**JÜLICH**
FORSCHUNGSZENTRUM

# Large-scale performance analysis of *PFLOTRAN* with Scalasca

2011-05-26  |  **Brian J. N. Wylie** & Markus Geimer
Jülich Supercomputing Centre

b.wylie@fz-juelich.de

# Overview

Dagstuhl challenge

- *PFLOTRAN* at large-scale on Cray XT5 & IBM BG/P

Scalasca performance analysis toolset

- Overview of architecture and usage
- Instrumentation of *PFLOTRAN* & PETSc
- Summary & trace measurement experiments
- Analysis report & trace exploration

Scalasca scalability improvement

- Unification of definition identifiers

Conclusions

# Dagstuhl challenge

Demonstrate performance measurement and analysis of *PFLOTRAN* using "2B" problem dataset and more than 10,000 MPI processes

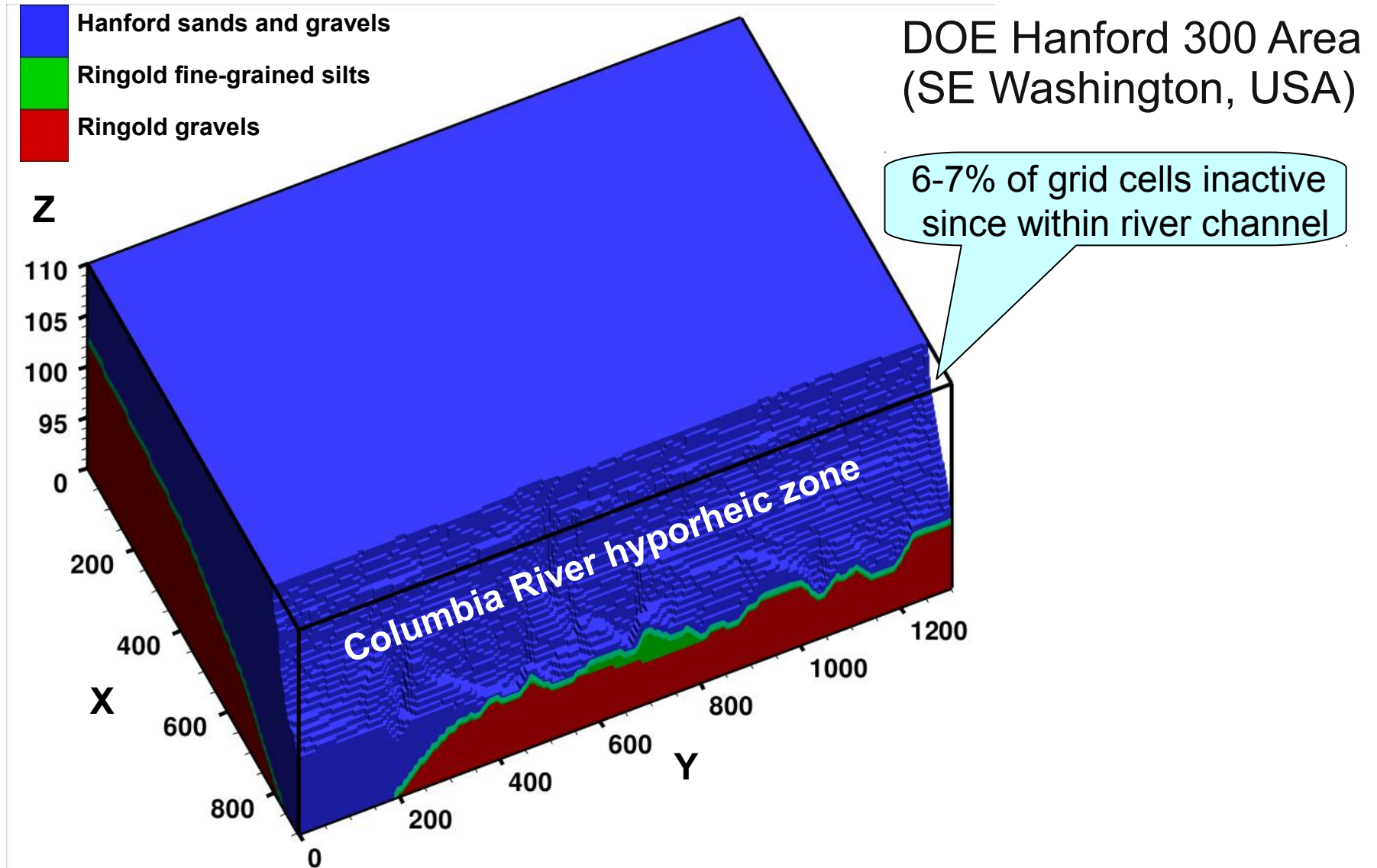Challenge issued for Dagstuhl seminar 10181 (3-7 May 2010) on "Program development for extreme-scale computing"

- two months notice with download/build/run instructions
- accounts/allocations provided by VI-HPS & INCITE/PEAC
    - *IBM BG/P (**jugene**.fz-juelich.de)*
    - *Cray XT5 (**jaguar**.nccs.ornl.gov)*
- numerous research groups and vendors presented their results and discussed issues encountered
    - *http://www.dagstuhl.de/10181*

# PFLOTRAN

3D reservoir simulator actively developed by LANL/ORNL/PNNL

- approx. 80,000 lines of Fortran9X, combining solvers for
  - *PFLOW non-isothermal, multi-phase groundwater flow*
  - *PTRAN reactive, multi-component contaminant transport*
- employs PETSc, LAPACK, BLAS & HDF5 I/O libraries
  - *87 PFLOTRAN + 789 PETSc source files*
  - *parallel I/O tuning via PERI active liaison*
- "2B" input dataset run for 10 timesteps
  - *uses 3-dimensional (non-MPI) PETSc Cartesian grid*
  - *TRAN(sport) step scales much better than FLOW step*
  - *FLOW step generally faster, but crossover at larger scales*

# *PFLOTRAN "2B" test case*



DOE Hanford 300 Area
(SE Washington, USA)

6-7% of grid cells inactive
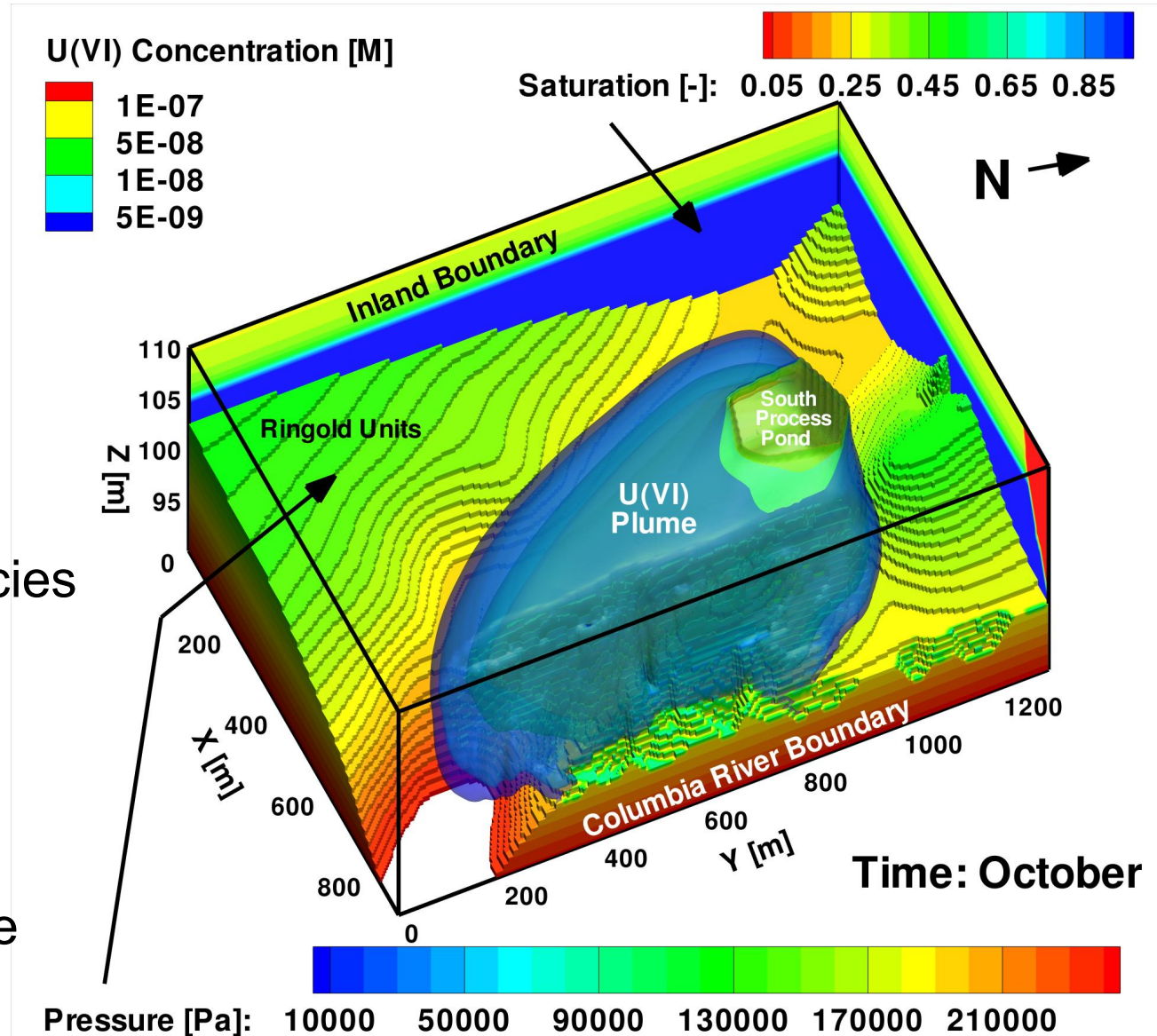since within river channel

# *PFLOTRAN* simulation of U(VI) migration

DOE Hanford 300 Area

Problem domain:
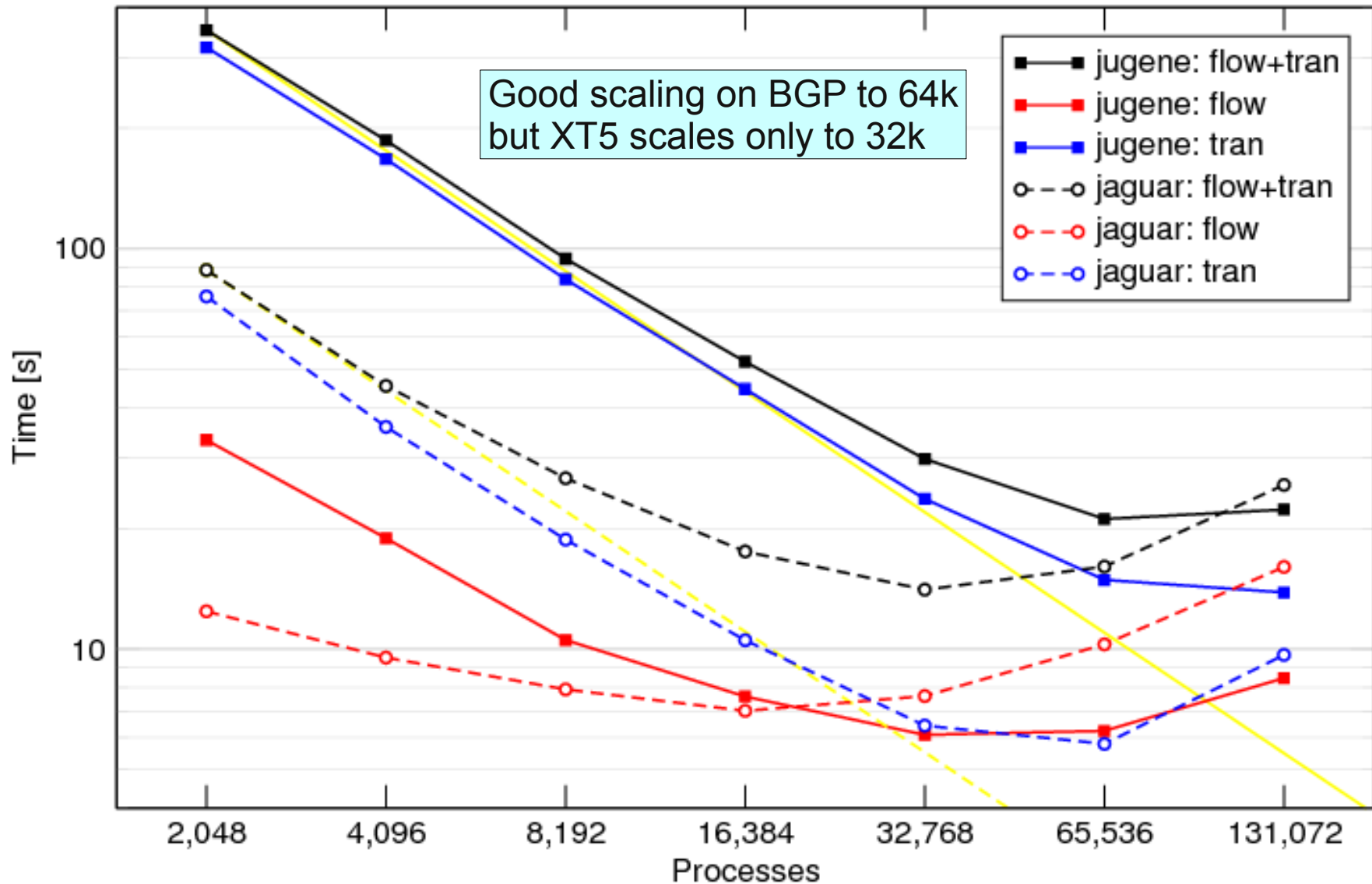
- 900x1300x20m
- $\Delta x/\Delta y$ = 5 m
- 1.87M grid cells
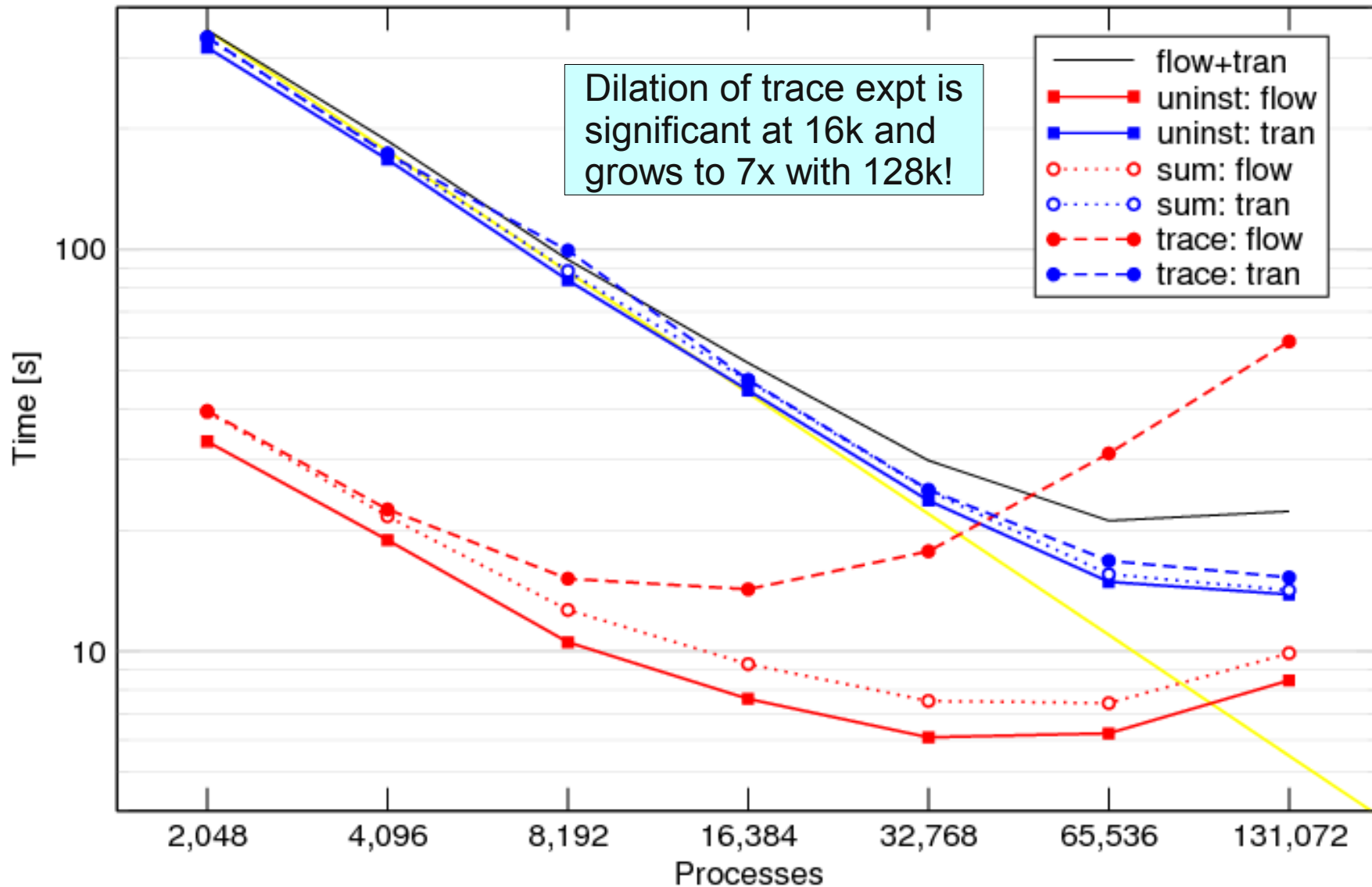- 15 chemical species
- 28M DoF total

1-year simulation:
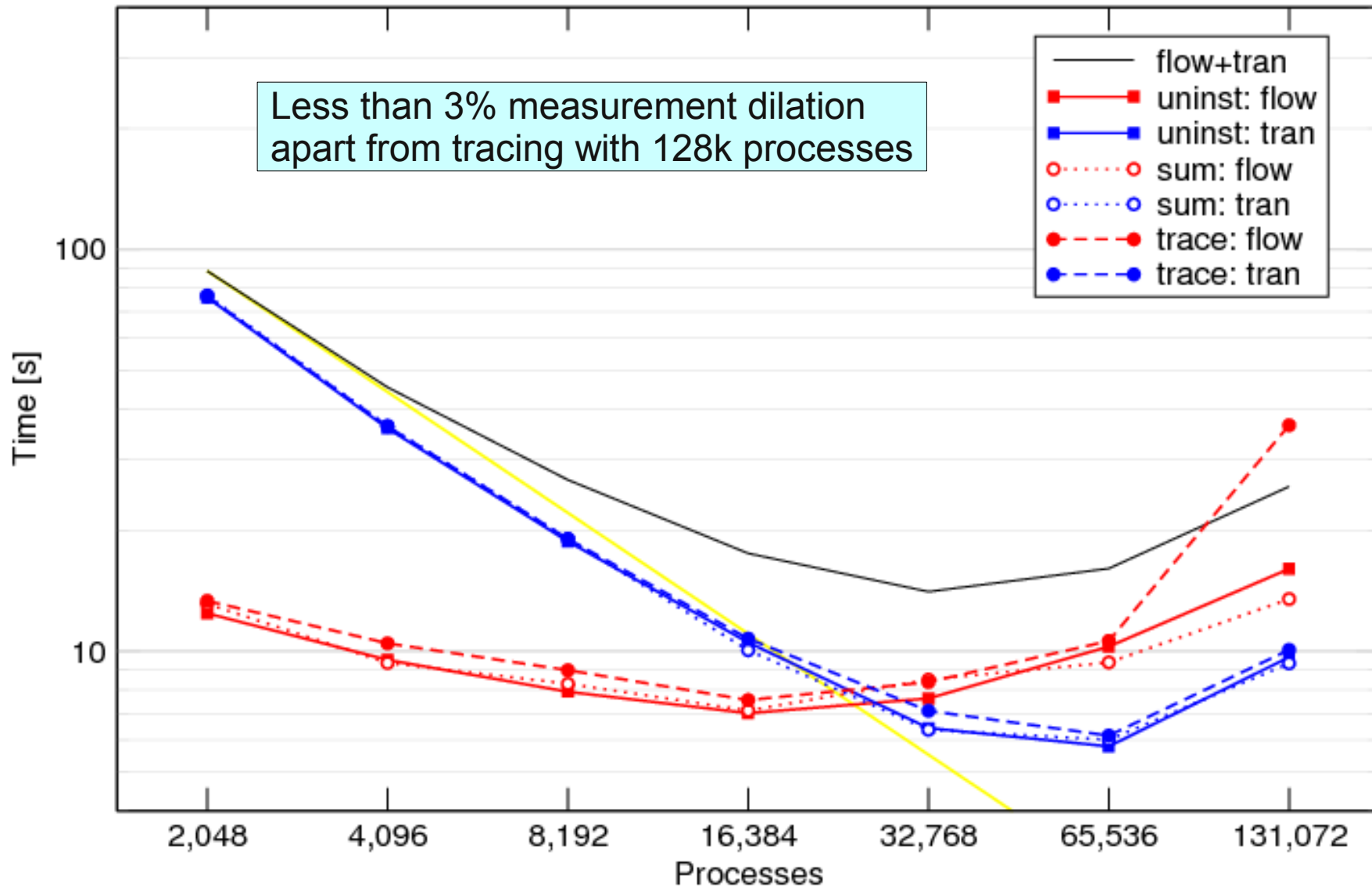
- $\Delta t$ = 1 hour
- 5-10 hour runtime
- 4096 XT5 cores

# *PFLOTRAN* "2B" strong scalability



Good scaling on BGP to 64k
but XT5 scales only to 32k

Legend:
- jugene: flow+tran
- jugene: flow
- jugene: tran
- jaguar: flow+tran
- jaguar: flow
- jaguar: tran

Y-axis: Time [s]

X-axis: Processes — 2,048 | 4,096 | 8,192 | 16,384 | 32,768 | 65,536 | 131,072

# Scalasca *PFLOTRAN* "2B" measurement dilation (BGP)



Dilation of trace expt is significant at 16k and grows to 7x with 128k!

Legend:
- flow+tran
- uninst: flow
- uninst: tran
- sum: flow
- sum: tran
- trace: flow
- trace: tran

Y-axis: Time [s]
X-axis: Processes — 2,048  4,096  8,192  16,384  32,768  65,536  131,072

# Scalasca *PFLOTRAN* "2B" measurement dilation (XT5)



Less than 3% measurement dilation apart from tracing with 128k processes

Legend:
- flow+tran
- uninst: flow
- uninst: tran
- sum: flow
- sum: tran
- trace: flow
- trace: tran

# Scalasca project

Overview

- Headed by Bernd Mohr (JSC) & Felix Wolf (GRS-Sim)
  - *Helmholtz Initiative & Networking Fund project started in 2006*
- 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**

Status

- Scalasca v1.3.3 released in March 2011
- Available for download from www.scalasca.org

# Scalasca features

Open source, New BSD license

Portable

- IBM BlueGene, IBM SP & blade clusters, Cray XT/XE, NEC SX, SGI Altix, SiCortex, Linux cluster® (SPARC, x86-64), ...

Supports typical HPC languages & parallel programming paradigms

- Fortran, C, C++
- MPI, OpenMP & hybrid MPI/OpenMP

Integrated instrumentation, measurement & analysis toolset

- Customizable automatic/manual instrumentation
- Runtime summarization (*aka* profiling)
- Automatic event trace analysis

# Scalasca components

- Automatic program instrumenter creates instrumented executable

- Unified measurement library supports both

  - *runtime summarization*

  - *trace file generation*

- Parallel, replay-based event trace analyzer invoked automatically on set of traces

- Common analysis report explorer & examination/processing tools

# Scalasca usage (commands)

1. Prepare application objects and executable for measurement:
   - **scalasca -instrument**  `cc -O3 -c …`
   - **scalasca -instrument**  `ftn -O3 -o pflotran.exe …`
     - *instrumented executable pflotran.exe produced*

2. Run application under control of measurement & analysis nexus:

   - **scalasca -analyze**     `aprun -N 12 -n 65536 pflotran.exe …`
     - *epik_pflotran_12p65536_sum*  experiment produced
   - **scalasca -analyze -t**  `aprun -N 12 -n 65536 pflotran.exe …`
     - *epik_pflotran_12p65536_trace*  experiment produced

   BATCH JOB

3. Interactively explore experiment analysis report:

   - **scalasca -examine**     `epik_pflotran_12p65536_trace`
     - *epik_pflotran_12p65536_trace/trace.cube.gz*  presented

# Measurement & analysis methodology

1. Run uninstrumented/optimized version (as reference for validation)
   - determine memory available for measurement
2. Run automatically-instrumented version collecting runtime summary
   - determine functions with excessive measurement overheads
       - *examine distortion and trace buffer capacity requirement*
   - if necessary, prepare filter file and repeat measurement
3. Reconfigure measurement to collect and automatically analyze traces

(optional) Customize analysis report and/or instrumentation, e.g.,
   - extract key code sections into specific analysis reports
   - annotate key code sections with *EPIK* instrumentation API macros

# Scalasca usage with *PFLOTRAN*

Automatic instrumentation

- both *PFLOTRAN* application (Fortran) & PETSc library (C)
- USR routines instrumented by IBM XL & PGI compilers
- MPI routine interposition with instrumented library (PMPI)

Initial summary measurements used to define filter files specifying all purely computational routines

Summary & trace experiments collected using filters

- parallel trace analysis initiated automatically on same partition

Post-processing of analysis reports

- cut to extract timestep loop; incorporation of application topology

Analysis report examination in GUI

# *PFLOTRAN* structural analysis

Determined by scoring summary experiment
using fully-instrumented application executable



- 29 MPI library routines used

- 1100+ PFLOTRAN & PETSc routines

  - *most not on a callpath to MPI, purely local calculation (USR)*

  - *~250 on callpaths to MPI, mixed calculation & comm. (COM)*

- Using measurement filter listing all USR routines

  - *maximum callpath depth 22 frames*

  - *~1750 unique callpaths (399 in FLOW, 375 in TRAN)*

  - *633 MPI callpaths (121 in FLOW, 114 in TRAN)*

- FLOW callpaths very similar to TRAN callpaths

# *PFLOTRAN* stepperrun callgraph

ParaProf view: nodes coloured by exclusive execution time

Many paths lead to MPI_Allreduce

# Scalasca analysis report: program call tree

# Scalasca analysis report: performance metric tree

# Scalasca analysis report: system tree/topology

# Exclusive Execution time



13% computational imbalance in various routines where a relativatively small number of processes on edge of grid are much faster
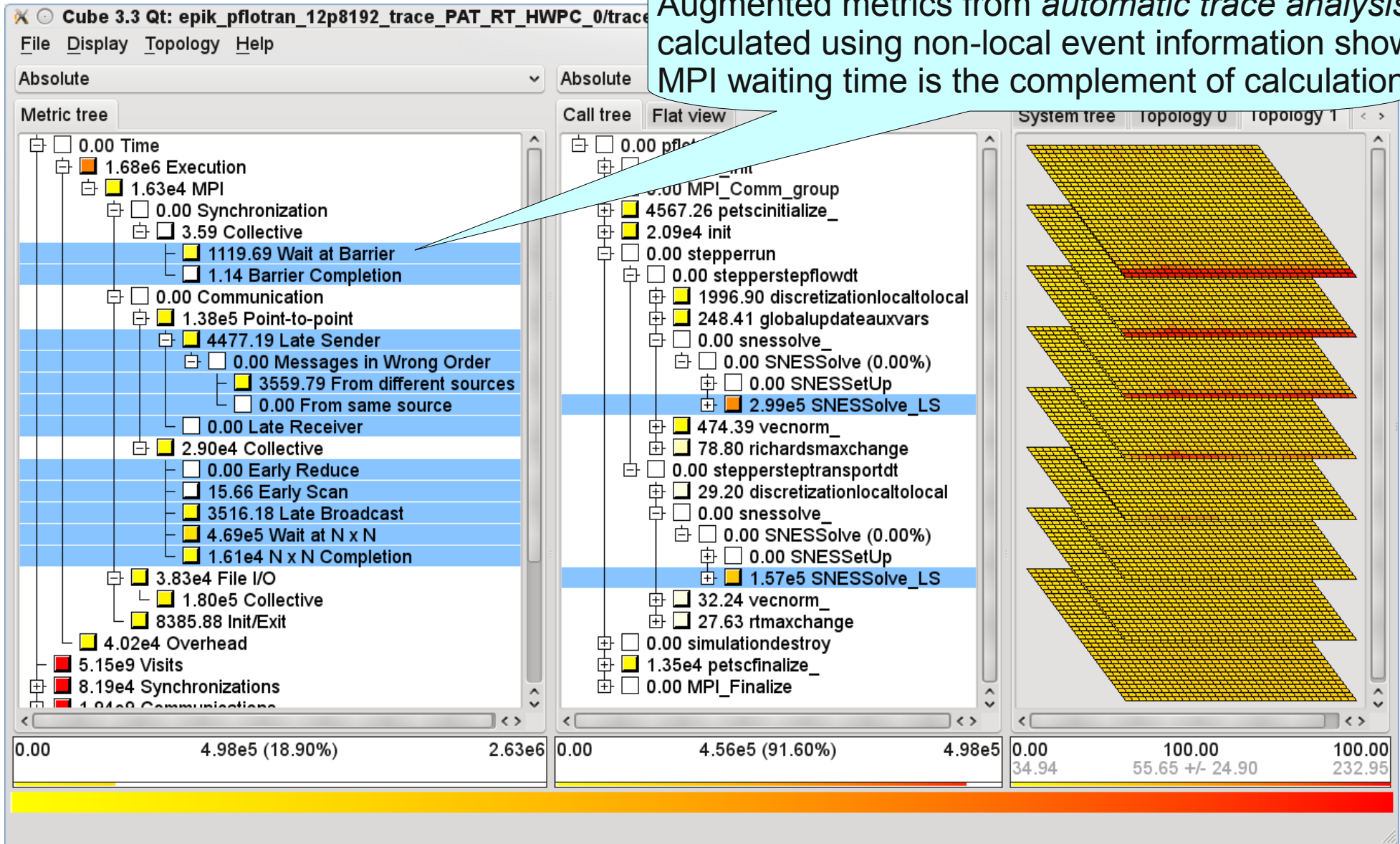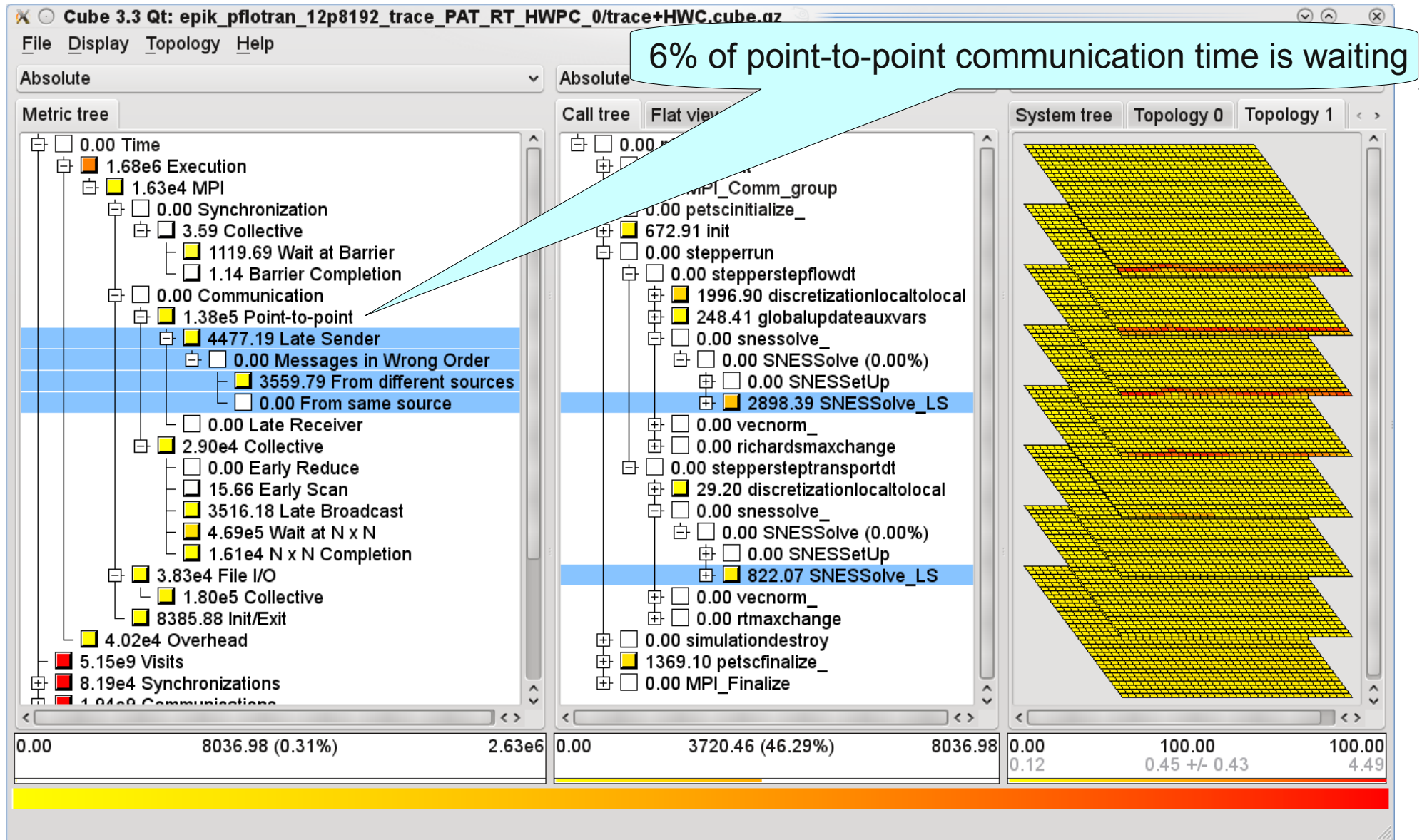
# Floating-point operations

# Scalasca trace analysis metrics (waiting time)



Augmented metrics from *automatic trace analysis* calculated using non-local event information show MPI waiting time is the complement of calculation

# Late Sender time (for early receives)



6% of point-to-point communication time is waiting

# Scalasca description for *Late Sender* metric



- Analysis report explorer GUI provides hyperlinked descriptions of performance properties

- Diagnosis hints suggest how to refine diagnosis of performance problems and possible remediation

# Computational imbalance: overload

# *PFLOTRAN* grid decomposition imbalance

850x1000x160 cells decomposed on 65536=64x64x16 process grid

- *x*-axis: 850/64=13 plus 18 extra cells
- *y*-axs: 1000/64=15 plus 40 extra cells
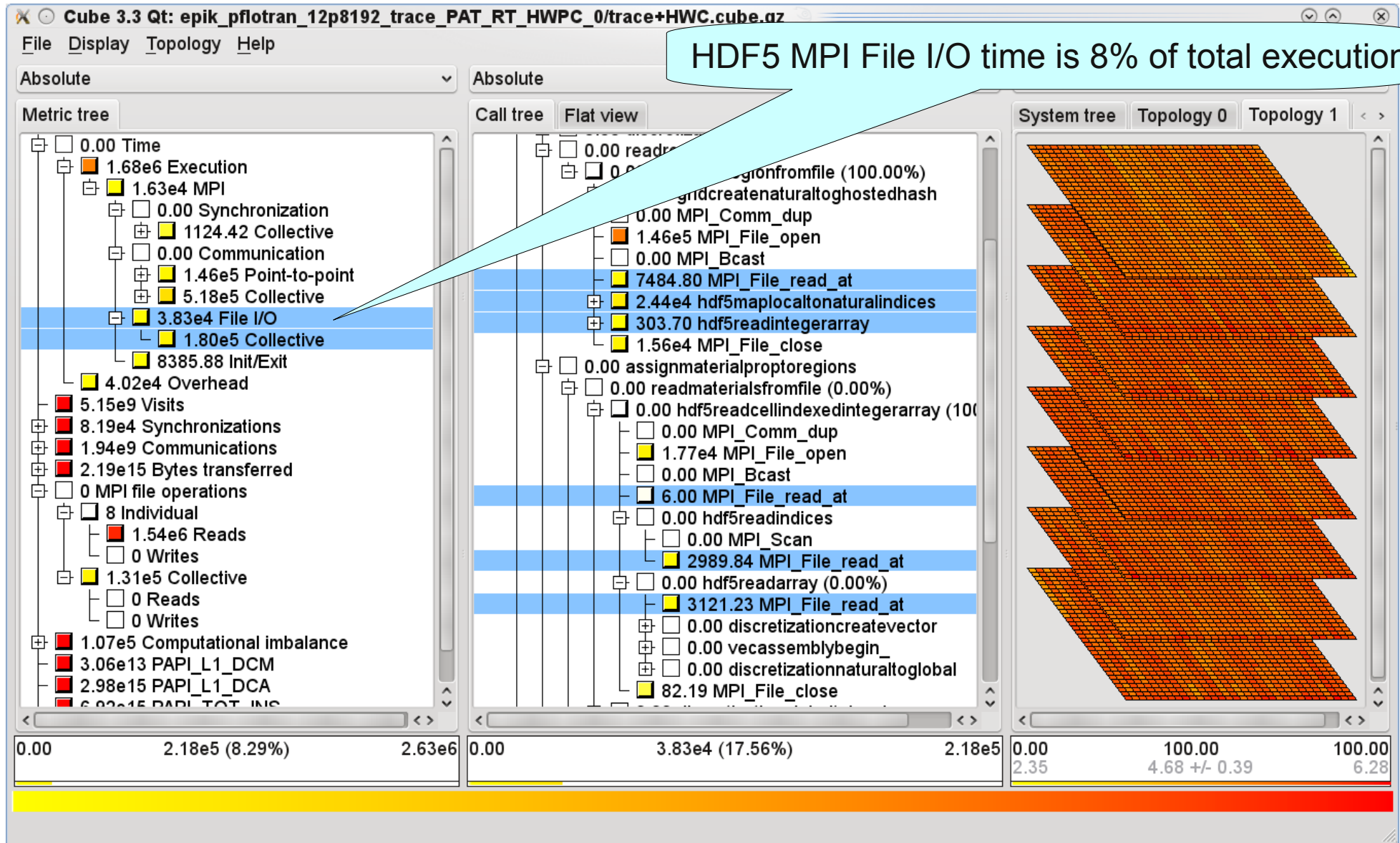- *z*-axis: 160/16=10

Perfect distribution would be 2075.2 cells per process

- but 20% of processes in each *z*-plane have 2240 cells
- 8% computation overload manifests as waiting times at the next communication/synchronization on the other processes
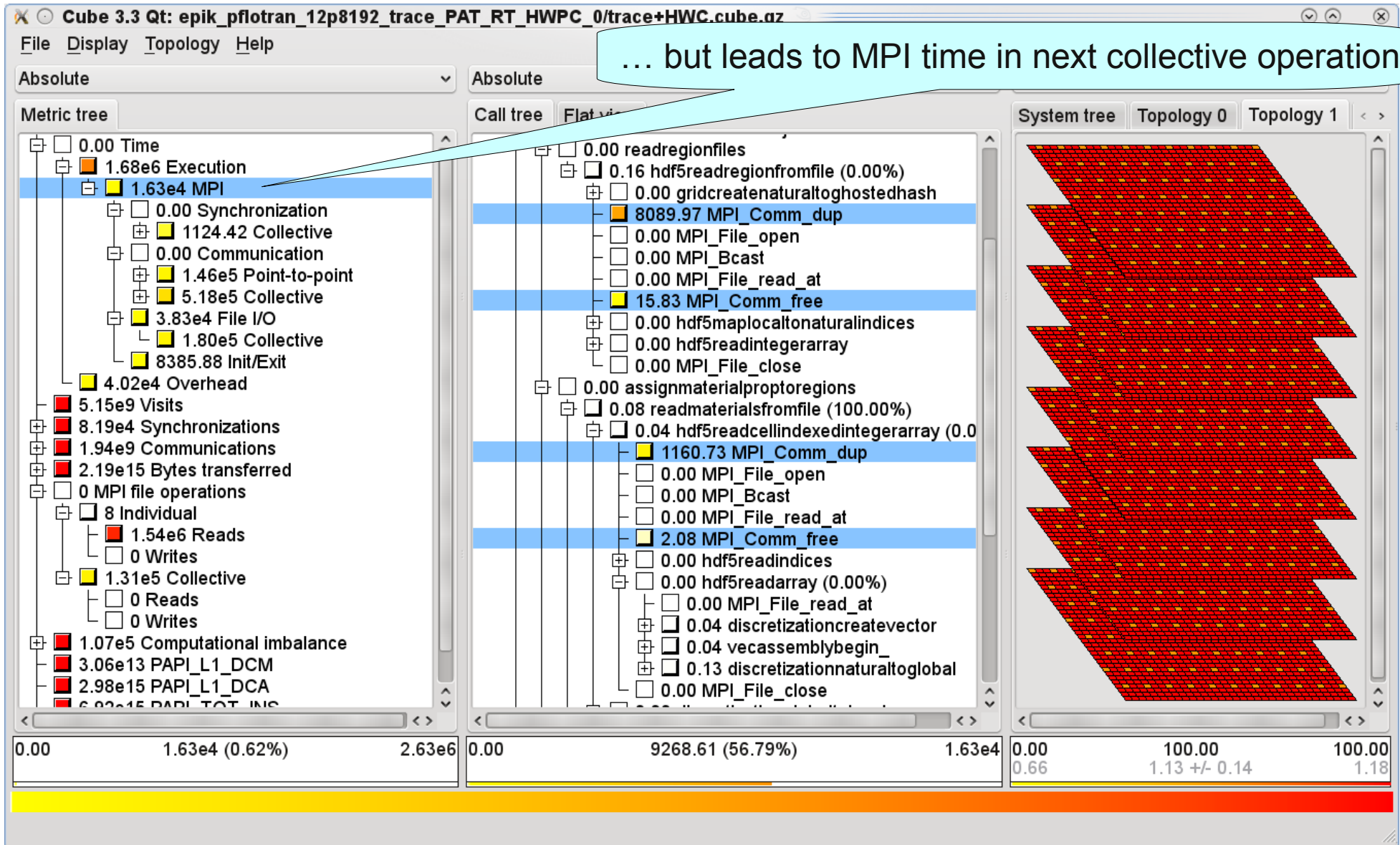
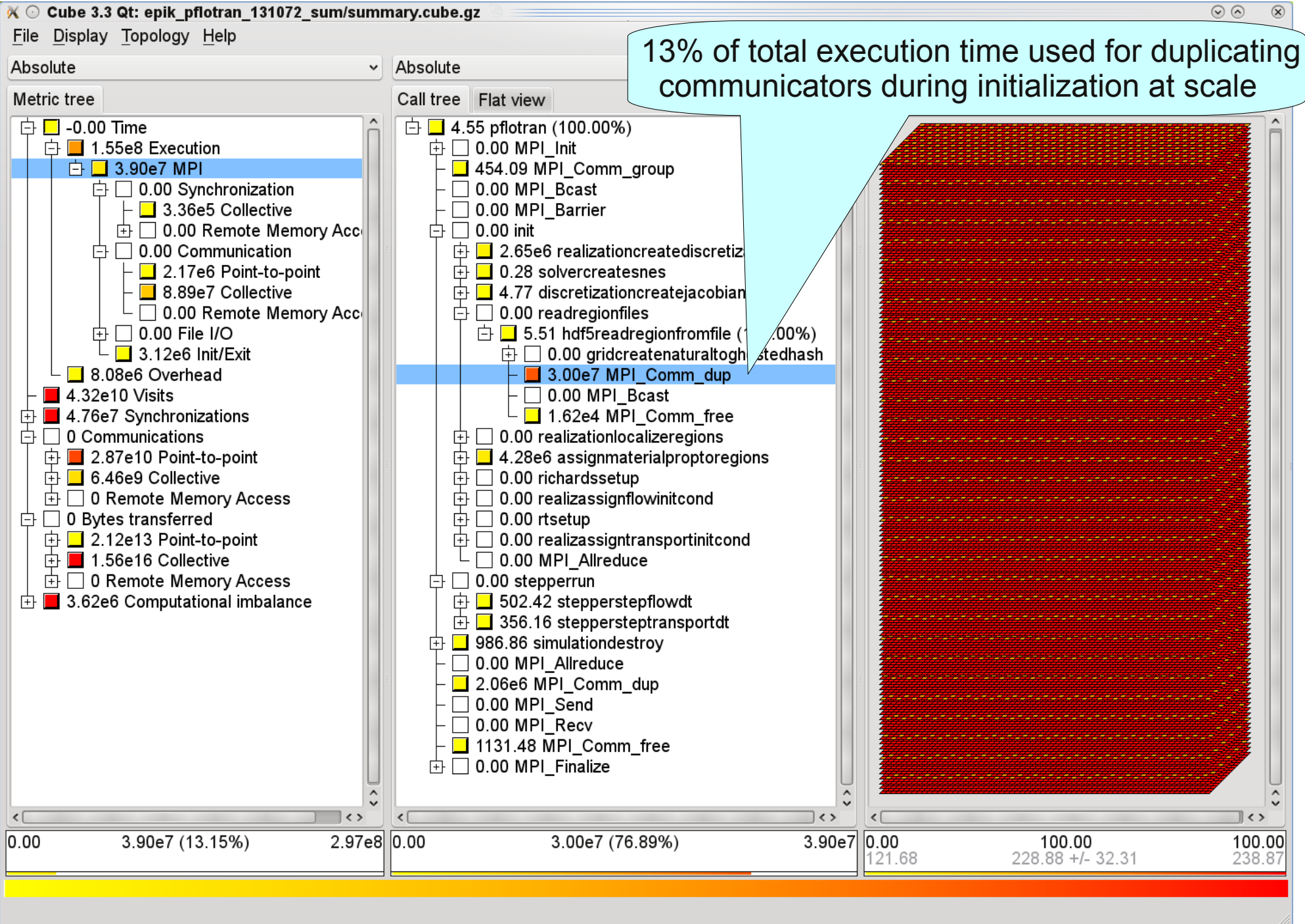The problem-specific localized imbalance in the river channel is minor

- reversing the assignments in the *x*-dimension won't help much since some of the *z*-planes have no inactive cells

# MPI File I/O time

# MPI Communicator duplication

13% of total execution time used for duplicating communicators during initialization at scale
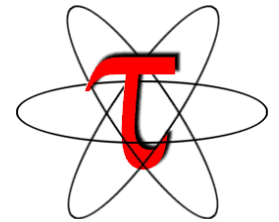
# Complementary analyses and visualizations

TAU/ParaProf can import Scalasca analysis reports

- part of portable open-source TAU performance system

- provides a callgraph display and various graphical profiles

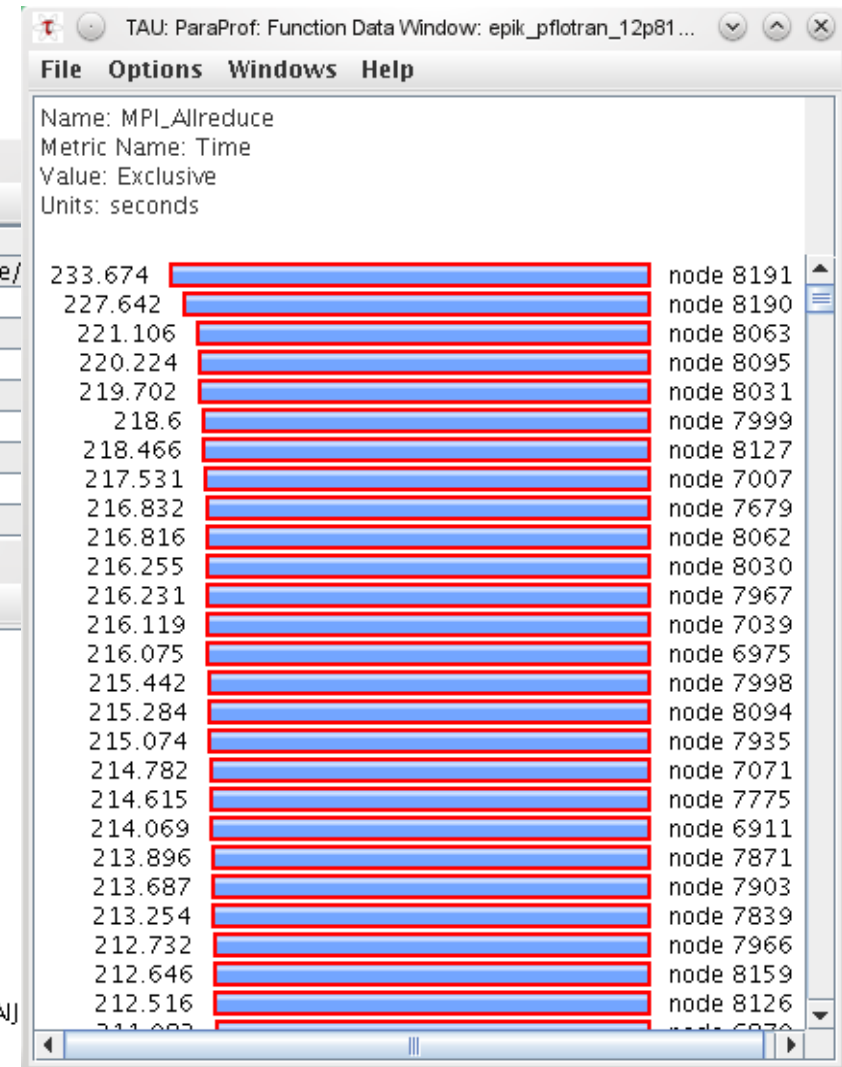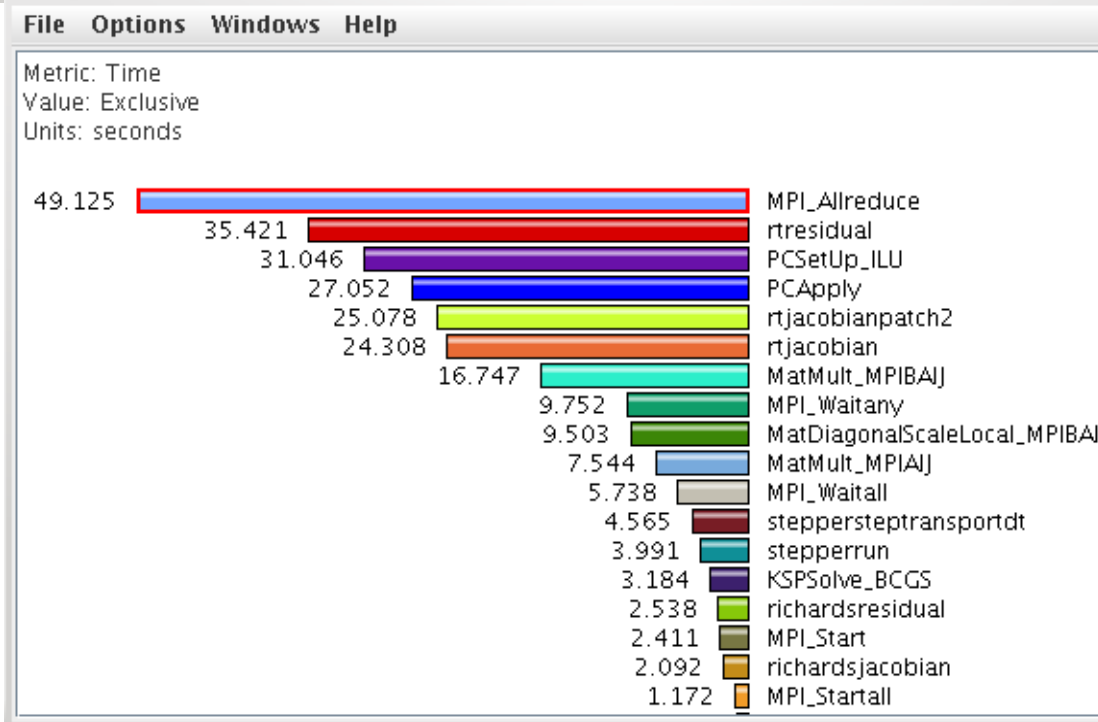- may need to extract part of report to fit available memory

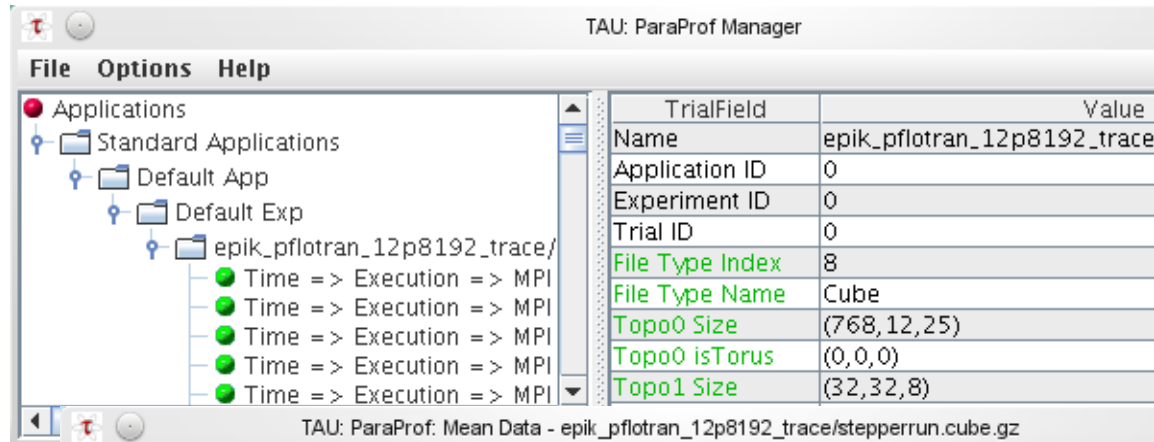Vampir 7 can visualize Scalasca distributed event traces

- part of commercially-licensed Vampir product suite

- provides interactive analysis & visualization of execution intervals

  - *powerful zoom and scroll to examine fine detail*

- avoids need for expensive trace merging and conversion

  - *required for visualization with Paraver & JumpShot*

  - *but often prohibitive for large traces*

# TAU/ParaProf graphical profiles



ParaProf views of Scalasca trace analysis report (extract featuring stepperrun subtree)

ParaProf callpath profile and process breakdown

# Vampir overview of *PFLOTRAN* execution (init + 10 steps)

# *PFLOTRAN* execution timestep 7 (flow + transport)

# *PFLOTRAN* execution flow/transport transition

# *PFLOTRAN* execution at end of flow phase

# *PFLOTRAN* execution at end of flow phase detail

# Scalasca/*PFLOTRAN* scalability issues

*PFLOTRAN* (via PETSc & HDF5) uses lots of MPI communicators

- 19 (MPI_COMM_WORLD) + 5xNPROCS (MPI_COMM_SELF)
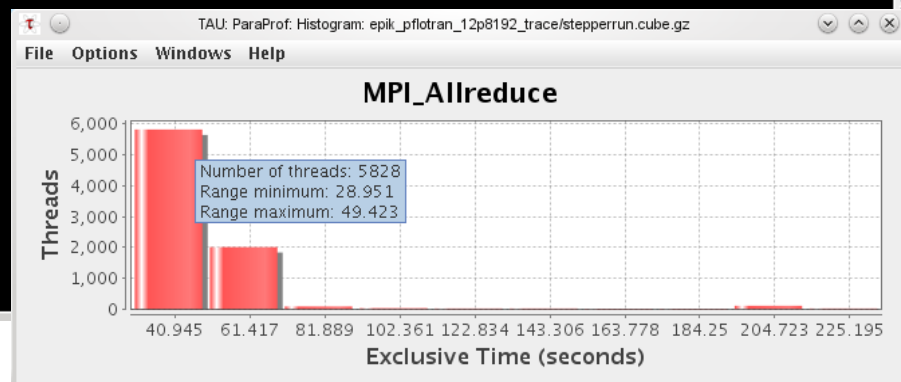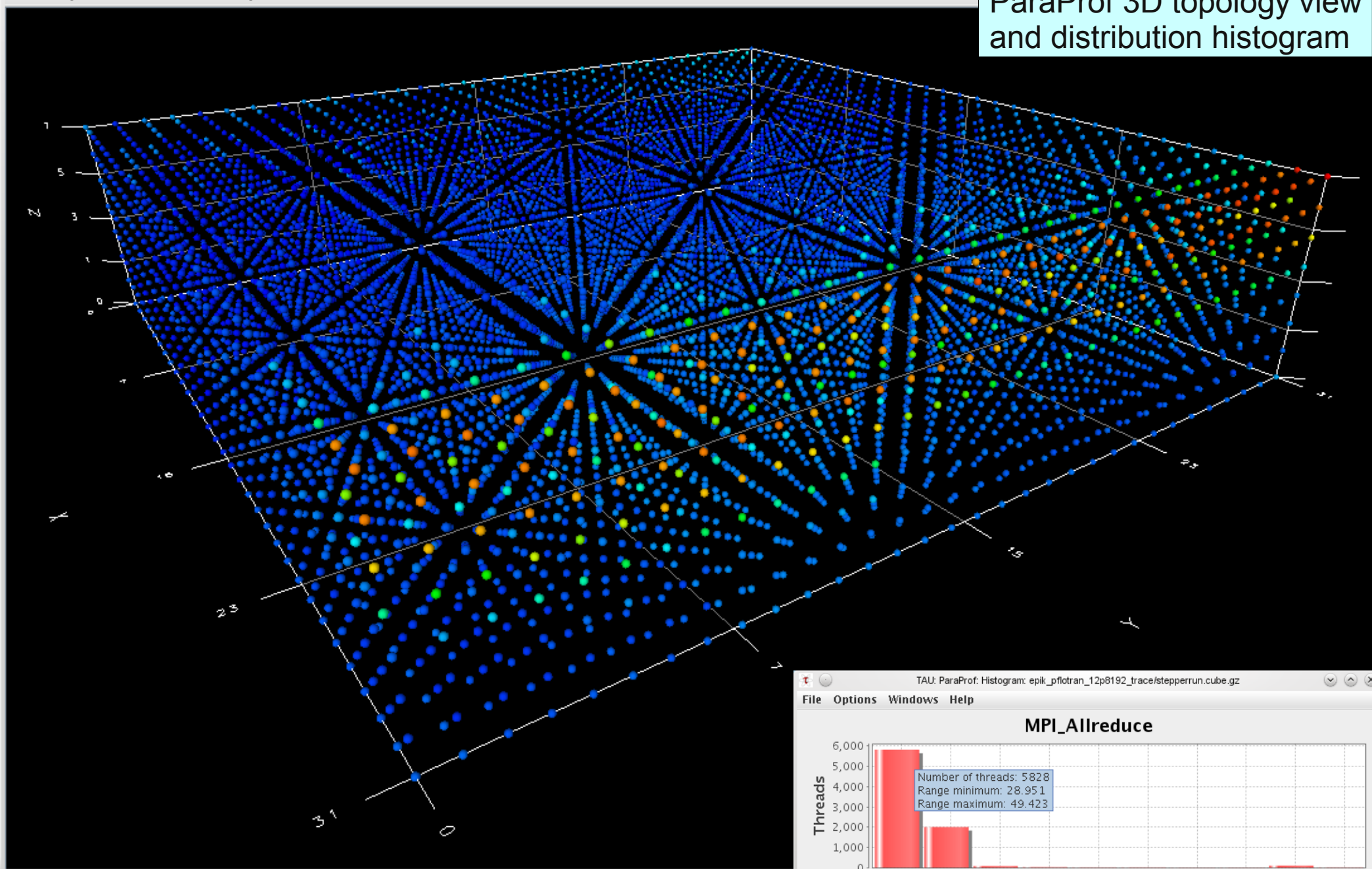  - *time shows up in collective MPI_Comm_dup*
- Not needed for summarization, but essential for trace replay
  - *definition storage & unification time grow accordingly*
- Original 1.3.1 version of Scalasca failed with 48k processes
  - *communicator definitions 1.42 GB, unification over 29 minutes*
- Revised 1.3.2 version resolved scalability bottleneck
  - *custom management of MPI_COMM_SELF communicators*
  - *hierarchical implementation of unification*

Scalasca trace collection and analysis costs increase with scale

# Improved unification of definition identifiers

# Trace collection & analysis scalability

# Conclusions

Very complex applications like *PFLOTRAN* provide significant challenges for performance analysis tools

Scalasca offers a range of instrumentation, measurement & analysis capabilities, with a simple GUI for interactive analysis report exploration

- works across Cray XT/XE and IBM BlueGene systems

- analysis reports and event traces can also be examined with complementary third-party tools such as TAU/ParaProf & Vampir

- convenient automatic instrumentation of applications and libraries must be moderated with selective measurement filtering

    - *analysis reports can still become awkwardly large*

Scalasca is continually improved in response to the evolving requirements of application developers and analysts

# Acknowledgments

The application and benchmark developers who generously provided their codes and/or measurement archives

The facilities who made their HPC resources available and associated support staff who helped us use them effectively

- ALCF, BSC, **CSC**, **CSCS**, **EPCC**, JSC, HLRN, **HLRS**, ICL, LRZ, NCAR, **NCCS**, **NICS**, RWTH, RZG, SARA, TACC, ZIH

- Access & usage supported by European Union, German and other national funding organizations

Cray personnel for answers to our detailed questions

Scalasca users who have provided valuable feedback and suggestions for improvements

# **Sc**alable performance **a**nalysis of **la**rge-**sc**ale parallel **a**pplications

- portable toolset for scalable performance measurement & analysis of MPI, OpenMP & hybrid OpenMP/MPI parallel applications

- supporting most popular HPC computer systems

- available under New BSD open-source license

- distributed on POINT/VI-HPS Parallel Productivity Tools Live-DVD

- sources, documentation & publications:

  - *http://www.scalasca.org*

  - *mailto: scalasca@fz-juelich.de*

# Scalasca training

Tutorials

- full- or part-day hands-on training using POINT/VI-HPS Live DVD
- *HP-SEE/LinkSCEEM/PRACE*, Athens, Greece: 15 July 2011
- *EU/US Summer School*, South Lake Tahoe, USA: 9 Aug 2011
- *PRACE Summer School*, Espoo/Helsinki, Finland: 1 Sep 2011

VI-HPS Tuning Workshops

- 5-day hands-on training with participants' own application codes, in context of Virtual Institute – High Productivity Supercomputing
- *8th VI-HPS Tuning Workshop*, Aachen, Germany: 5-9 Sept 2011

Visit www.vi-hps.org/training for further details & announcements

# Scalasca functionality in latest releases

Generic improvements

- Hybrid OpenMP/MPI measurement & analysis
- Selective routine instrumentation with *PDToolkit* instrumenter
- Runtime filtering of routines instrumented by PGI compilers
- Clock synchronization and event timestamp correction
- *SIONlib* mapping of task-local files into physical files
- Separate *CUBE3* GUI package for shared/distributed installations

Enhancements for Cray XT

- Support for multiple programming environments (PrgEnvs)
  - *PGI, GNU, CCE, Intel, Pathscale*
- Capture of mappings of processes to processors in machine

# POINT/VI-HPS Linux Live-DVD/ISO

Bootable Linux installation on DVD or USB memory stick

Includes everything needed to try out parallel tools
on an x86-architecture notebook computer

- VI-HPS tools: KCachegrind, Marmot, PAPI, Periscope, Scalasca, TAU, VT/Vampir*

- Also: Eclipse/PTP, TotalView*, etc.

  - *time/capability-limited evaluation licenses provided*

- GCC (w/OpenMP), OpenMPI

- User Guides, exercises & examples

http://www.vi-hps.org/training/material/

**Parallel Productivity Tools**
**Live DVD**
**Also includes:** TotalView, DyninstAPI, PDT, Eclipse PTP, Berkeley UPC, ptoolsrte, Chapel, and much more...

**Partners:**
ParaTools, Inc.
University of Florida
University of Oregon
Totalview Technologies
RWTH Aachen University
HLRS / University of Stuttgart
Jülich Supercomputing Centre
Technische Universität Dresden
Technische Universität München
University of Wisconsin at Madison
Pittsburgh Supercomputing Center
University of Tennessee at Knoxville
National Center for Supercomputing Applications

September 2010

**POINT VI-HPS**

http://nic.uoregon.edu/point    http://www.vi-hps.org