# Performance analysis of Sweep3D on Blue Gene/P with Scalasca

2010-04-23 | **Brian J. N. Wylie**, David Böhme, Bernd Mohr, Zoltán Szebenyi & Felix Wolf
Jülich Supercomputing Centre

b.wylie@fz-juelich.de

# Overview

Introduction

- Scalasca performance analysis toolset

- Blue Gene/P supercomputer

- Sweep3D compact application

Measurements & Analyses

- Default & improved Sweep3D configurations

- Refined instrumentation & analyses

Scalasca scalability

Conclusions

# Scalasca project

Overview

- Helmholtz Initiative & Networking Fund project started in 2006
    - *Headed by Prof. Felix Wolf (RWTH Aachen, GRS & JSC)*
- 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 released in March 2010
- Available for download from www.scalasca.org

# Supercomputers at Jülich Supercomputing Centre

| Period | Machine | #Cores |
|--------|---------|--------|
| 2003-8 | "jump" IBM p690 cluster | 1,312 |
| 2006-8 | "jubl" IBM Blue Gene/L | 16,384 |
| 2009- | "juropa" Bull/Sun/Intel blades | 26,304 |
| 2009- | "jugene" IBM Blue Gene/P | 294,912 |

*[jugene is the largest system in the Top500 list, Jun/Nov 2009]*

# JUGENE - IBM Blue Gene/P system



- 72 racks with 32 nodecards of 32 compute nodes (multiple networks)
  - *Quad-core 850 MHz PowerPC 450 processors*
  - *2 Gbytes main memory (aggregate 144 TB)*
- Proprietary Linux microkernel, MPI, XL compilers, GPFS filesystem

# Large-scale parallel applications (* analyzed by Scalasca)

Selected applications from Jülich Blue Gene Scaling Workshops:

- 294,912 (288k) cores

    - *GENE (MPI-RZG/D) gyrokinetic turbulence\**
    - *KKRnano (FZJ/D) electronic structure\**
    - *MUPHY/PEBS (Harvard/USA) MD/evolutionary biology*
    - *OSIRIS (IST/PT) 3D relativistic plasma\**
    - *PHASTA (RPI/USA) unsteady fluid dynamics\**
    - *UNIC (ANL/USA) neutron transport in fast reactors*
    - *(various) lattice QCD*

- 262,144 (256k) cores

    - *MP2C (FZJ/D) mesoscale particle dynamics\**
    - *NEK5000 (ANL/USA) CFD\**

# Sweep3D

Ubiquitous ASCI benchmark code from Los Alamos National Laboratory

- 3-dimensional neutron transport simulation
- direct order solve uses diagonal wavefront sweeps over grid cells combined with pipelining of blocks of $k$-planes and octants
- execution performance extensively modeled & analyzed

MPI parallel version using 2D domain decomposition

- ~2,000 lines of code (12 source modules), mostly Fortran77
- very portable, and highly scalable
- tunable via input deck, e.g., number of $k$-planes in blocks (MK)
- benchmark configuration does 12 iterations
    - *flux correction 'fixups' applied after 7th iteration*

# Scalasca features

Open source, New BSD license

Portable

- IBM BlueGene/L, BlueGene/P, SP & blade clusters, Cray XT4/5, NEC SX, SGI Altix, SiCortex, Linux cluster® (SPARC, x86-64), ...

Supports typical HPC languages & parallel programming paradigms
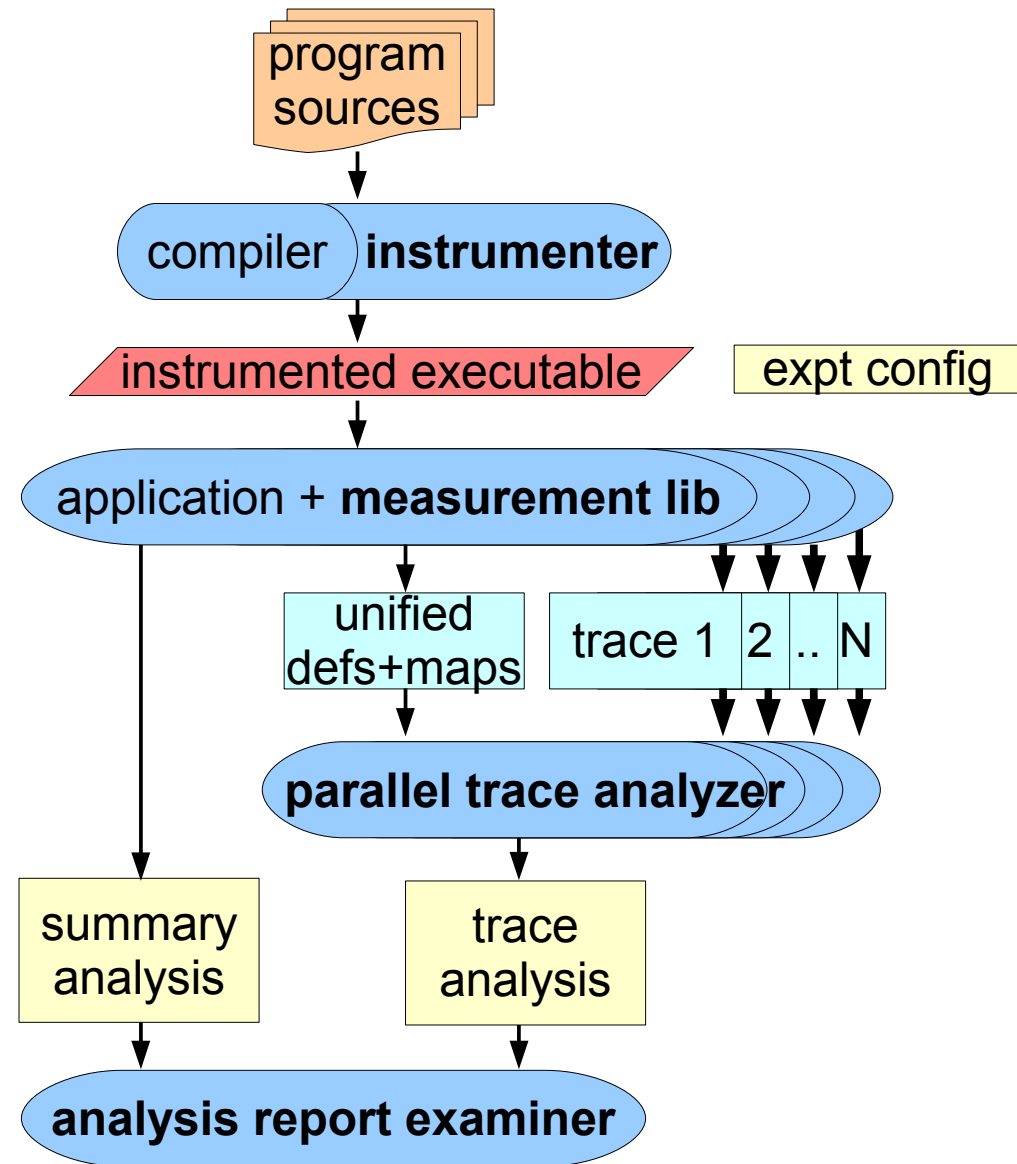
- Fortran, C, C++
- MPI, OpenMP & hybrid MPI/OpenMP

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

program sources → compiler / **instrumenter** → instrumented executable / expt config → application + **measurement lib** → unified defs+maps / trace 1 2 .. N → **parallel trace analyzer** → summary analysis / trace analysis → **analysis report examiner**

# Scalasca usage

1. Prepare application objects and executable for measurement:

   - **scalasca -instrument**  `mpixlf77 -O3 -qarch=450 -qtune=450 …`

2. Run application under control of measurement nexus:

   - **scalasca -analyze**  `mpirun -mode VN -np 294912 sweep3d`

     - *epik_sweep3d_vn294912_sum*  experiment produced
   - **scalasca -analyze -t**  `mpirun -mode VN -np 294912 sweep3d`

     - *epik_sweep3d_vn294912_trace*  experiment produced

3. Interactively explore measurement analysis report

   - **scalasca -examine**  `epik_sweep3d_vn294912_trace`

     - *epik_sweep3d_vn294912_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 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

4. Refine instrumentation by manually annotating key code sections

   - use EPIK instrumentation API macros


Compare original MK=10 and alternative MK=1 Sweep3D configurations

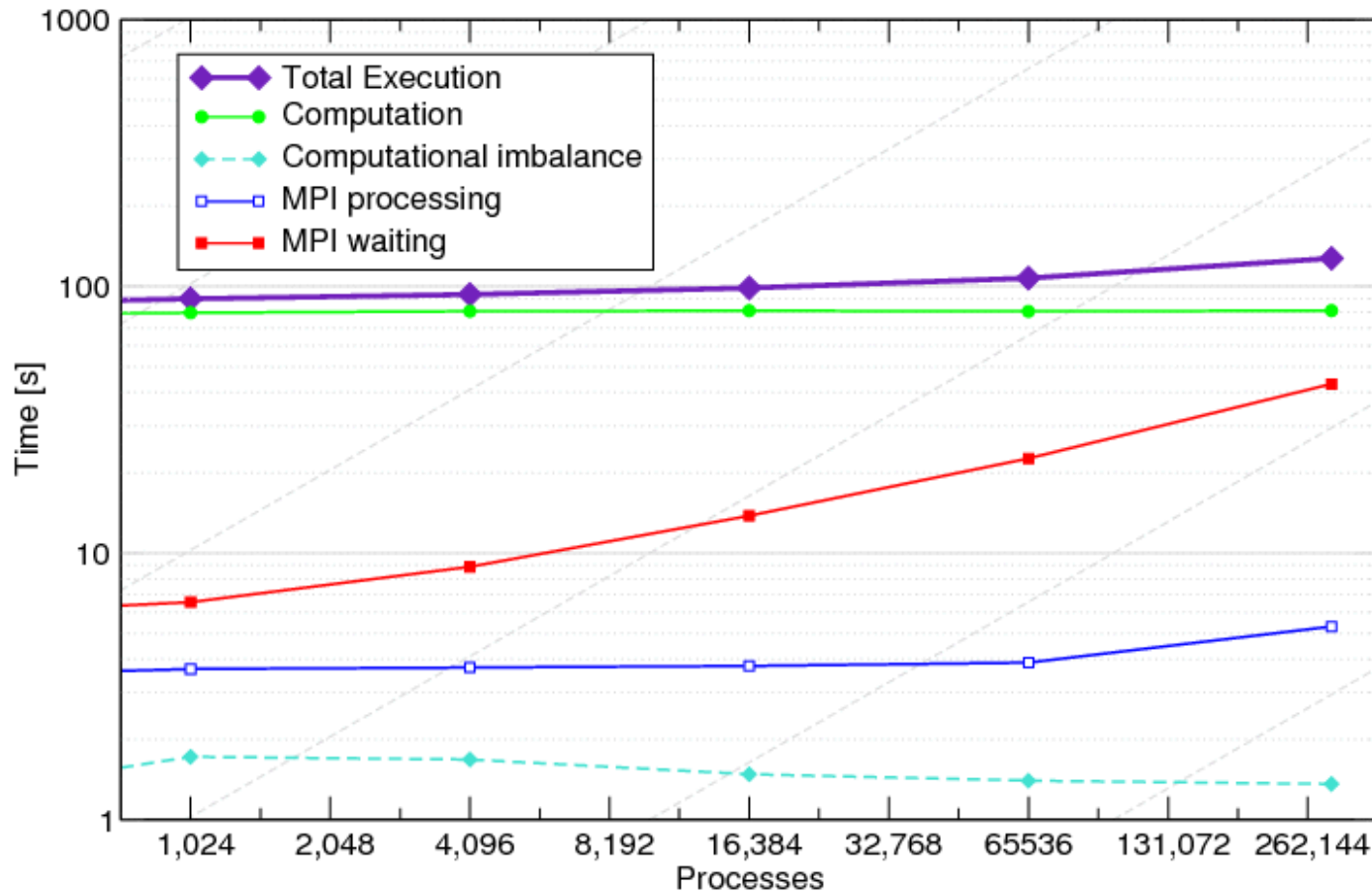   - varies computation wavefront pipelining concurrency

# Sweep3D execution time scaling: original MK=10



- Default input deck settings

- *Weak scaling* due to fixed problem size per process

- Reasonable scalability(?)

- Constant time for *computation*

- Rapid growth in *MPI* time, which is almost all *waiting* time

[Replication of scaling characteristics reported by others]

# Sweep3D execution time scaling: improved MK=1



- Much improved performance and scalability

- *Computation* 20% faster (better caching?)

- *MPI processing* time increased (10x messages)

- *MPI waiting* time significantly reduced, though still growing markedly
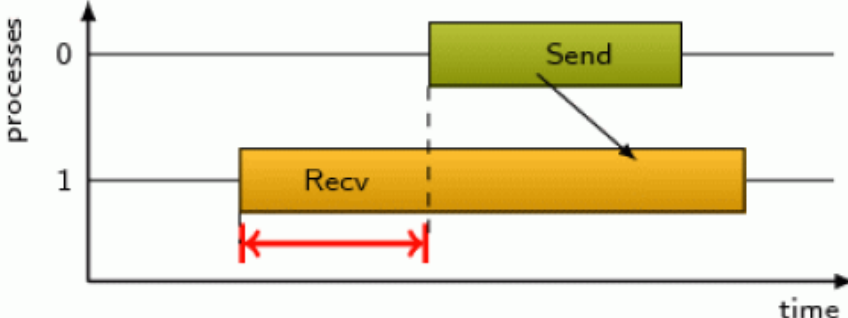
[Single *k*-planes rather than default batches of 10]

# Scalasca GUI description for *Late Sender* metric



- Analysis report explorer GUI provides hyperlinked online descriptions of metrics

- Diagnosis hints suggest how to refine diagnosis of performance problems and possible remediation
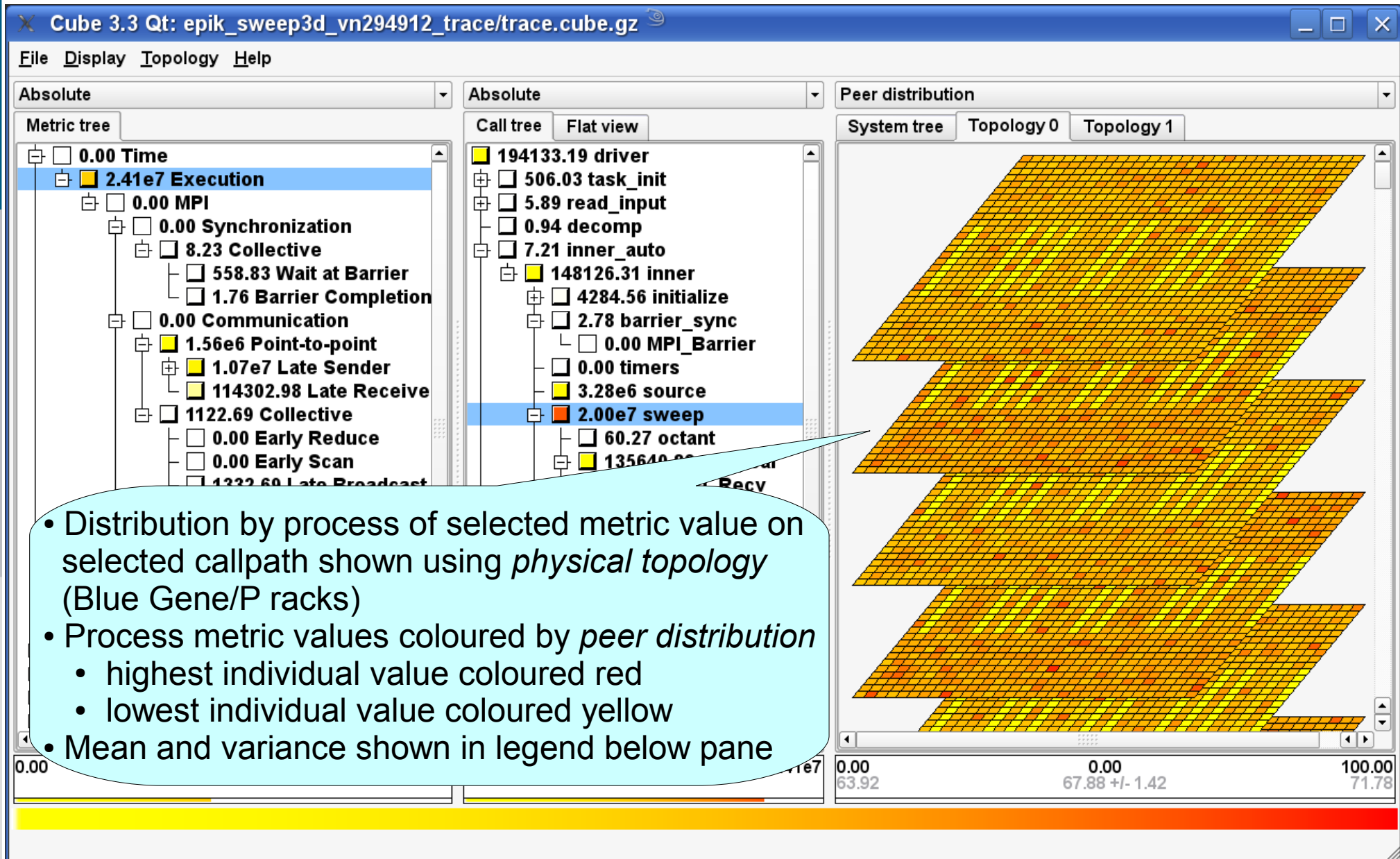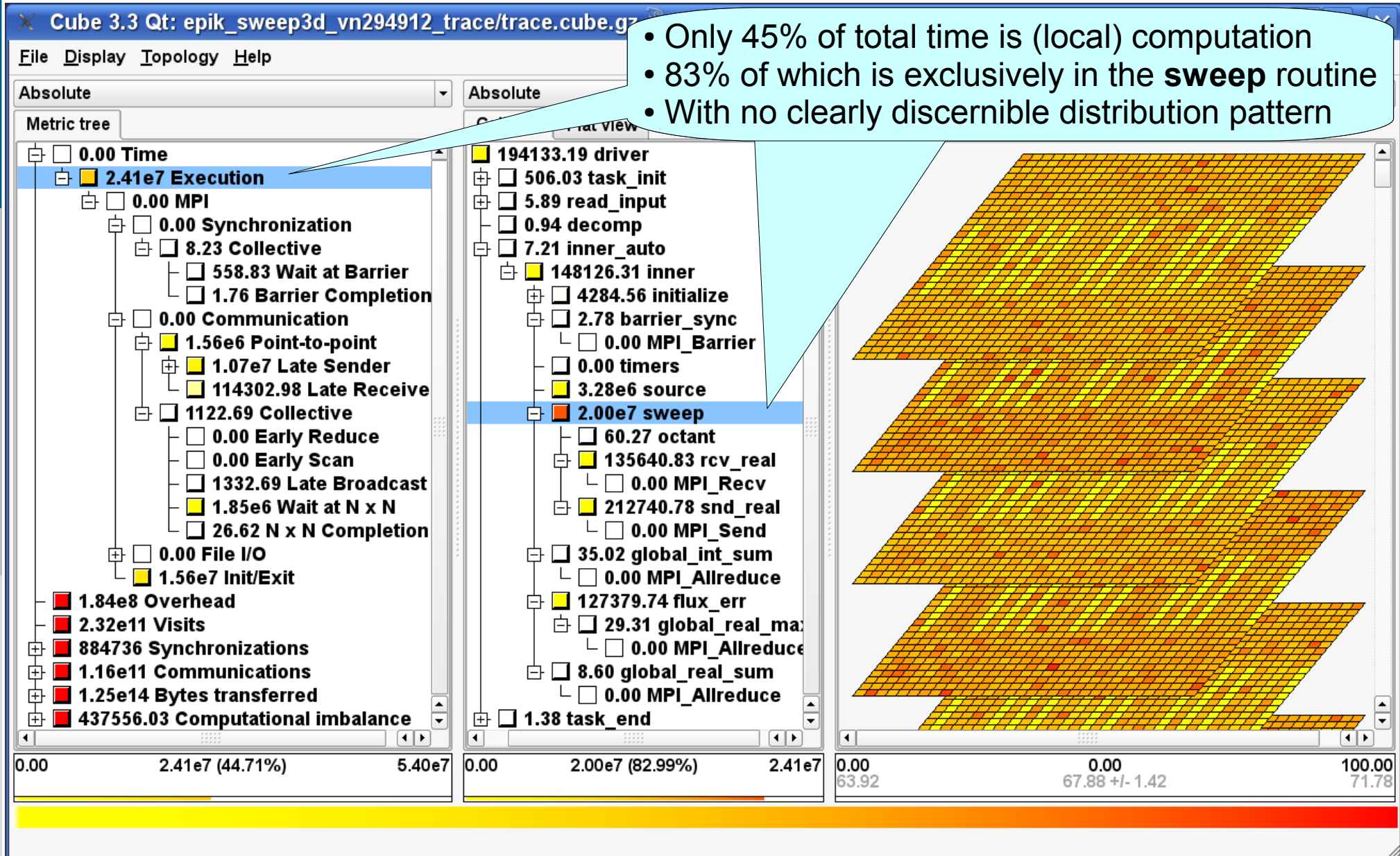
# Scalasca analysis report: performance metric tree



**Cube 3.3 Qt: epik_sweep3d_vn294912_trace/trace.cube.gz**

File  Display  Topology  Help

Absolute | Absolute | Peer distribution

**Metric tree** | Call tree | Flat view | System tree | Topology 0 | Topology 1

- □ 0.00 Time
  - ■ **2.41e7 Execution**
    - □ 0.00 MPI
      - □ 0.00 Synchronization
        - □ 8.23 Collective
          - □ 558.83 Wait at Barrier
          - □ 1.76 Barrier Completion
      - □ 0.00 Communication
        - ■ 1.56e6 Point-to-point
          - ■ 1.07e7 Late Sender
          - ■ 114302.98 Late Receive
        - □ 1122.69 Collective
          - □ 0.00 Early Reduce
          - □ 0.00 Early Scan
          - □ 1332.69 Late Broadcast
          - ■ 1.85e6 Wait at N x N
          - □ 26.62 N x N Completion
      - □ 0.00 File I/O
      - ■ 1.56e7 Init/Exit
  - ■ 1.84e8 Overhead
  - ■ 2.32e11 Visits
  - ■ 884736 Synchronizations
  - ■ 1.16e11 Communications
  - ■ 1.25e14 Bytes transferred
  - ■ 437556.03 Computational imbalance

Call tree:
- ■ 194133.19 driver
  - □ 506.03 task_init
  - □ ... input
  - ...
  - □ 0.00 MPI_Send
  - □ 0.00 MPI_Allreduce
  - □ 8.60 global_real_sum
    - □ 0.00 MPI_Allreduce
  - □ 1.38 task_end

0.00    2.41e7 (44.71%)    5.40e7
0.00    2.00e7 (82.99%)    2.41e7
0.00    0.00    100.00
63.92   67.88 +/- 1.42    71.78

- Metrics arranged hierarchically in trees
  - *inclusive* values shown for closed nodes
  - *exclusive* values shown for open nodes
- Metrics values aggregated for all processes & callpaths
- Colour-coded according to scale at window base
- Selected metric value projected to panes on right
- Most metrics calculated by local *runtime summarization*

- Augmented with metrics from *automatic trace analysis* calculated using non-local event information

# Scalasca analysis report: program call tree



Graphical presentation of callpath profile
- metric selected from pane on left
- shown projected on program call tree

Tree selectively expanded following high metric values
- *inclusive* value when node closed
- *exclusive* value when node opened

Colour-coded according to scale at base

# Scalasca analysis report: system tree/topology



- Distribution by process of selected metric value on selected callpath shown using *physical topology* (Blue Gene/P racks)
- Process metric values coloured by *peer distribution*
  - highest individual value coloured red
  - lowest individual value coloured yellow
- Mean and variance shown in legend below pane

# Scalasca analysis report: sweep computation time

# Scalasca analysis report: *Late Sender* time



- 55% of total time is *MPI* messaging overhead
- 23% of total time is *Point-to-point communication*
- 20% of total time is due to *Late Sender* situations

# Sweep3D grid topology view

- Application's 2D topology of 576x512 procs

- Reveals clear pattern in the distribution of *Late Sender* metric values

- Arises from superimposing octant sweeps with imbalanced computation

# Sweep3D computation time by iteration (16,384 processes)

# Iteration execution time breakdown (16,384 processes)



- Initial 7 (non-fixup) iterations faster than later 5 (fixup) iterations
- *MPI waiting* time dominates original; amplified by fixup imbalance
- *MPI processing* time a fixed cost in each iteration

# Sweep3D computation time distribution evolution

| its=7 | its=8 | its=9 | its=10 |
|-------|-------|-------|--------|



[Values coloured on scale 5 to 10 seconds]

- Initial 7 (non-fixup) iterations are very well balanced

- Fixups introduced after iteration 7 are computationally imbalanced

- Iteration 8 clearly shows central rectangular imbalance feature

- Subsequent iterations have less imbalance in the rectangular core, however, pronounced symmetric oblique lines of excess computation

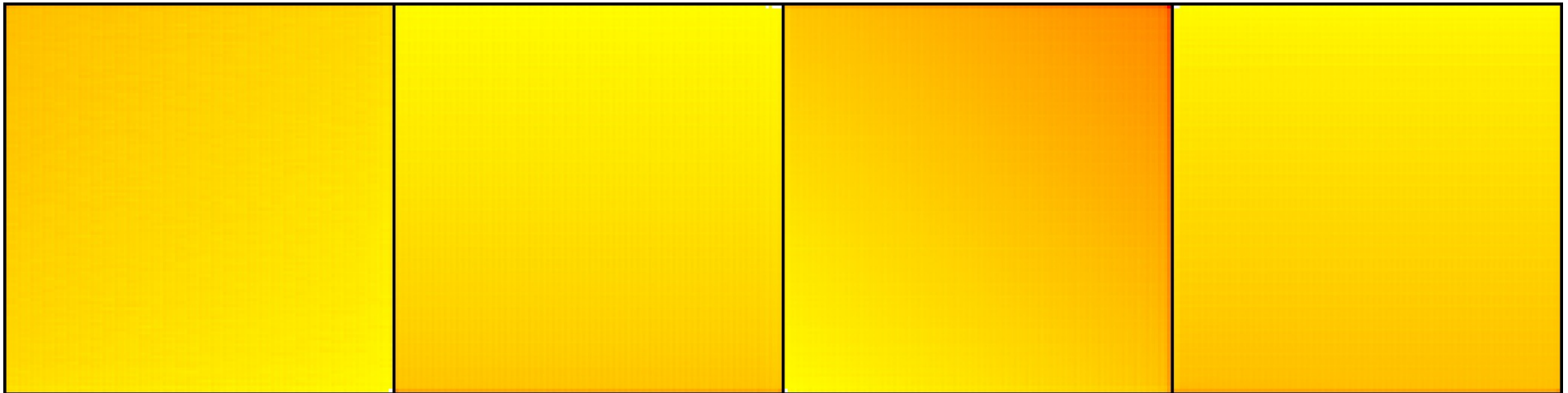# Scalasca source code browser



- Source file **sweep.f** contains 625-line **sweep** flow routine

- Computationally imbalanced *i*-line section where corrective 'fixups' for negative fluxes are recursively applied in *i*, *j*, *k* grid directions

- Hotspot also typically identified by profilers based on sampling

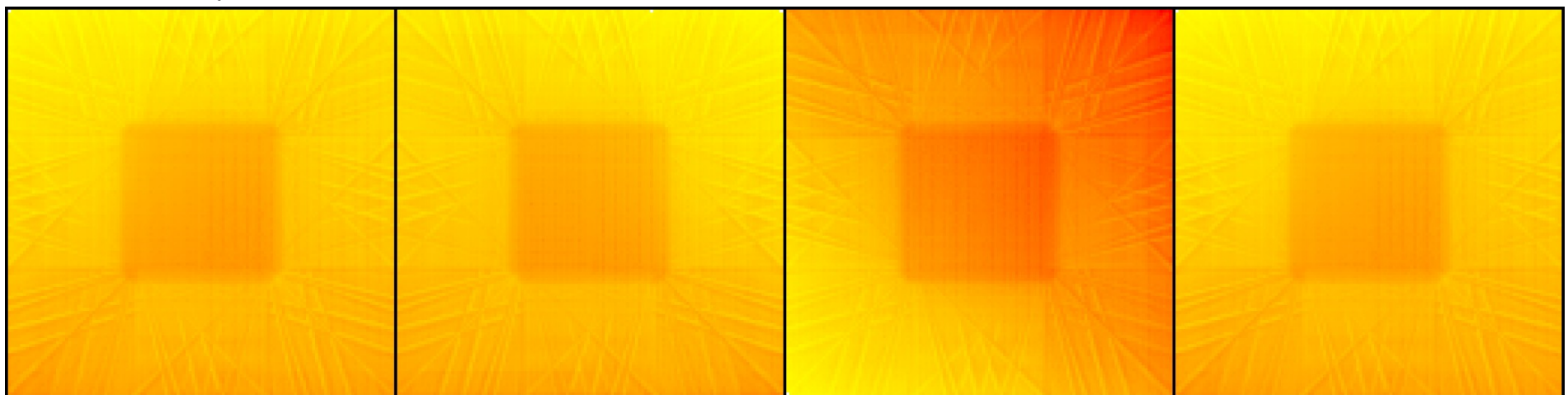# MPI *Late Sender* time variation by sweep octant

# MPI *Late Sender* time distributions for octant pairs

Iterations without fixups:



Iterations with fixups:
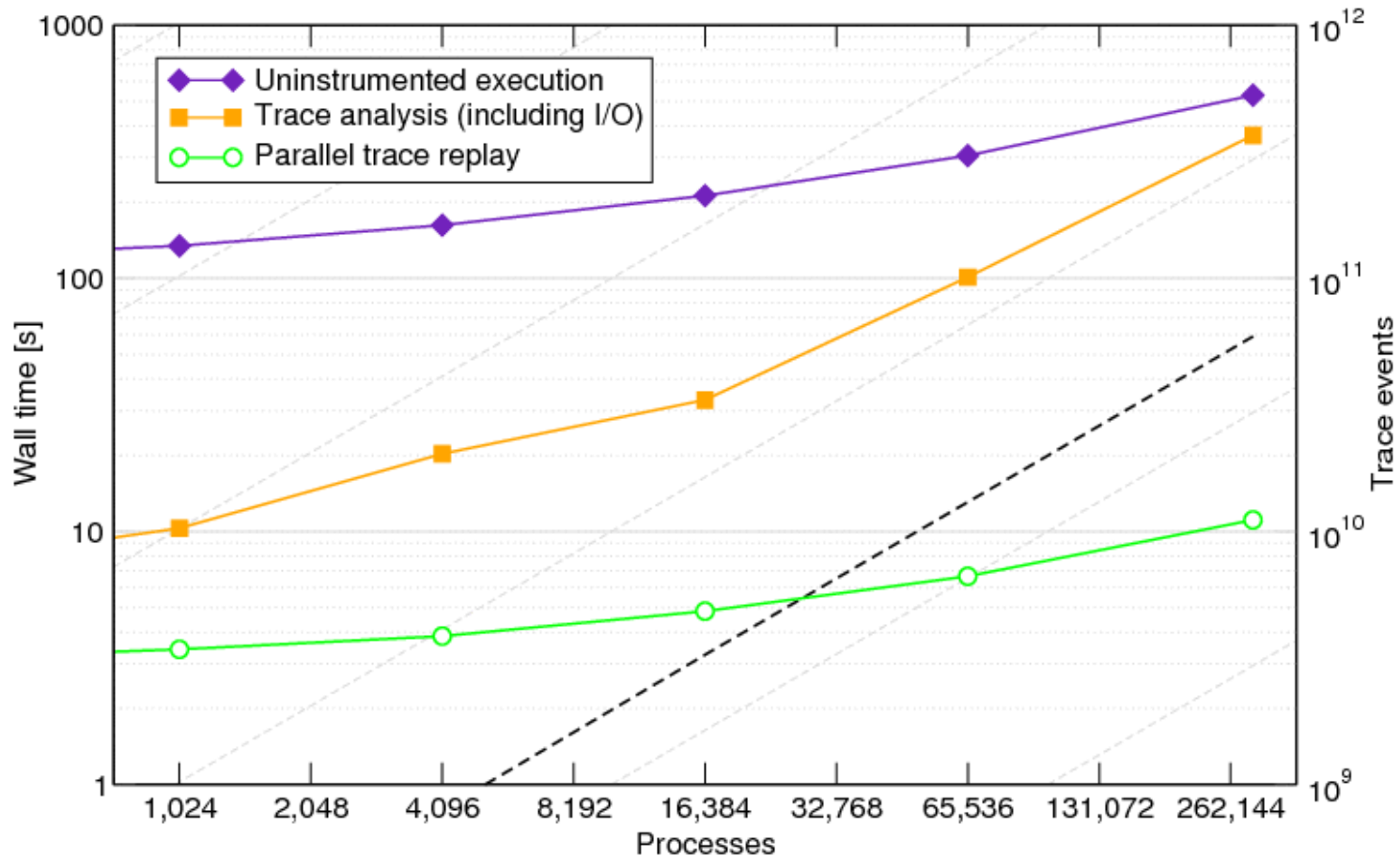


| octant=1+2 | octant=3+4 | octant=5+6 | octant=7+8 |

# Scalasca experiment statistics: sweep3d_vn294912_trace

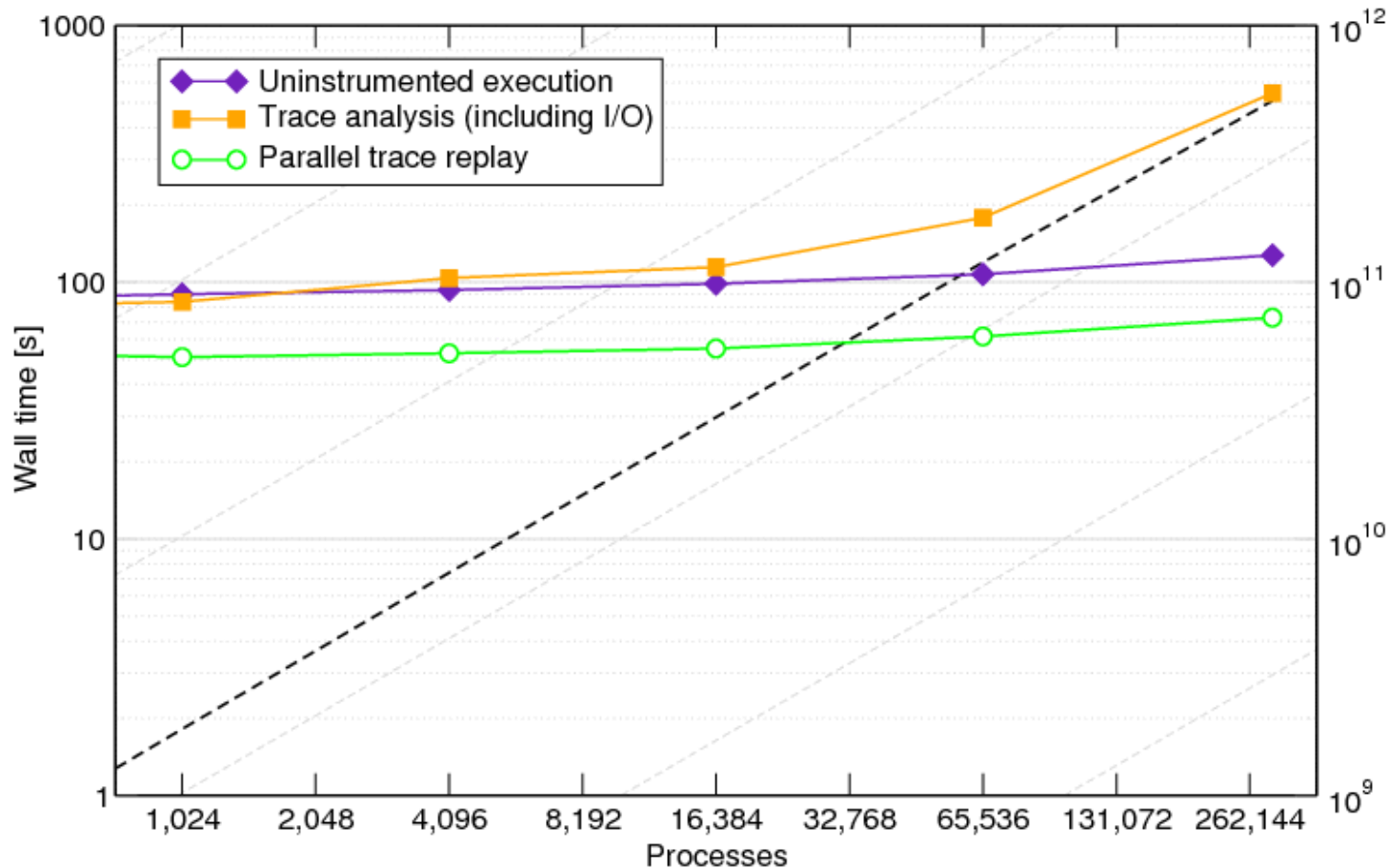| Scalasca | v1.2 | v1.3+ |
|---|---|---|
| Sweep3d size of *k*-blocks MK | 10 | 1 |
| Sweep3d elapsed time [s] | 505 | 129 |
| Measurement dilation [%] | 5 | 3 |
| | | |
| Unification time [mins] | 43 | 41 |
| | | |
| Trace event records [G] | 59 | 510 |
| Trace buffer content [MB] | 2.75 | 27 |
| Trace total size [TB] | 0.79 | 7.6 |
| Trace (physical) files [#] | 294912 | 576 |
| Trace open/create time [mins] | 86 | 10 |
| Trace flush rate [GB/s] | 3.7 | 19.1 |
| | | |
| Trace analysis time [s] | 368 | 546 |
| - Trace analysis collation [s] | 162 | 91 |
| - Trace analysis replay [s] | 11 | 74 |

# Scalasca trace analysis time scaling: original MK=10



- 2.75MB of trace event records per process (determined from summary)

- Total trace size (---) increases linearly to 790GB for 59G events

- Trace replay time scales with application execution time

[Scalasca v1.2]

# Scalasca trace analysis time scaling: improved MK=1



- 10x larger traces due to 10x messages

- 27MB of trace event records per process

- Total trace size (---) increased to 7.6TB for 510G events

- 10x longer trace replay time still scales with application execution time

[Scalasca v1.3 using SIONlib & prototype binary-format report]

# Scalasca scalability issues/optimizations

Time for unification of identifiers grows linearly with number of processes

- use *hierarchical* unification scheme

Creating individual traces files for each process is prohibitive

- use SIONlib *multi-files* (e.g., one per Blue Gene/P IONode)

Analysis reports are large and slow to collate

- *binary format* for metric value data plus *separate* XML metadata

Analysis presentation is slow and requires lots of memory

- store metrics as *inclusive* values and load them *incrementally*

Full analysis presentation requires very large, high-resolution screens

- :-)

# Conclusions

Blue Gene/P efficiently supports application executions at extreme scale

Sweep3D is scalable to 294,912 MPI processes on Blue Gene/P

- appropriate input configuration/tuning is necessary
- 'fixups' are computationally imbalanced and amplify waiting times
  - *they therefore need to be used sparingly*
- there might be further opportunities for performance improvement
  - *e.g., with non-default mapping of processes to cores*

Scalasca has demonstrated performance measurement/analysis at unprecedented scale

- however, further significant improvements are in progress

# **Sc**alable performance **a**nalysis
## of *extremely* **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*