?
- How is my code performing relative to peak performance?
 - if my code is not performing well, what is the nature of its problems?



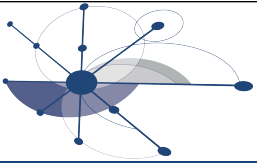
Performance Tool Requirements

- Cope with complex application characteristics
 - large, multi-lingual programs
 - fully optimized code: loop optimization, templates, inlining
 - binary-only libraries, sometimes partially stripped
 - hybrid programs: MPI + OpenMP
- Cope with complex execution environments
 - static or dynamic binaries
 - batch jobs
- Provide effective performance analysis
 - pinpoint and quantify problems
 - yield actionable results
- Scale to leadership computing platforms

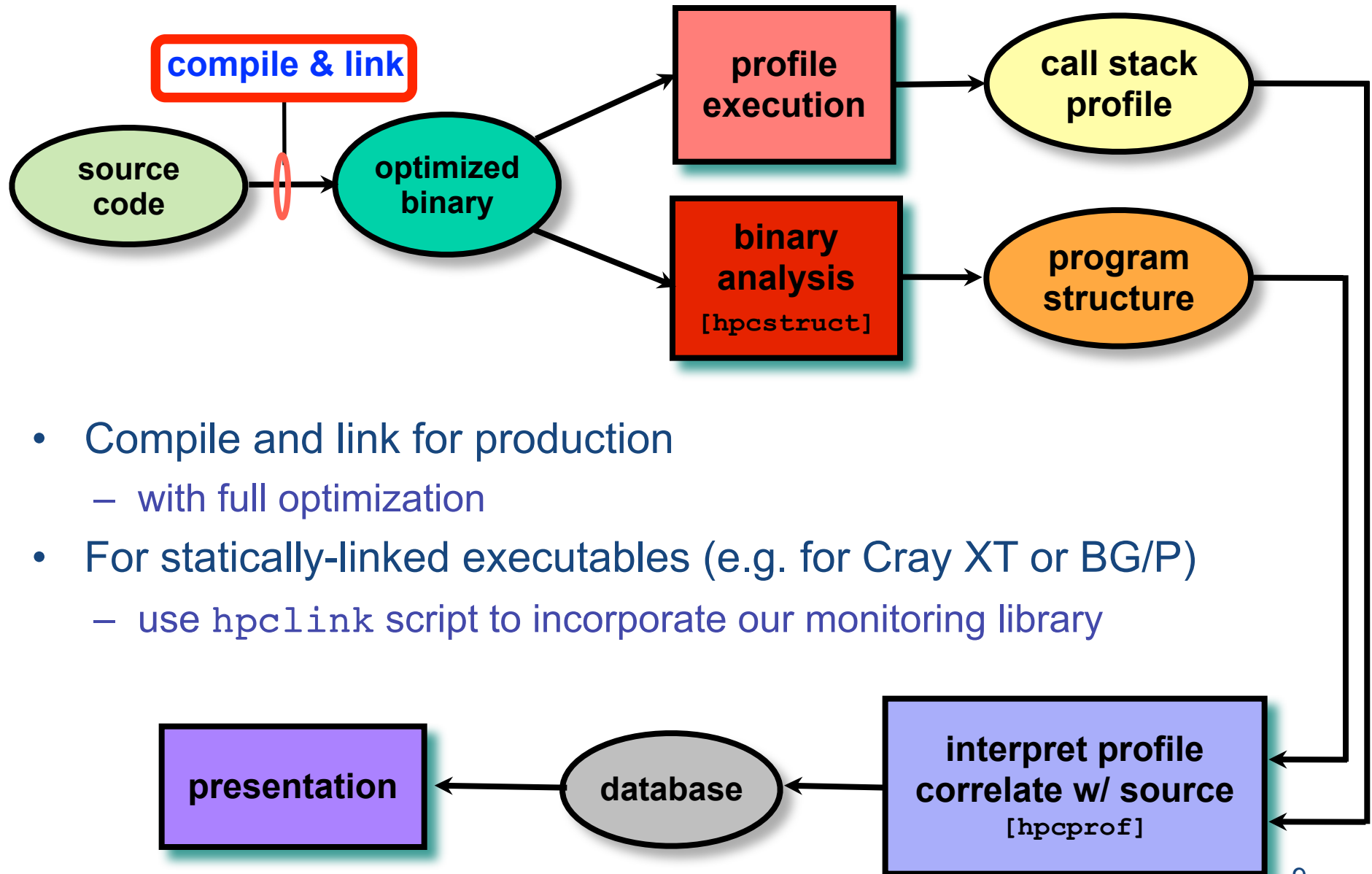


Outline

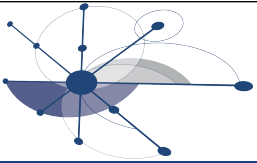
- Introduction to HPCToolkit
- Five new approaches for analyzing parallel program performance
 - scalability analysis using call path profiles [SC09]
 - blame shifting to analyze lock contention in threaded codes [PPoPP10]
 - pinpointing load imbalance in parallel codes [SC10]
 - understanding temporal dynamics of parallel codes
 - data centric analysis of program performance
- Conclusions



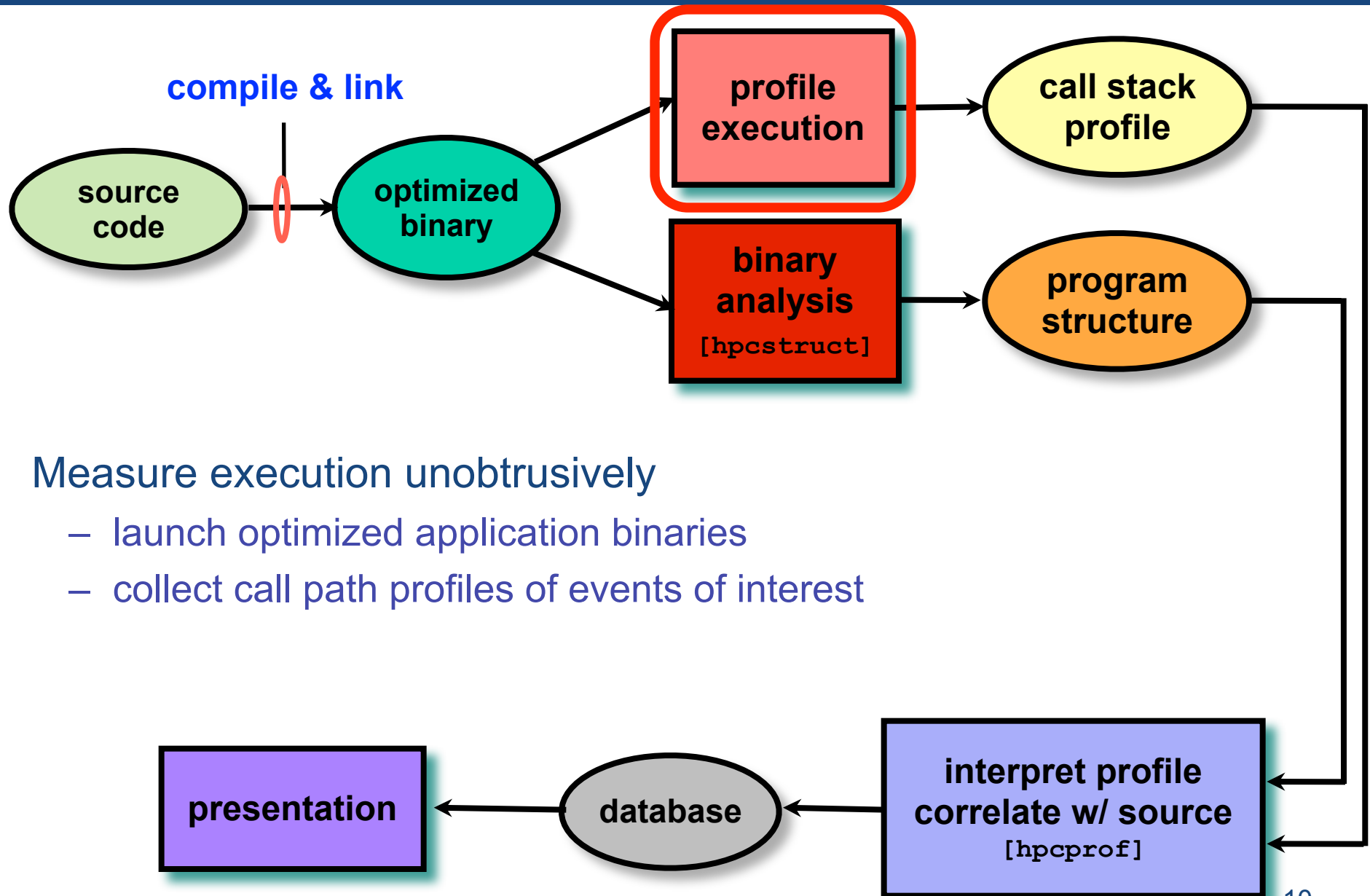
HPCToolkit Performance Tools



- Compile and link for production
 - with full optimization
- For statically-linked executables (e.g. for Cray XT or BG/P)
 - use `hpcLink` script to incorporate our monitoring library



HPCToolkit Performance Tools

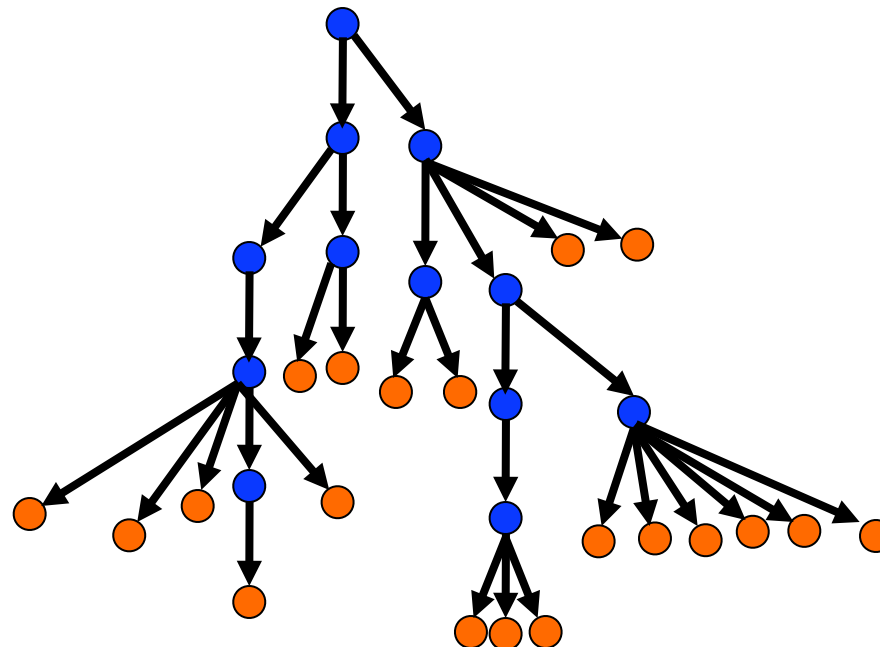
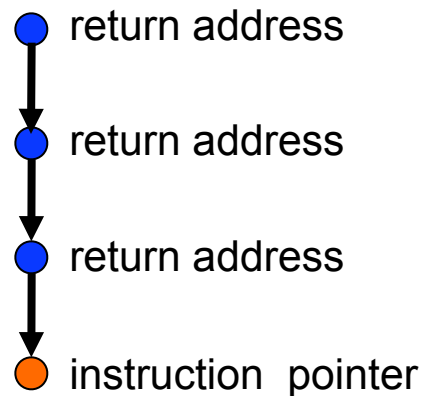


Measure execution unobtrusively

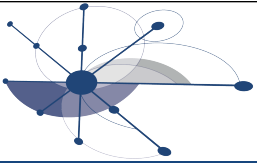
- launch optimized application binaries
- collect call path profiles of events of interest



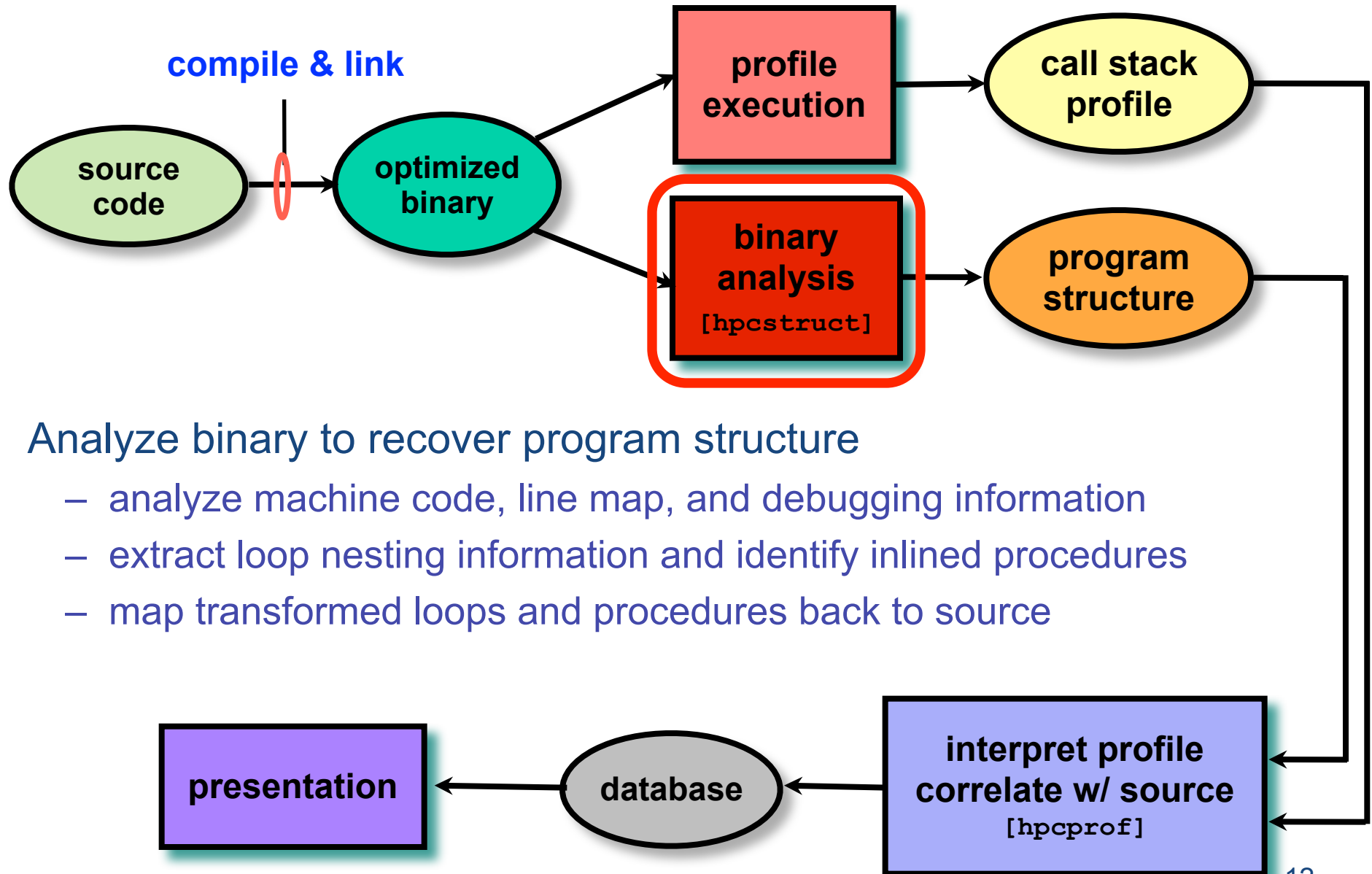
- Sample timer or hardware counter overflows
- Gather calling context using stack unwinding



11

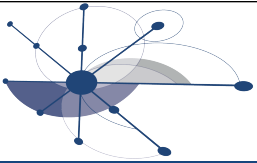


HPCToolkit Performance Tools

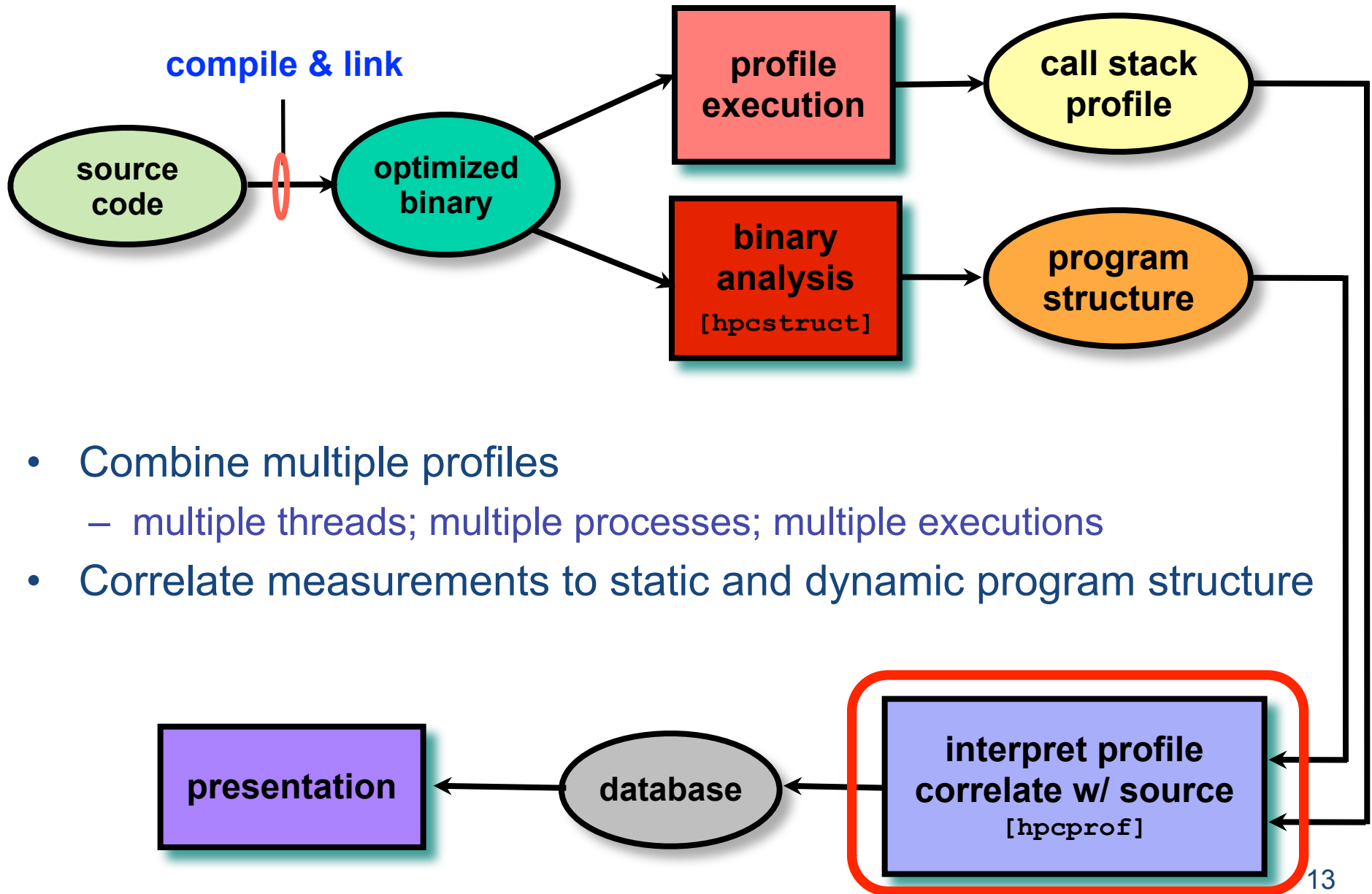


Analyze binary to recover program structure

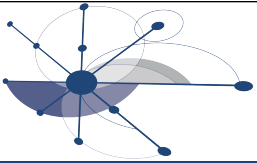
- analyze machine code, line map, and debugging information
- extract loop nesting information and identify inlined procedures
- map transformed loops and procedures back to source



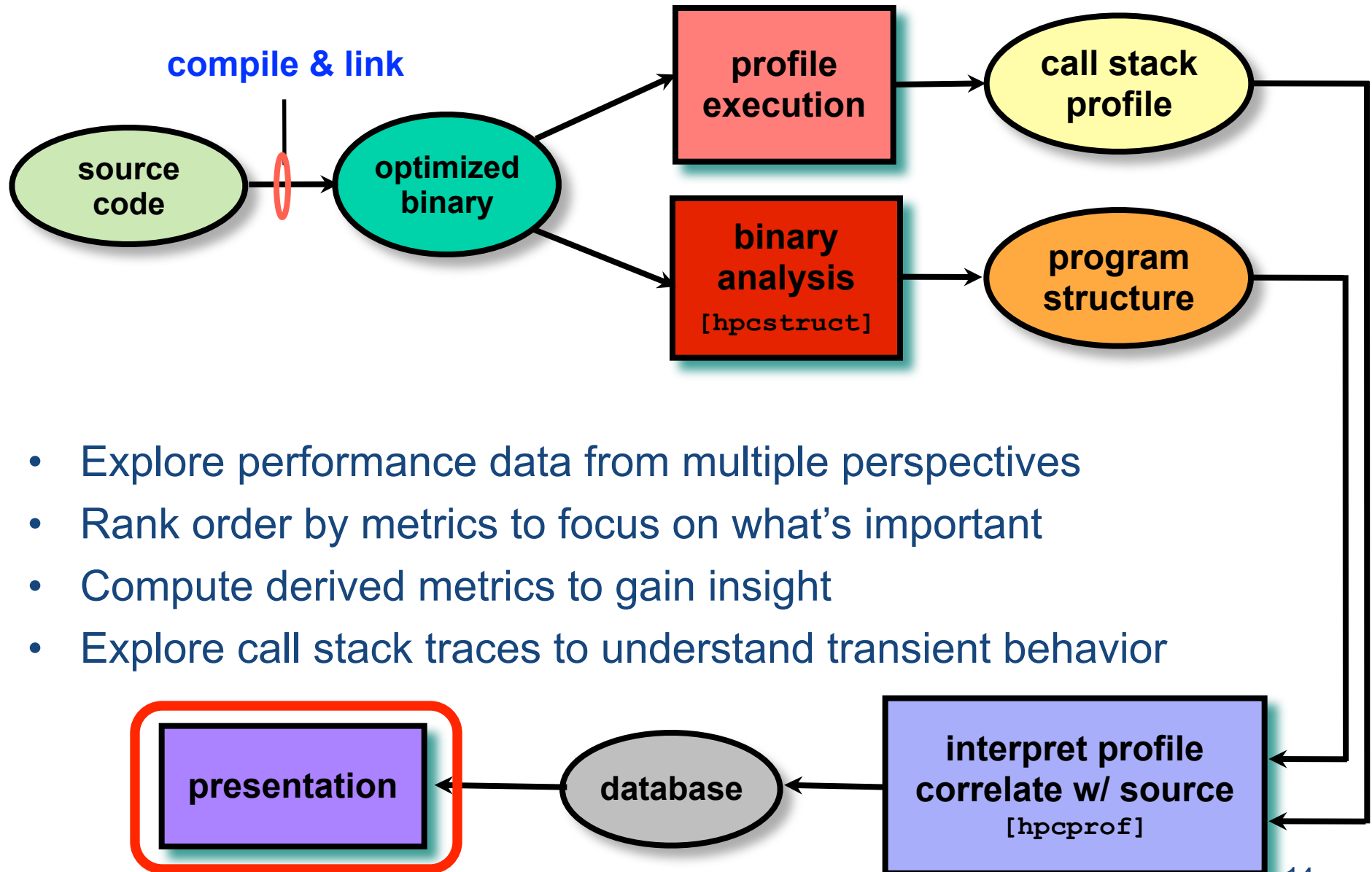
HPCToolkit Performance Tools



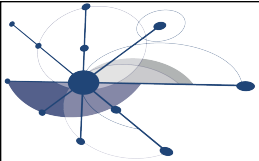
- Combine multiple profiles
 - multiple threads; multiple processes; multiple executions
- Correlate measurements to static and dynamic program structure



HPCToolkit Performance Tools



- Explore performance data from multiple perspectives
- Rank order by metrics to focus on what's important
- Compute derived metrics to gain insight
- Explore call stack traces to understand transient behavior



MOAB Mesh Library from ITAPS

hpcviewer: MOAB: mbperf_iMesh 200 B (Barcelona 2360 SE)

calling context view

```
mbperf_iMesh.cpp  TypeSequenceManager.hpp  stl_tree.h
```

```
22  * Define less-than comparison for EntitySequence pointers as a comparison
23  * of the entity handles in the pointed-to EntitySequences.
24  */
25  class SequenceCompare {
26  public: bool operator()( const EntitySequence* a, const EntitySequence* b ) const
27  { return a->end_handle() < b->start_handle(); }
28  };
```

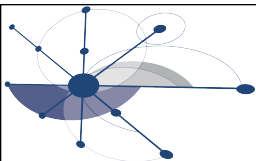
costs for

- inlined procedures
- loops
- function calls in full context

Calling Context View Callers View Flat View

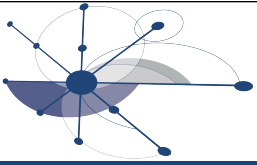
Scope PAPI_L1_DCM (I) PAPI_TOT_CYC (I) P

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 39.8%	3.37e+10 29.9%	
imesh_getvtxarrcoords_	3.20e+08 37.1%	2.18e+10 19.3%	
MBCore::get_coords(unsigned long const*, int, double*) c	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%	

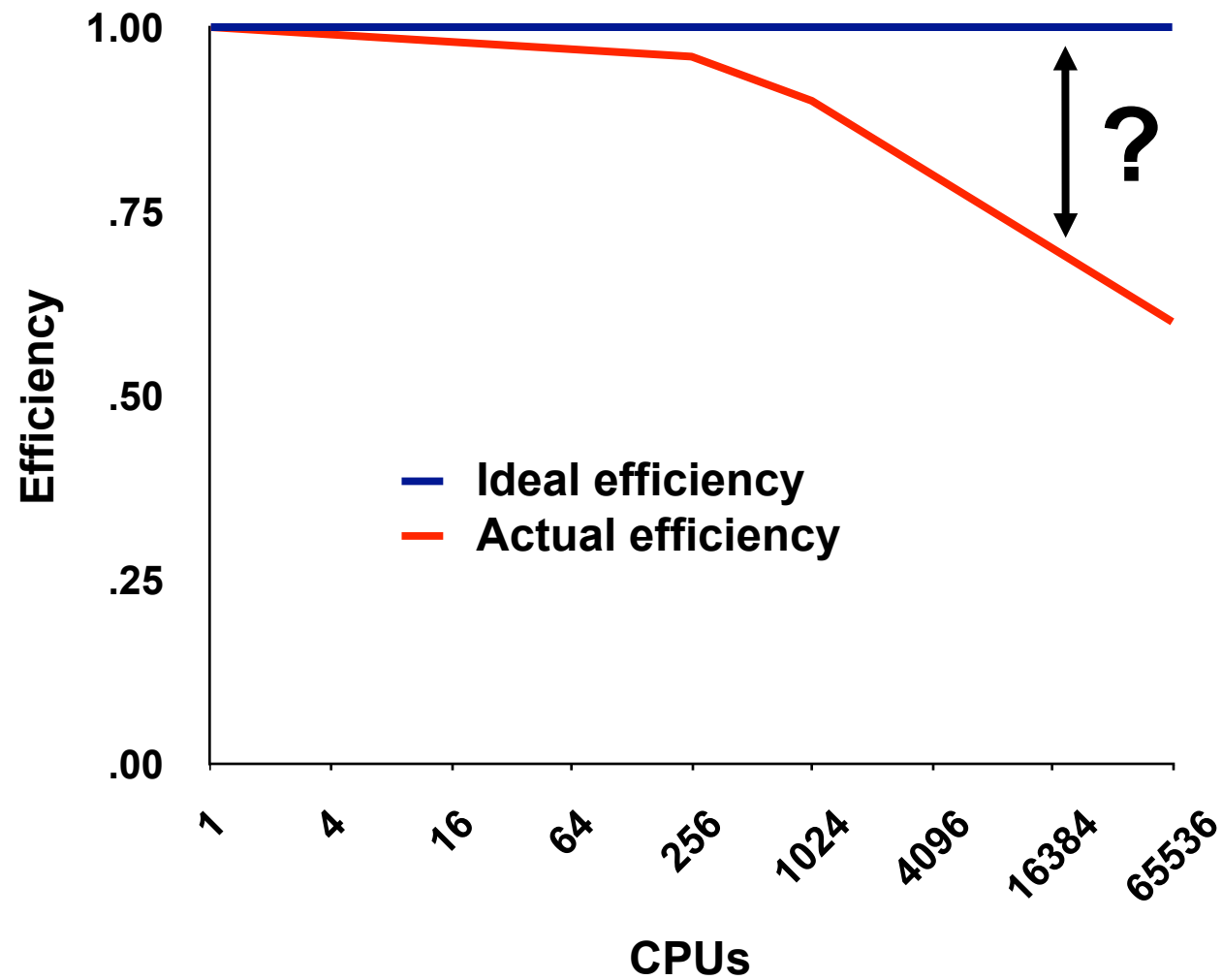


Outline

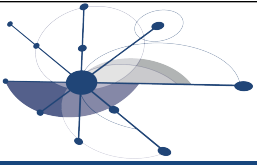
- Introduction to HPCToolkit
- Five new approaches for analyzing parallel program performance
 - scalability analysis using call path profiles [SC09]
 - blame shifting to analyze lock contention in threaded codes [PPoPP10]
 - pinpointing load imbalance in parallel codes [SC10]
 - understanding temporal dynamics of parallel codes
 - data centric analysis of program performance
- Conclusions



Pinpointing Scalability Bottlenecks



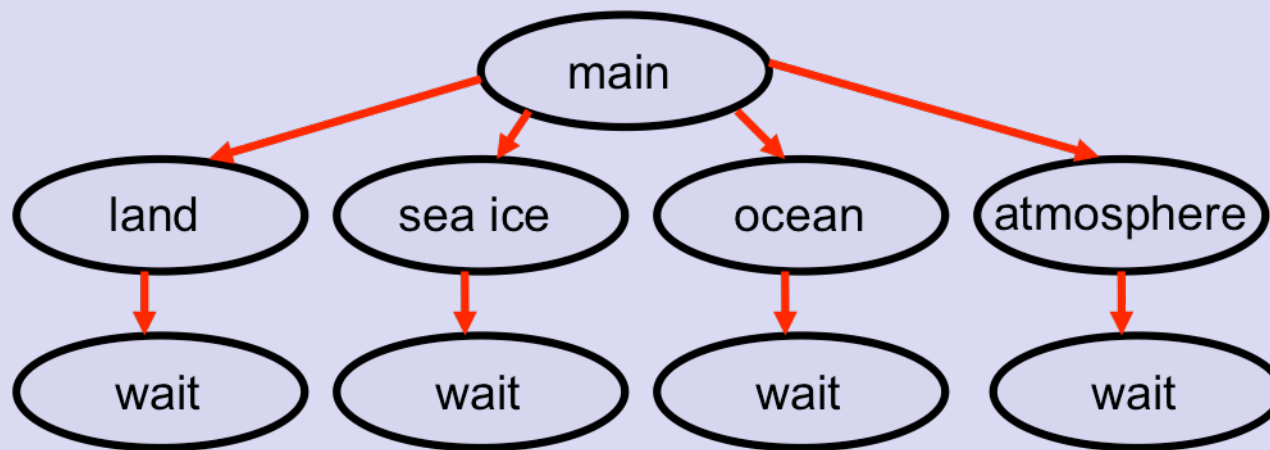
Note: higher is better

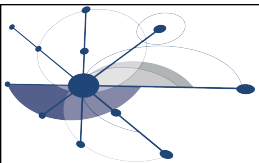


Scalability Analysis Challenges

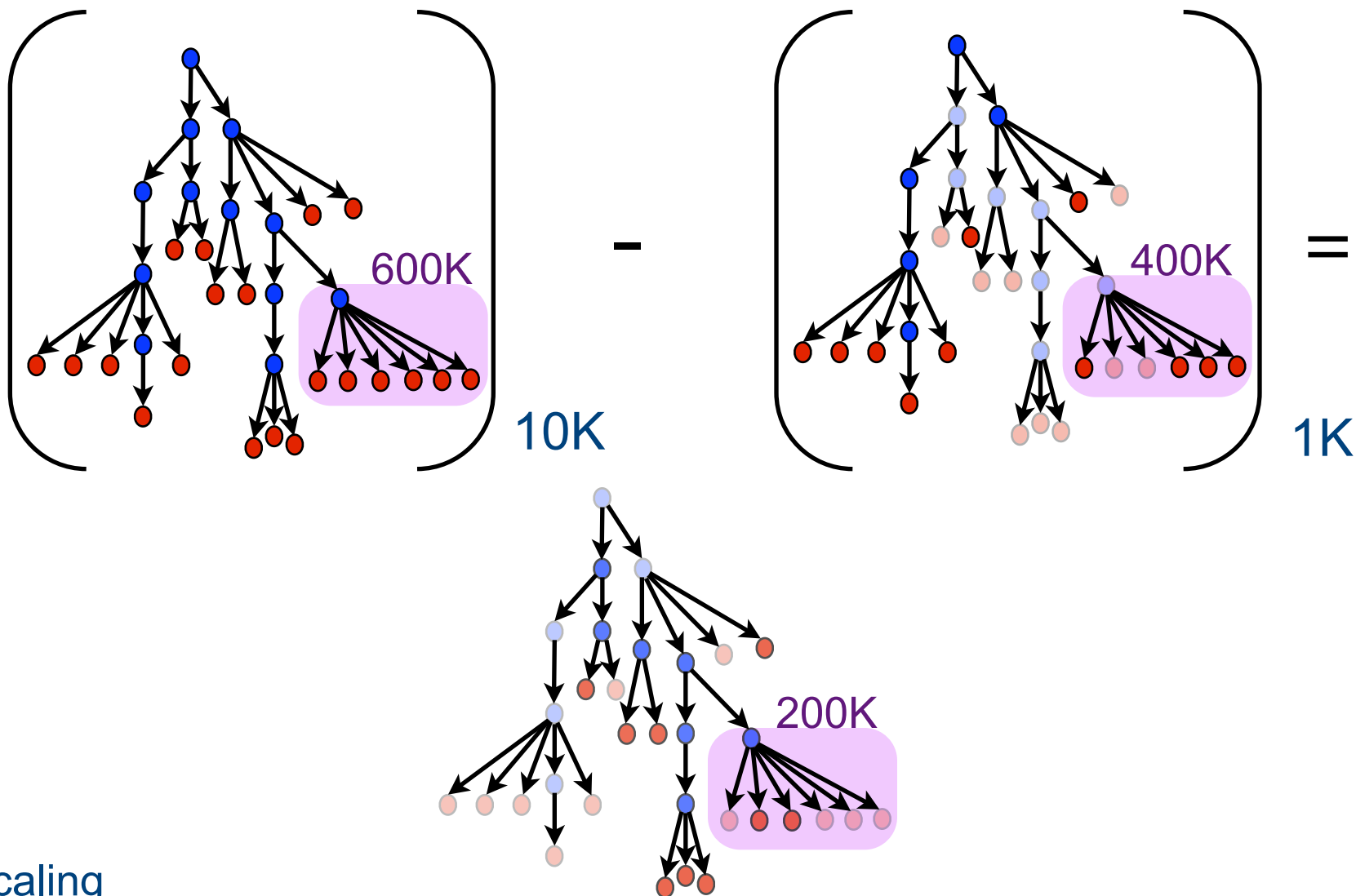
- Parallel applications
 - modern software uses layers of libraries
 - performance is often context dependent
- Monitoring
 - bottleneck nature: computation, data movement, synchronization?
 - pragmatics: need low data volume and low perturbation

Example climate code skeleton

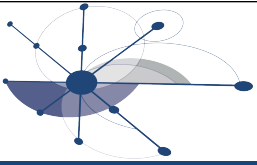




Analyzing Weak Scaling: 1K to 10K processors

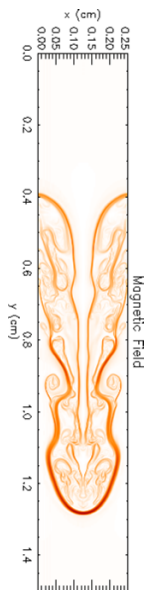


Weak scaling

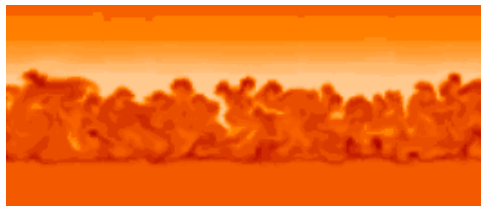


Scalability Analysis Demo: FLASH

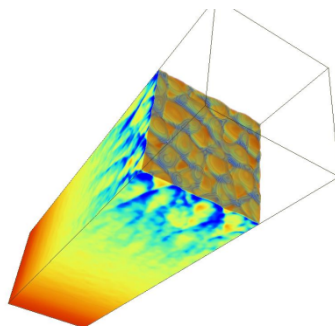
Code:	University of Chicago FLASH
Simulation:	white dwarf collapse
Platform:	Blue Gene/P
Experiment:	8192 vs. 256 processors
Scaling type:	weak



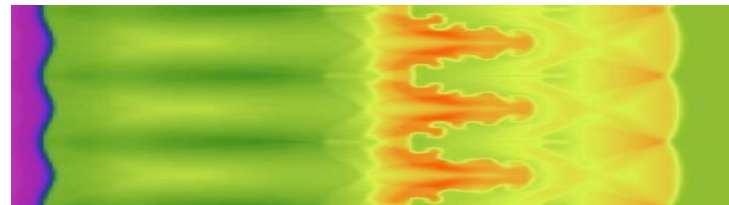
*Magnetic
Rayleigh-Taylor*



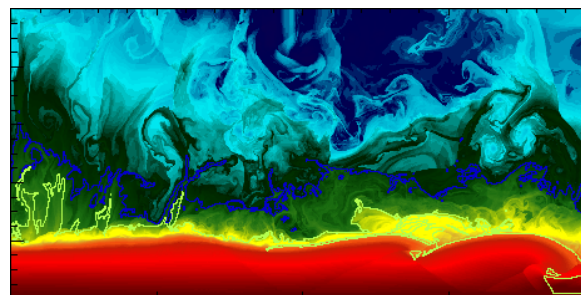
Nova outbursts on white dwarfs



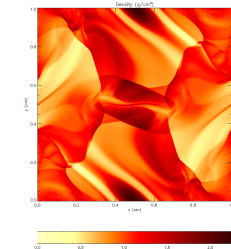
Cellular detonation



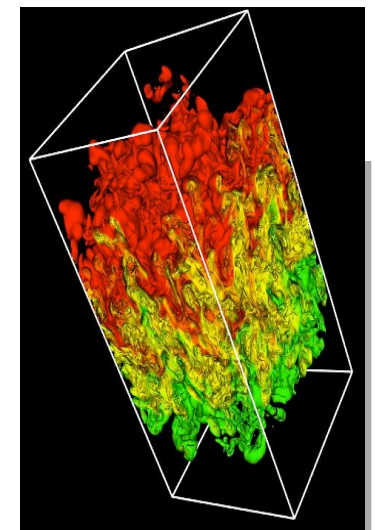
Laser-driven shock instabilities



Helium burning on neutron stars

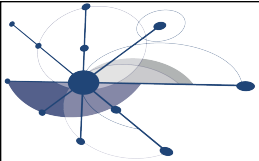


*Orzag/Tang MHD
vortex*



Rayleigh-Taylor instability

Figures courtesy of FLASH Team, University of Chicago



Pinpointing a Scalability Loss in Flash

hpcviewer: FLASH/white dwarf: IBM BG/P, weak 256->8192

Flash.F90 Driver_initFlash.F90 local_tree_build.F90

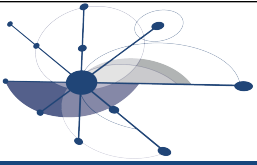
```
210 lnblocks_old = lnblocks
211 nproc = mype
212 !-----Loop through all processors
213 Do iproc = 0, nprocs-1
214
215 If (iproc == 0) Then
216   off_proc = .False.
217 Else
218   off_proc = .True.
219 End If
220
```

21% of the program's scaling loss is due to a loop over all processors in the adaptive mesh refinement setup called during program initialization

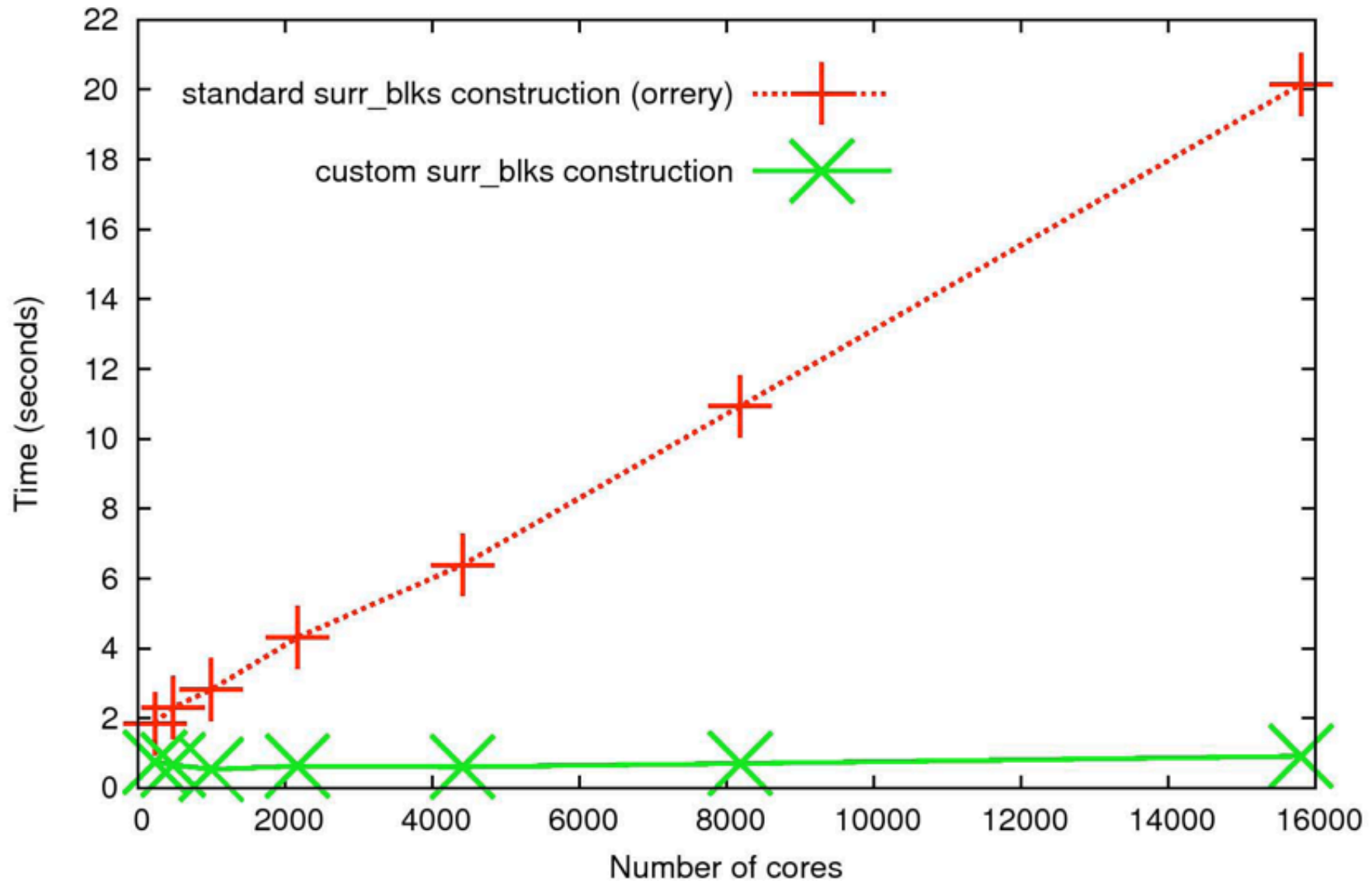
Calling Context View Callers View Flat View

Scope	% scalability loss	256/WALLCLOCK (us) (I)	8192/WALLCLOCK (us) (I)
Experiment Aggregate Metrics	2.46e+01 100 %	5.07e+08 100 %	6.71e+08 100 %
flash	2.46e+01 100 %	5.07e+08 100 %	6.71e+08 100 %
driver_evolveflash	1.41e+01 57.5%	4.46e+08 88.1%	5.41e+08 80.6%
driver_initflash	1.04e+01 42.5%	6.02e+07 11.9%	1.30e+08 19.4%
grid_initdomain	8.58e+00 34.9%	3.45e+07 6.8%	9.21e+07 13.7%
gr_expanddomain	8.58e+00 34.9%	3.45e+07 6.8%	9.21e+07 13.7%
loop at gr_expandDomain.F90: 119	6.85e+00 27.9%	3.42e+07 6.7%	8.02e+07 11.9%
amr_refine_derefine	5.56e+00 22.6%	2.87e+06 0.6%	4.02e+07 6.0%
amr_morton_process	5.45e+00 22.2%	9.75e+05 0.2%	3.76e+07 5.6%
find_surrblks	5.18e+00 21.1%	8.40e+05 0.2%	3.56e+07 5.3%
local_tree_build	5.18e+00 21.1%	8.25e+05 0.2%	3.56e+07 5.3%
loop at local_tree_build.F90: 211	5.18e+00 21.1%	8.25e+05 0.2%	3.56e+07 5.3%
loop at local_tree_build.F90: 216	5.18e+00 21.1%	8.25e+05 0.2%	3.56e+07 5.3%
pmpi_allreduce	1.49e-03 0.0%		1.00e+04 0.0%
pmpi_allreduce	7.45e-04 0.0%		5.00e+03 0.0%
free_local_tree	7.45e-04 0.0%		5.00e+03 0.0%

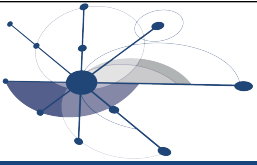
54M of 433M



Improved Flash Scaling of AMR Setup

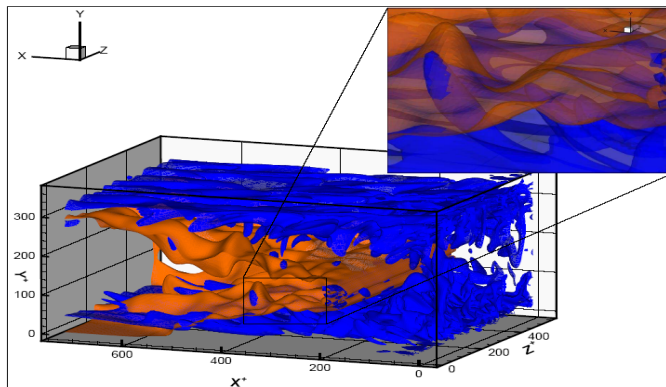


Graph courtesy of Anshu Dubey, U Chicago

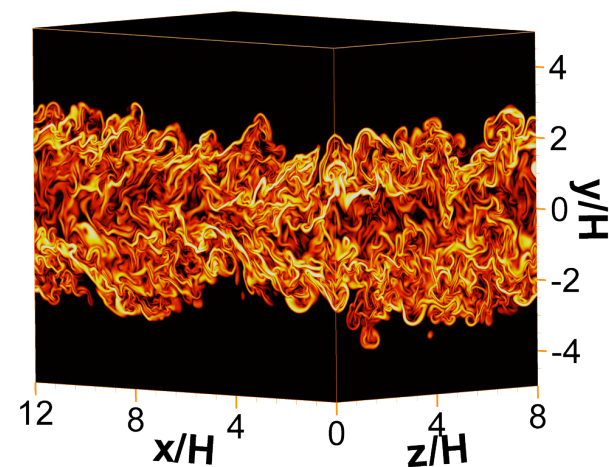


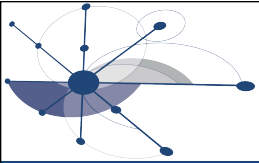
S3D - DNS Solver

- Solves compressible reacting Navier-Stokes equations
- High fidelity numerical methods
 - 8th order finite-difference
 - 4th order explicit RK integrator
- Hierarchy of molecular transport models
- Detailed chemistry
- Multi-physics (sprays, radiation and soot)
 - from SciDAC-TSTC (Terascale Simulation of Turbulent Combustion)

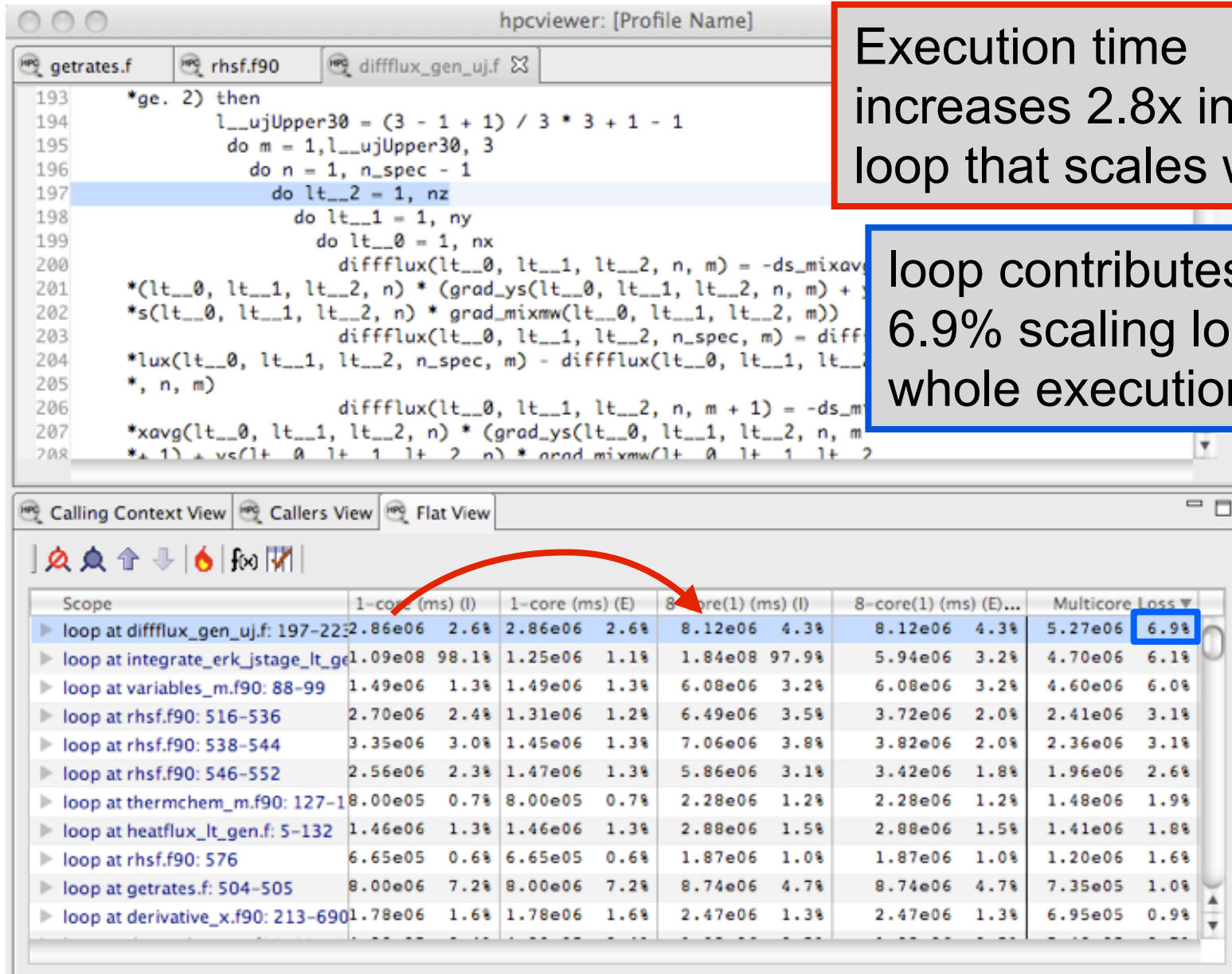


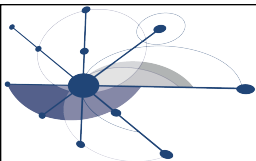
Text and figures courtesy of Jacqueline H. Chen, SNL





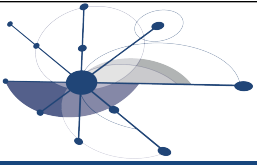
S3D: Multicore Losses at the Loop Level





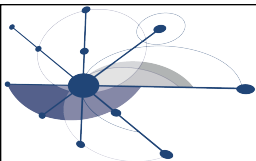
Outline

- Introduction to HPCToolkit
- Five new approaches for analyzing parallel program performance
 - scalability analysis using call path profiles [SC09]
 - blame shifting to analyze lock contention in threaded codes [PPoPP10]
 - pinpointing load imbalance in parallel codes [SC10]
 - understanding temporal dynamics of parallel codes
 - data centric analysis of program performance
- Conclusions



Understanding Lock Contention in Threaded Code

- Lock contention => idleness
 - explicitly threaded programs (Pthreads, etc)
 - implicitly threaded programs (critical sections in OpenMP, ...)
- Strategy: “blame-shifting” of contention 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 associated with a 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

Quantum chemistry; MPI + pthreads

- **65M distinct locks**
- **max. of 340K live locks**
- **30K lock acquisitions/sec/thread**

1-5% overhead

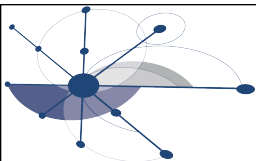
16 cores; 1 thread/core (4 x Barcelona)

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

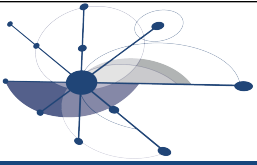
- Introduction to HPCToolkit
- Five new approaches for analyzing parallel program performance
 - scalability analysis using call path profiles [SC09]
 - blame shifting to analyze lock contention in threaded codes [PPoPP10]
 - pinpointing load imbalance in parallel codes [SC10]
 - understanding temporal dynamics of parallel codes
 - data centric analysis of program performance
- Conclusions



2. Identify balance points (procedures or loops that cannot contribute to imbalance)

3. Blame imbalance on the computation subtree in which it originates

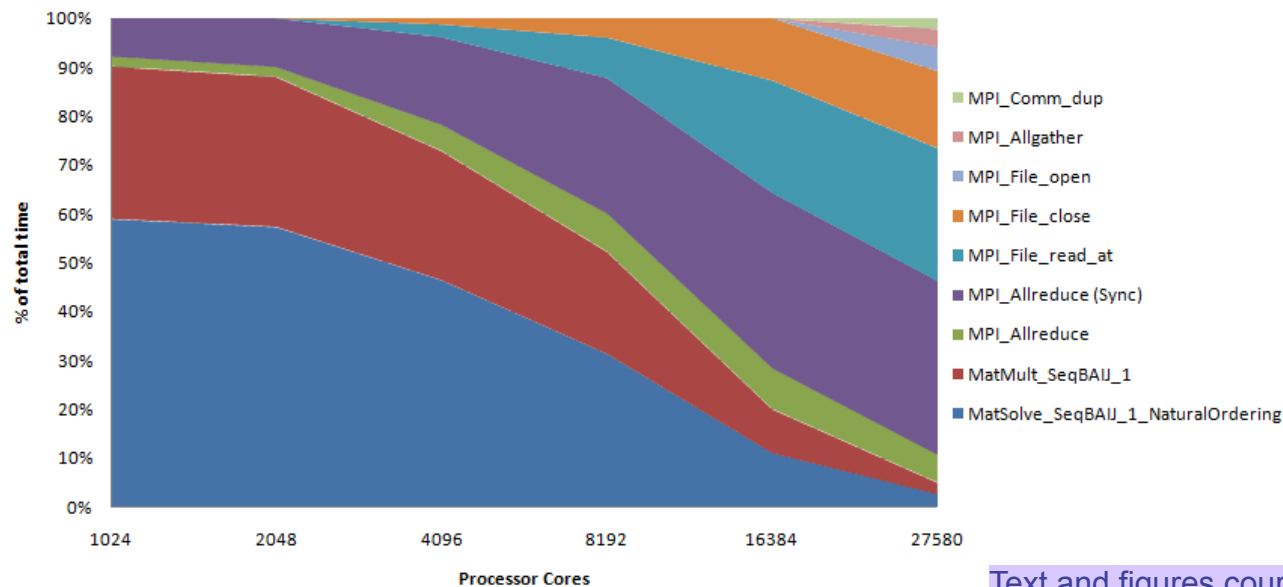
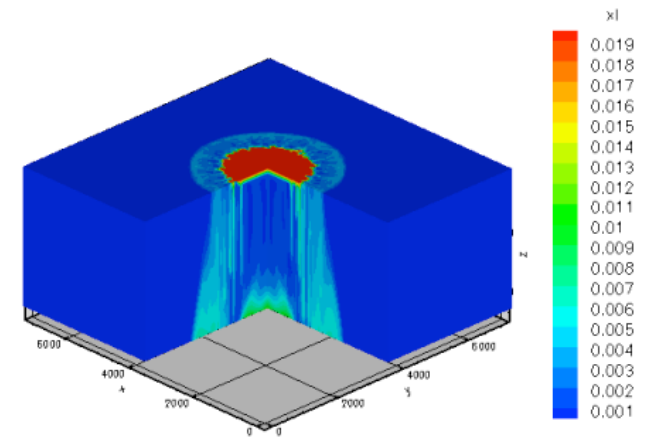
28



Load Imbalance Analysis Example

PFLOTRAN: modeling multi-scale, multiphase, multi-component subsurface reactive flows

Example use: modeling sequestration of CO₂ in deep geologic formations, where resolving density-driven fingering patterns is necessary to accurately describe the rate of dissipation of the CO₂ plume



Strong scaling
study on Cray XT

PFLOTRAN

8K cores, Cray XT5

1. Drill down 'hot path' to loop (a balance point)

2. Notice top two call sites...

3. Plot the per-process values:

Early finishers...

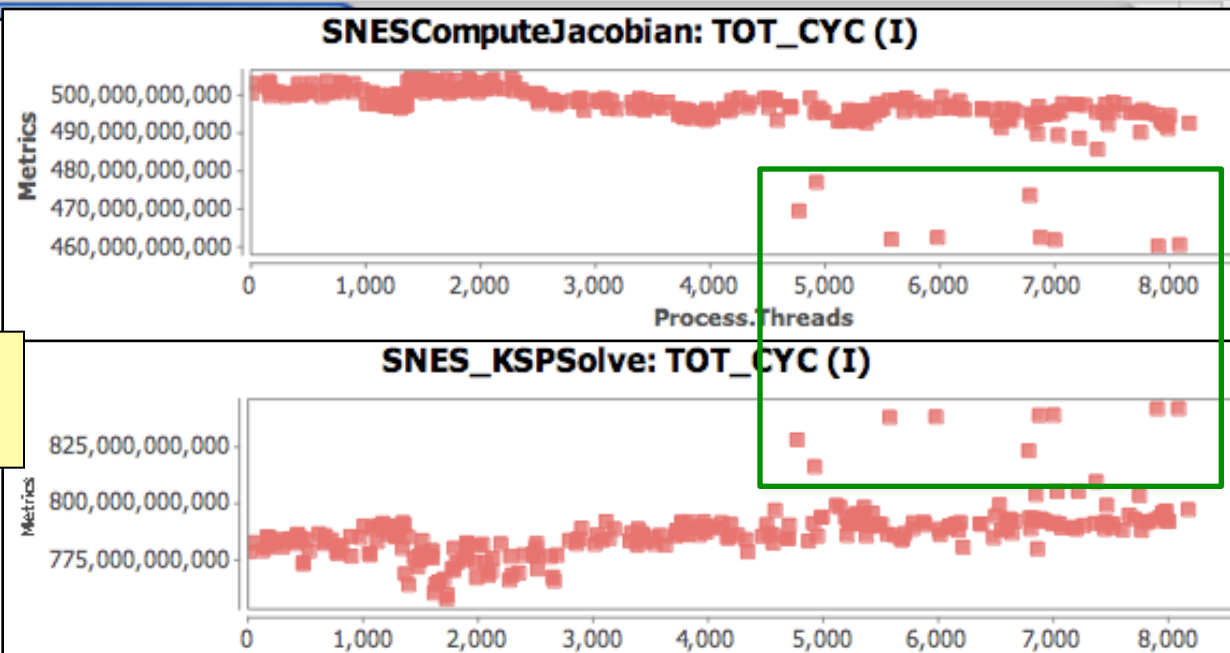
... become early arrivers at Allreduce

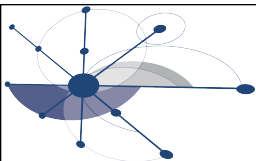
	imbalance (I)	TOT_CYC:Sum (I)	
pflotran	5.28e+15	1.85e+16	100 %
timestepper_module_stepperrun_	5.17e+15	1.82e+16	98.3%
loop at timestepper.F90: 384	5.17e+15	1.82e+16	98.2%
timestepper_module_steppersteptransportdt_	2.22e+15	1.33e+16	72.0%
loop at timestepper.F90: 1230	2.22e+15	1.33e+16	72.0%
loop at timestepper.F90: 1254	2.22e+15	1.32e+16	71.3%
snessolve_	2.22e+15	1.30e+16	70.4%
SNESSolve	2.22e+15	1.30e+16	70.4%
SNESSolve_LS	2.22e+15	1.30e+16	70.4%
loop at ls.c: 181	2.15e+15	1.27e+16	68.8%
SNES_KSPSolve	1.19e+15	6.44e+15	34.8%
SNESComputeJacob	6.21e+14	4.07e+15	22.0%

```

89 ierr = SNESComputeJacob(snes,X,&snes->jacobian,&snes->jacobian_pre,&
190 ierr = KSPSetOperators(snes->ksp,snes->jacobian,snes->jacobian_pre,flg)
191 ierr = SNES_KSPSolve(snes,snes->ksp,F,Y);CHKERRQ(ierr);

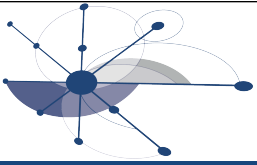
```





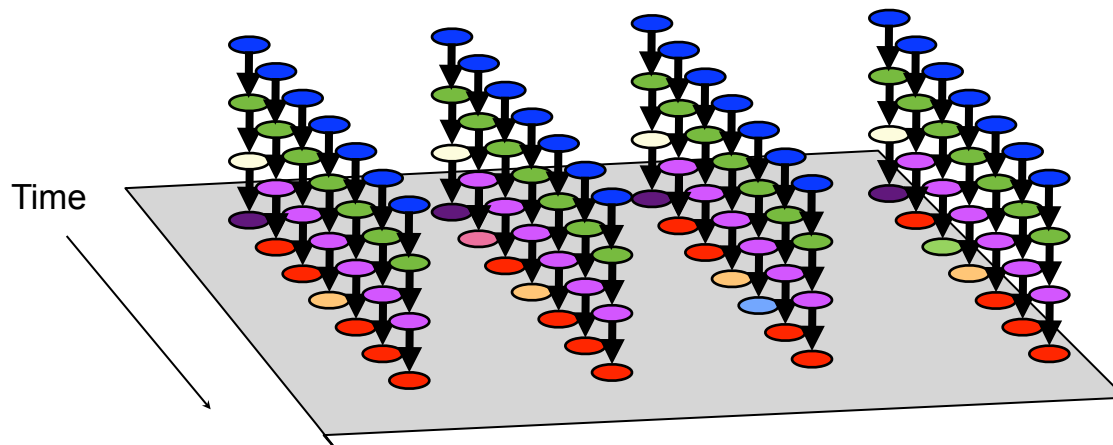
Outline

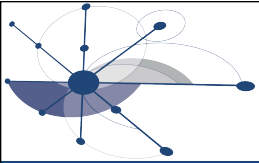
- Introduction to HPCToolkit
- Five new approaches for analyzing parallel program performance
 - scalability analysis using call path profiles [SC09]
 - blame shifting to analyze lock contention in threaded codes [PPoPP10]
 - pinpointing load imbalance in parallel codes [SC10]
 - understanding temporal dynamics of parallel codes
 - data centric analysis of program performance
- Conclusions



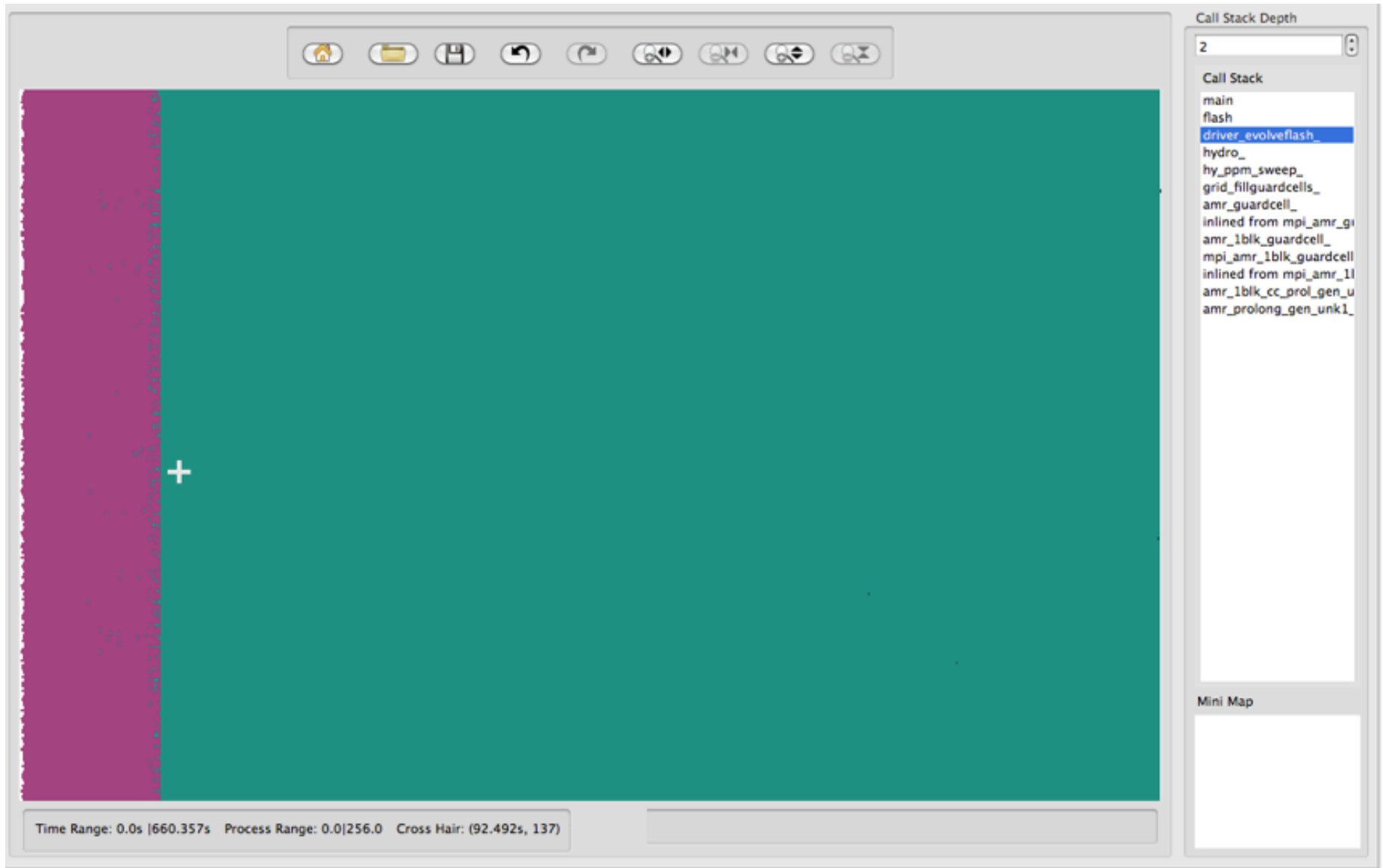
Understanding Temporal Behavior

- Profiling compresses out the temporal dimension
 - that's why serialization is 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 execution with a depth slice

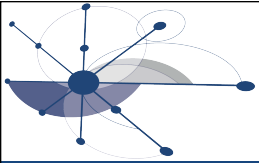




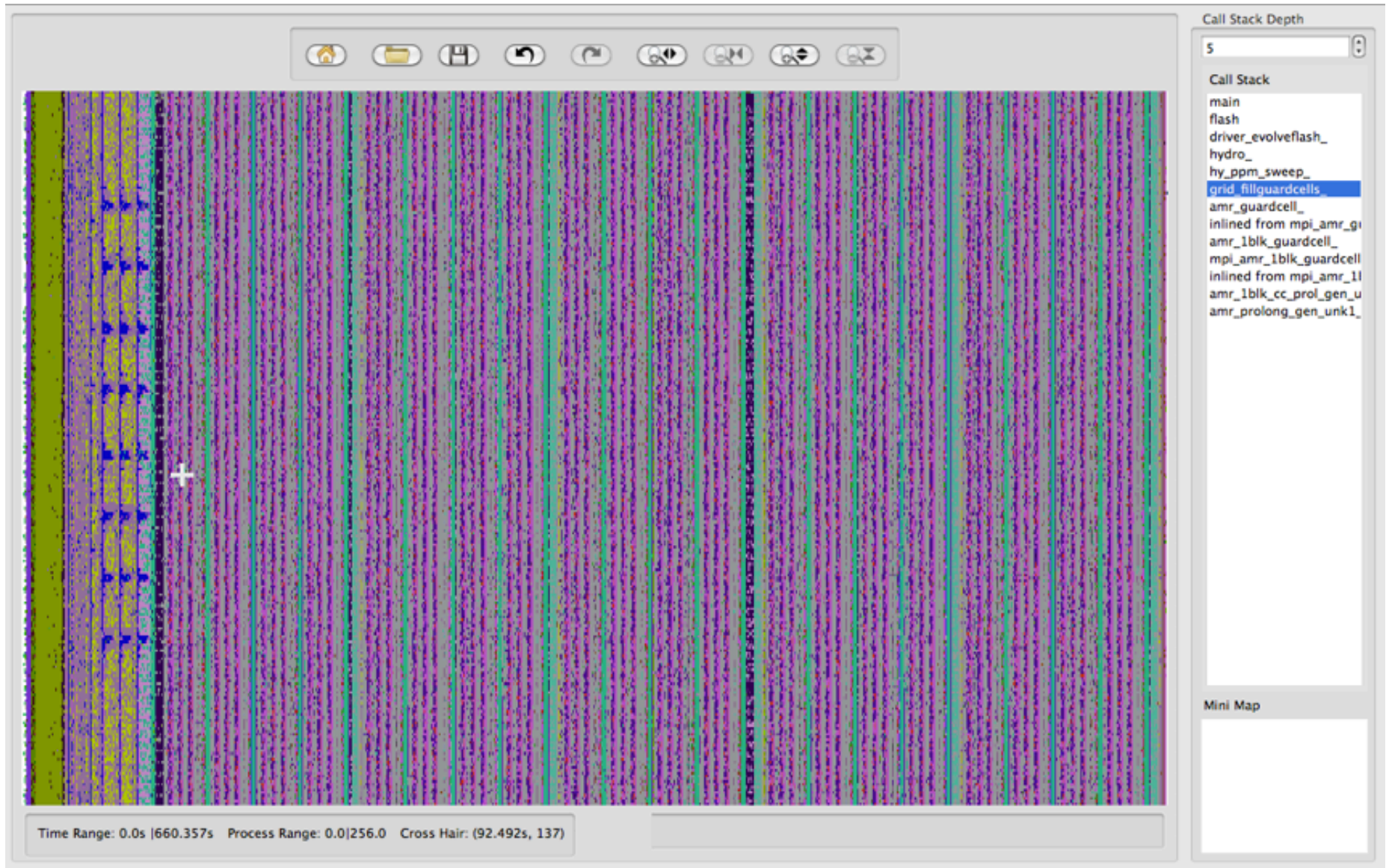
Flash White Dwarf Collapse on 256 Cores



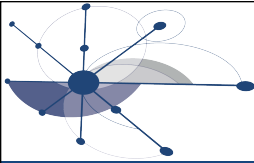
Full execution at call stack depth 2



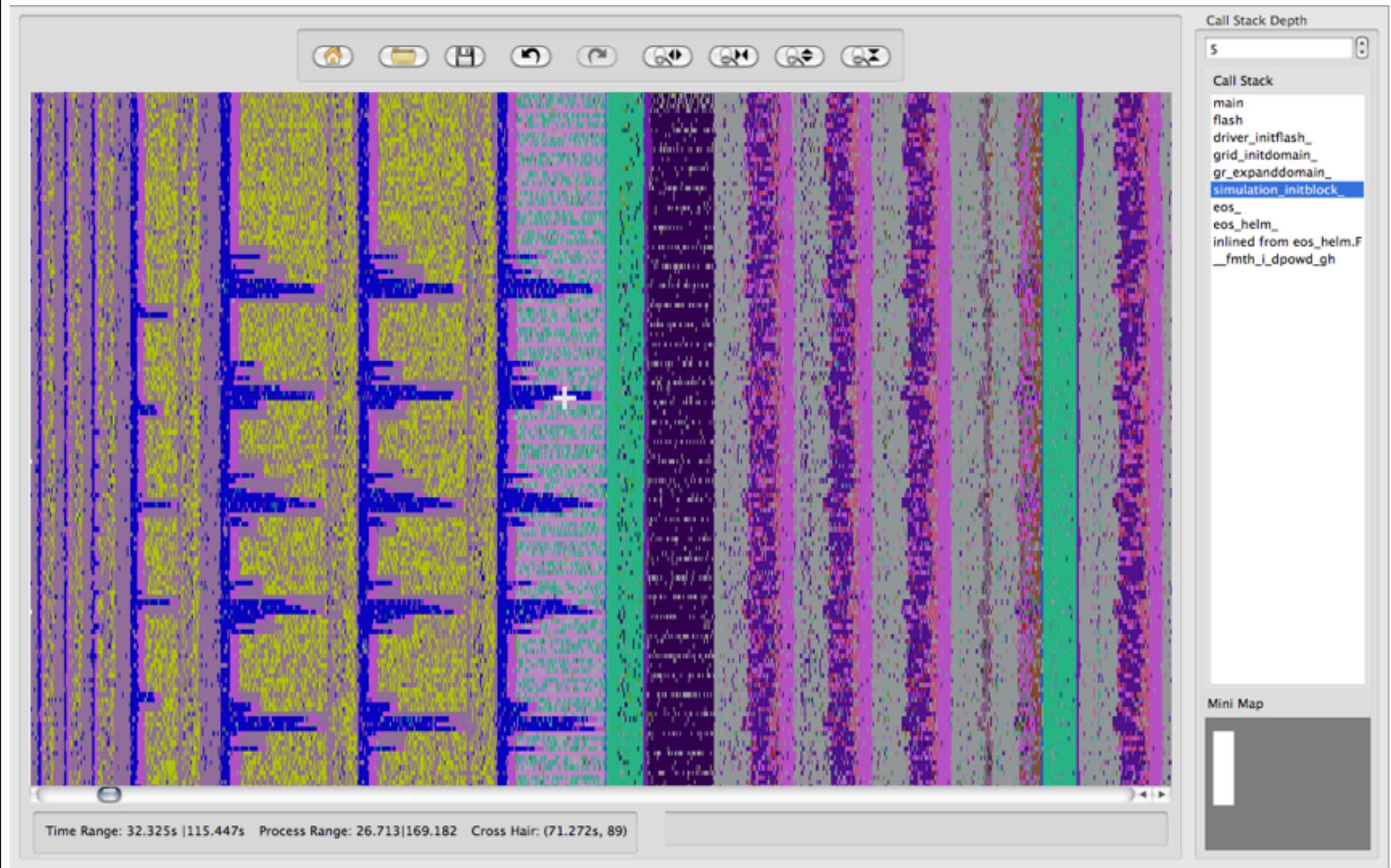
Flash White Dwarf Collapse on 256 Cores



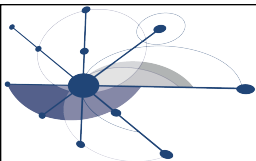
Full execution at call stack depth 5



Flash White Dwarf Collapse on 256 Cores

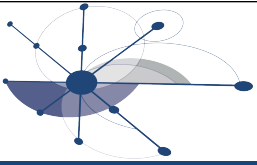


Execution detail at call stack depth 5



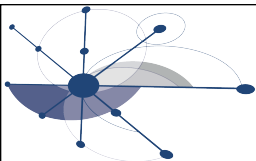
Outline

- Introduction to HPCToolkit
- Five new approaches for analyzing parallel program performance
 - scalability analysis using call path profiles [SC09]
 - blame shifting to analyze lock contention in threaded codes [PPoPP10]
 - pinpointing load imbalance in parallel codes [SC10]
 - understanding temporal dynamics of parallel codes
 - data centric analysis of program performance
- Conclusions

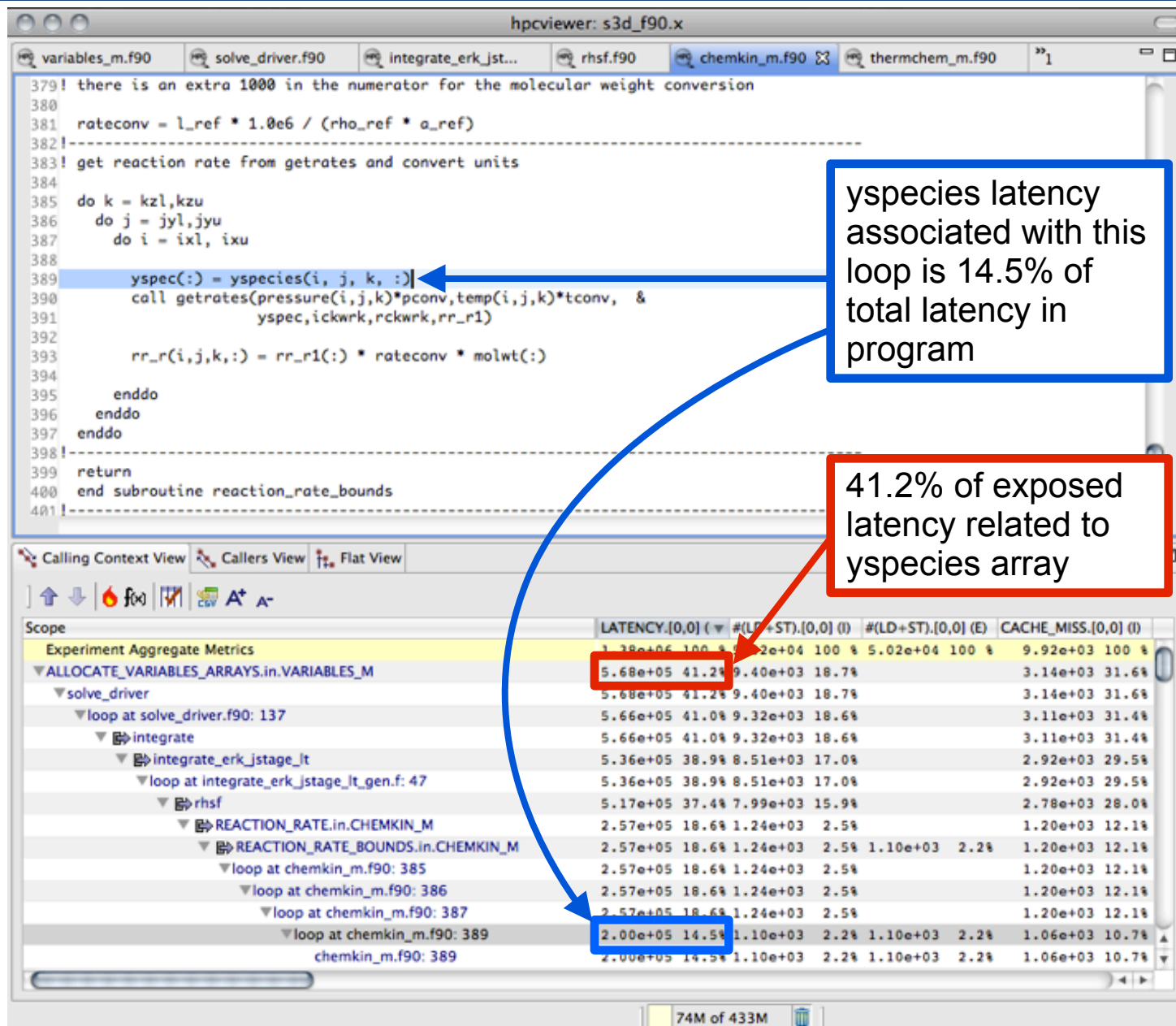


Data Centric Analysis

- Goal: associate memory hierarchy locality problems with particular data structures
- Approach
 - intercept memory allocations to associate data range with allocation
 - associate latency with data structures using “instruction based sampling” capability of 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 results in hpcviewer



Data Centric Analysis of S3D





Conclusions

- Obtain insight, accuracy & precision by combining call path profiling, binary analysis, and blame shifting
- Show surprisingly effective measurement and source-level attribution for fully optimized code (1-3% overhead)
 - statements in their full static and dynamic context
 - project low-level measurements to much higher levels
- Sampling-based measurements can deliver insight into a range of phenomena
 - scalability bottlenecks
 - sources of lock contention
 - load imbalance
 - temporal dynamics
 - problematic data structures



Some Challenges Ahead

- Data management for scalable measurement and analysis
- Moving from descriptive to prescriptive feedback
- Increasing importance of threading as core counts increase
- Heterogeneous architectures, e.g. GPU accelerators