

VI-HPS



scalasca

Scalable performance analysis of large-scale parallel applications

Brian Wylie & Markus Geimer
Jülich Supercomputing Centre

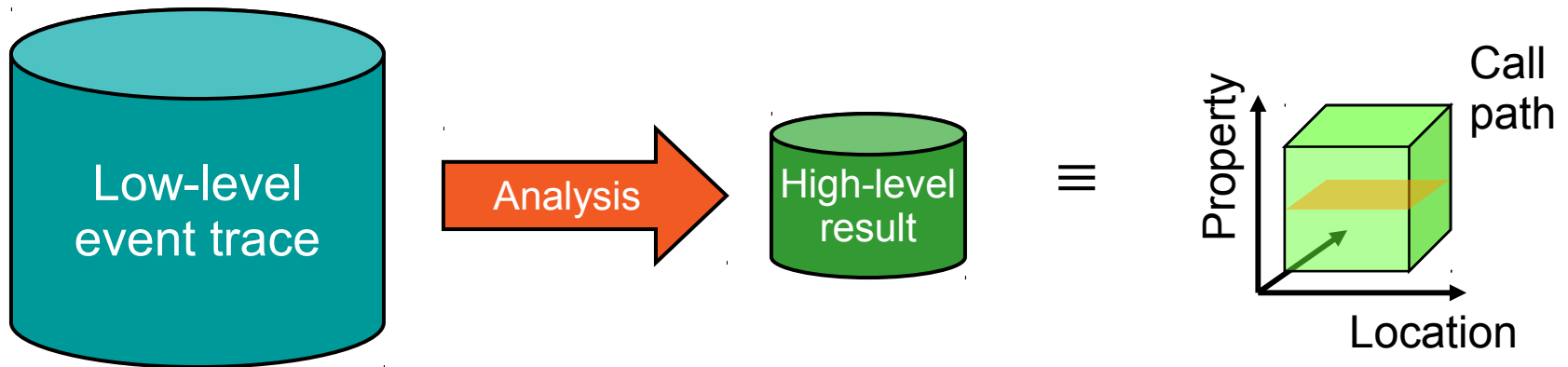
scalasca@fz-juelich.de

September 2011

- Profile analysis
 - Summary of aggregated metrics
 - ▶ per function/callpath and/or per process/thread
 - Most tools (can) generate and/or present such profiles
 - ▶ but they do so in *very* different ways, often from event traces!
 - e.g., gprof, mpiP, ompP, **Scalasca**, TAU, Vampir, ...
- Time-line analysis
 - Visual representation of the space/time sequence of events
 - Requires an execution trace
 - e.g., Vampir, Paraver, JumpShot, Intel TAC, Sun Studio, ...
- Pattern analysis
 - Search for event sequences characteristic of inefficiencies
 - Can be done manually, e.g., via visual time-line analysis
 - or automatically, e.g., KOJAK, **Scalasca**, Periscope, ...

- Idea

- Automatic search for patterns of inefficient behaviour
- Classification of behaviour & quantification of significance



- Guaranteed to cover the entire event trace
- Quicker than manual/visual trace analysis
- Parallel replay analysis exploits memory & processors to deliver scalability

- Overview
 - Helmholtz Initiative & Networking Fund project started in 2006
 - Headed by Bernd Mohr (JSC) & Felix Wolf (GRS)
 - Follow-up to pioneering KOJAK project (started 1998)
 - ▶ Automatic pattern-based trace analysis
- Objective
 - Development of a **scalable** performance analysis toolset
 - Specifically targeting **large-scale** parallel applications
 - ▶ such as those running on BlueGene/P or Cray XT with 10,000s to 100,000s of processes
- Latest release March 2011: Scalasca v1.3.3
 - Download from www.scalasca.org
 - Available on POINT/VI-HPS Parallel Productivity Tools DVD

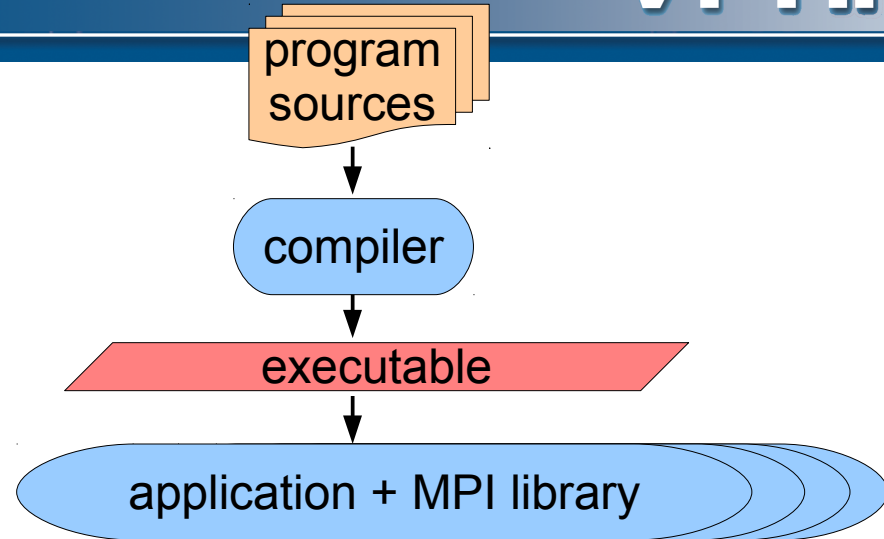
- Open source, New BSD license
- Portable
 - IBM BlueGene, IBM SP & blade clusters, Cray XT, NEC SX, SGI Altix, SiCortex, Solaris & Linux clusters, ...
- Supports parallel programming paradigms & languages
 - MPI, OpenMP & hybrid OpenMP/MPI
 - Fortran, C, C++
- Integrated instrumentation, measurement & analysis toolset
 - Automatic and/or manual customizable instrumentation
 - Runtime summarization (aka profiling)
 - Automatic event trace analysis
 - Analysis report exploration & manipulation

- MPI 2.2 apart from dynamic process creation
 - C++ interface deprecated with MPI 2.2
- OpenMP 2.5 apart from nested thread teams
 - partial support for dynamically-sized/conditional thread teams*
 - no support for OpenMP used in macros or included files
- Hybrid OpenMP+MPI
 - partial support for non-uniform thread teams*
 - no support for MPI_THREAD_MULTIPLE

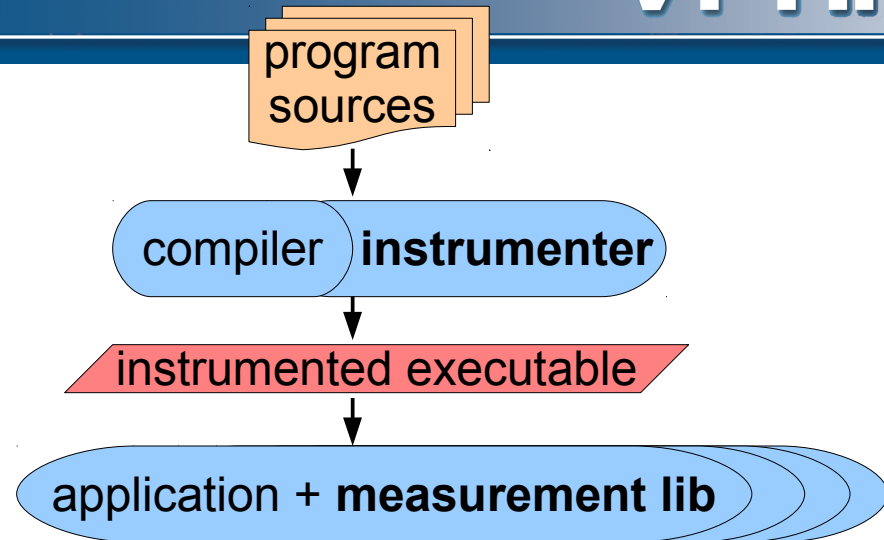
* Summary & trace measurements are possible, and traces may be analyzed with Vampir or other trace visualizers

- automatic trace analysis currently not supported

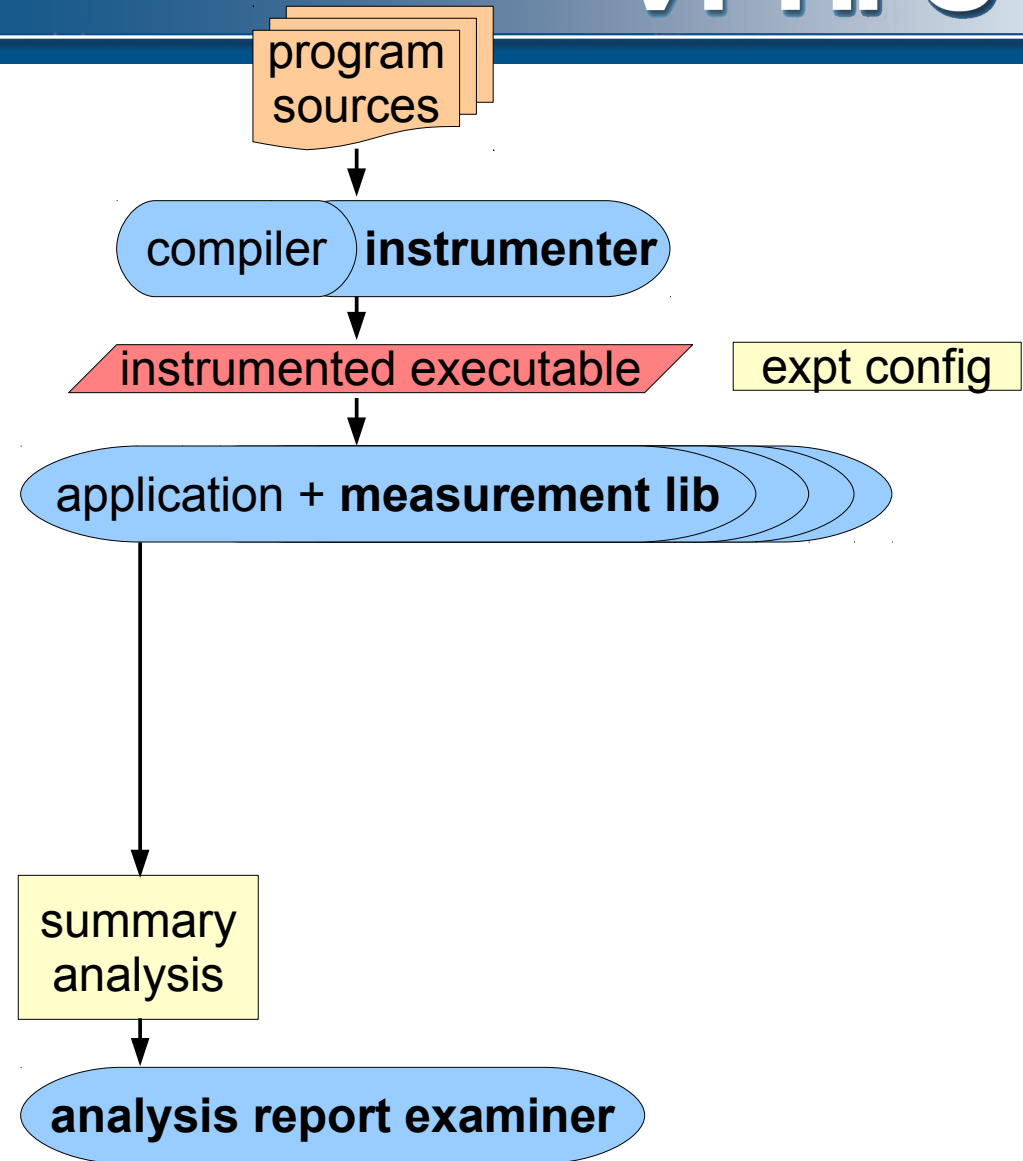
- Application code compiled & linked into executable using MPICC/CXX/FC
- Launched with MPIEXEC
- Application processes interact via MPI library



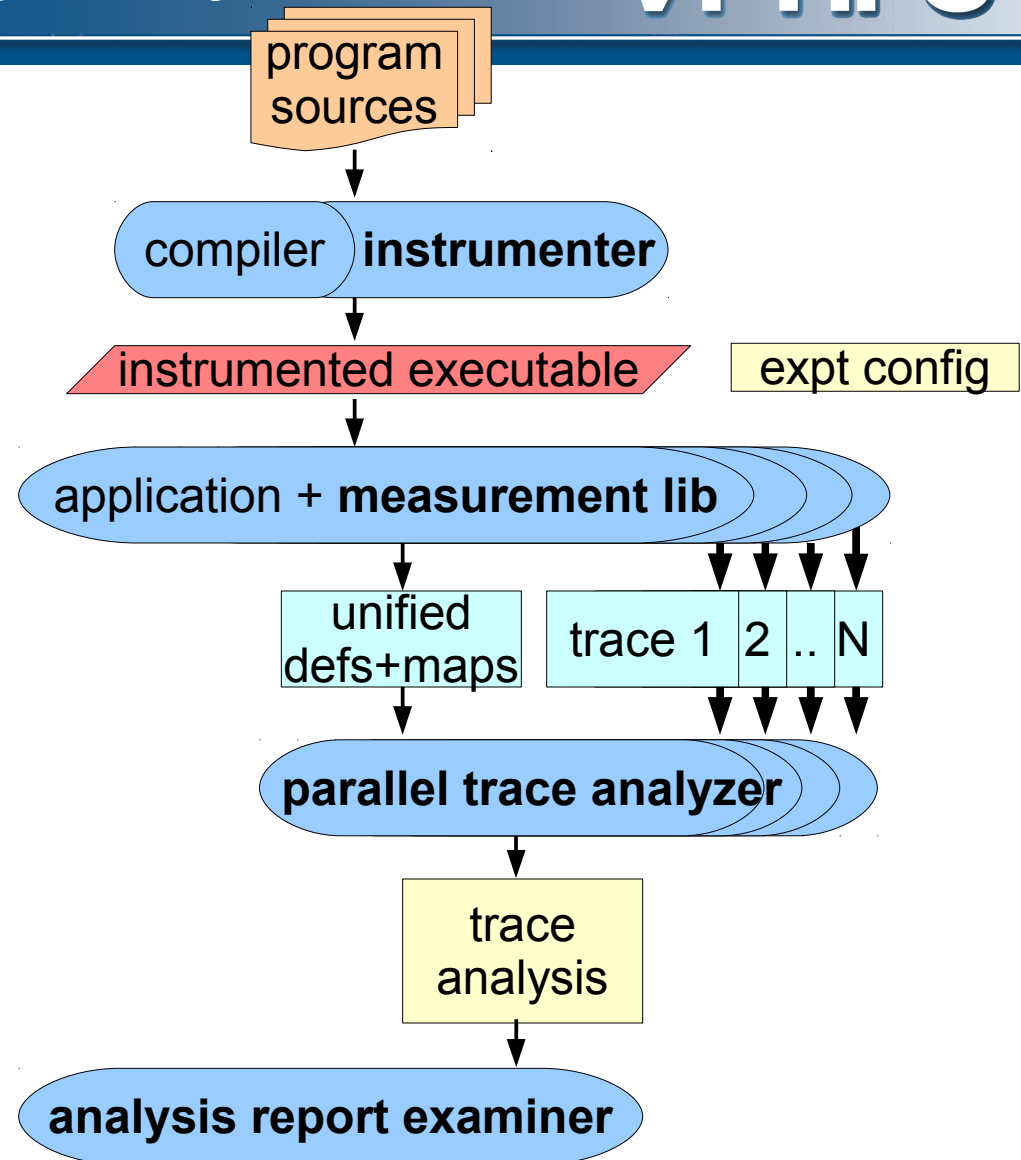
- Automatic/manual code instrumenter
- Program sources processed to add instrumentation and measurement library into application executable
- Exploits MPI standard profiling interface (PMPI) to acquire MPI events



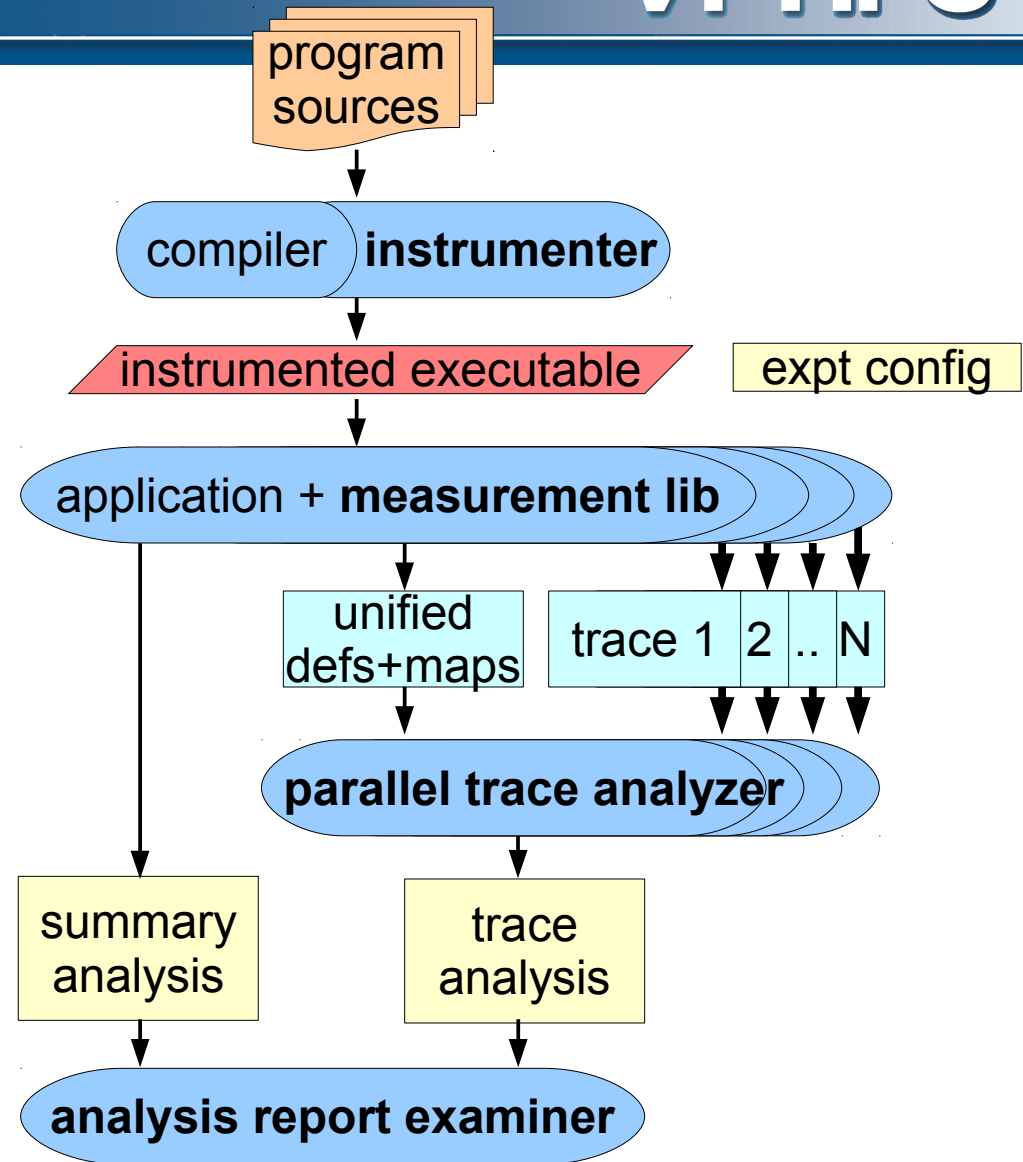
- Measurement library manages threads & events produced by instrumentation
- Measurements summarized by thread & call-path during execution
- Analysis report unified & collated at finalization
- Presentation of summary analysis



- During measurement time-stamped events buffered for each thread
- Flushed to files along with unified definitions & maps at finalization
- Follow-up analysis replays events and produces extended analysis report
- Presentation of analysis report



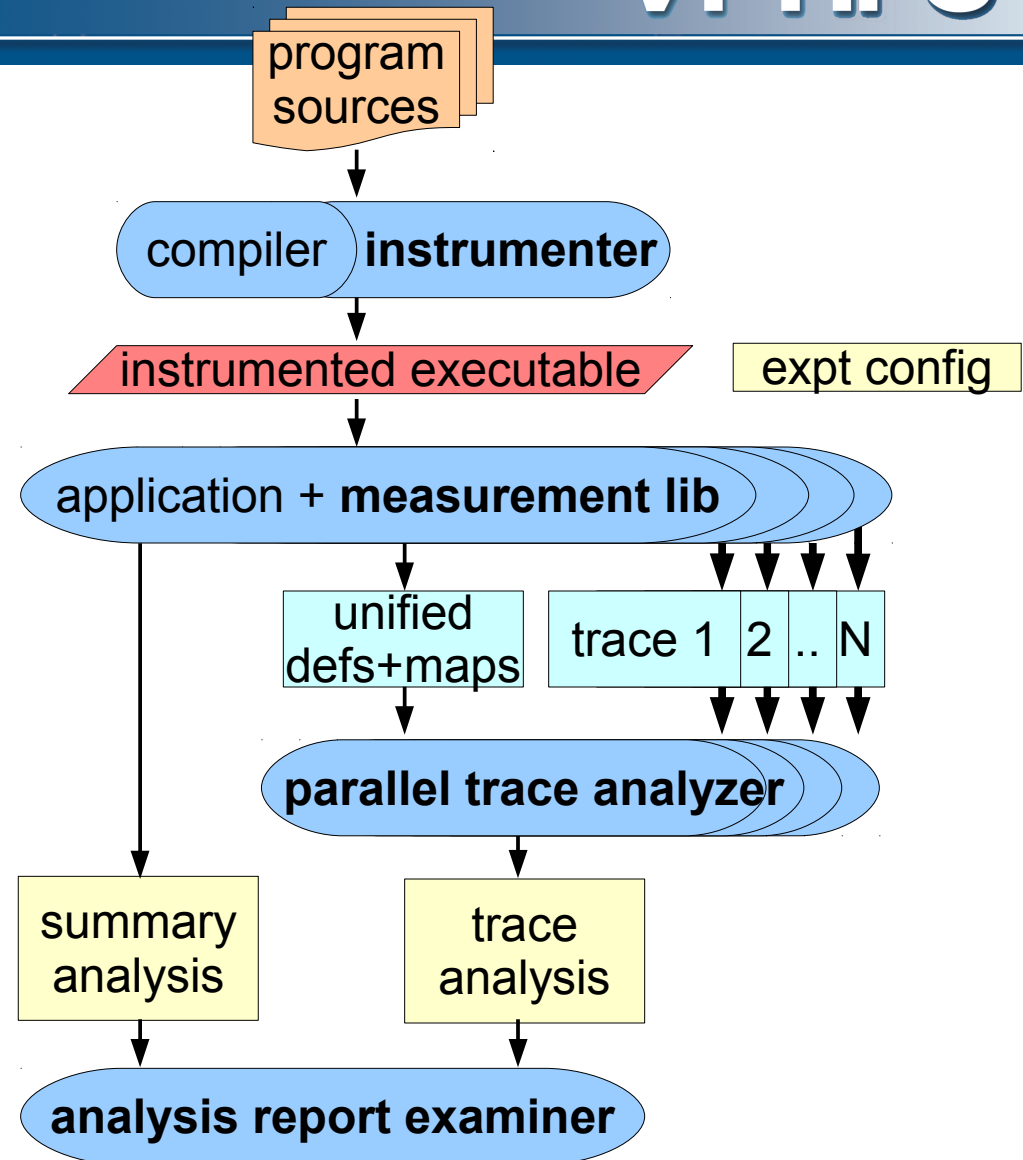
- Automatic/manual code instrumenter
- Measurement library for runtime summary & event tracing
- Parallel (and/or serial) event trace analysis when desired
- Analysis report examiner for interactive exploration of measured execution performance properties



- Scalasca instrumenter = SKIN

- Scalasca measurement collector & analyzer = SCAN

- Scalasca analysis report examiner = SQUARE



- One command for everything

% **scalasca**

Scalasca 1.3

Toolset for scalable performance analysis of large-scale apps

usage: scalasca [-v][-n] {action}

1. prepare application objects and executable for measurement:

scalasca *-instrument* <compile-or-link-command> # **skin**

2. run application under control of measurement system:

scalasca *-analyze* <application-launch-command> # **scan**

3. post-process & explore measurement analysis report:

scalasca *-examine* <experiment-archive|report> # **square**

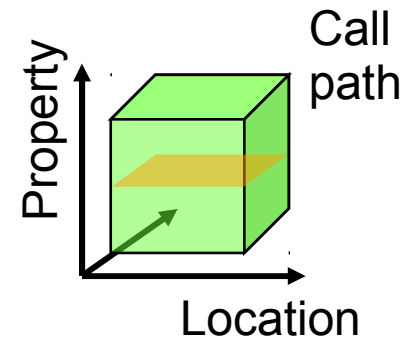
[-h] show quick reference guide (only)

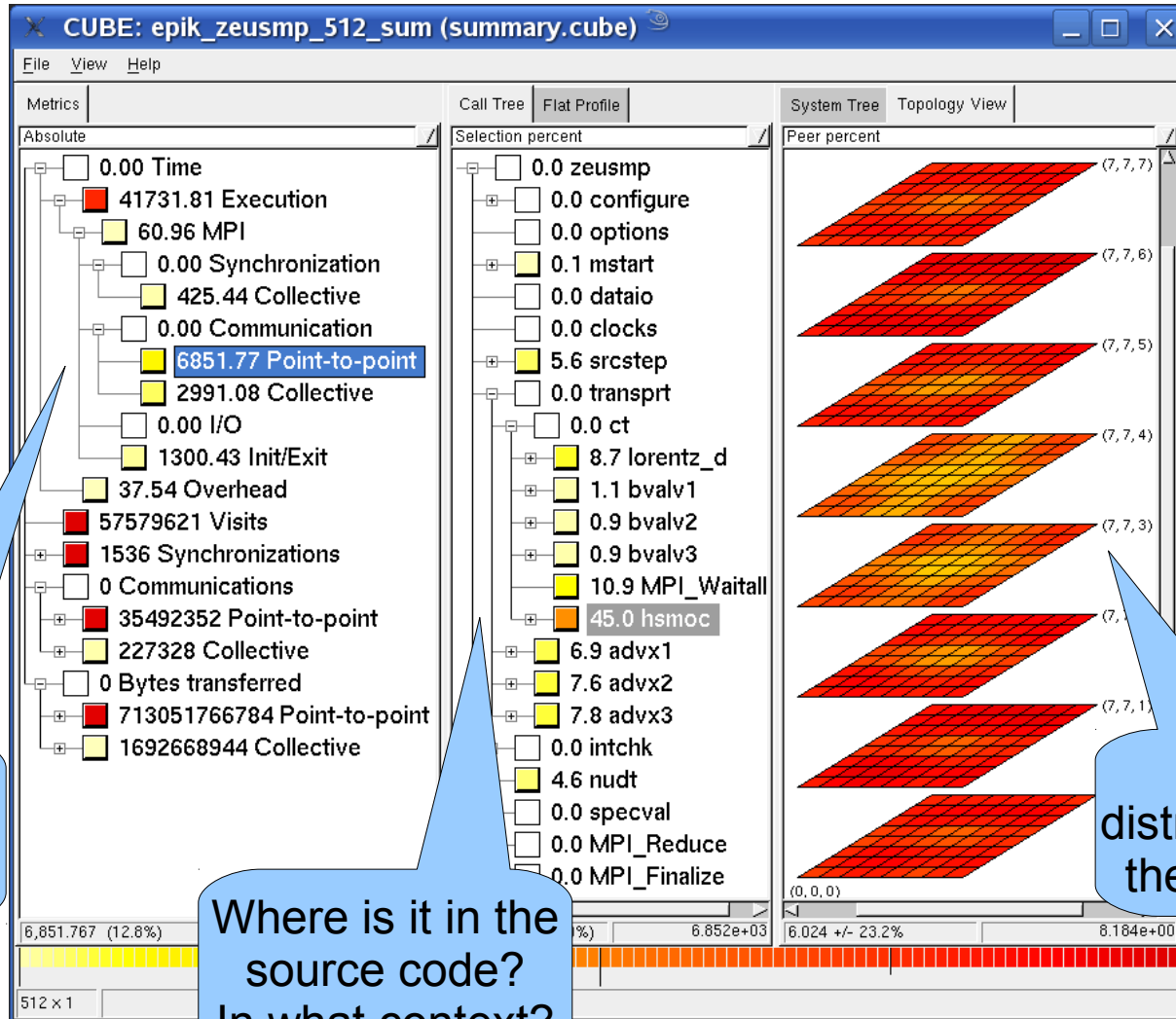
- Measurement & analysis runtime system
 - Manages runtime configuration and parallel execution
 - Configuration specified via EPIK.CONF file or environment
 - ▶ epik_conf reports current measurement configuration
 - Creates experiment archive (directory): **epik_<title>**
 - Optional runtime summarization report
 - Optional event trace generation (for later analysis)
 - Optional filtering of (compiler instrumentation) events
 - Optional incorporation of HWC measurements with events
 - ▶ via PAPI library, using PAPI preset or native counter names
- Experiment archive directory
 - Contains (single) measurement & associated files (e.g., logs)
 - Contains (subsequent) analysis reports

- Automatic instrumentation of OpenMP & POMP directives via source pre-processor
 - Parallel regions, worksharing, synchronization
 - Currently limited to OpenMP 2.5
 - ▶ No special handling of guards, dynamic or nested thread teams
 - ▶ Support for OpenMP 3.0 tasks currently in development
 - Configurable to disable instrumentation of locks, etc.
 - Typically invoked internally by instrumentation tools
- Used by Scalasca/Kojak, ompP, TAU, VampirTrace, etc.
 - Provided with Scalasca, but also available separately
 - ▶ OPARI 1.1 (October 2001)
 - ▶ OPARI 2.0 (July 2011)

- Parallel program analysis report exploration tools
 - Libraries for XML report reading & writing
 - Algebra utilities for report processing
 - GUI for interactive analysis exploration
 - ▶ requires Qt4 or wxGTK widgets library
 - ▶ can be installed independently of Scalasca instrumenter and measurement collector/analyzer, e.g., on laptop or desktop
- Used by Scalasca/Kojak, Marmot, ompP, PerfSuite, etc.
 - Analysis reports can also be viewed/stored/analyzed with TAU Paraprof & PerfExplorer
 - Provided with Scalasca, but also available separately
 - ▶ CUBE 3.3.2 (March 2011)

- Representation of values (severity matrix) on three hierarchical axes
 - Performance property (metric)
 - Call-tree path (program location)
 - System location (process/thread)
- Three coupled tree browsers
- CUBE3 displays severities
 - As value: for precise comparison
 - As colour: for easy identification of hotspots
 - Inclusive value when closed & exclusive value when expanded
 - Customizable via display mode





What kind of performance problem?

Where is it in the source code? In what context?

How is it distributed across the processes?

CUBE: epik_zeusmp_512_trace (trace.cube)

File View Help

Metrics

Absolute

- 0.00 Time
- 41826.08 Execution
- 92.79 MPI
- 0.00 Synchronization
- 429.92 Collective
- 0.00 Communication
- 3037.73 Point-to-point
- 4034.67 Late Sender
- 0.00 Late Receiver
- 3089.32 Collective
- 0.00 File I/O
- 1488.13 Init/Exit
- 1849.98 Overhead
- 57579621 Visits
- 1536 Synchronizations
- 0 Communications
- 35492352 Point-to-point
- 227328 Collective
- 0 Bytes transferred
- 713051766784 Point-to-point
- 1692668944 Collective
- 2042.71 Computational imbalance

Selection percent

- 0.0 zeusmp
- 0.0 configure
- 0.0 options
- 0.1 mstart
- 0.0 dataio
- 0.0 clocks
- 11.8 srcstep
- 0.0 transprt
- 0.0 ct
- 11.3 lorentz_d
- 2.4 bvalv1
- 2.1 bvalv2
- 2.0 bvalv3
- 1.5 MPI_Waitall
- 14.9 hsmoc
- 13.6 advx1
- 15.2 advx2
- 16.0 advx3
- 0.0 intchk
- 9.1 nudt
- 0.0 specval
- 0.0 MPI_Reduce
- 0.0 MPI_Finalize

3,037.735 (5.4%) 5.585e+04 453.671 (14.9%) 3.038e+03 0.886 +/- 19.3% 1.260e+00

512 x 1

CUBE metric description <@j36>

Late Sender Time

Description:
Refers to the time lost waiting caused by a blocking receive operation (e.g., MPI_Recv() or MPI_Wait()) that is posted earlier than the corresponding send operation.

location

time

<< Close >>

(7,7,4)

(7,7,3)

(7,7,2)

(7,7,1)

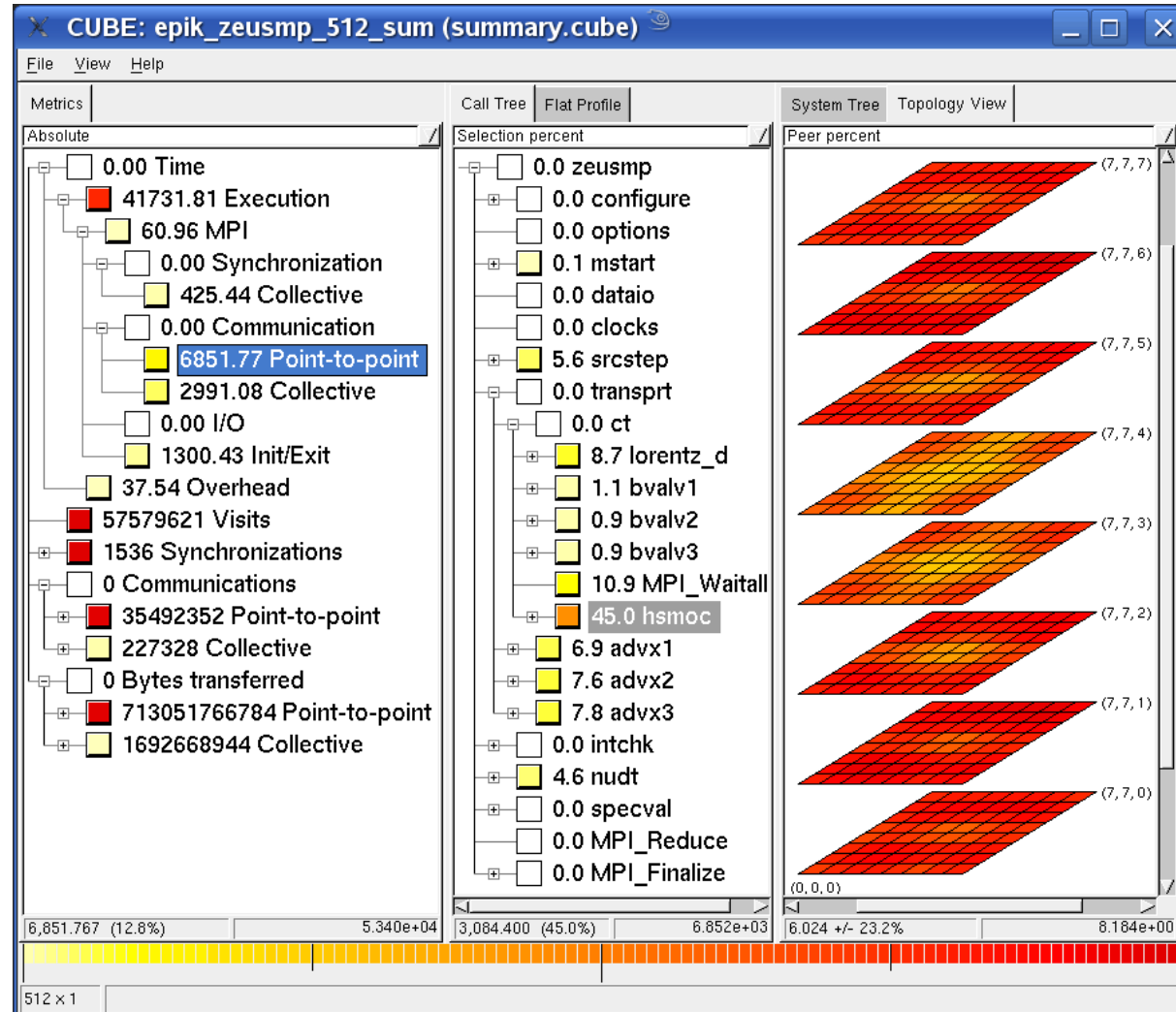
(7,7,0)

(0,0,0)

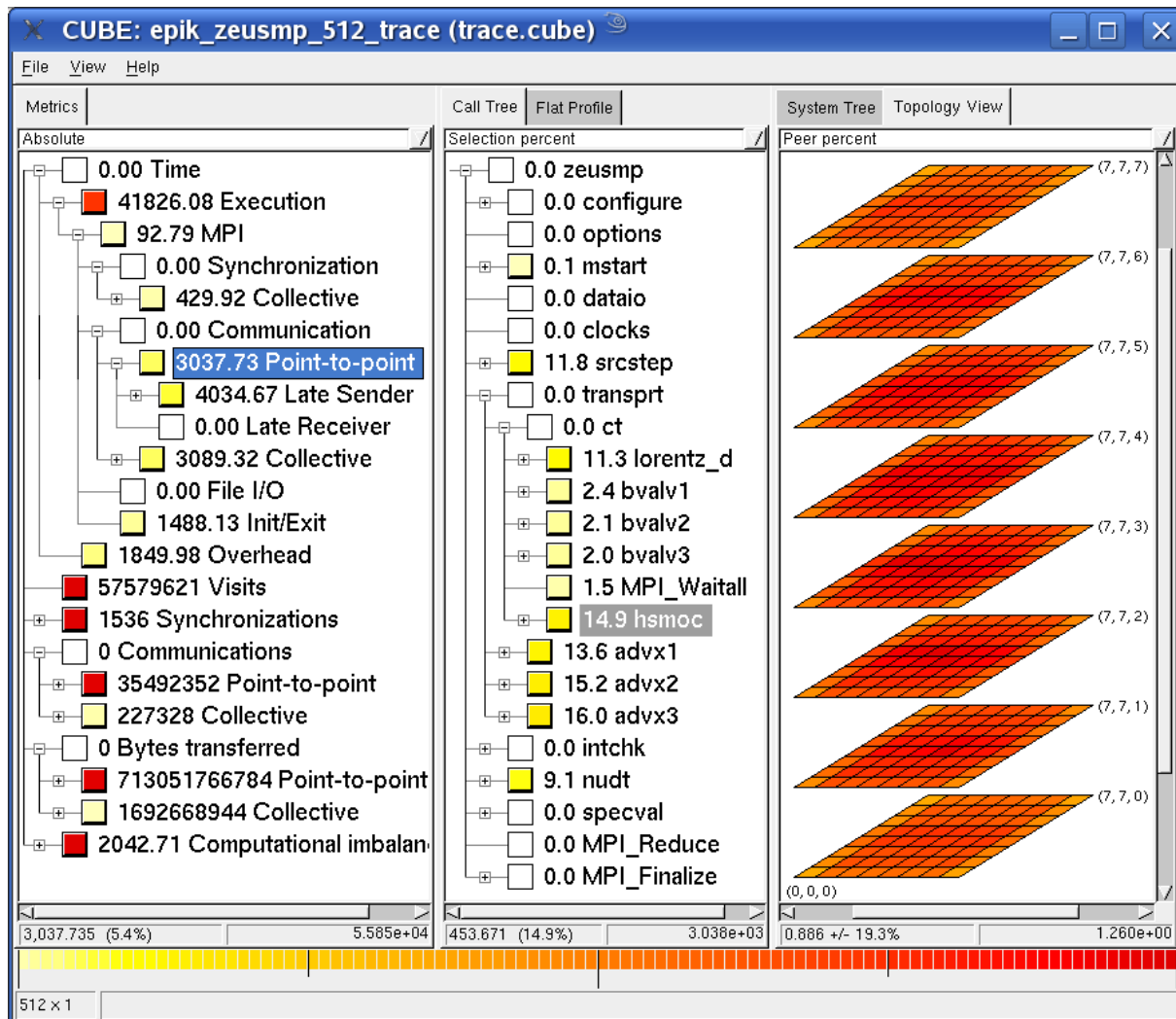
Additional metrics determined from trace

- Computational astrophysics
 - (magneto-)hydrodynamic simulations on 1-, 2- & 3-D grids
 - part of SPEC MPI2007 1.0 benchmark suite (132.zeusmp2)
 - developed by UCSD/LLNL
 - >44,000 lines Fortran90 (in 106 source modules)
 - provided configuration scales to 512 MPI processes
- Run with 512 processes on JUMP
 - IBM p690+ eServer cluster with HPS at JSC
- Scalasca summary and trace measurements
 - ~5% measurement dilation (full instrumentation, no filtering)
 - 2GB trace analysis in 19 seconds
 - application's 8x8x8 grid topology automatically captured from MPI Cartesian

- 12.8% of time spent in MPI point-to-point communication
- 45.0% of which is on program callpath transprt/ct/hsmoc
- With 23.2% std dev over 512 processes
- Lowest values in 3rd and 4th planes of the Cartesian grid



- MPI point-to-point communication time separated into transport and Late Sender fractions
- Late Sender situations dominate (57%)
- Distribution of transport time (43%) indicates congestion in interior of grid



- Automatic function instrumentation (and filtering)
 - CCE, GCC, IBM, Intel, PathScale & PGI compilers
 - optional PDTToolkit selective instrumentation (when available) and manual instrumentation macros/pragmas/directives
- MPI measurement & analyses
 - scalable runtime summarization & event tracing
 - only requires application executable re-linking
 - P2P, collective, RMA & File I/O operation analyses
- OpenMP measurement & analysis
 - requires (automatic) application source instrumentation
 - thread management, synchronization & idleness analyses
- Hybrid OpenMP/MPI measurement & analysis
 - combined requirements/capabilities
 - parallel trace analysis requires uniform thread teams

- Improved configure/installation
- Support for using PDTToolkit to instrument sources
 - selective instrumentation of source files and routines
- Consistent instrumentation selection
 - automatic (compiler/pdt) and/or manual (pomp/user)
- Measurement configuration of MPI event wrappers
 - specify desired categories of events, e.g., P2P, COLL, RMA
- MPI RMA (one-sided communication) analysis
- Improved OpenMP (and hybrid) measurement & analysis
 - specify desired number of threads: ESD_MAX_THREADS
 - consistent automatic analyses of traces
- Improved documentation of analysis reports

- Instrumentation
 - Separate OpenMP instrumenter (OPARI) distribution
 - Scalasca source instrumentation via TAU/PDTToolkit
 - Adapter for VT manual instrumentation macros
 - TAU instrumentation with Scalasca measurement libraries
- Trace utilities
 - Trace conversion utilities for VT/OTF, Paraver, JumpShot
 - Vampir visualization of Scalasca traces (without conversion)
- Analysis report utilities
 - Separate report generation/manipulation library and GUI (CUBE3) distribution
 - Alternative presentation with TAU Paraprof/PerfExplorer
- Part of Uniform Integrated Tool Environment (UNITE)