

VI-HPS



scalasca

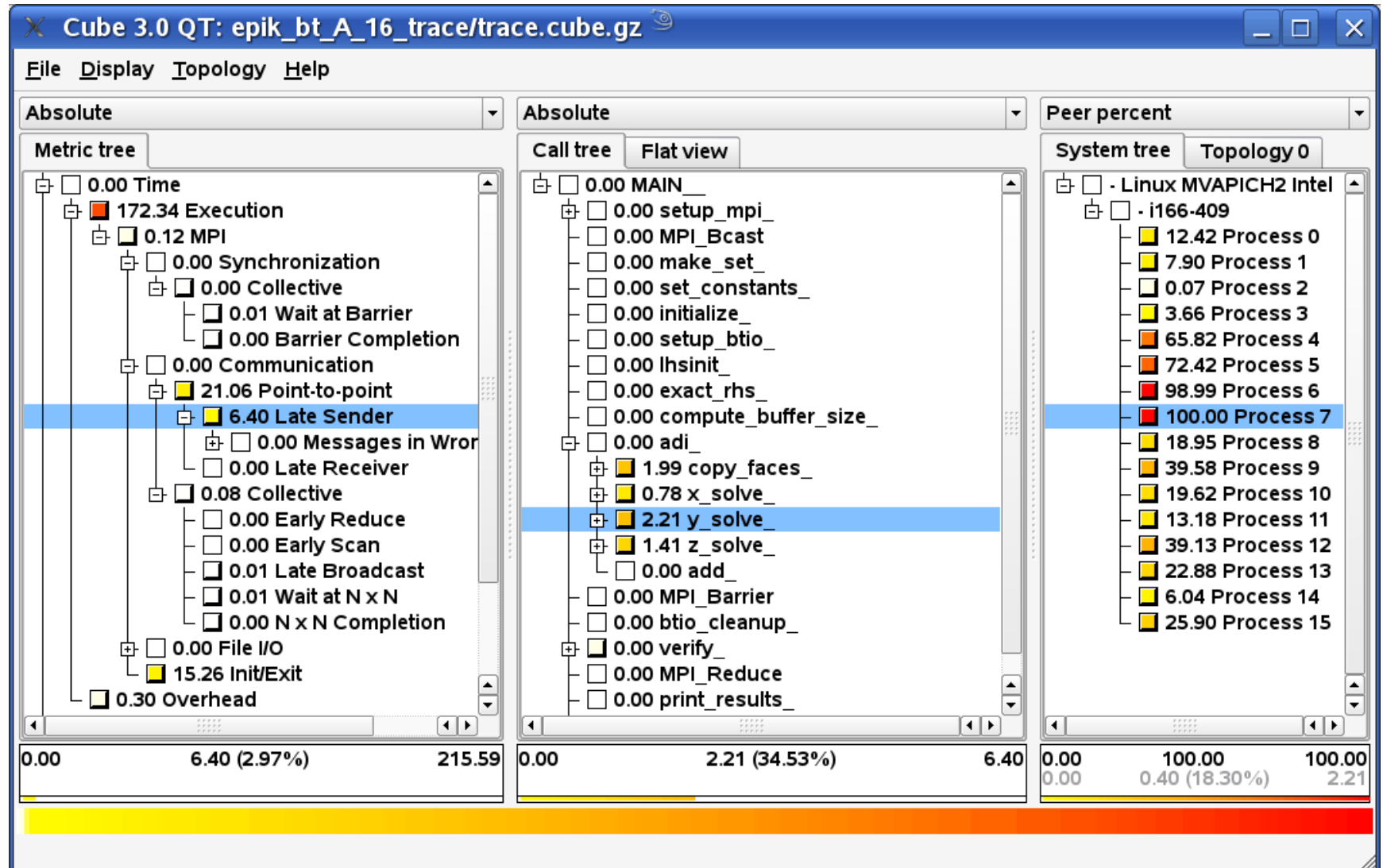
Performance analysis & tuning case studies

Brian Wylie & Markus Geimer
Jülich Supercomputing Centre
scalasca@fz-juelich.de
April 2012

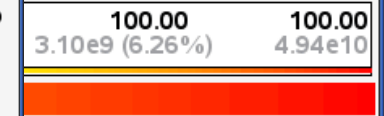
- Example experiment archives provided for examination:
 - jugene_sweep3d
 - ▶ 294,912 & 65,536 MPI processes on BG/P (trace)
 - jump_zeusmp2
 - ▶ 512 MPI processes on p690 cluster (summary & trace)
 - marenosturm_wrf-nmm
 - ▶ 1600 MPI processes on JS21 blade cluster, solver extract
 - ▶ summary analysis with 8 PowerPC hardware counters
 - ▶ trace analysis showing NxN completion problem on some blades
 - neptun_jacobi
 - ▶ 12 MPI processes, or 12 OpenMP threads, or 4x3 hybrid parallelizations implemented in C, C++ & Fortran on SGI Altix
 - ranger_smg2000
 - ▶ 12,288 MPI processes on Sun Constellation cluster, solve extract

- Comparison of NPB-BT class A in various configurations run on a single dedicated 16-core cluster compute node
 - 16 MPI processes
 - ▶ optionally built using MPI File I/O (e.g., SUBTYPE=full)
 - ▶ optionally including PAPI counter metrics in measurement (e.g., EPK_METRICS=PAPI_FP_OPS:DISPATCH_STALLS)
 - 16 OpenMP threads
 - 4 MPI processes each with 4 OpenMP threads (MZ-MPI)
- NPB-BT-MZ class B on Cray XT5 (8-core compute nodes)
 - 32 MPI processes with OMP_NUM_THREADS=8
 - ▶ More threads created on some processes (and fewer on others) as application attempts to balance work distribution
- NPB-MPI-BT on BlueGene/P with 144k processes
 - 1536x1536x1536 gridpoints distributed on 384x384 processes

VI-HPS



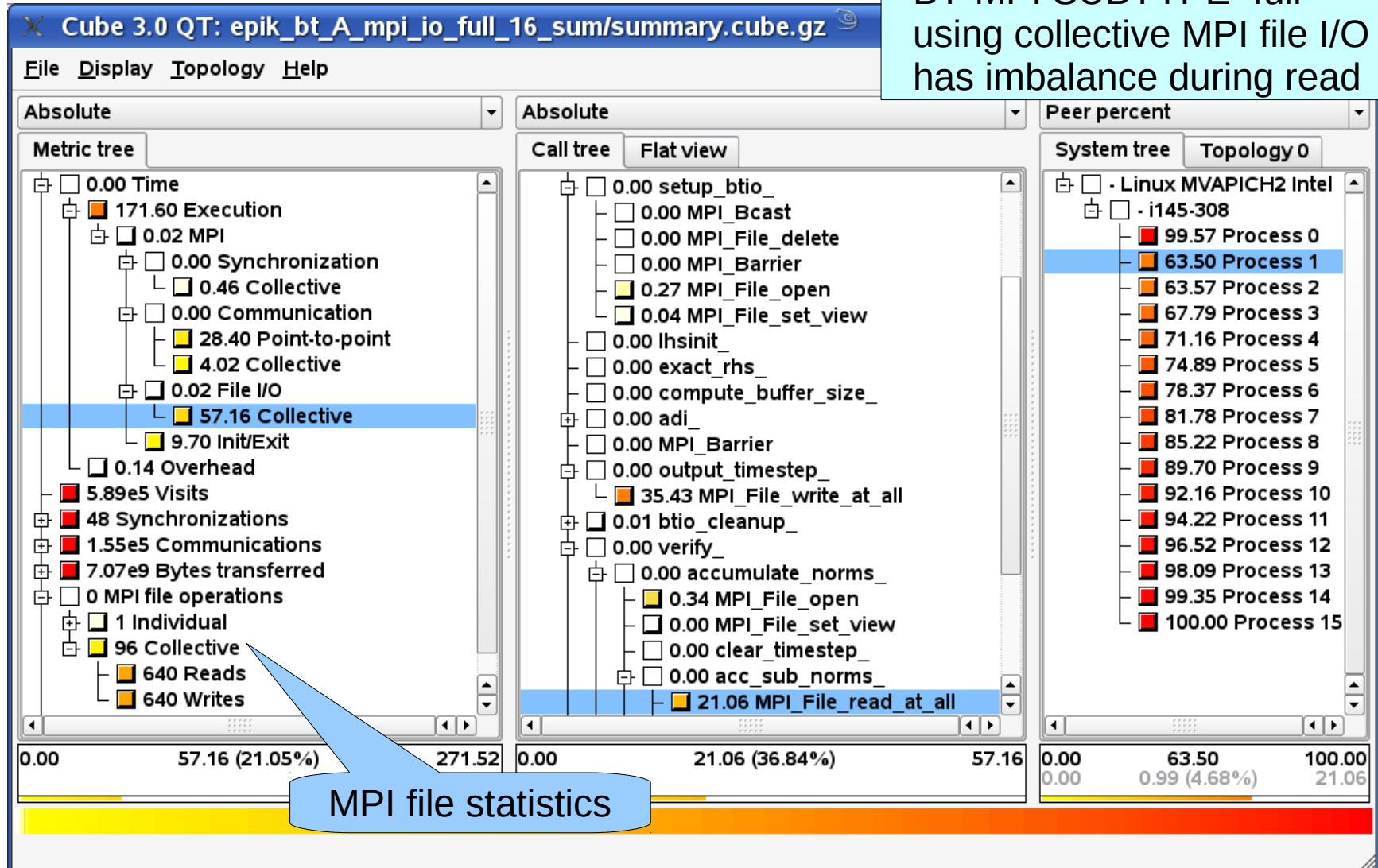
VI-HPS



16-process summary analysis: MPI File I/O time



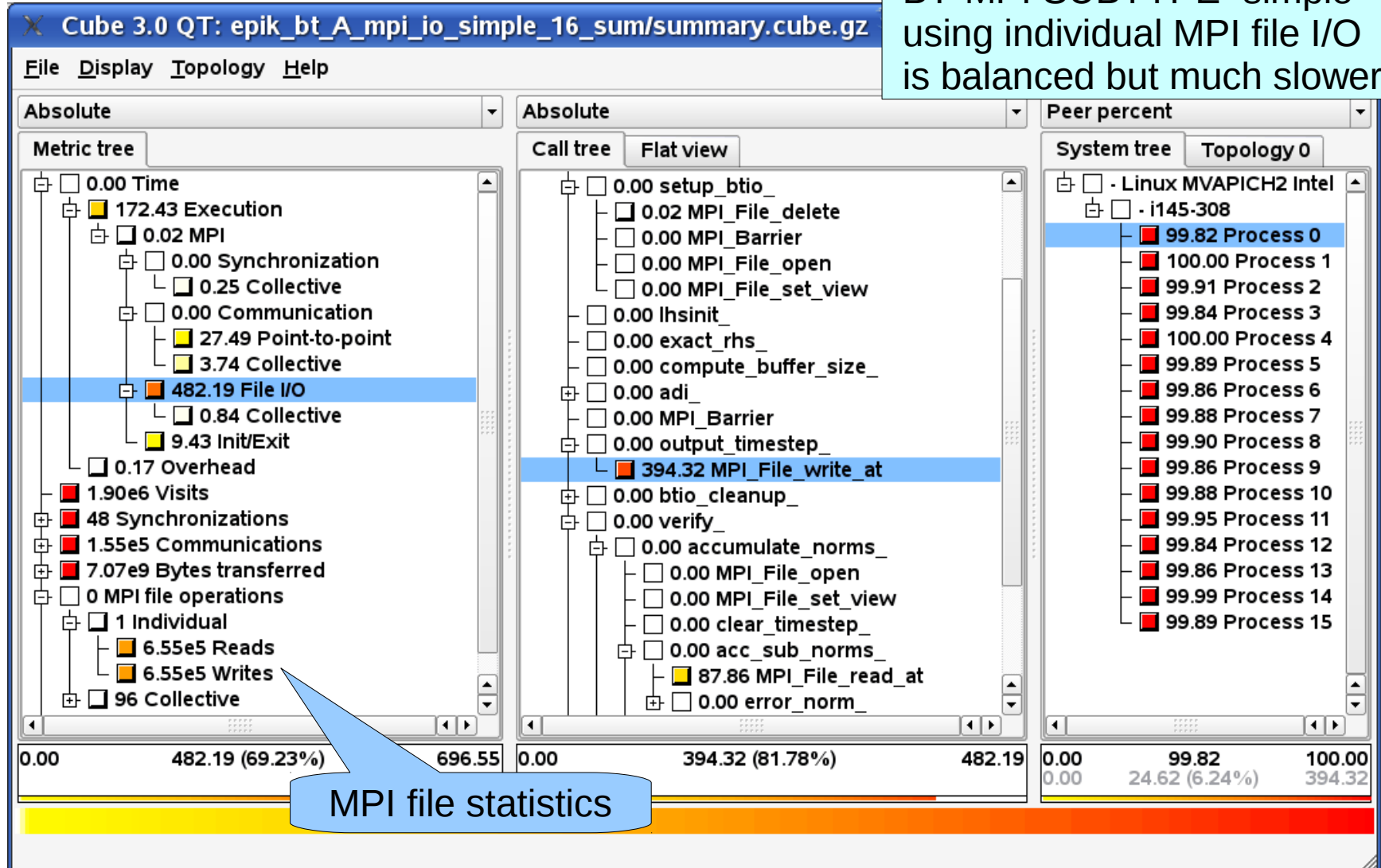
BT-MPI SUBTYPE=full
using collective MPI file I/O
has imbalance during read



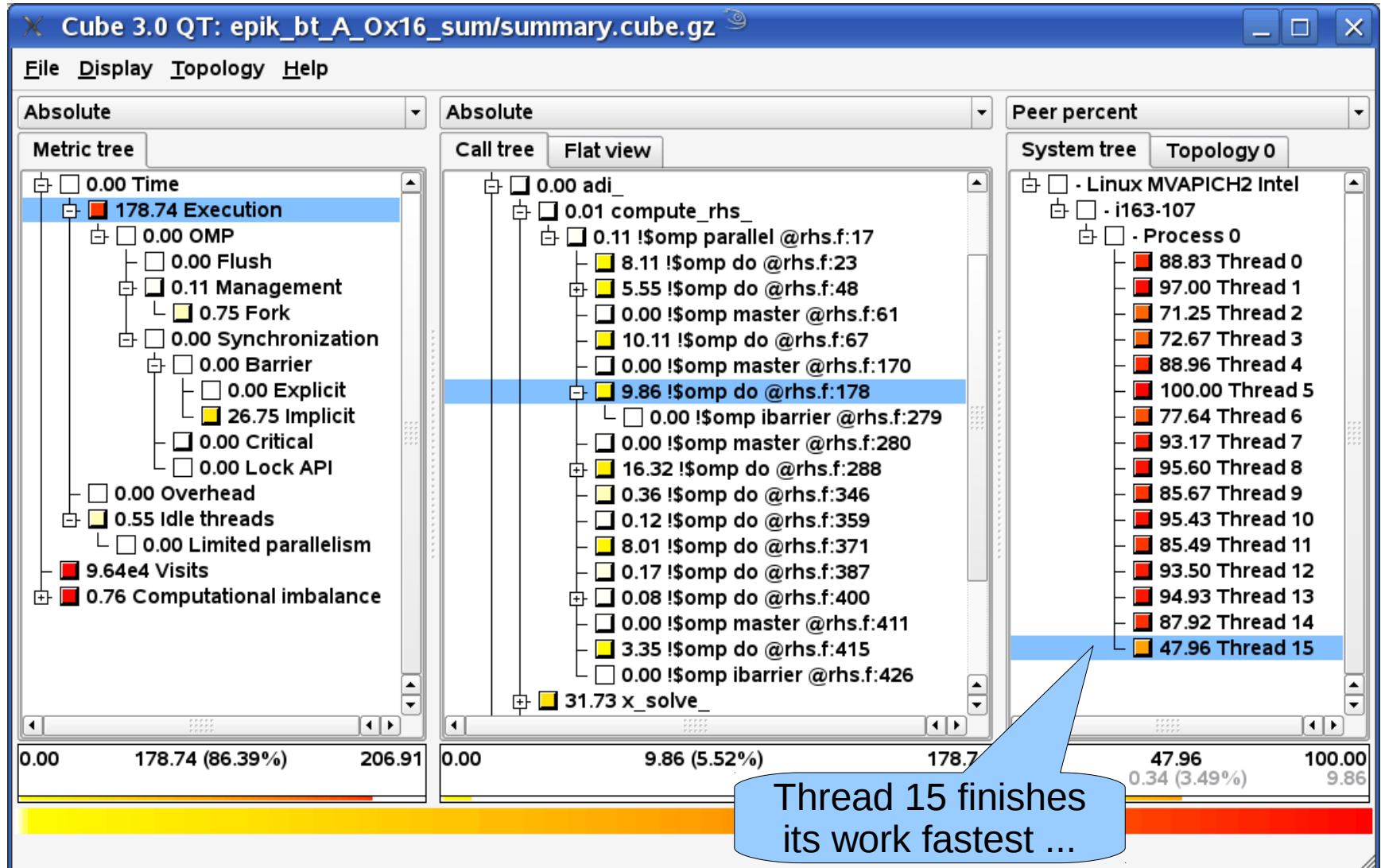
16-process summary analysis: MPI File I/O time



BT-MPI SUBTYPE=simple
using individual MPI file I/O
is balanced but much slower

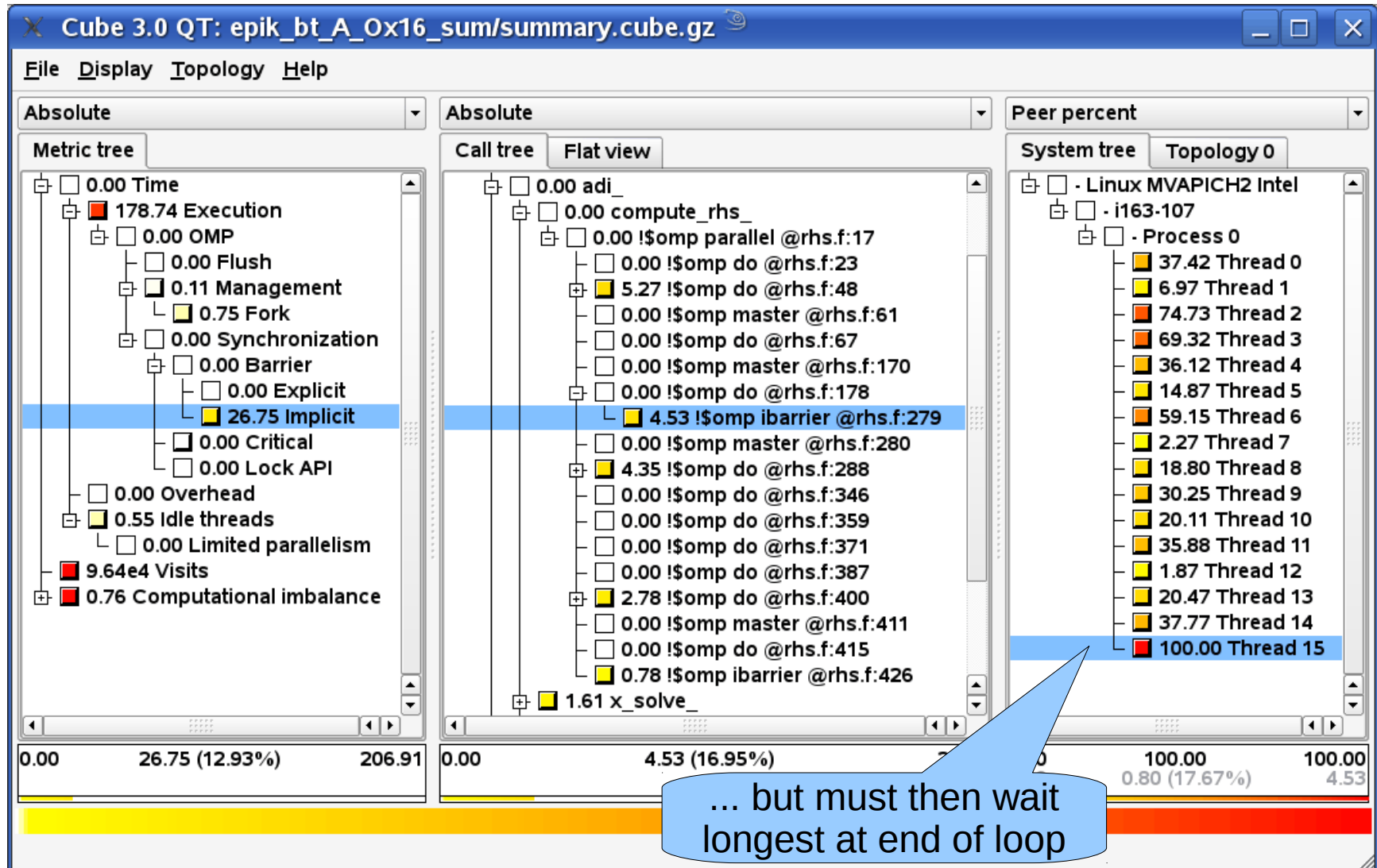


16-thread summary analysis: Execution time

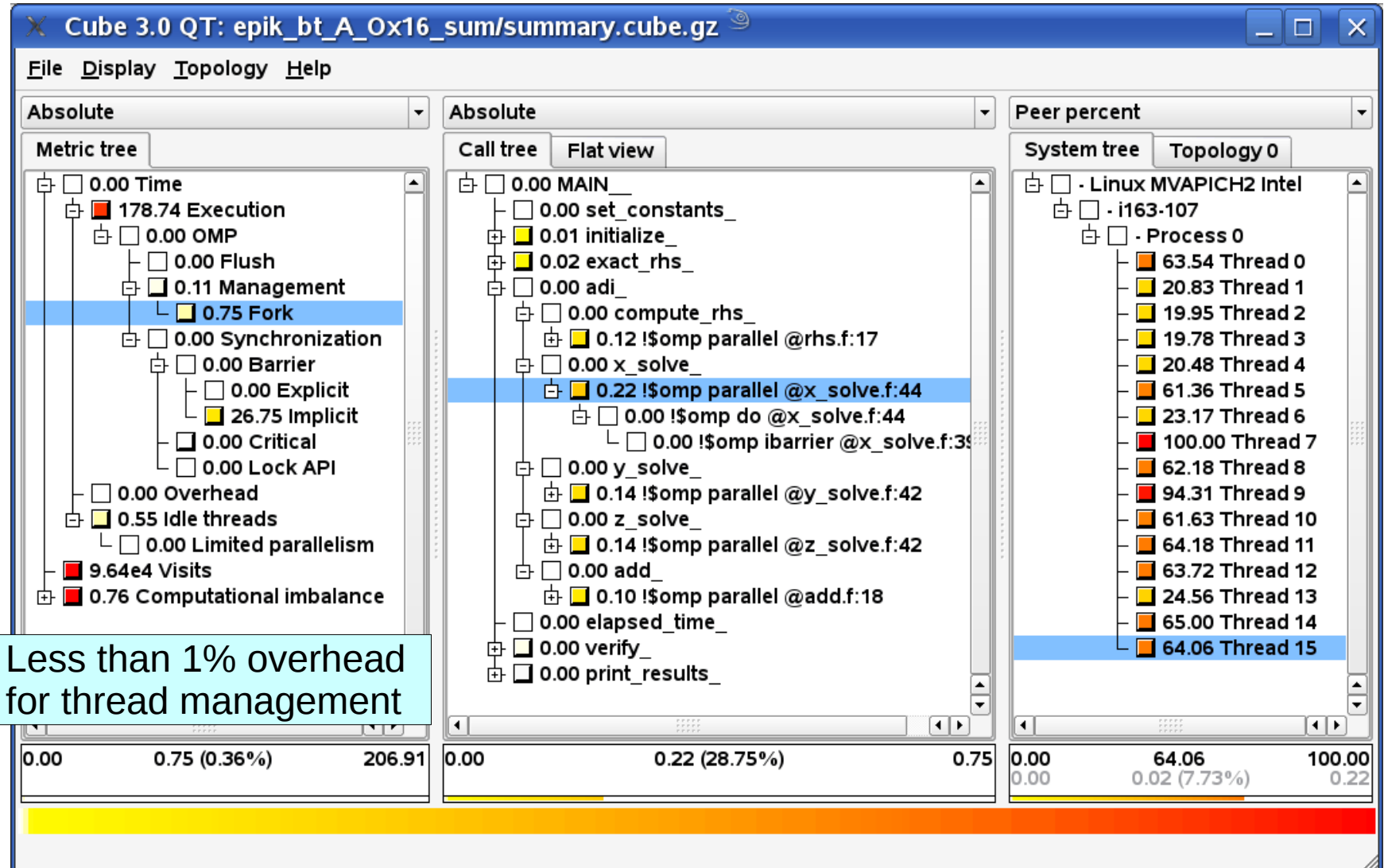


16-thread summary analysis: Implicit barrier time

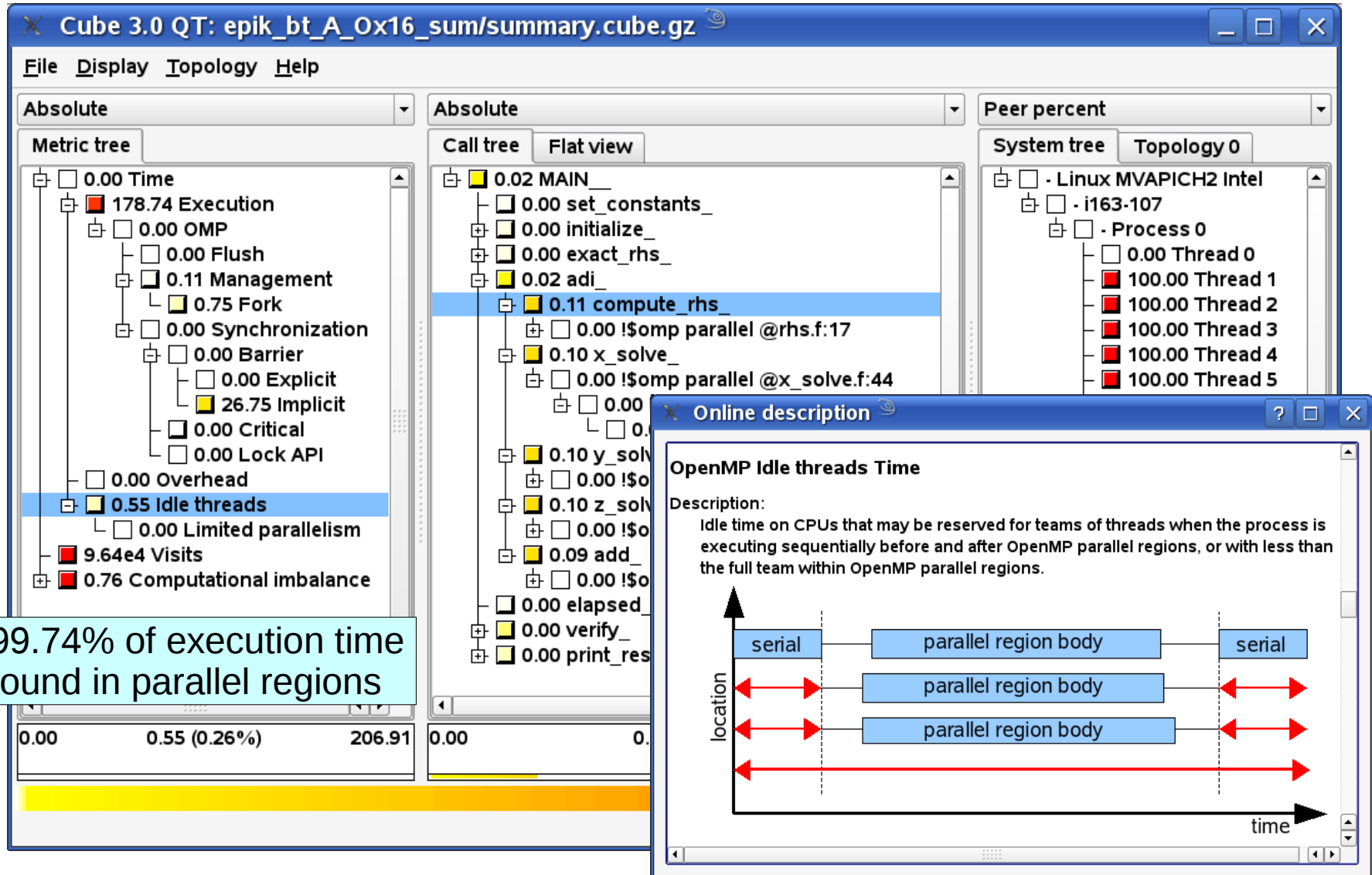
VI-HPS



16-thread summary analysis: Thread fork time

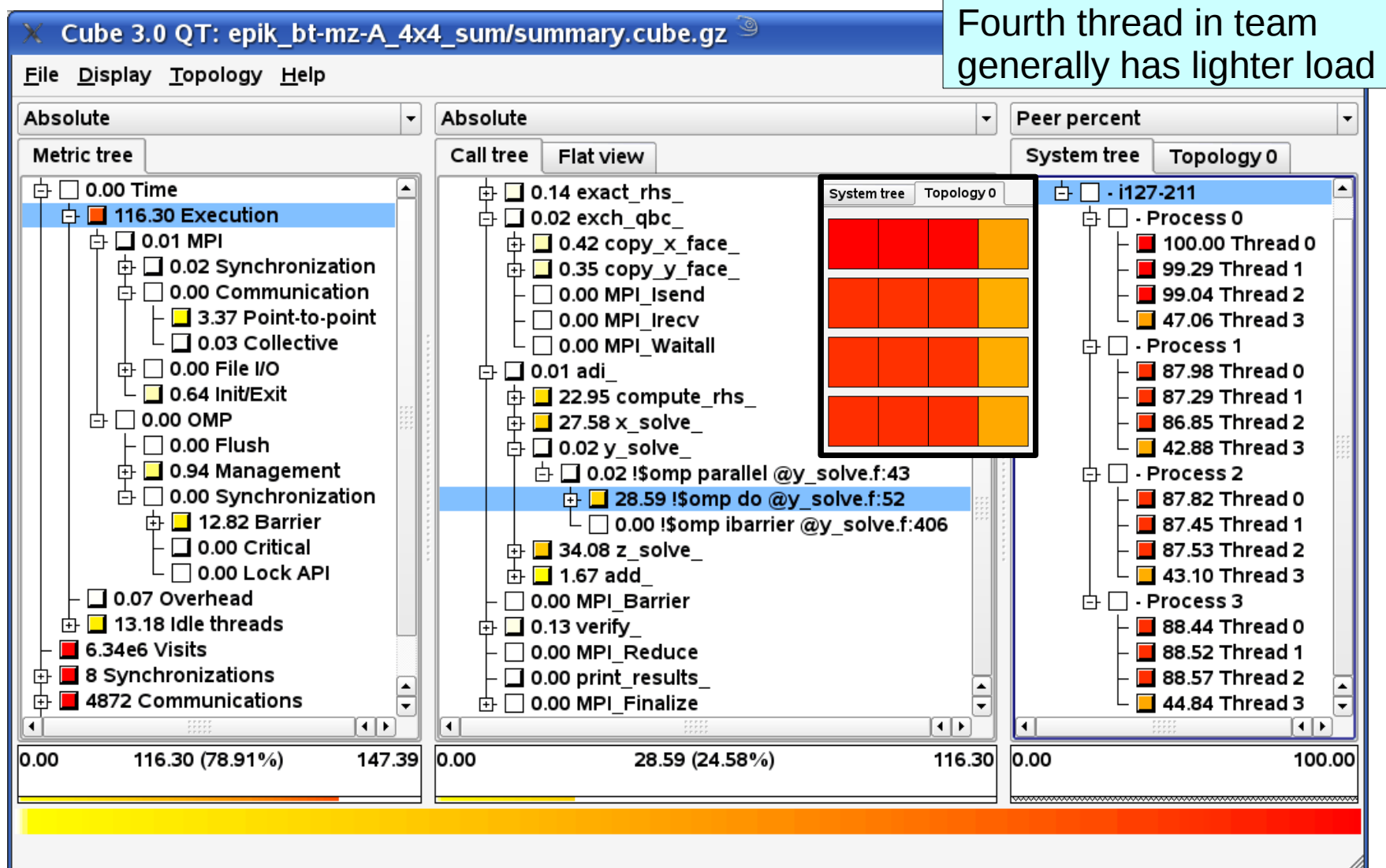


16-thread summary analysis: Idle threads time



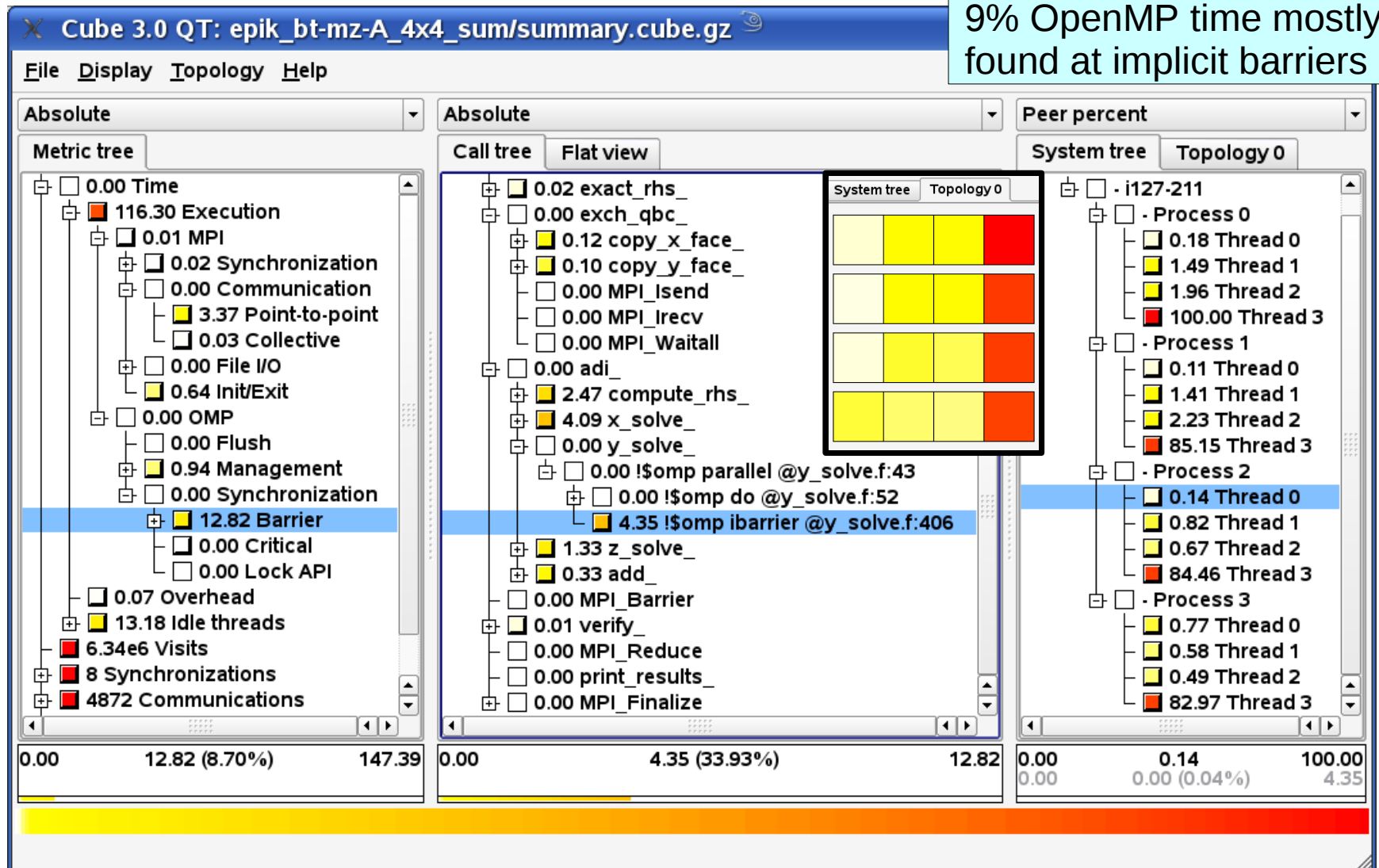
99.74% of execution time found in parallel regions

4x4 summary analysis: Execution time

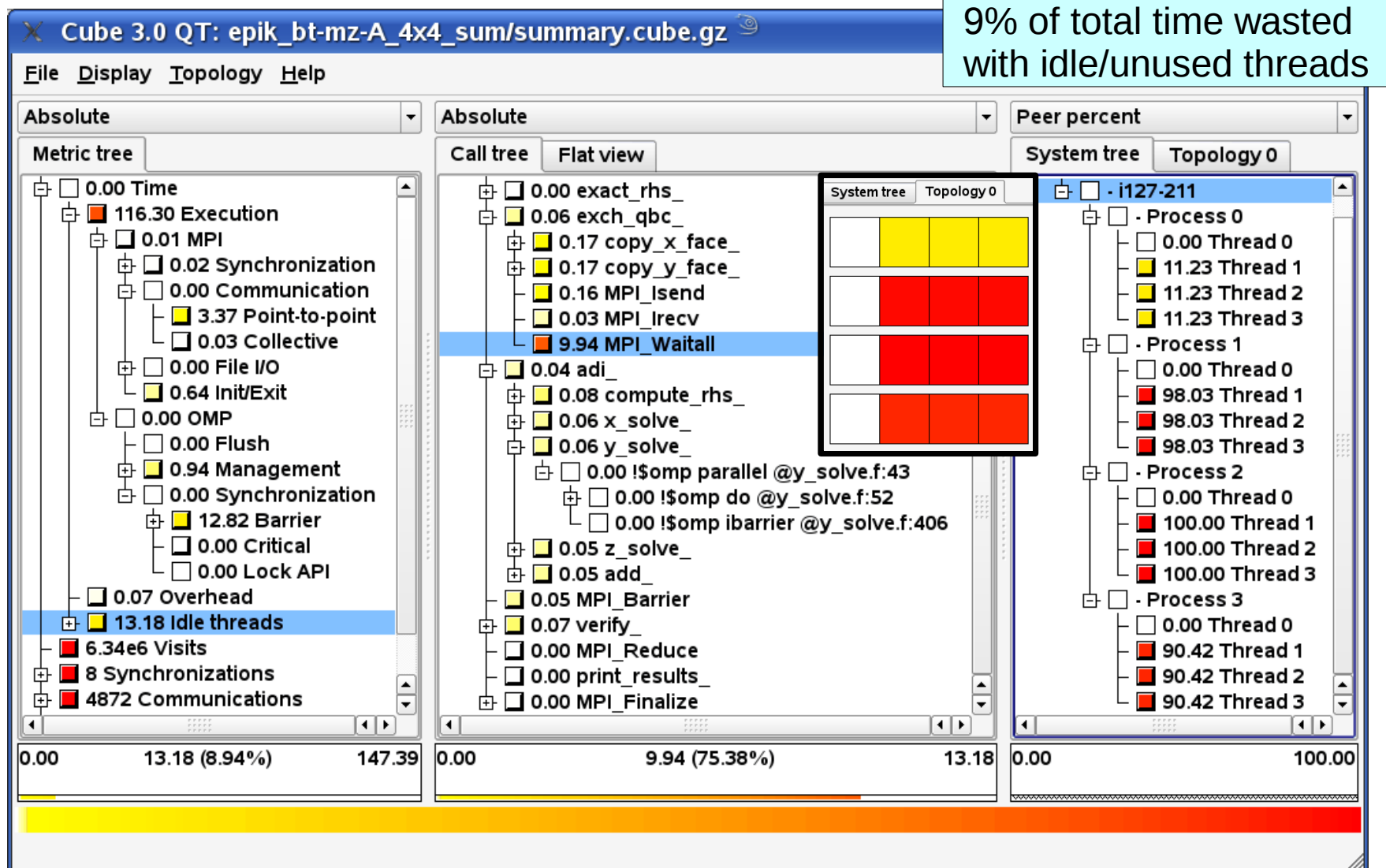


4x4 summary analysis: OpenMP time

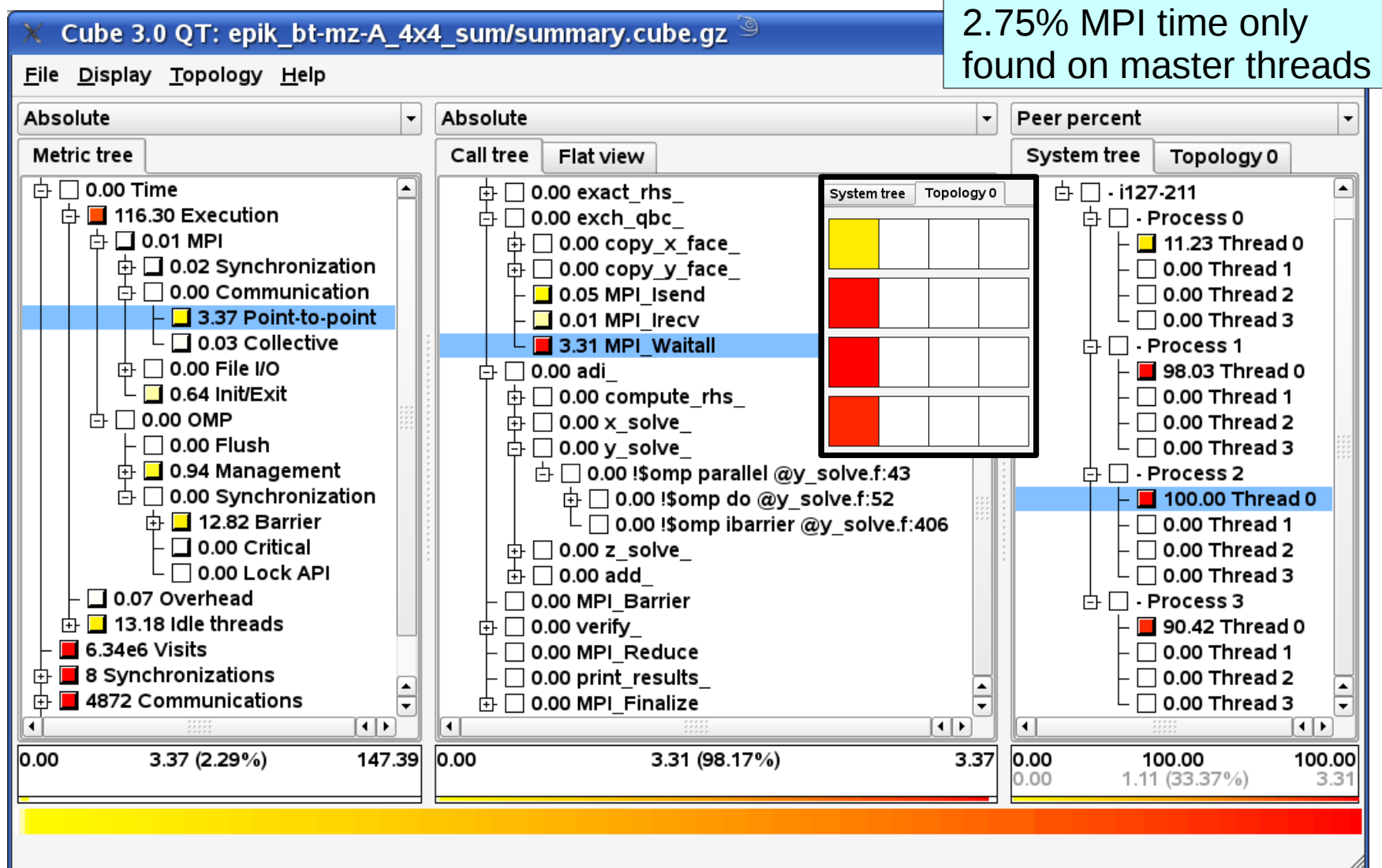
9% OpenMP time mostly found at implicit barriers



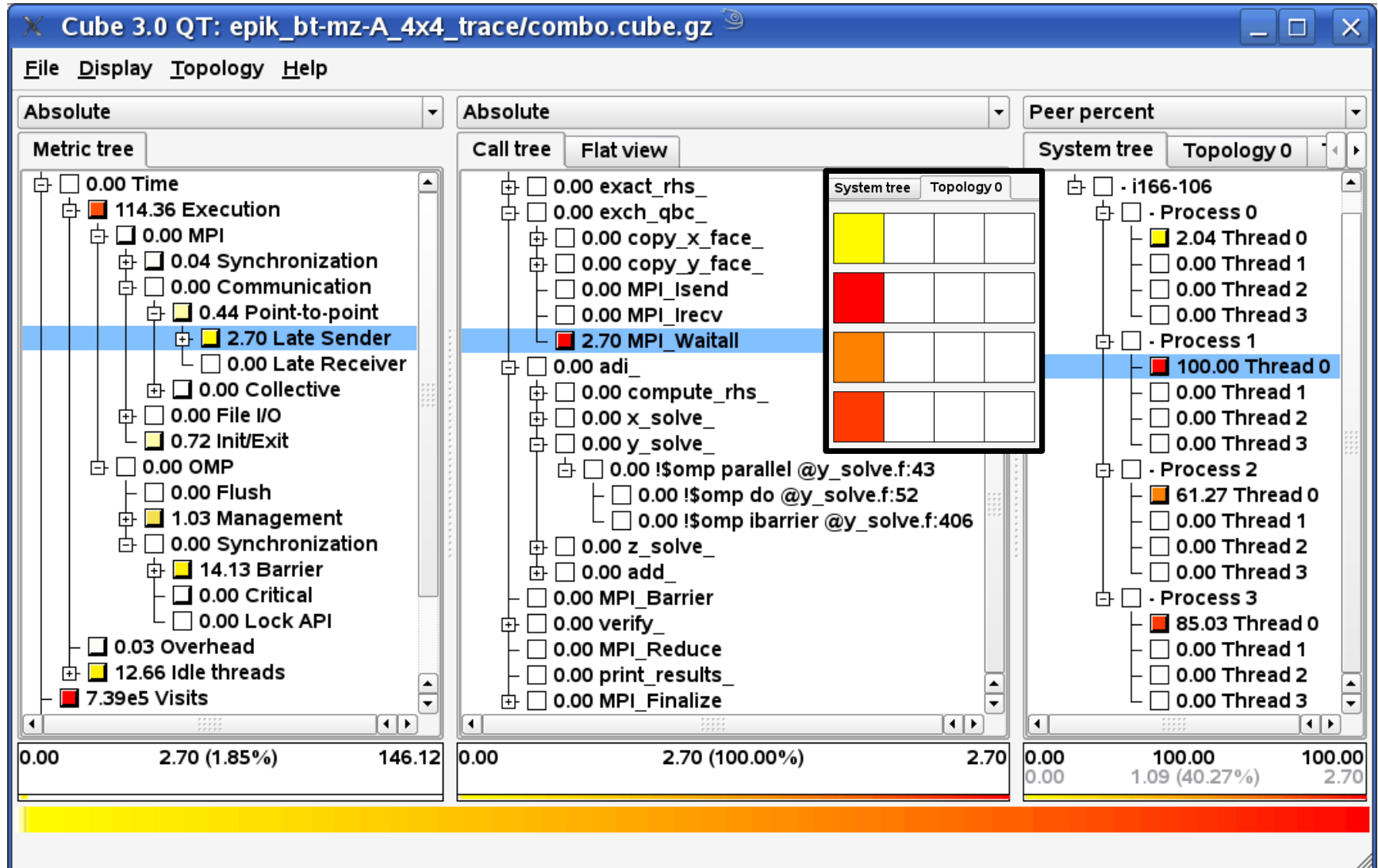
4x4 summary analysis: Idle threads time



4x4 summary analysis: MPI time



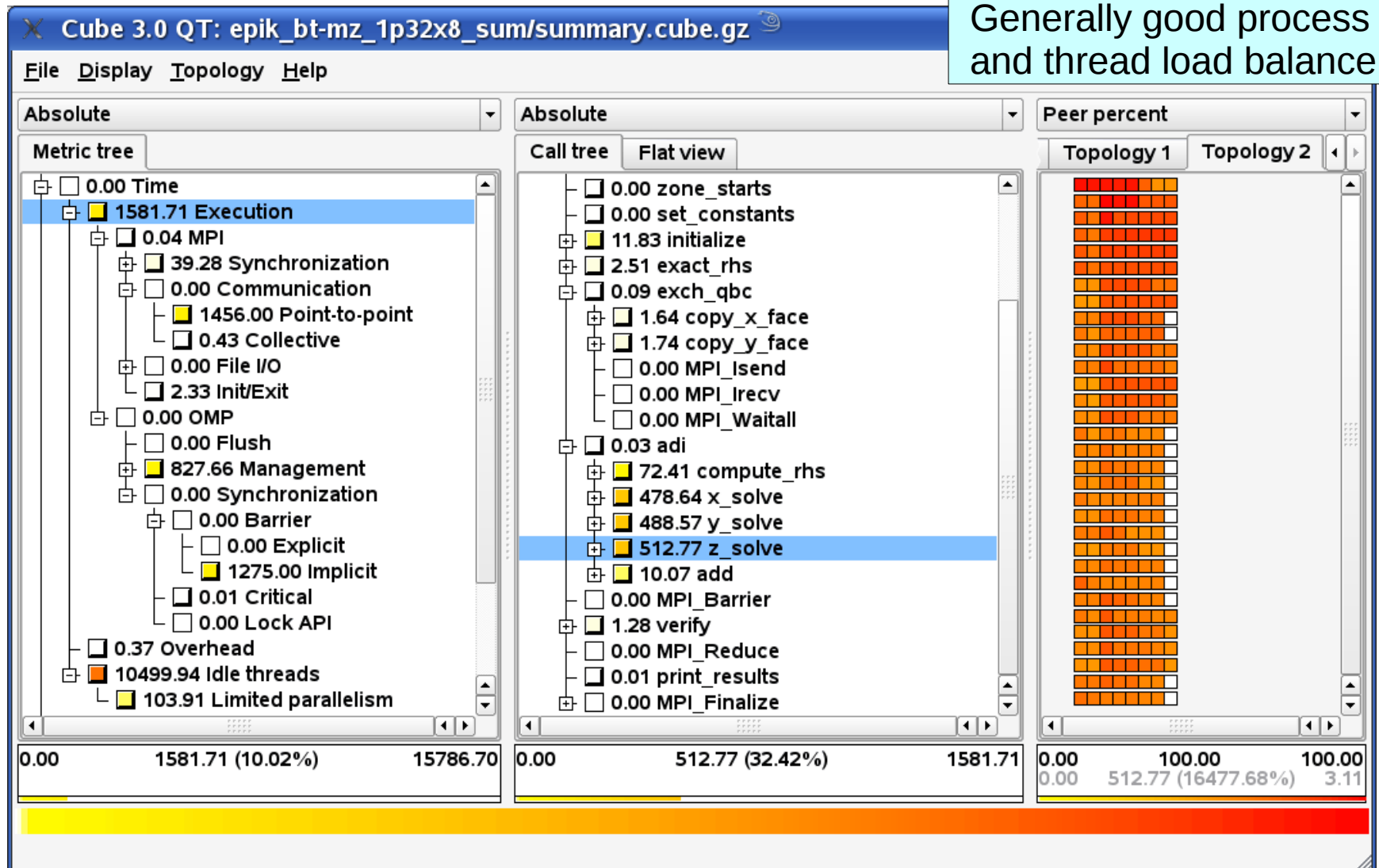
4x4 combined summary & trace analysis



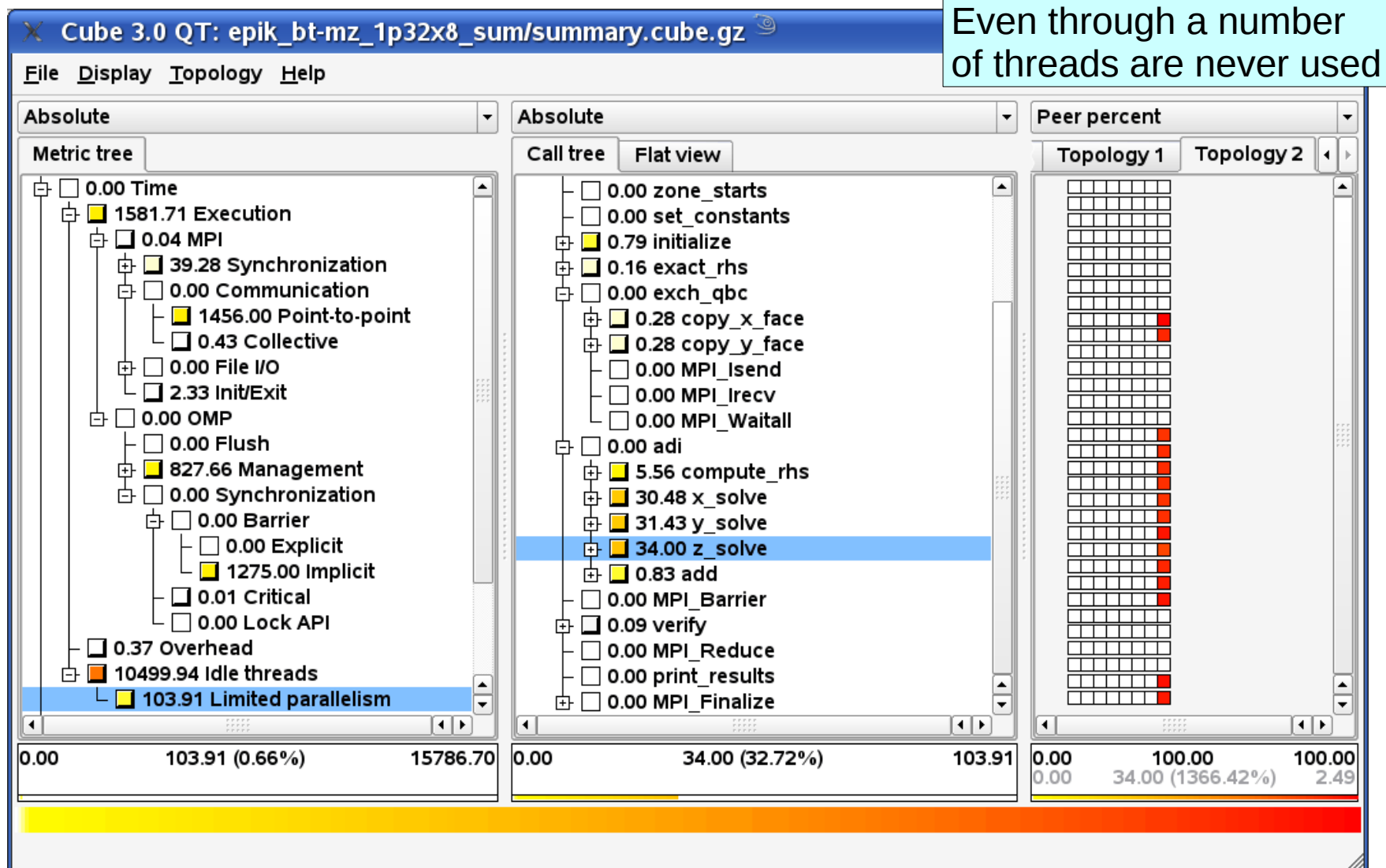
32x8 summary analysis: Excl. execution time



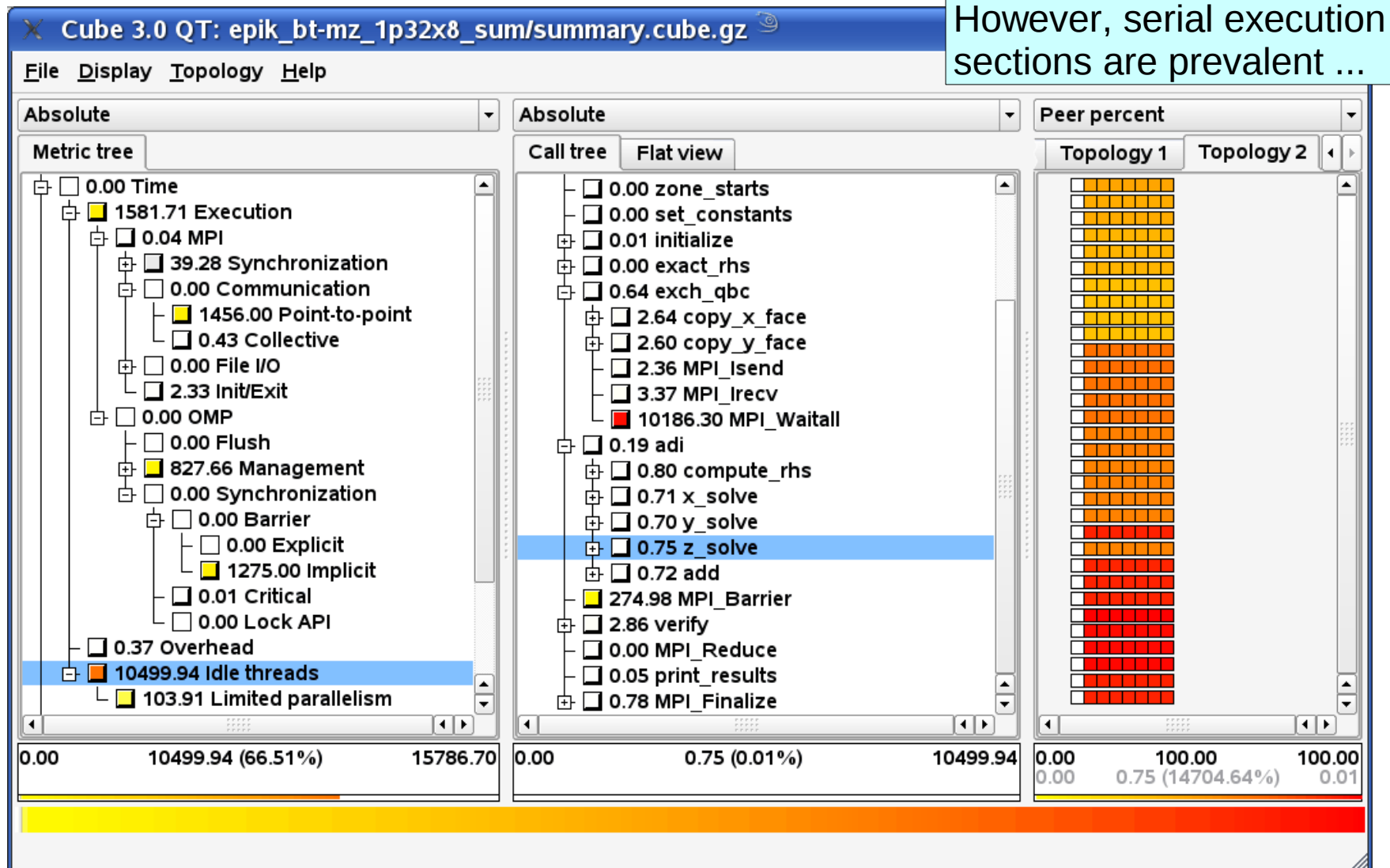
Generally good process and thread load balance



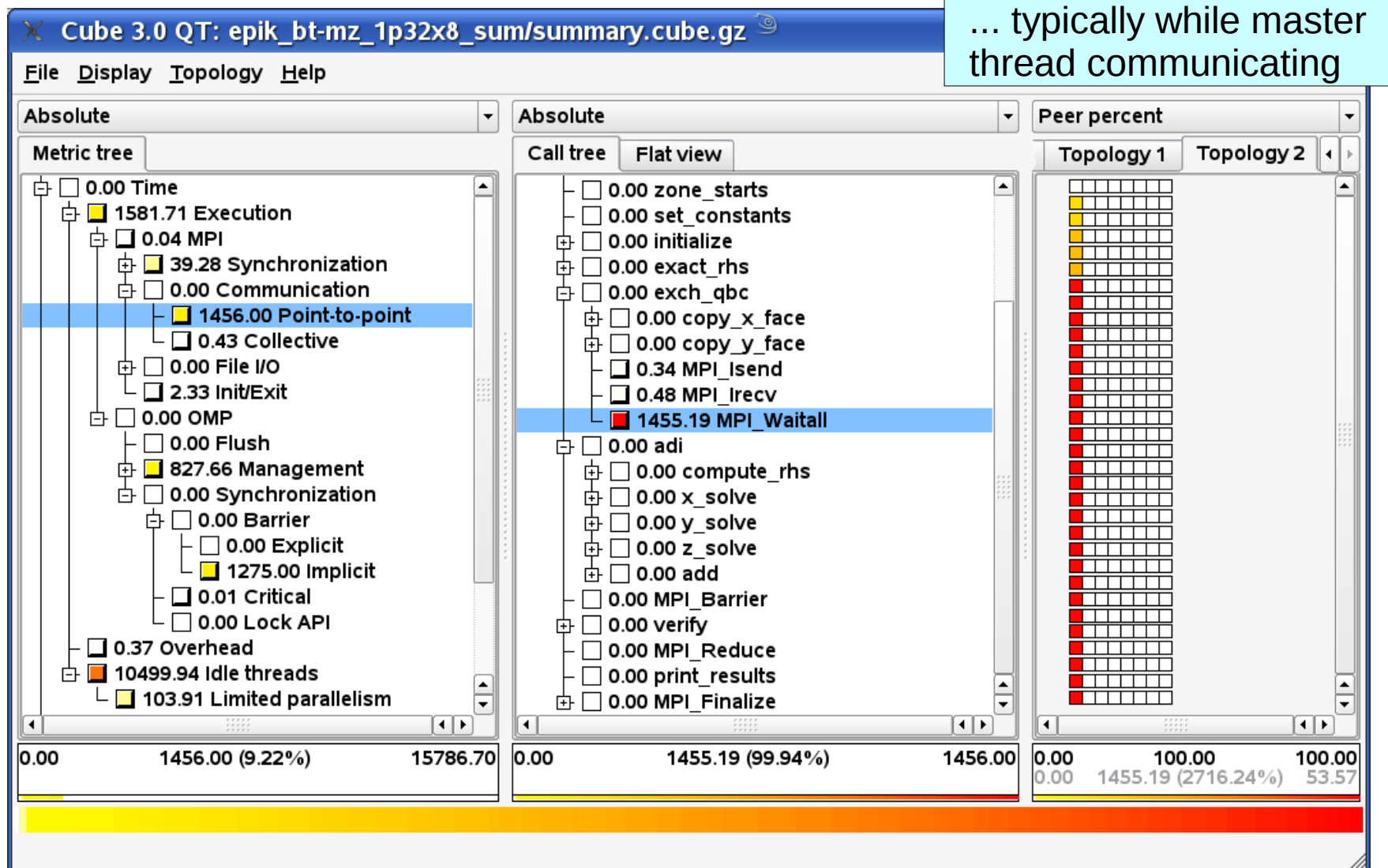
32x8 summary analysis: Limited parallelism



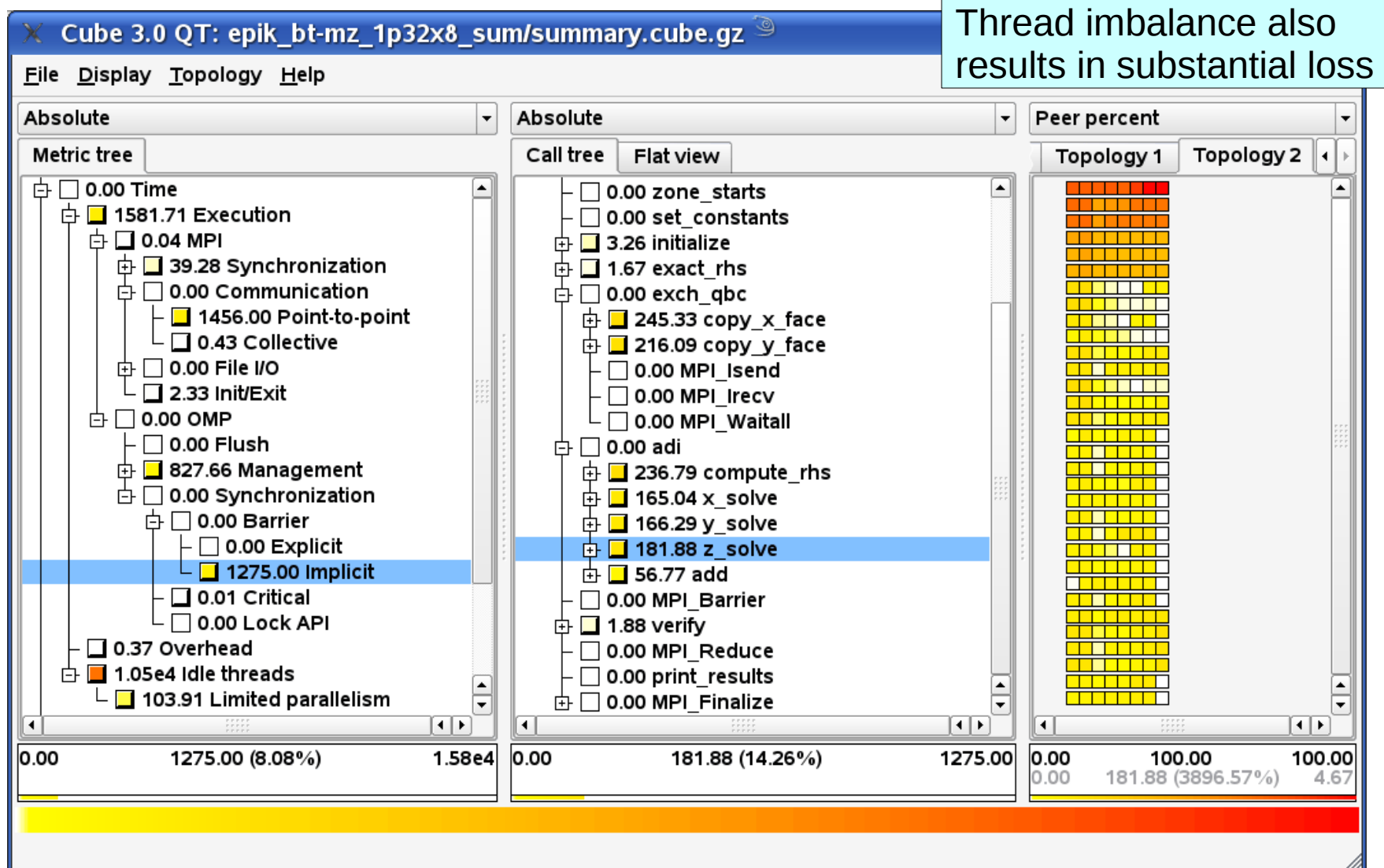
32x8 summary analysis: Idle threads time



32x8 summary analysis: MPI communication time

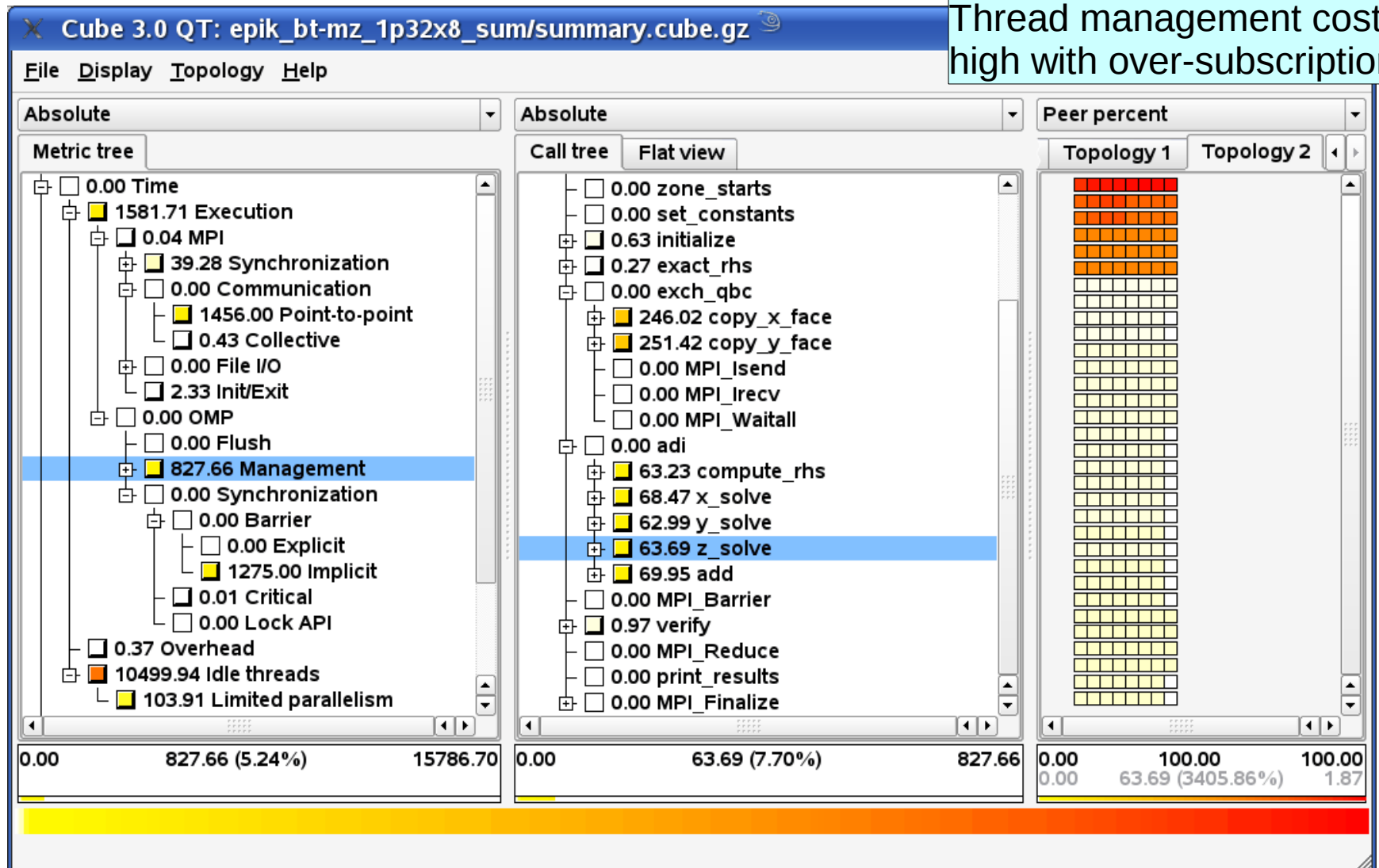


32x8 summary analysis: Implicit barrier time



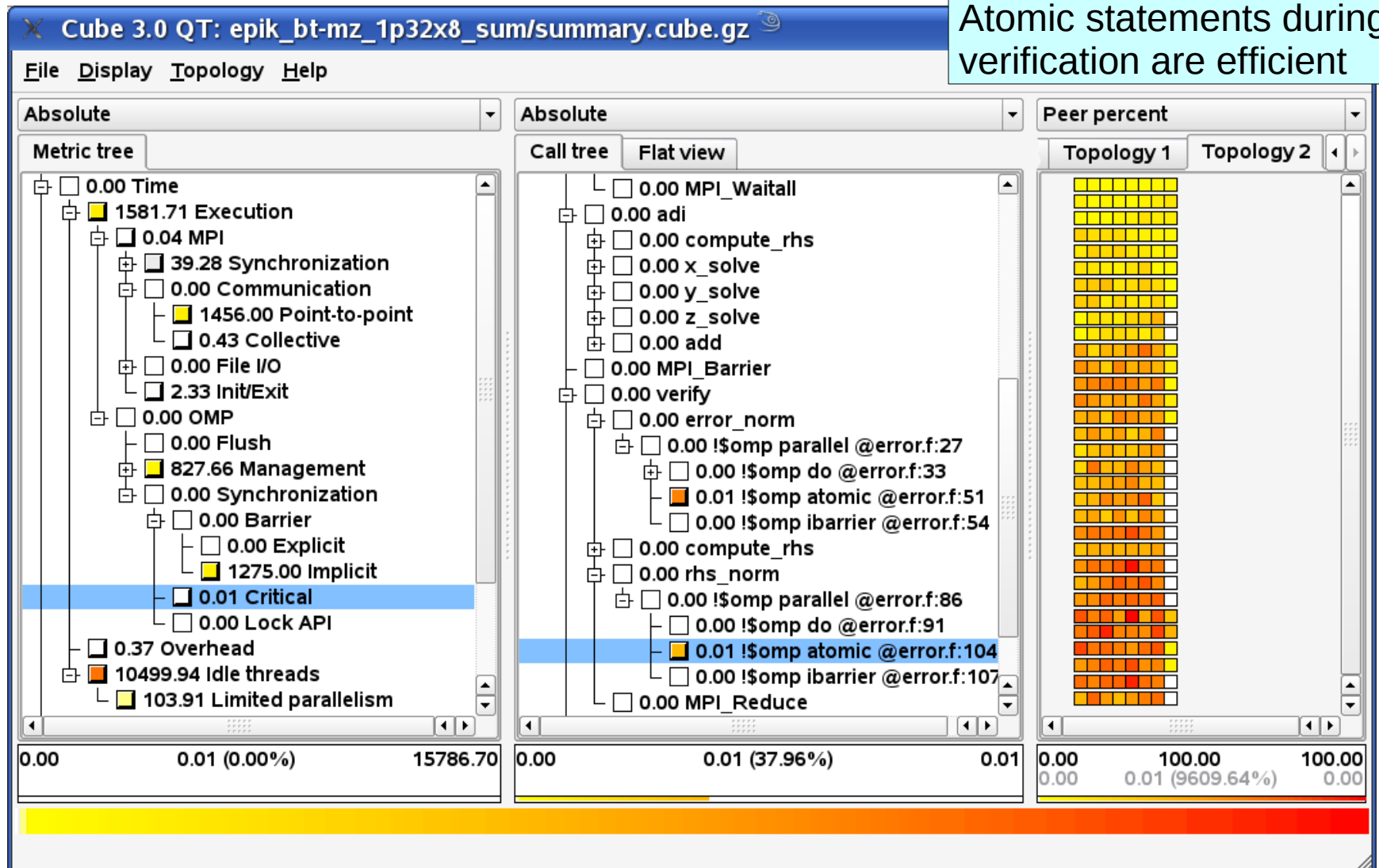
32x8 summary analysis: Thread management

Thread management cost high with over-subscription



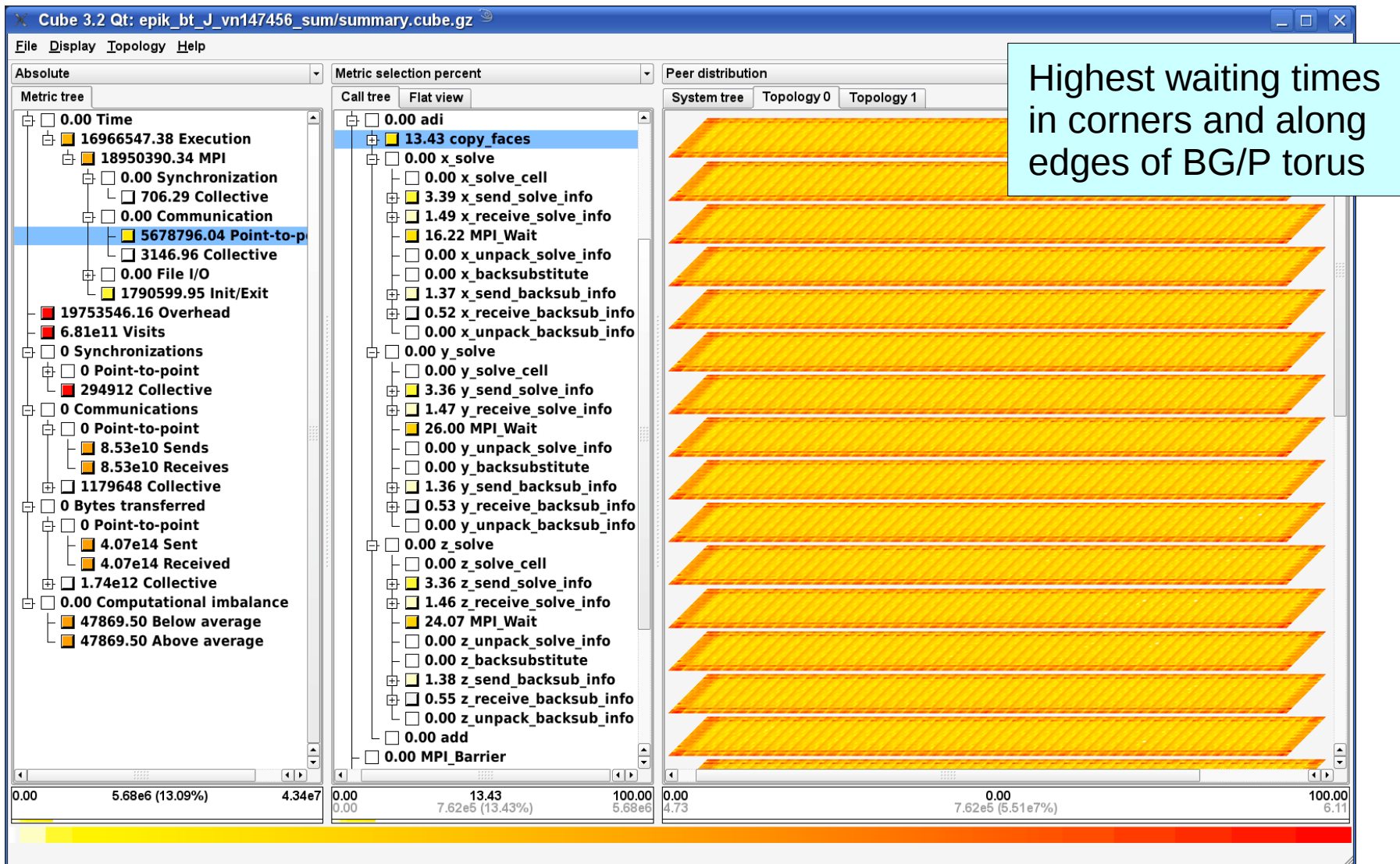
32x8 summary analysis: Critical section time

Atomic statements during verification are efficient

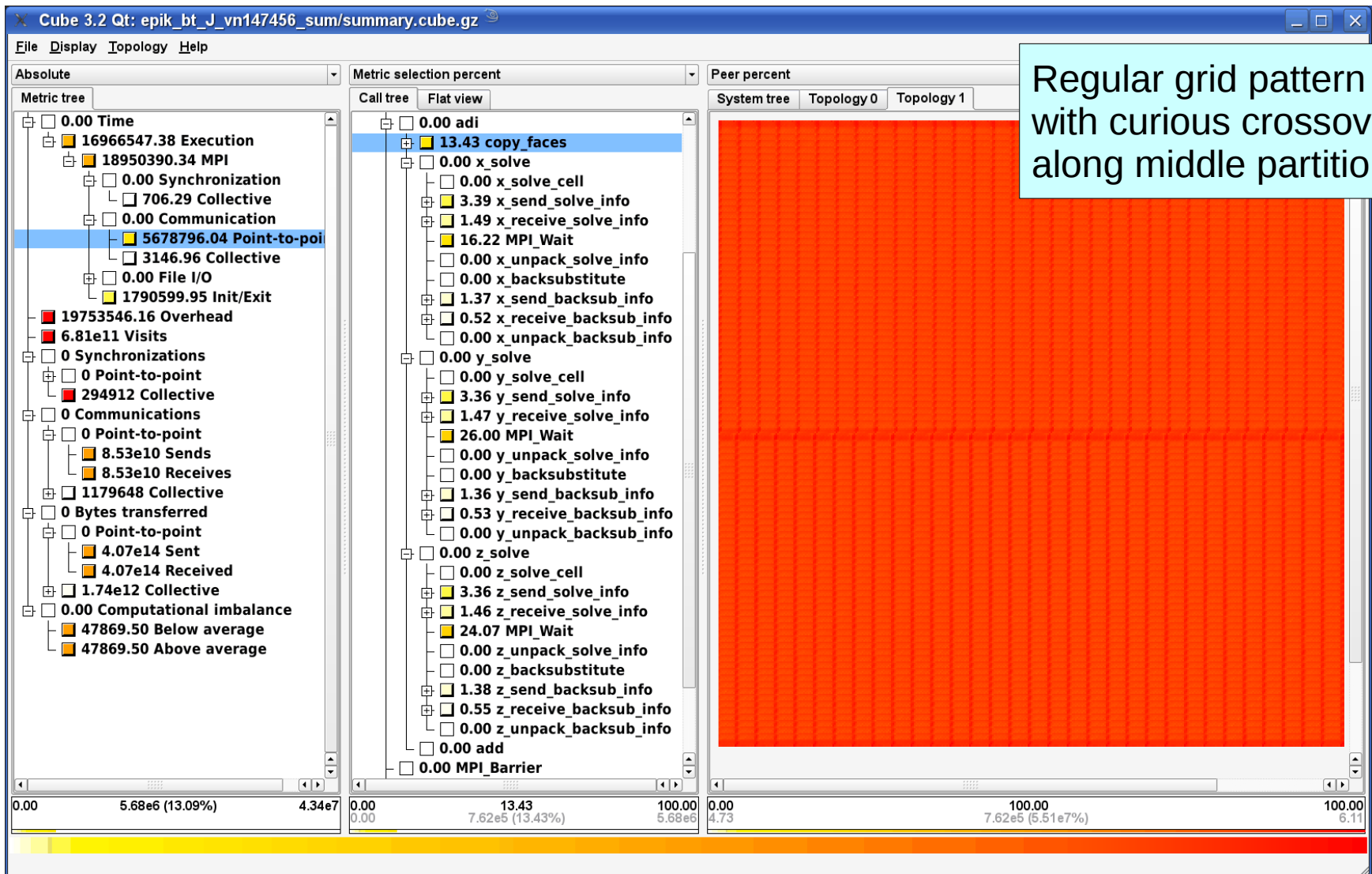


- 3D solution of unsteady, compressible Navier-Stokes eqs
 - NASA NAS parallel benchmark suite Block-Tridiagonal solver
 - series of ADI solve steps in X, Y & Z dimensions
 - ~9,500 lines (20 source modules), mostly Fortran77
- Run on IBM BlueGene/P in VN mode with 144k processes
 - Good scaling when problem size matched to architecture
 - ▶ 1536x1536x1536 gridpoints mapped onto 384x384 processes
 - Measurement collection took 53 minutes
 - 38% dilation for summarization measurement compared to uninstrumented execution (using 10 function filter)
 - MPI trace size would be 18.6TB
 - 25% of time in ADI is point-to-point communication time
 - ▶ 13% copy_faces, 23% x_solve, 33% y_solve, 31% z_solve
 - 128s for a single MPI_Comm_split during setup!

NPB-MPI-BT on jugene@144k summary analysis VI-HPS



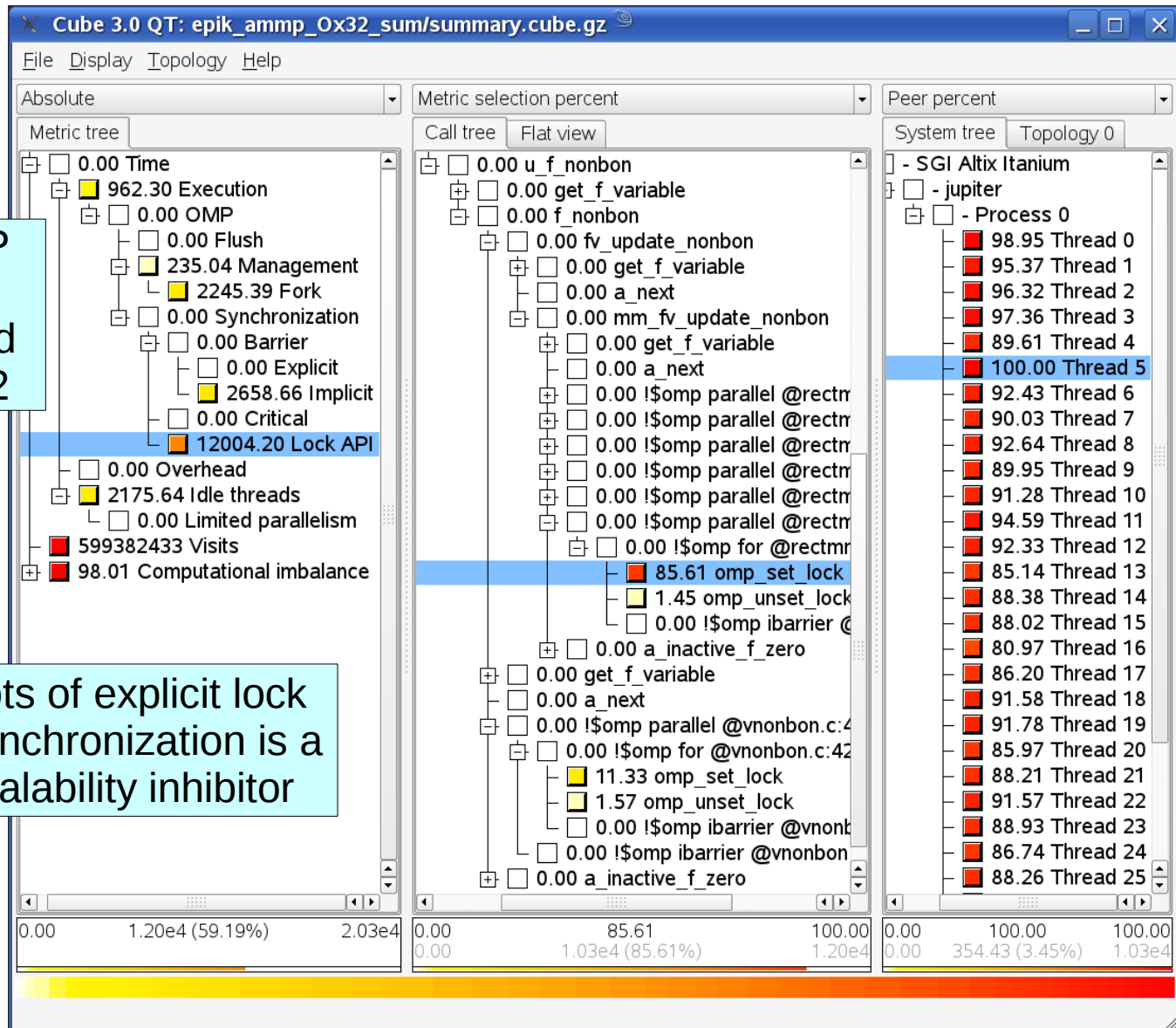
NPB-MPI-BT on jugene@144k summary analysis VI-HPS



- Molecular mechanics simulation
 - original version developed by Robert W. Harrison
- SPEC OMP benchmark parallel version
 - ~14,000 lines (in 28 source modules): 100% C
- Run with 32 threads on SGI Altix 4700 at TUD-ZIH
 - Built with Intel compilers
 - 333 simulation timesteps for 9,582 atoms
- Scalasca summary measurement
 - Minimal measurement dilation
 - 60% of total time lost in synchronization with lock API
 - 12% thread management overhead

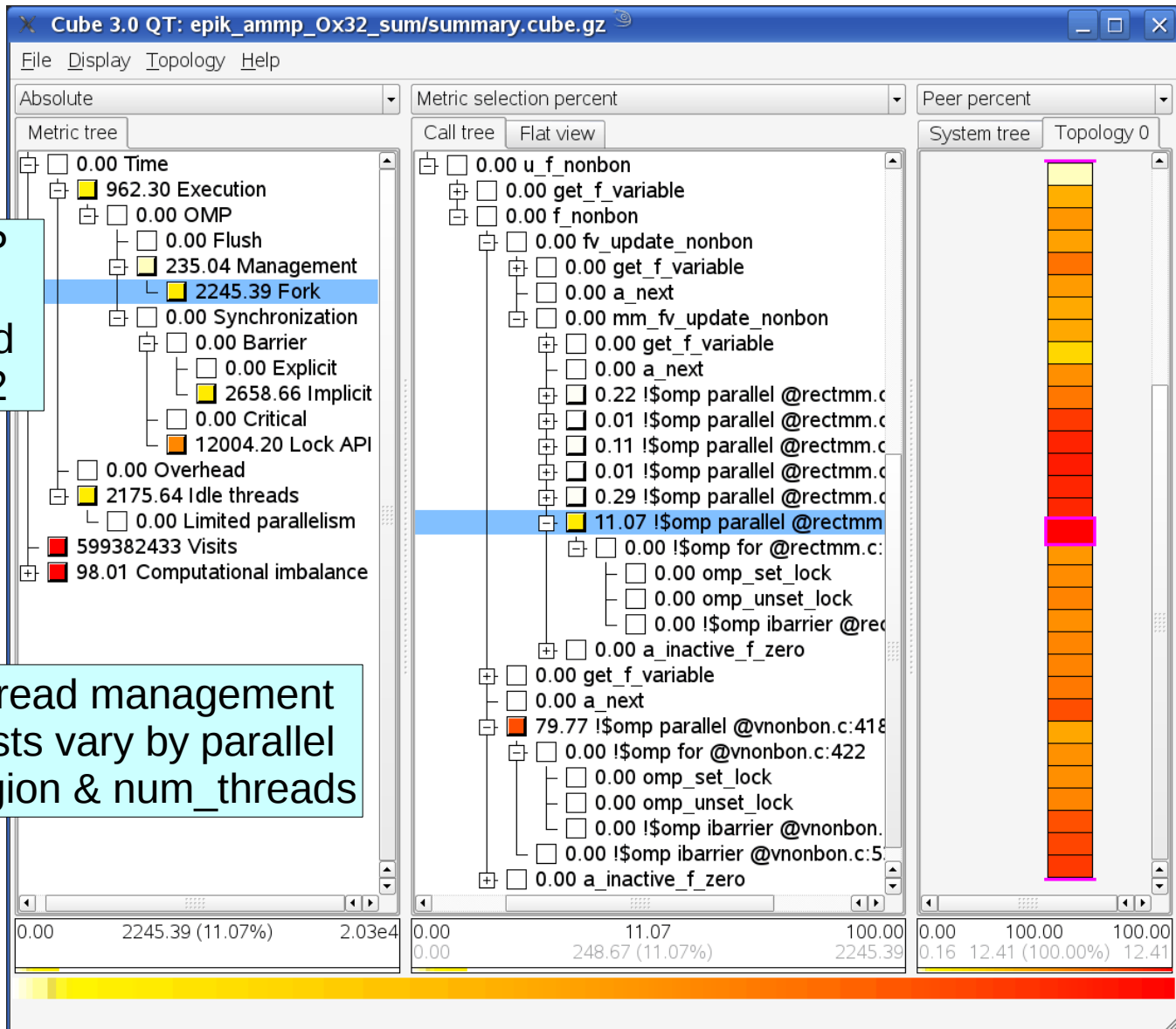
OpenMP
metrics
reworked
with v1.2

Lots of explicit lock
synchronization is a
scalability inhibitor

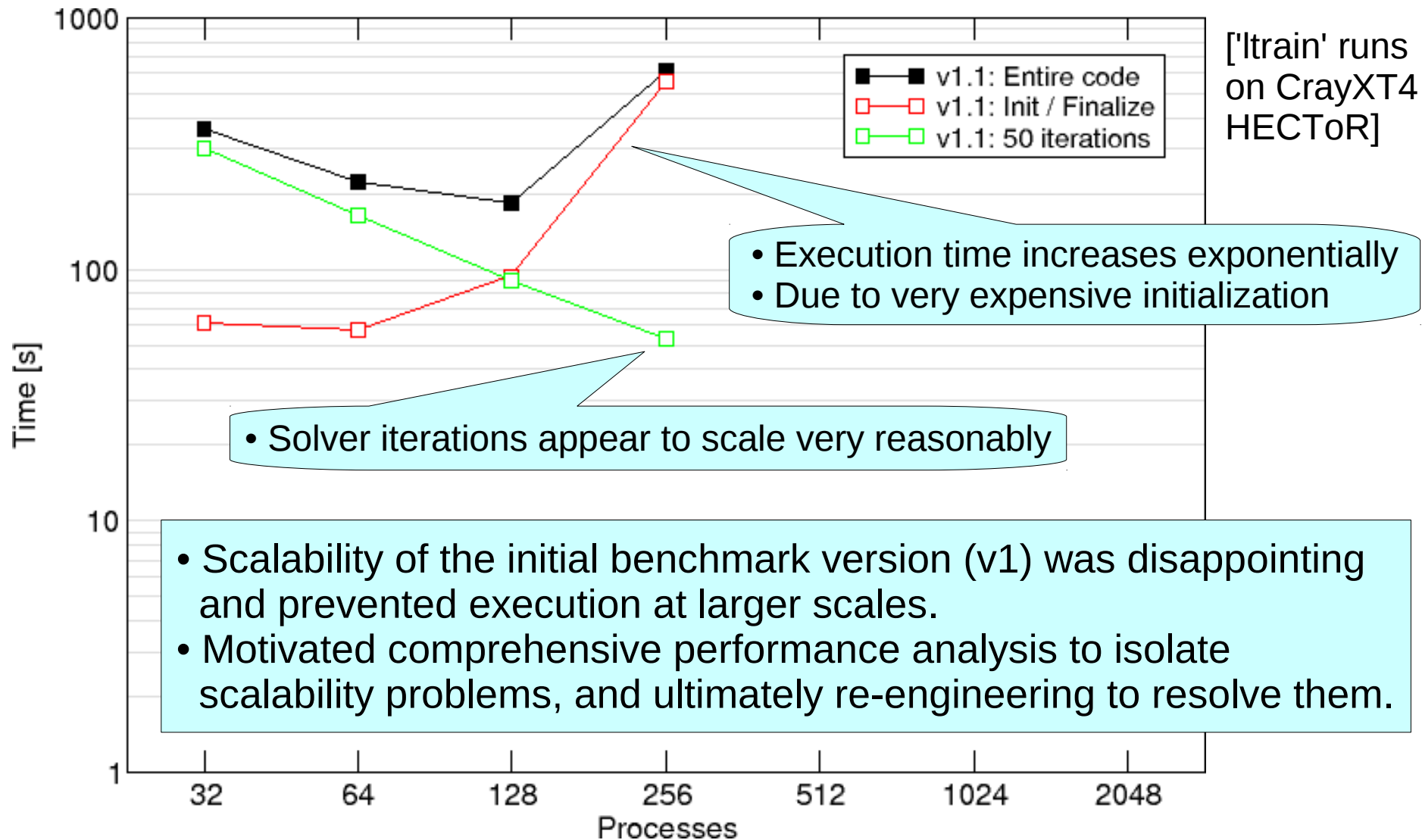


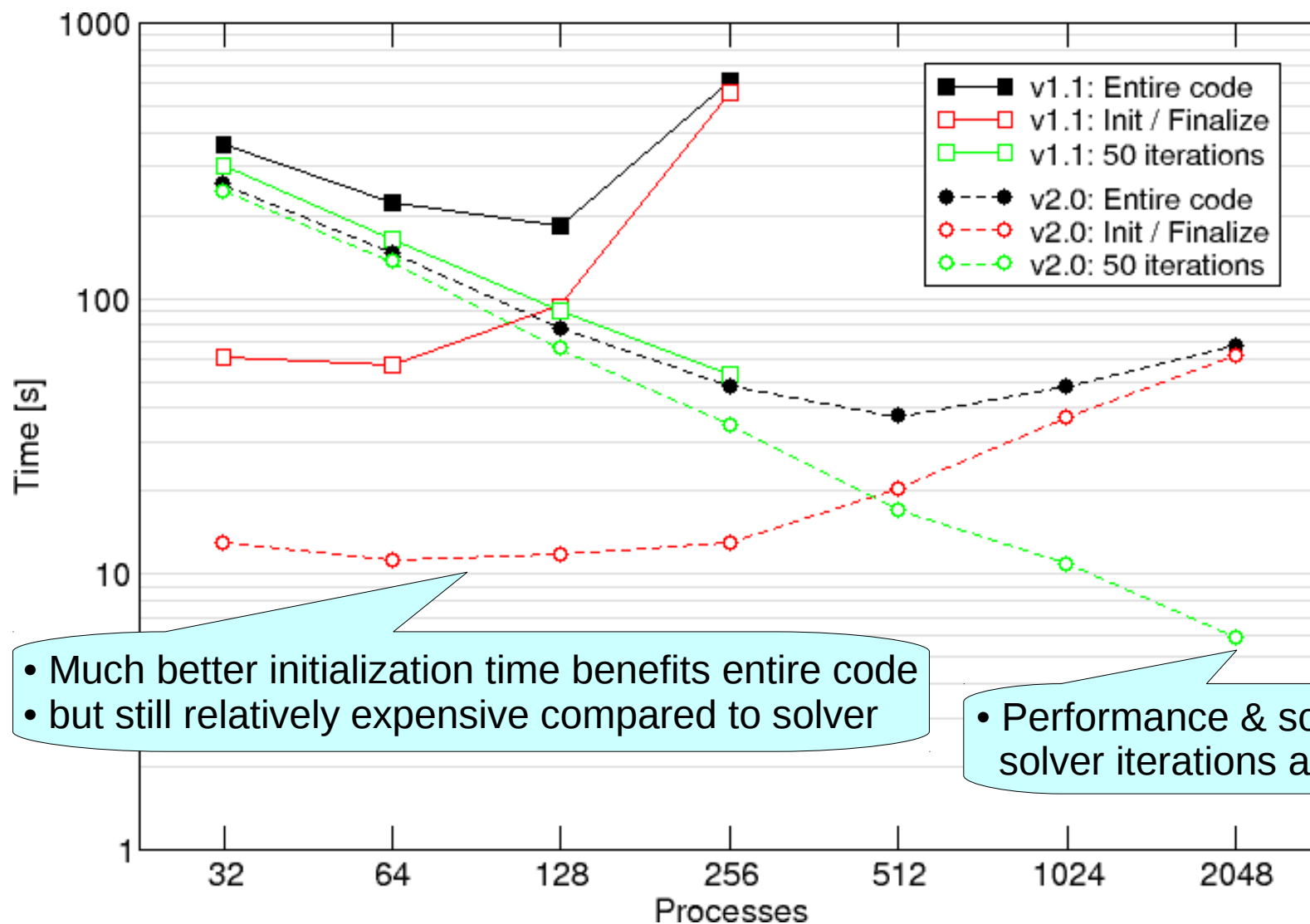
OpenMP
metrics
reworked
with v1.2

Thread management
costs vary by parallel
region & num_threads



- Computational electromagnetics solver
 - originates from KTH General ElectroMagnetics Solvers project
 - finite-difference time-domain method for Maxwell equations
- MPI parallel versions in SPEC MPI2007 benchmark suite
 - original **v1.1** (113.GemsFDTD) “medium” size
 - revised **v2.0** (145.lGemsFDTD) “large” size
 - built with PGI 9.0.4 Fortran90 compiler (21k lines of code)
 - ▶ typical benchmark optimization: `-fastsse -O3 -Mipa=fast,inline`
- Both run on 'hector' Cray XT4 at EPCC
 - using “ltrain” dataset from v2.0 benchmark (50 timesteps)
 - default Scalasca instrumentation for measurements
 - ▶ 9 of 90 application user-level source routines specified in filter determined by scoring initial summary experiment





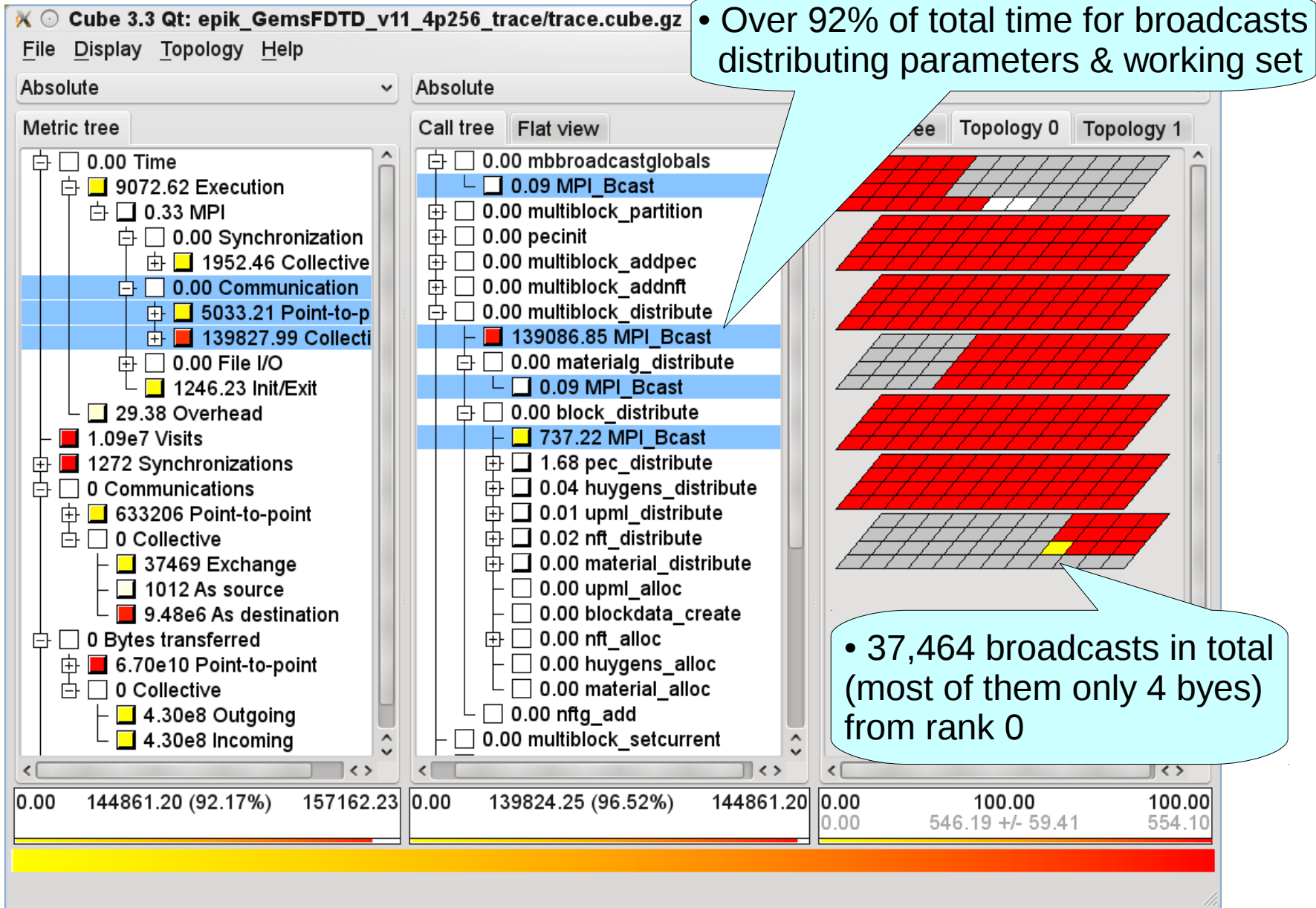
['ltrain' runs on CrayXT4 HECToR]

- Much better initialization time benefits entire code
- but still relatively expensive compared to solver

- Performance & scalability of solver iterations also improved

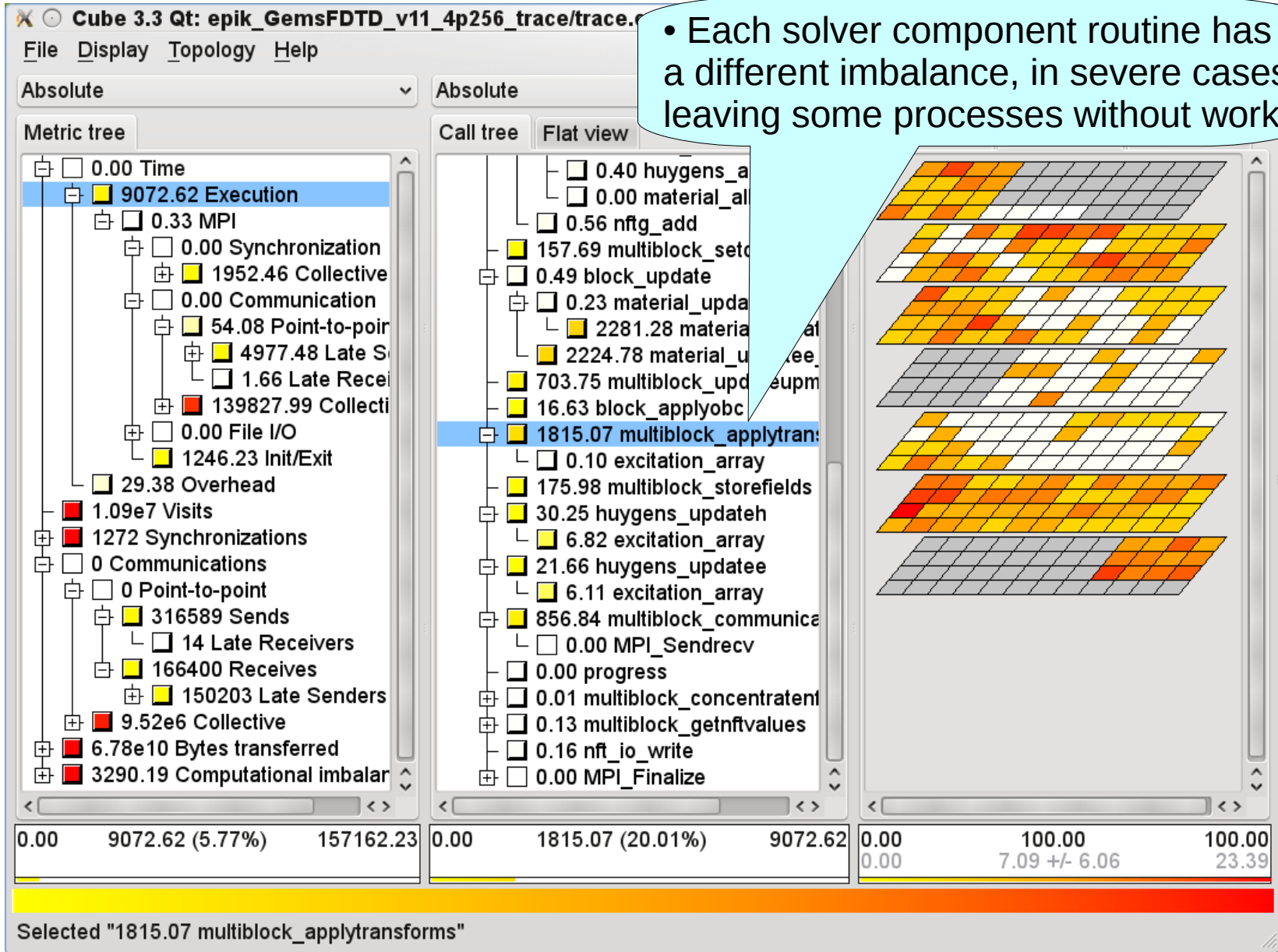
Time for initialization broadcasts (v1.1)

VI-HPS



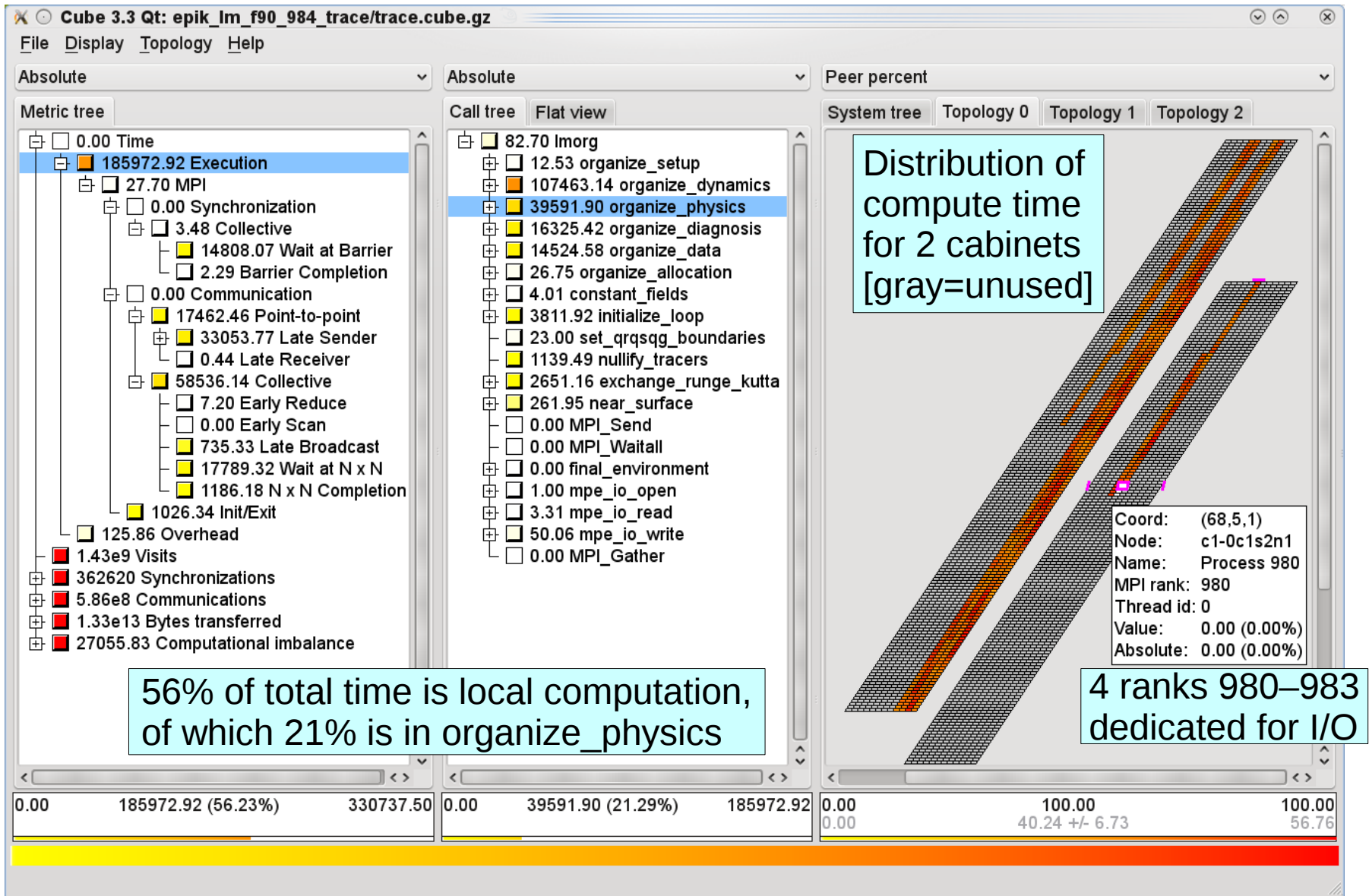
Computation time in solver transforms (v1.1)

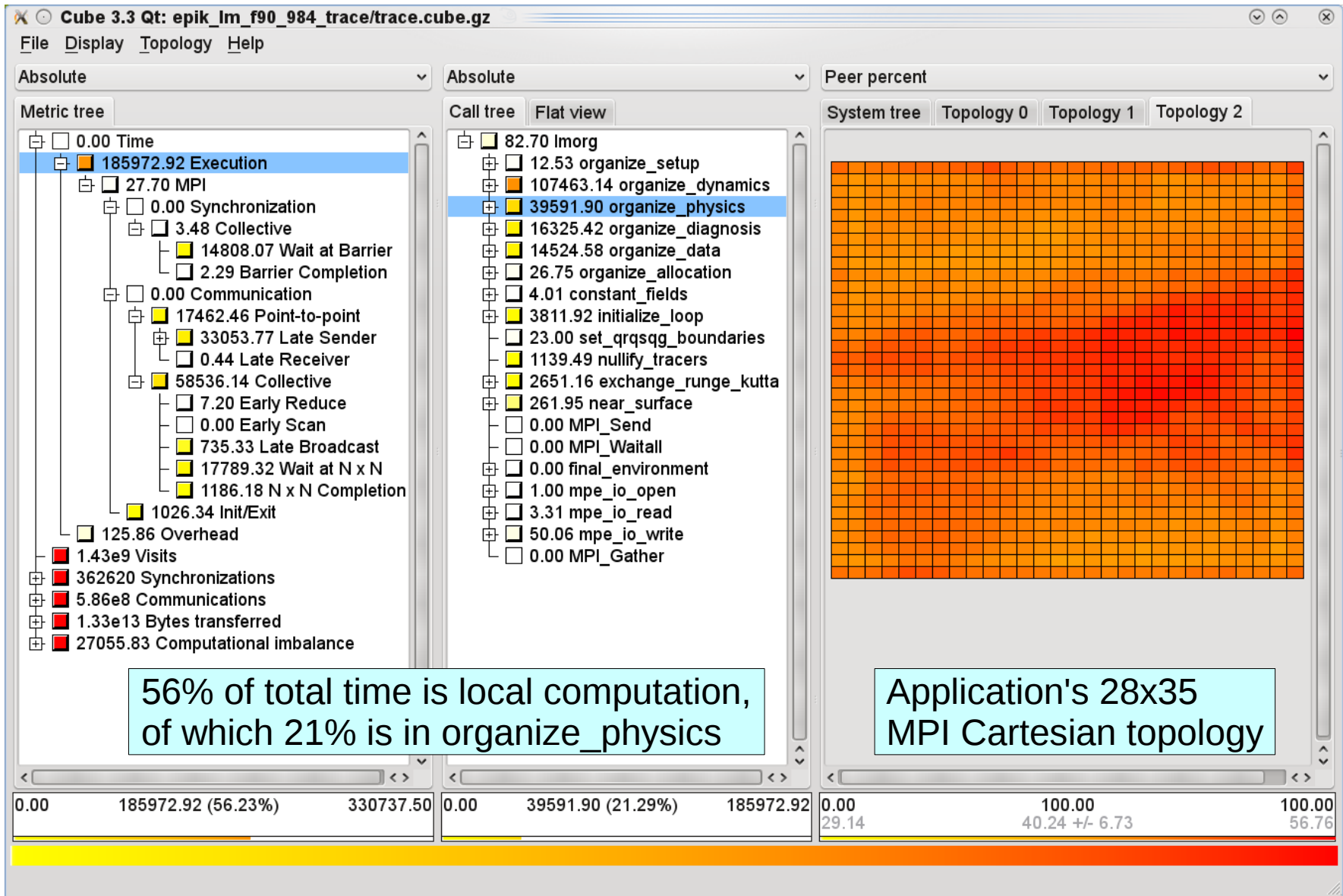
- Each solver component routine has a different imbalance, in severe cases leaving some processes without work

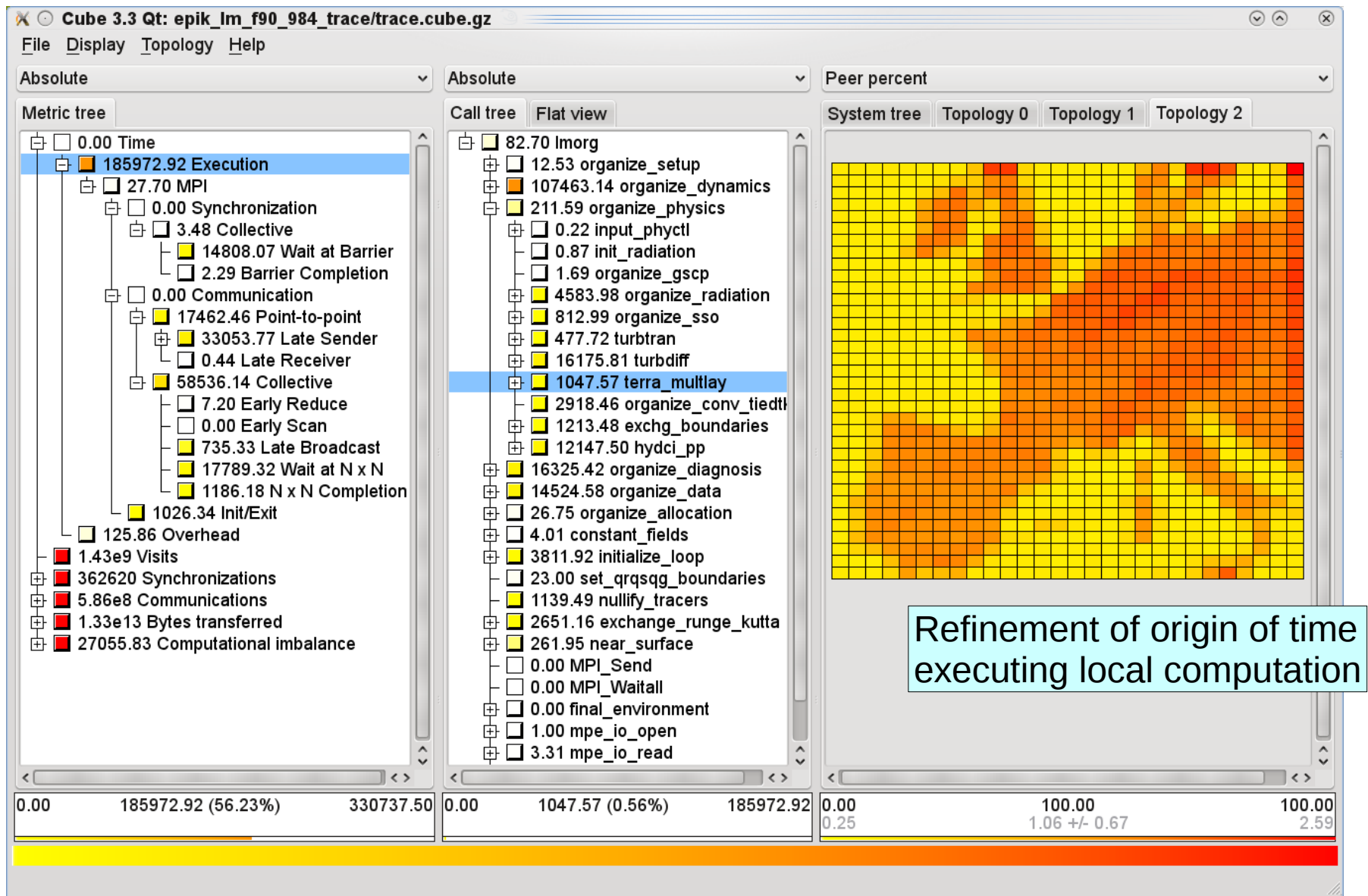


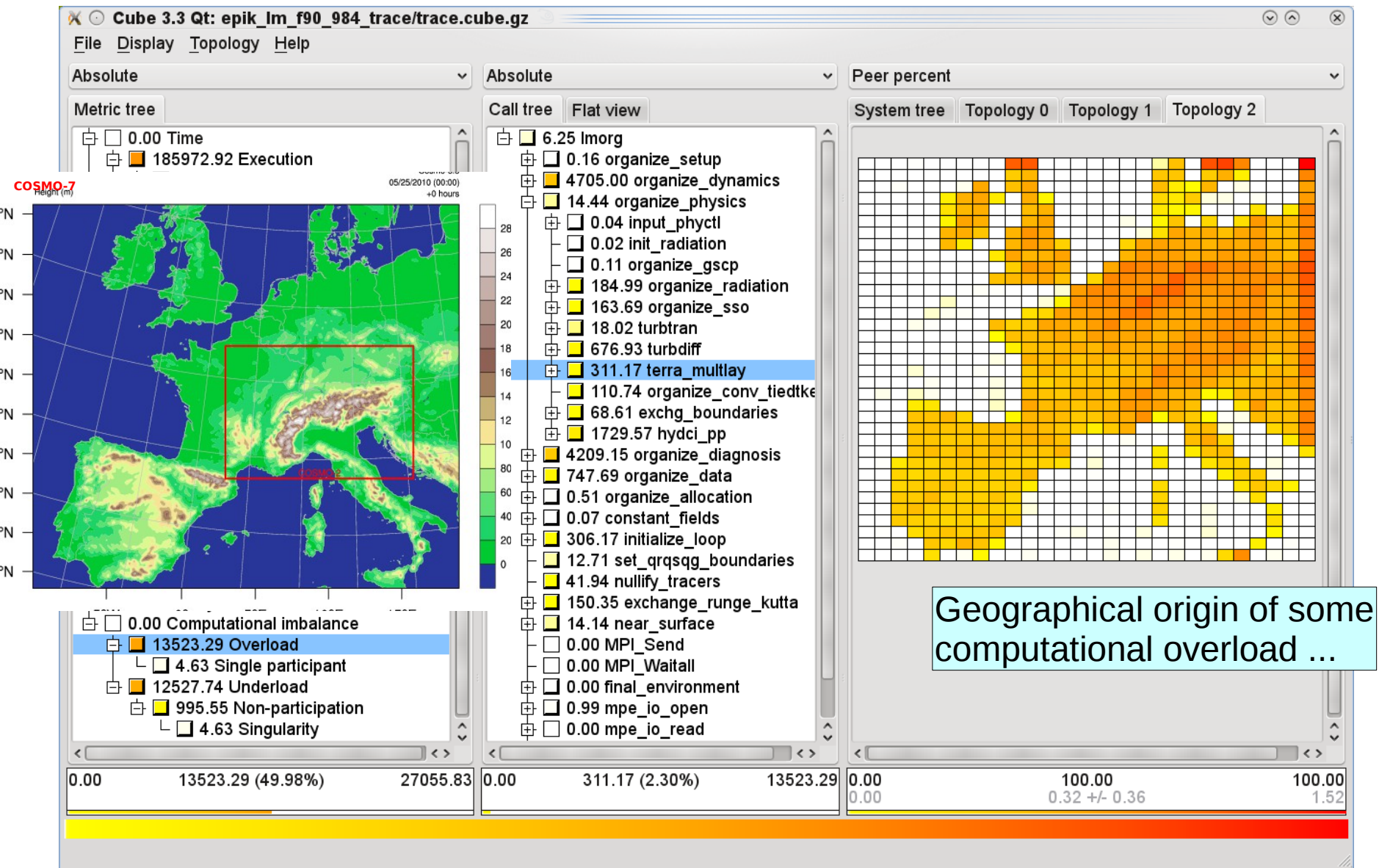
- Initialization originally dominated by numerous broadcasts and expensive serial multiblock partition by rank 0
 - Re-engineered implementation of scalable partition routine, aggregation of multiple data values into larger messages, and postpones allocations until all block information in broadcast
 - ▶ Initialization time reduced to less than 2% of total time
- Solver iterations using blocking communication manifests as *Late Sender* waiting originating from imbalance in local computation time (due to different computations)
 - Re-engineered implementation uses non-blocking comms and re-uses communication pattern used to exchange blocks (as well as 2 of 256 processes unintentionally idled throughout)
 - ▶ computation & communication time both improved more than 25%
- Scalability improved from 128 processes to more than 1024

- Regional climate and weather model
 - developed by Consortium for Small-scale Modeling (COSMO)
 - ▶ DWD, MeteoSwiss and others
 - non-hydrostatic limited-area atmospheric model (6.6km grid)
- MPI parallel version 4.12 (Jan-2011)
 - built with PGI 10.9 Fortran90 compiler (222k lines of code)
- MeteoSwiss operational 24-hour forecast of 06-Dec-2010
 - Western Europe 393x338x60 resolution, 1440 timesteps
- Run with 984 processes on 'palu' Cray XE6 at CSCS
 - 28x35 compute grid + 4 dedicated I/O processes
 - used 41 Opteron compute nodes each with 24 cores
 - Scalasca trace measurement with 19 of 178 routines filtered
 - 44GB trace written in 23s and analyzed in 82s

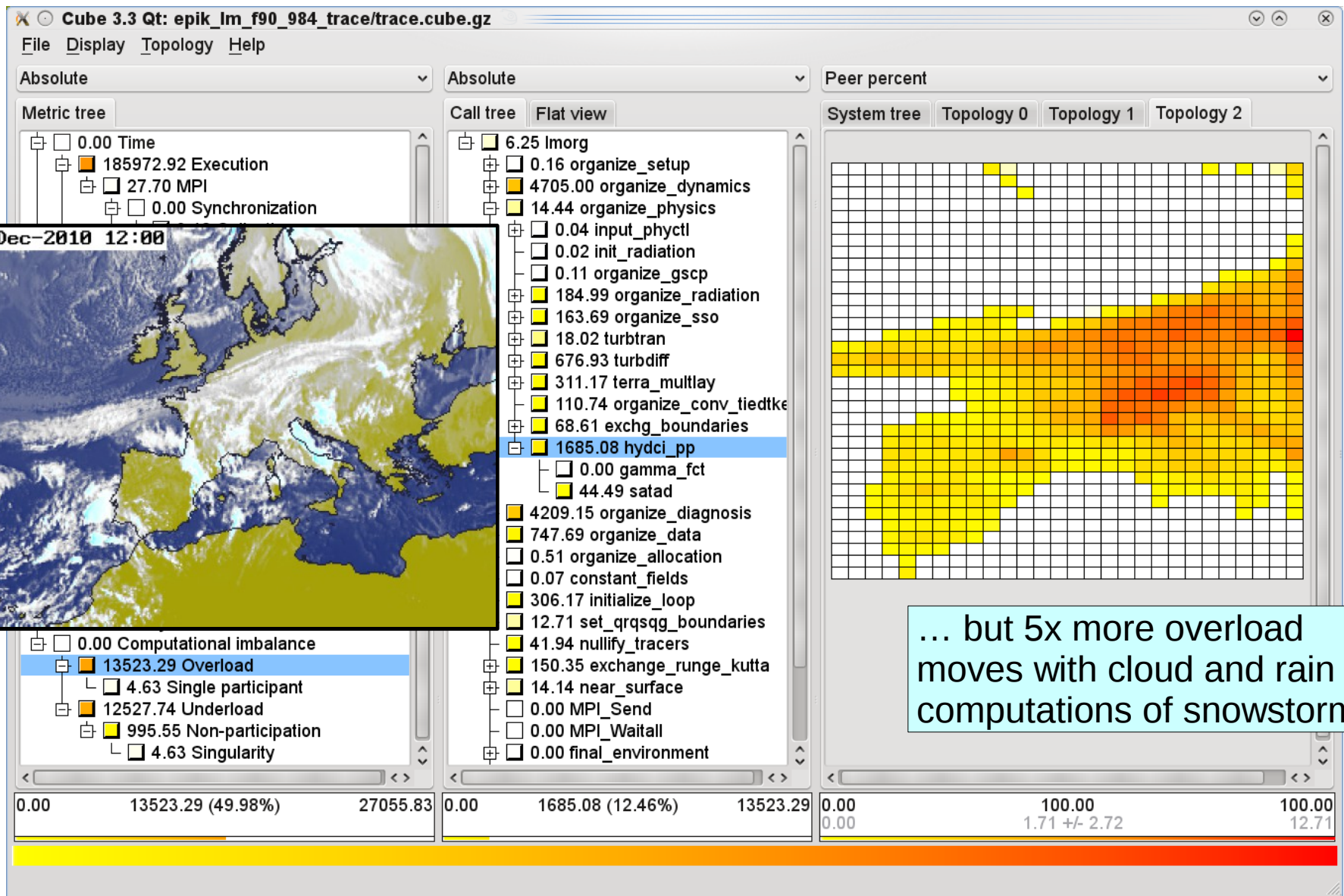




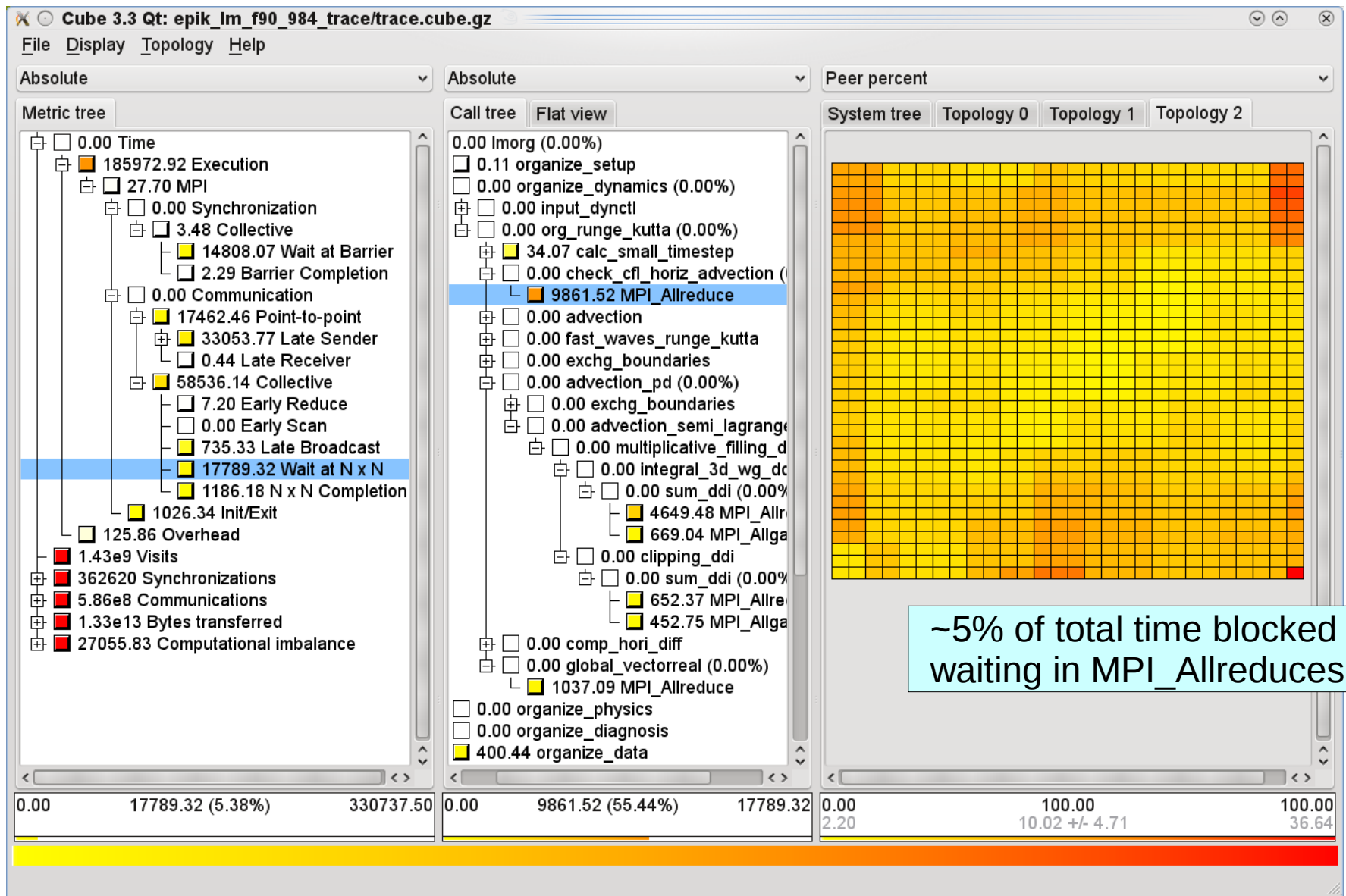




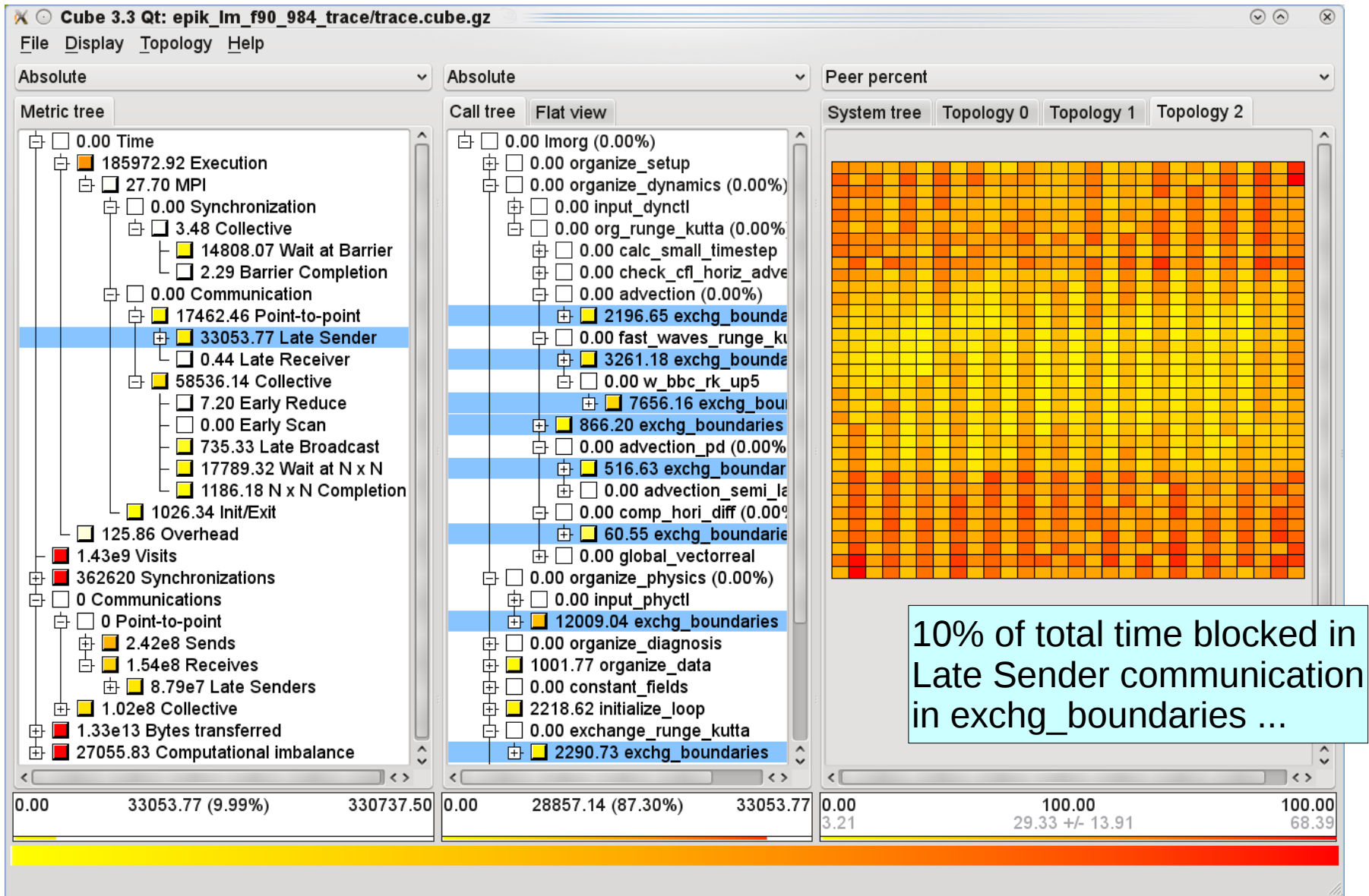
COSMO/XE6 computational overload (hydro)

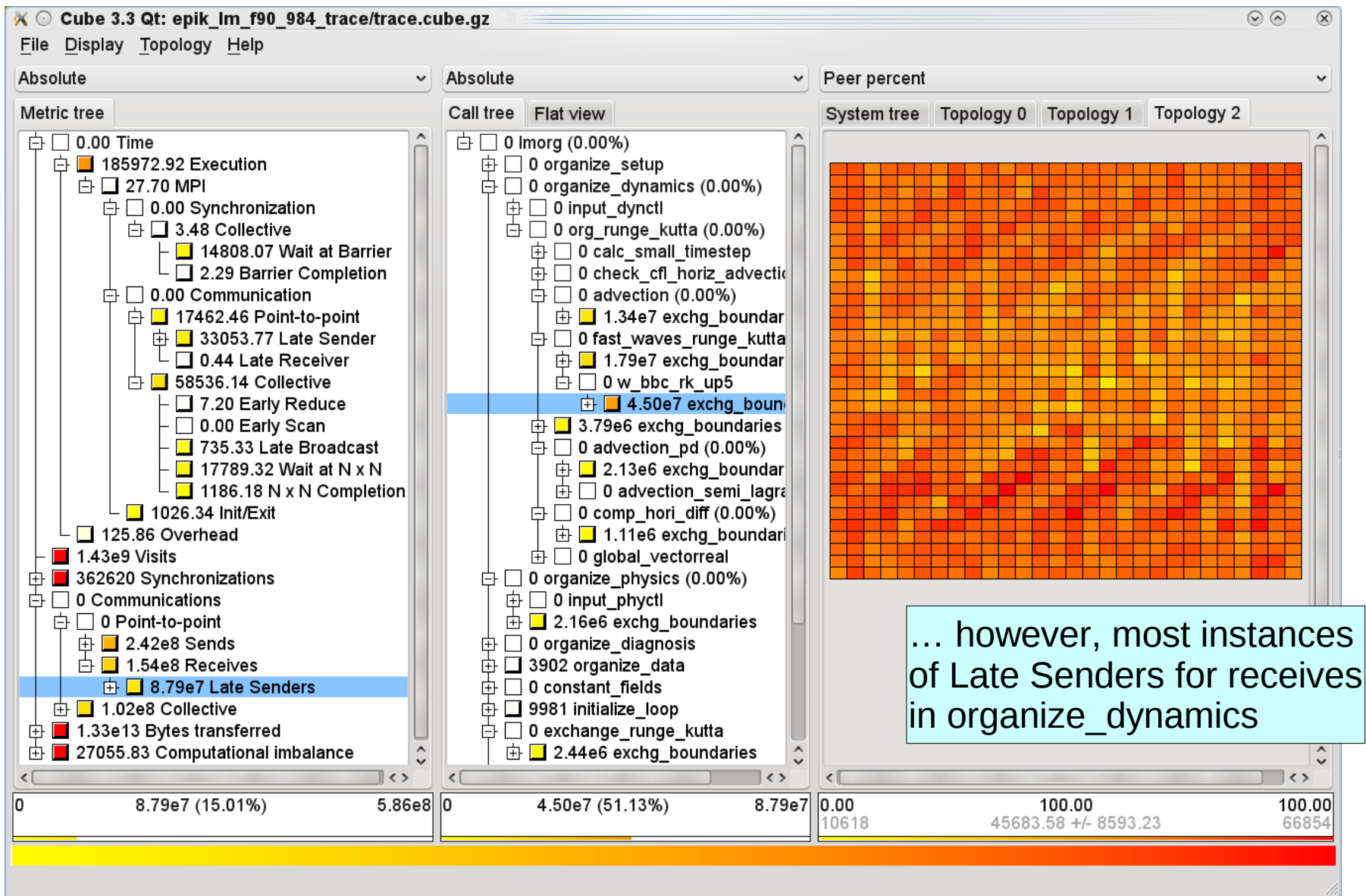


COSMO/XE6 collective wait at N x N time



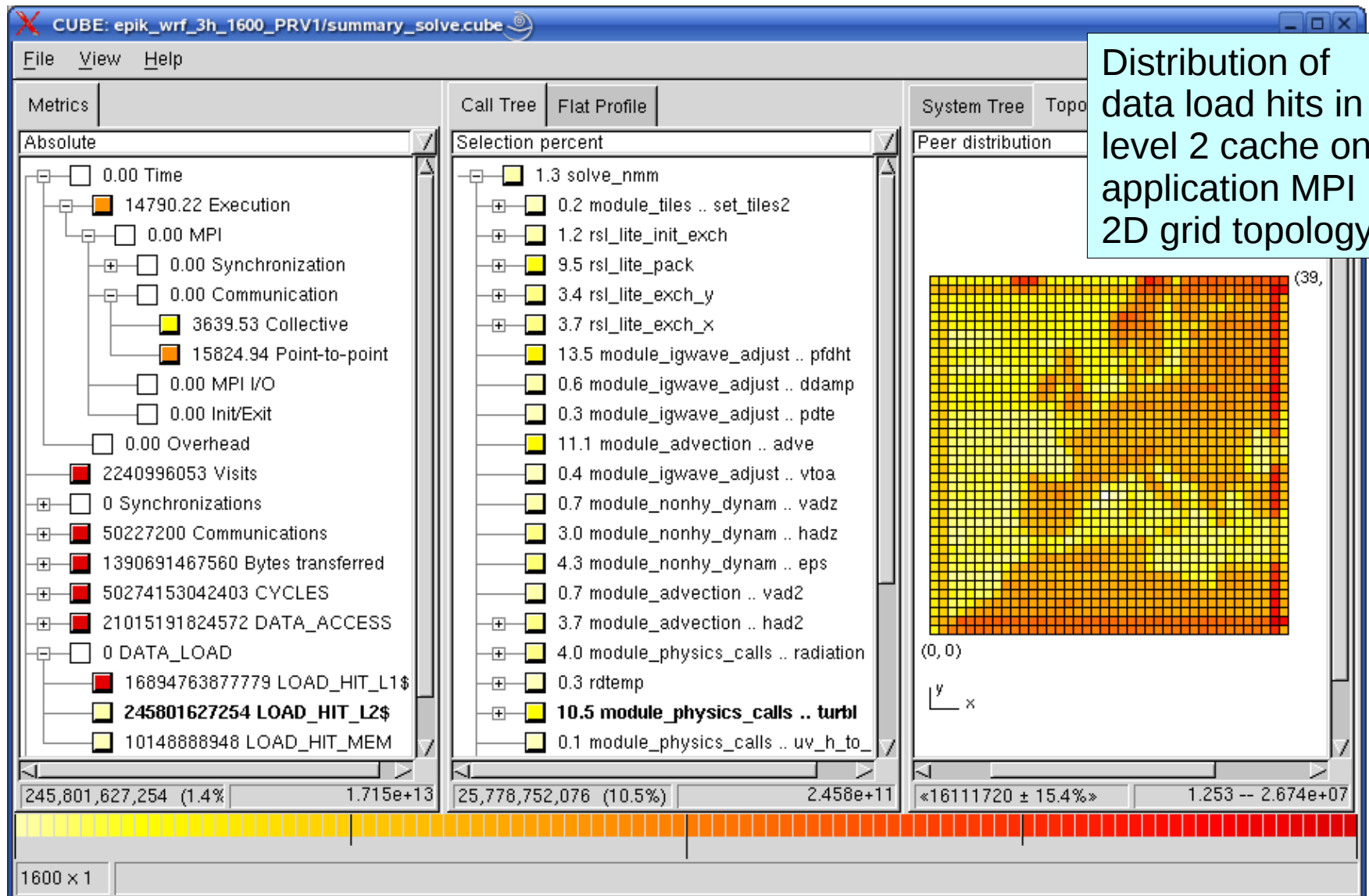
COSMO/XE6 late sender waiting time



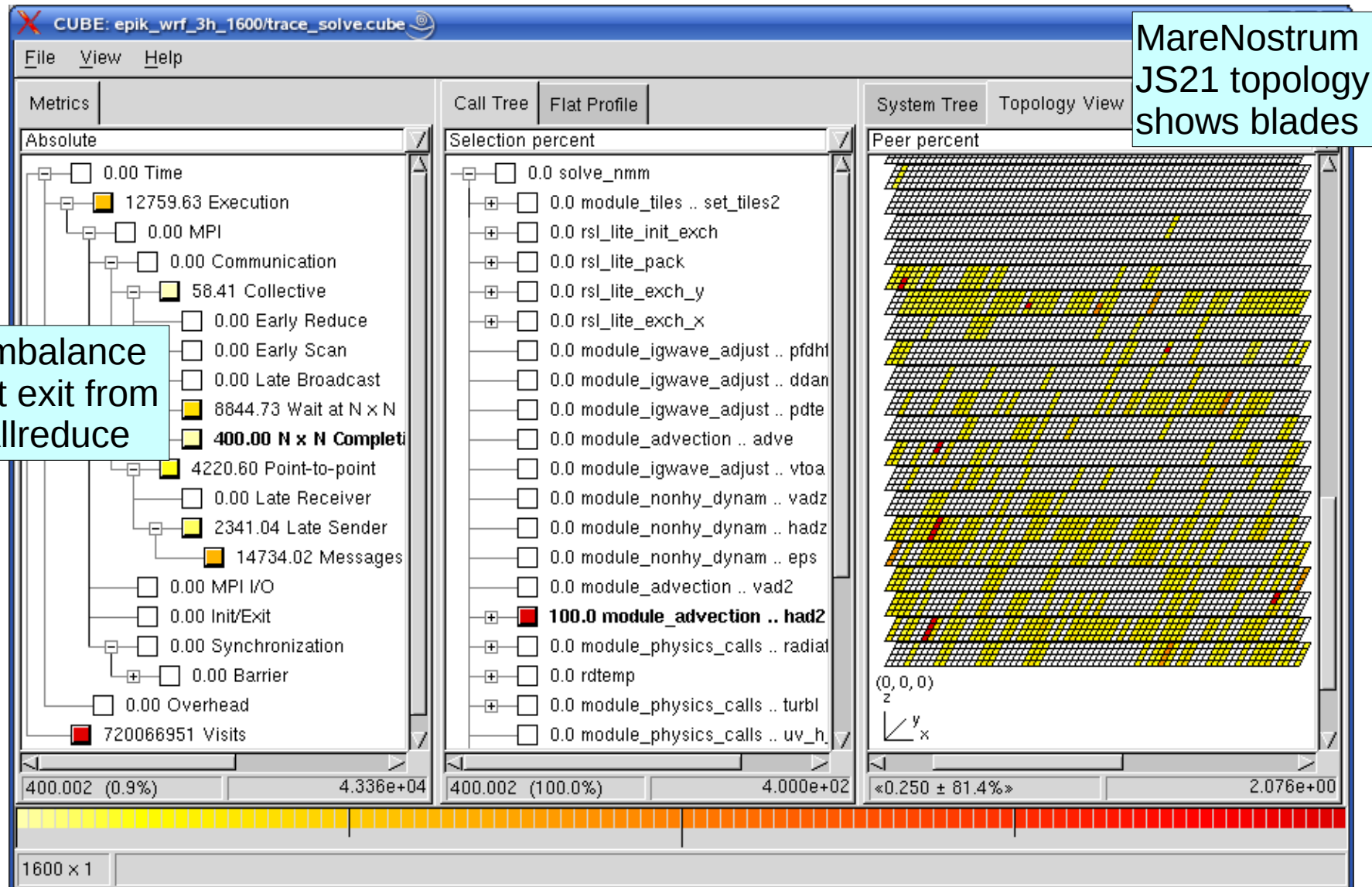


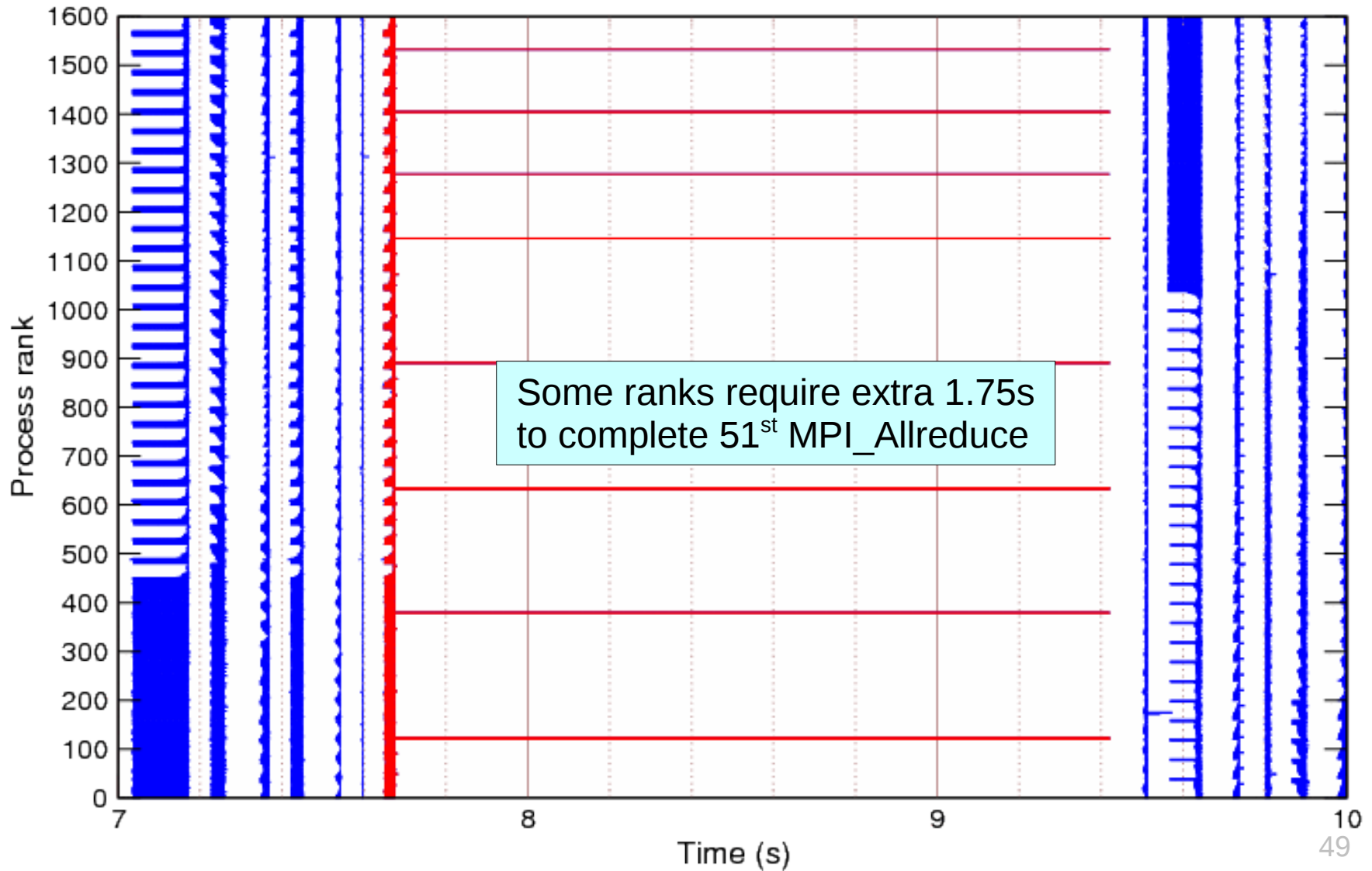
- 56% of total time in local computation
 - 32% in dynamics which is quite well balanced (11% std.dev)
 - 12% in physics is rather less well balanced (17% std.dev)
 - much of the imbalance is inherently physical/geographical
- 44% of total time in MPI
 - 5% collective synchronization (92% output_data)
 - 24% collective communication
 - ▶ 14% for MPI_Gather operations in output_data
 - ▶ 5% “Wait at NxN” mostly in dynamics check_cfl_horiz_advection
 - 15% point-to-point communication (91% exchg_boundaries)
 - ▶ 10% “Late Sender” time (44% dynamics, 36% physics)
 - ▶ 36% of receives are for “Late Senders” (95% in dynamics)
- Communication associated with file I/O was a major factor
 - the 4 dedicated I/O processes idle 95% of the time

- Numerical weather prediction
 - public domain code developed by US NOAA
 - flexible, state-of-the-art atmospheric simulation
 - Non-hydrostatic Mesoscale Model (NMM)
- MPI parallel version 2.1.2 (Jan-2006)
 - >315,000 lines (in 480 source modules): 75% Fortran, 25% C
- Eur-12km dataset configuration
 - 3-hour forecast (360 timesteps) with checkpointing disabled
- Run with 1600 processes on MareNostrum
 - IBM BladeCenter cluster at BSC
- Scalasca summary and trace measurements
 - 15% measurement dilation with 8 hardware counters
 - 23GB trace analysis in 5 mins



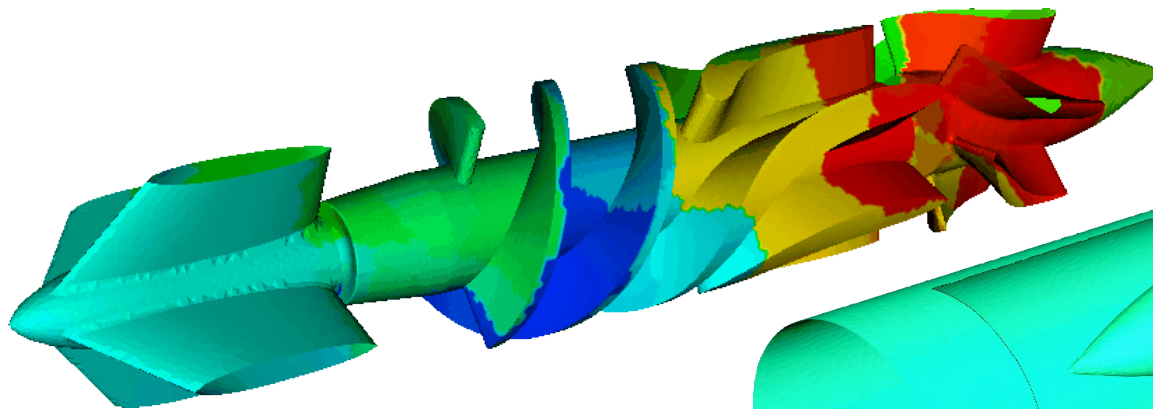
WRF on MareNostrum@1600 trace analysis



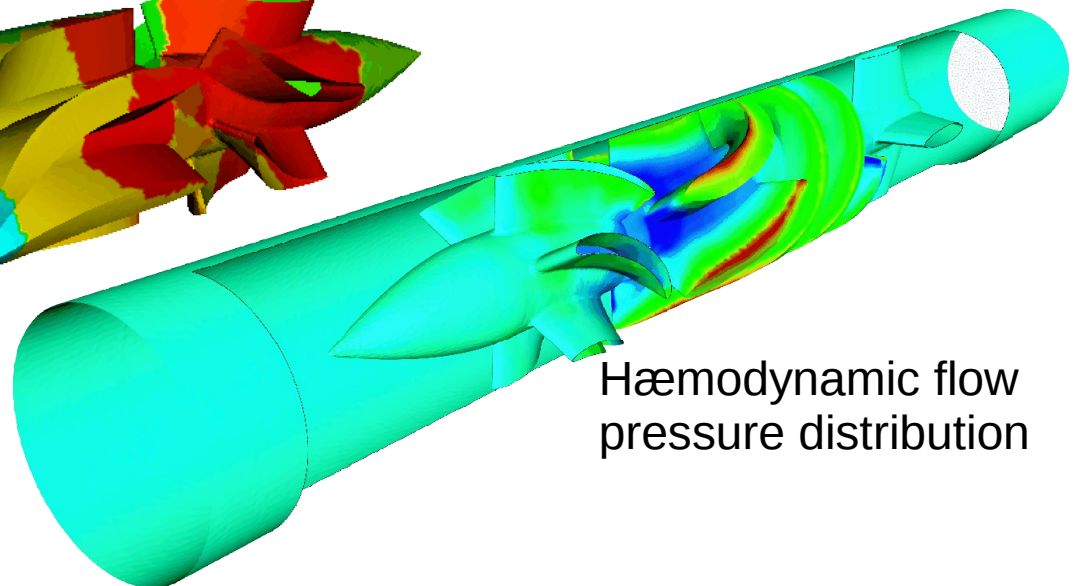


- Limited system I/O requires careful management
 - Selective instrumentation and measurement filtering
- PowerPC hardware counter metrics included in summary
 - Memory/cache data access hierarchy constructed
- Automated trace analysis quantified impact of imbalanced exit from MPI_Allreduce in “NxN completion time” metric
 - Intermittent but serious MPI library/system problem, that restricts application scalability
 - Only a few processes directly impacted, however, communication partners also quickly blocked
- Presentation using logical and physical topologies
 - MPI Cartesian topology provides application insight
 - Hardware topology helps localize system problems

- CFD simulation of unsteady flows
 - developed by RWTH CATS group of Marek Behr
 - exploits finite-element techniques, unstructured 3D meshes, iterative solution strategies
- MPI parallel version (Dec-2006)
 - >40,000 lines of Fortran & C
 - DeBaKey blood-pump dataset (3,714,611 elements)



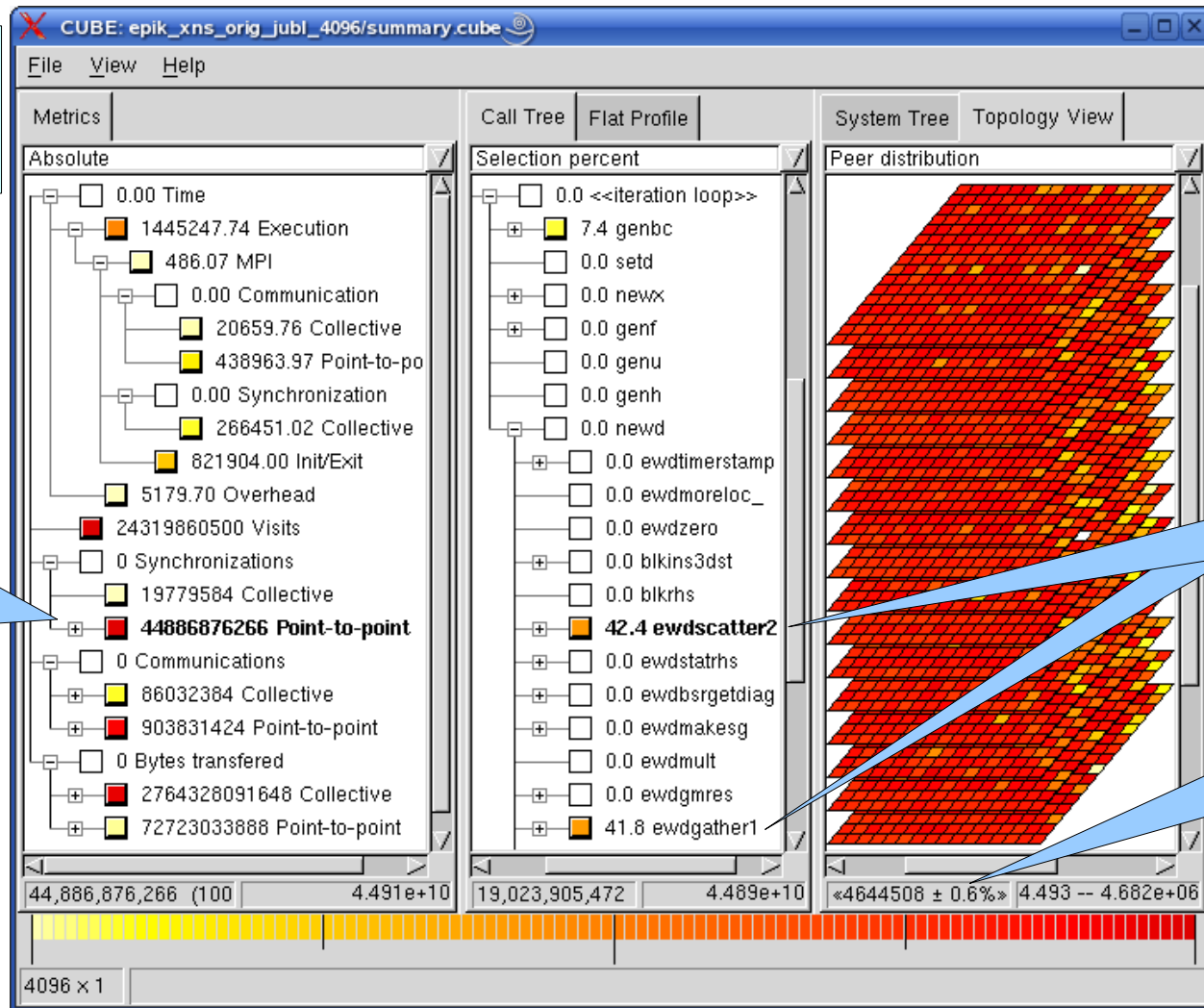
Partitioned finite-element mesh



Hæmodynamic flow
pressure distribution

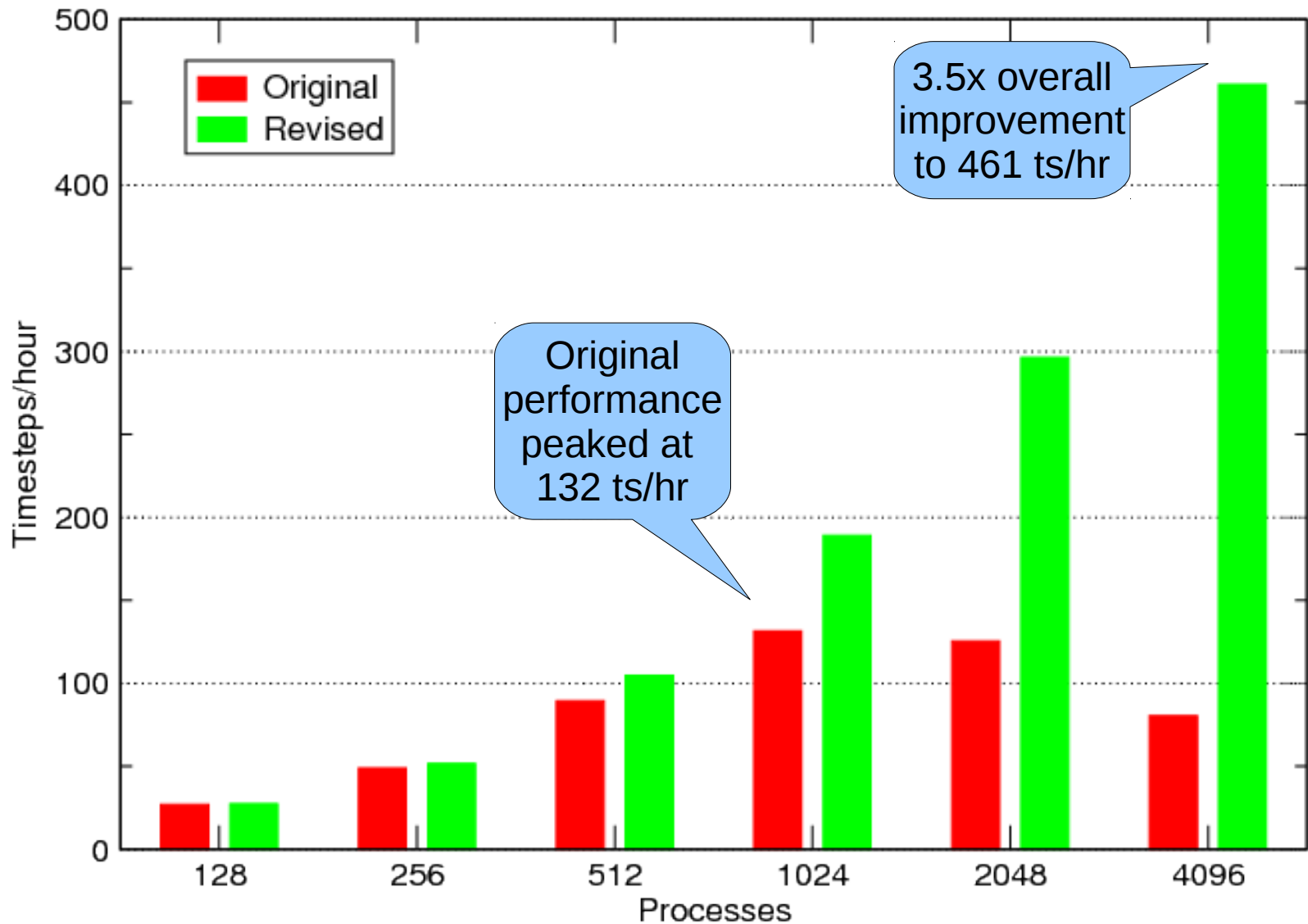
Point-to-point msgs
w/o data

Masses of
P2P synch
operations



Primarily
in scatter
& gather

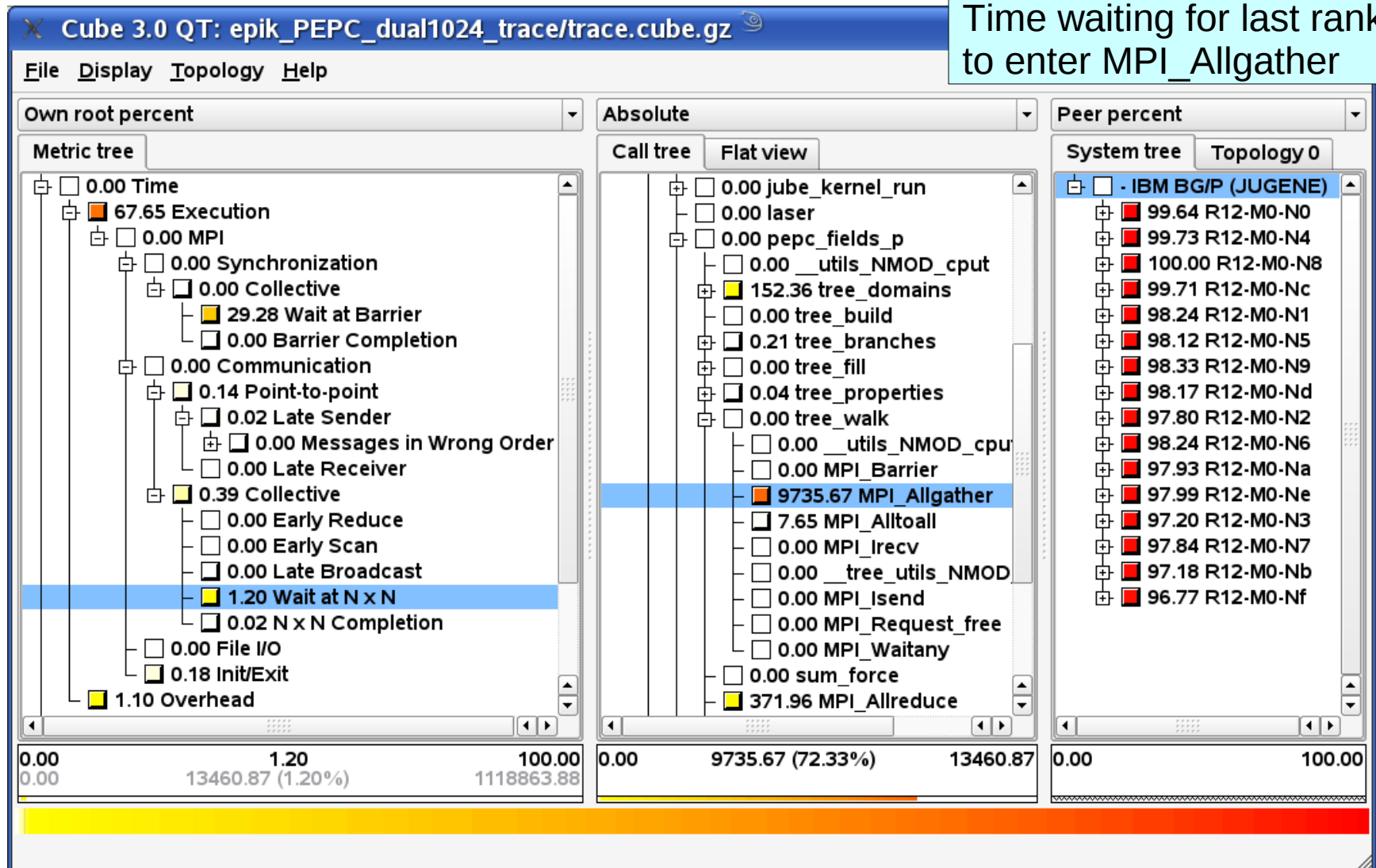
Processes
all equally
responsible



- Globally synchronized high-resolution clock facilitates efficient measurement & analysis
- Restricted compute node memory limits trace buffer size and analyzable trace size
- Summarization identified bottleneck due to unintended P2P synchronizations (messages with zero-sized payload)
- 4x solver speedup after replacing MPI_Sendrecv operations with size-dependant separate MPI_Send and MPI_Recv
- Significant communication imbalance remains due to mesh partitioning and mapping onto processors
- MPI_Scan implementation found to contain implicit barrier
 - responsible for 6% of total time with 4096 processes
 - decimated when substituted with simultaneous binomial tree

- Coulomb solver used for laser-plasma simulations
 - Developed by Paul Gibbon (JSC)
 - Tree-based particle storage with dynamic load-balancing
- MPI version
 - PRACE benchmark configuration, including file I/O
- Run on BlueGene/P in dual mode with 1024 processes
 - 2 processes per quad-core PowerPC node, 1100 seconds
 - IBM XL compilers, MPI library and torus/tree interconnect
- Run on Cray XT in VN (4p) mode with 1024 processes
 - 4 processes per quad-core Opteron node, 360 seconds
 - PGI compilers and Cray MPI, CNL, SeaStar interconnect

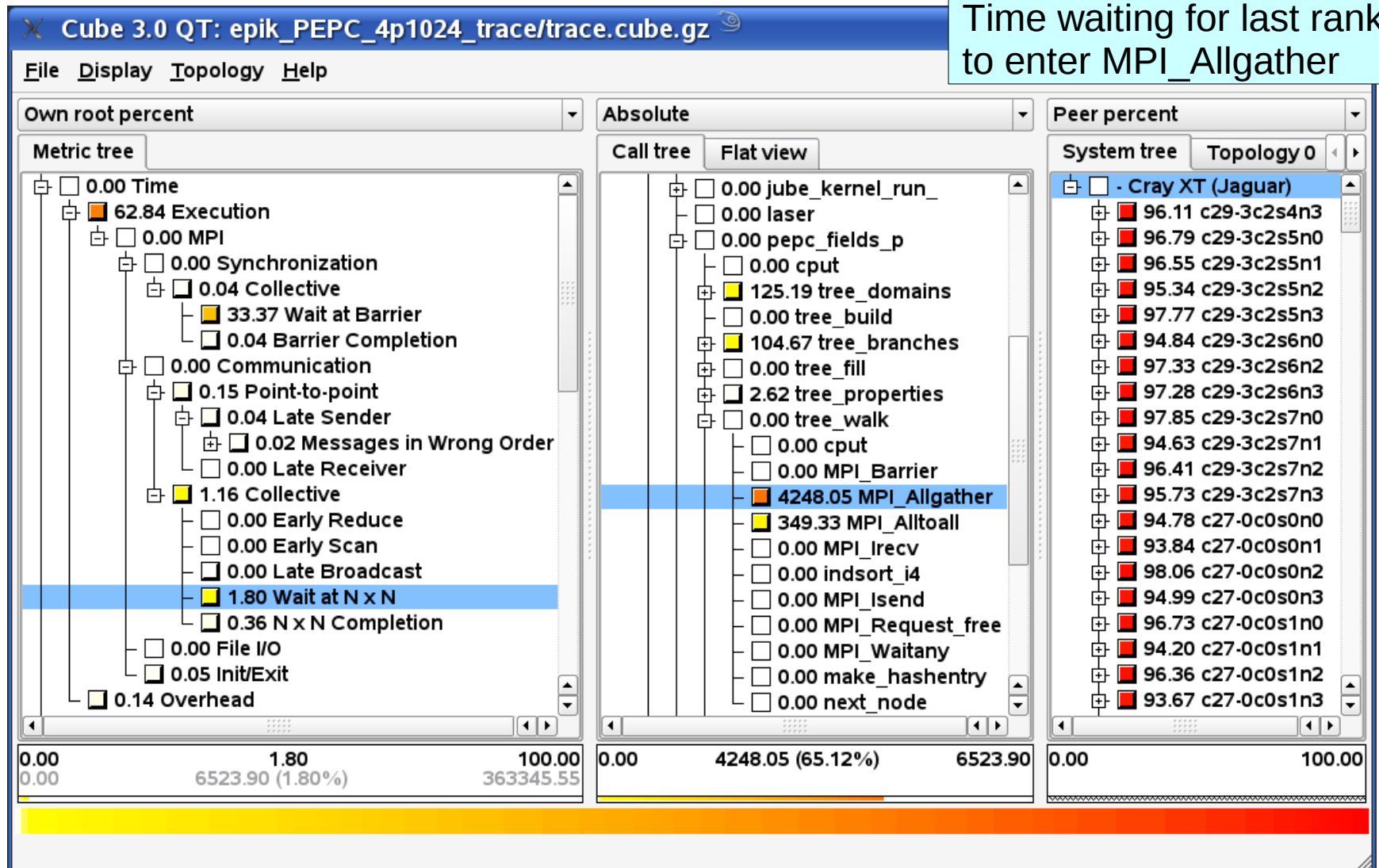
Time waiting for last rank to enter MPI_Allgather



PEPC@1024 on Cray XT4: Wait at NxN time



Time waiting for last rank
to enter MPI_Allgather

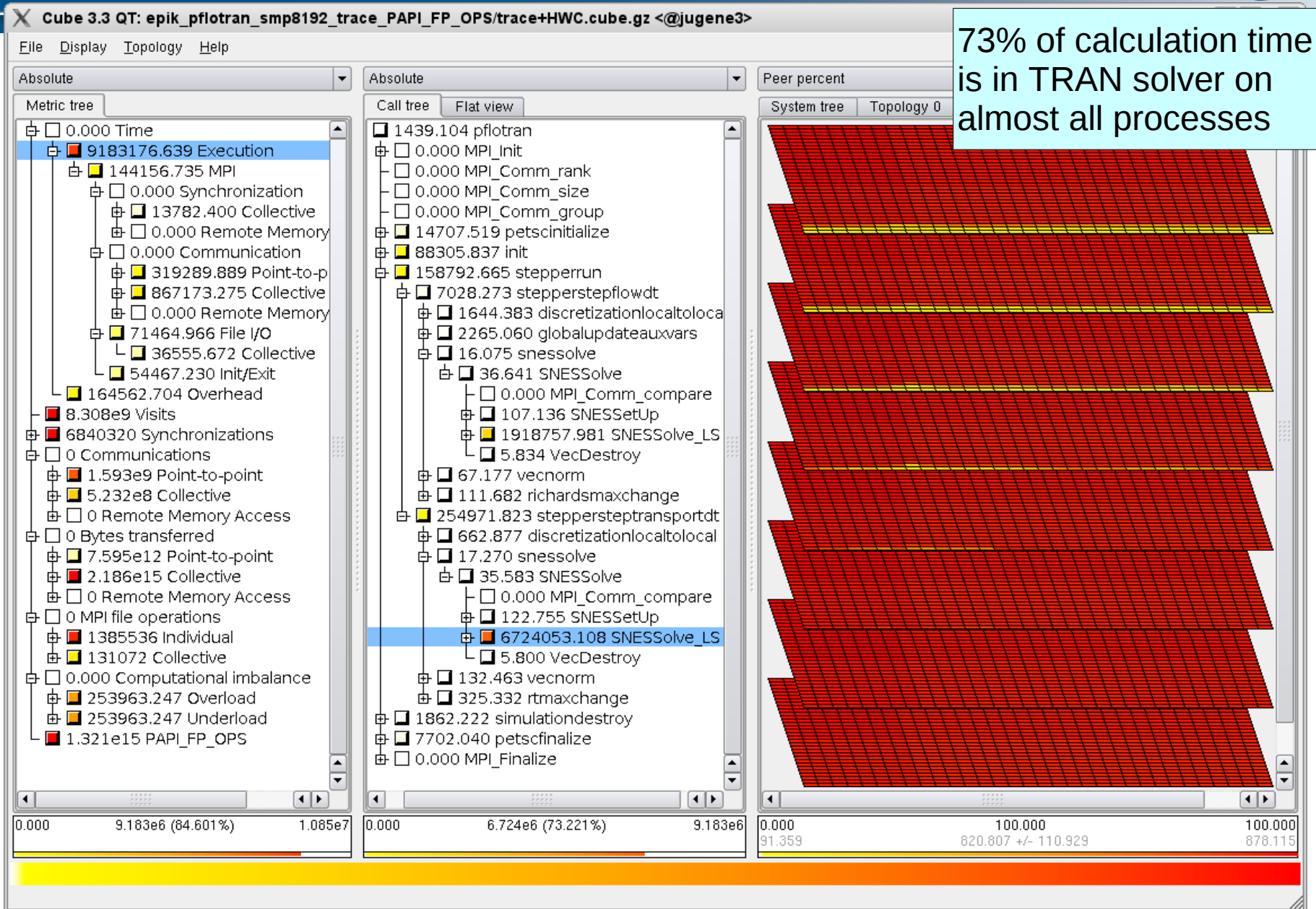


- Despite very different processor and network performance, measurements and analyses can be easily compared
 - different compilers affect function naming & in-lining
- Both spend roughly two-thirds of time in computation
 - tree_walk has expensive computation & communication
- Both waste 30% of time waiting to enter MPI_Barrier
 - not localized to particular processes, since particles are regularly redistributed
- Most of collective communication time is also time waiting for last ranks to enter MPI_Allgather & MPI_Alltoall
 - imbalance for MPI_Allgather twice as severe on BlueGene/P, however, almost 50x less for MPI_Alltoall
 - collective completion times also notably longer on Cray XT

- 3D reservoir simulator combining alternating
 - PFLOW non-isothermal, multiphase groundwater flow
 - PTRAN reactive, multi-component contaminant transport
 - developed by LANL/ORNL/PNNL
- MPI with PETSc, LAPACK, BLAS & HDF5 I/O libraries
 - ~80,000 lines (97 source files) Fortran9X
 - PFLOTRAN & PETSc fully instrumented by IBM XL compilers
 - ▶ filter produced listing 856 USR routines (leaving 291 COM)
 - ▶ 1732 unique callpaths (399 in FLOW, 375 in TRAN)
 - ▶ 633 MPI callpaths (121 in FLOW, 114 in TRAN)
 - 29 distinct MPI routines recorded (excludes 15 misc. routines)
- Run on IBM BlueGene/P with '2B' input dataset (10 steps)
 - Scalasca summary & trace measurements (some with PAPI)
 - 22% dilation of FLOW, 10% dilation of TRAN [8k summary]

PFLOTRAN jugene@smp8192 trace analysis

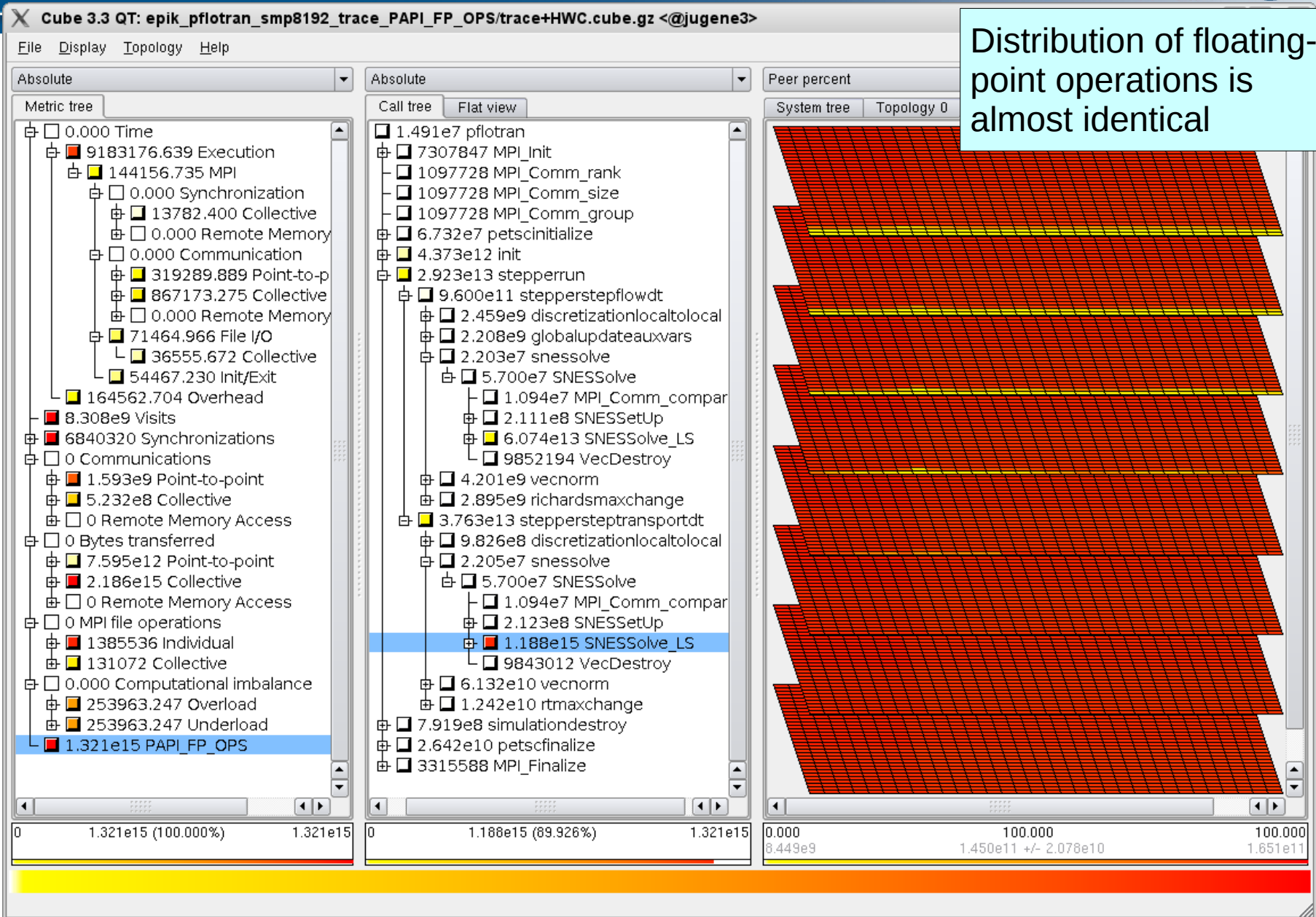
VI-HPS



PFLOTRAN jugene@smp8192 trace analysis

VI-HPS

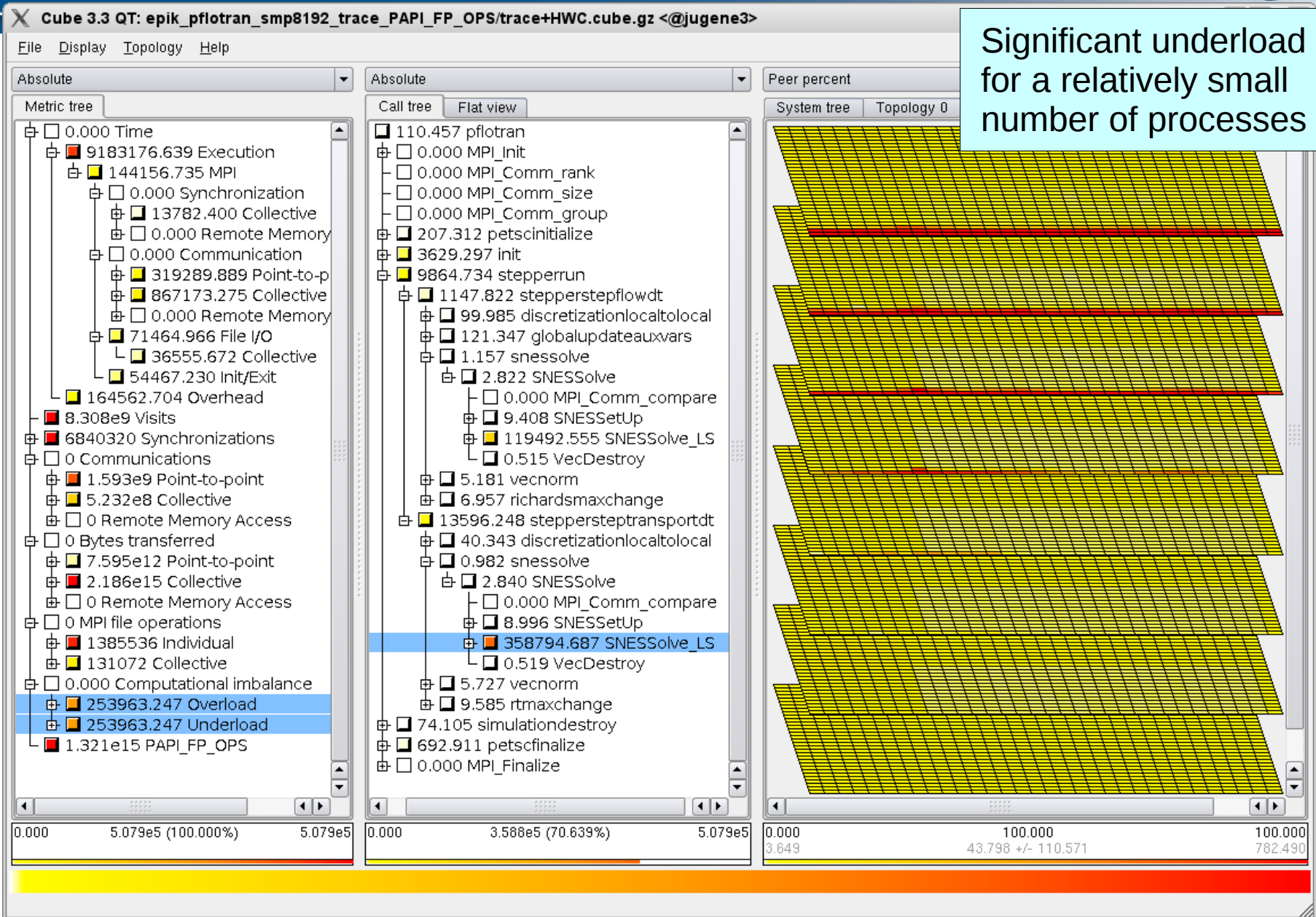
Distribution of floating-point operations is almost identical



PFLOTRAN jugene@smp8192 trace analysis



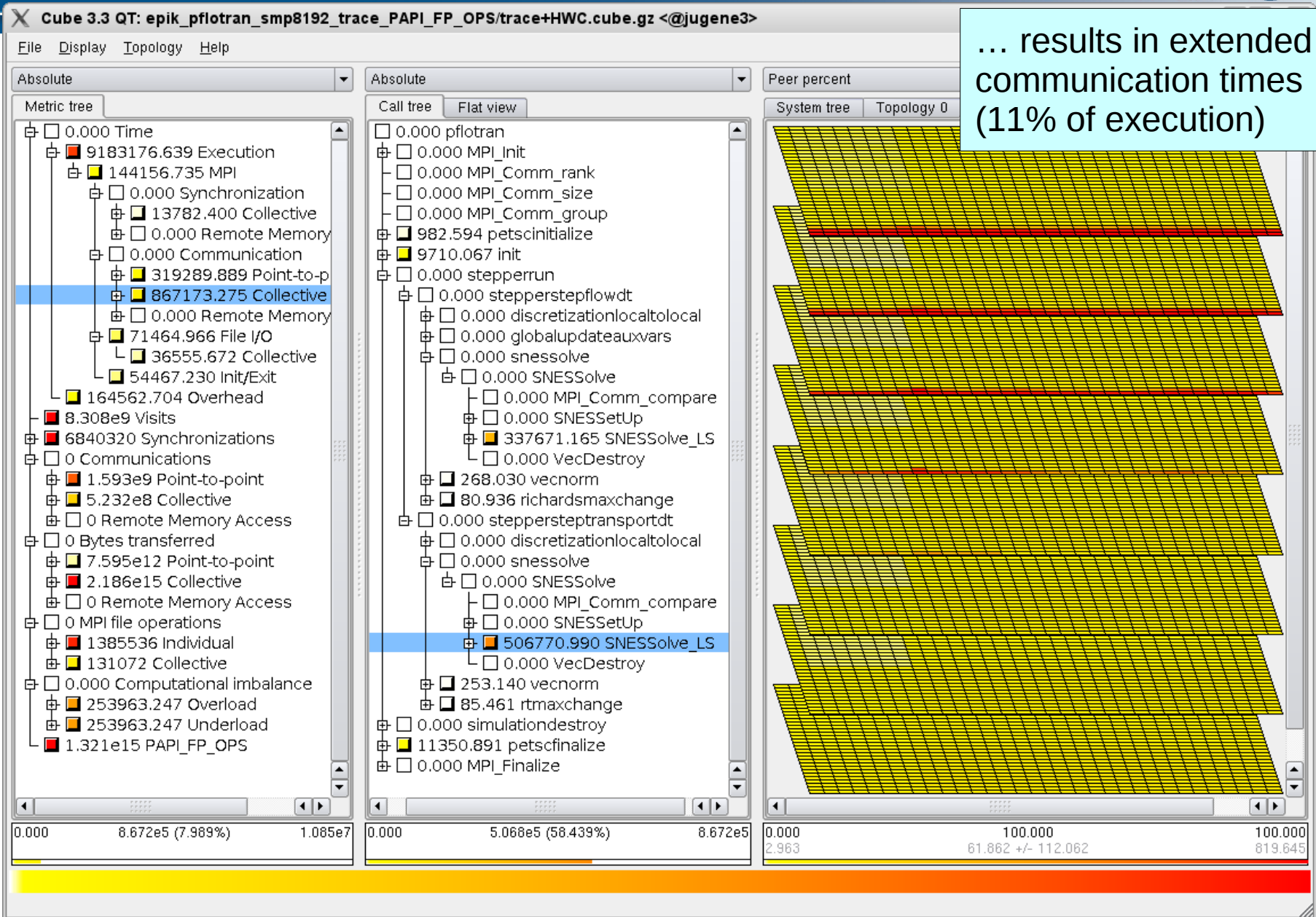
Significant underload
for a relatively small
number of processes



PFLOTRAN jugene@smp8192 trace analysis

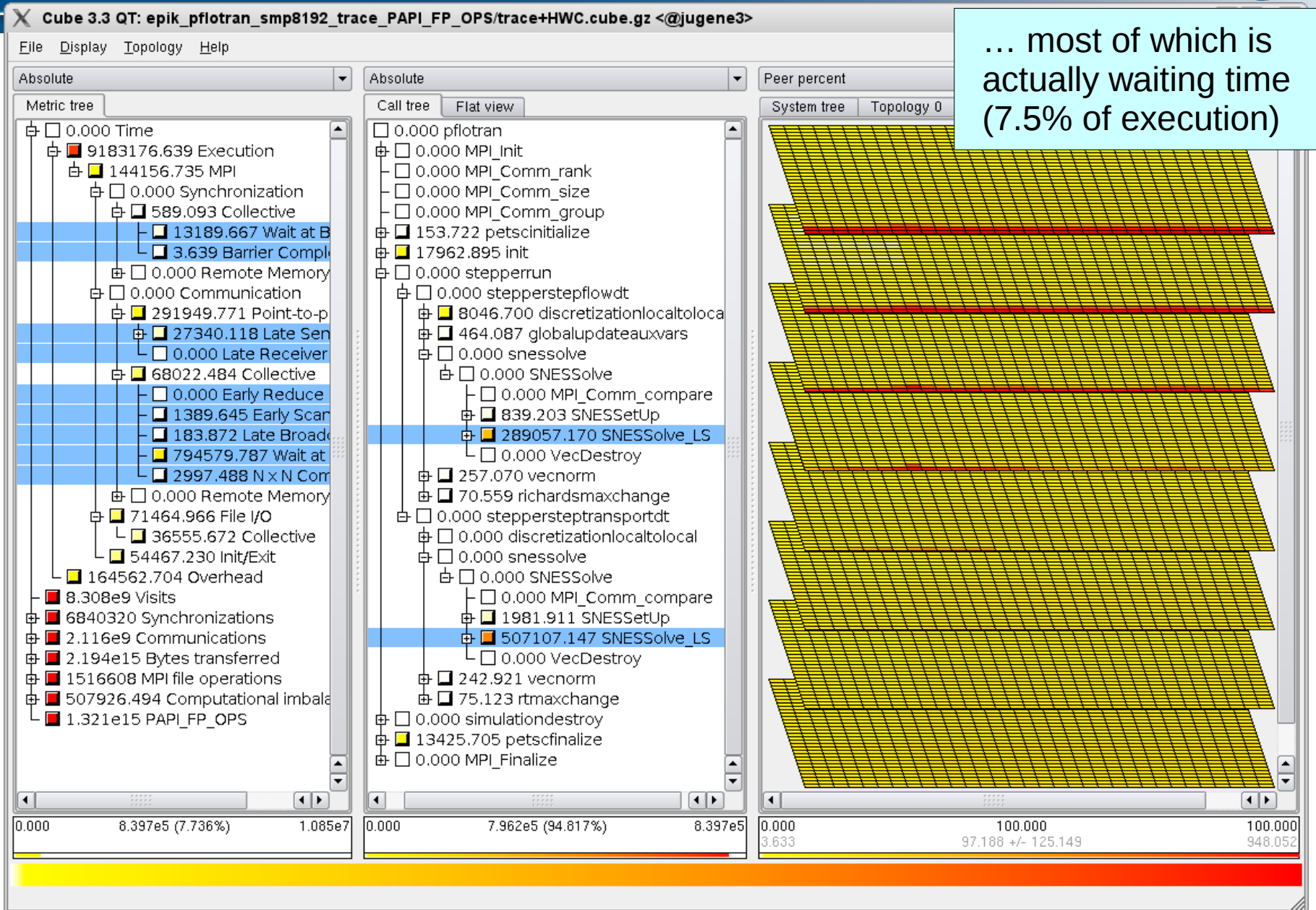


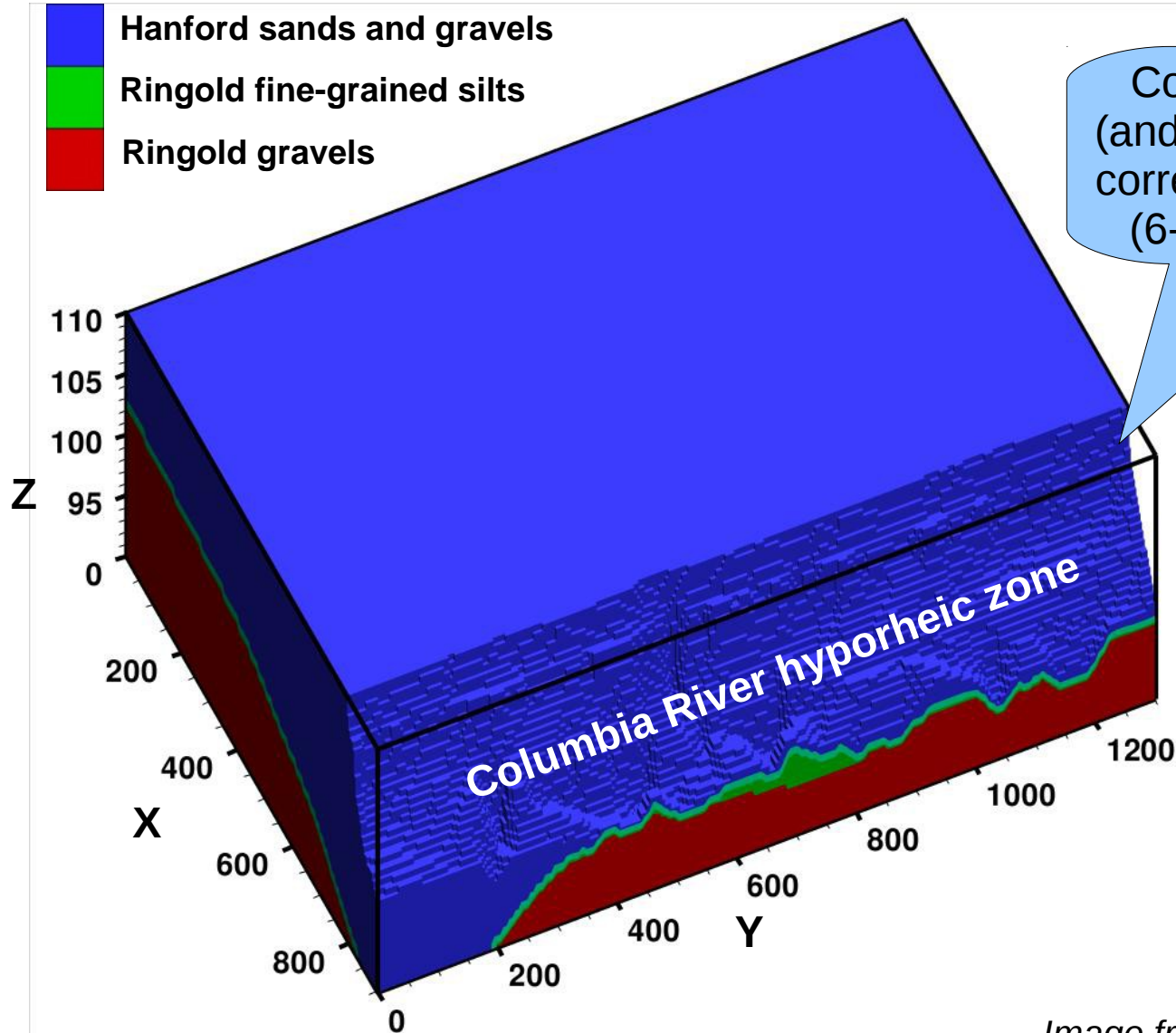
... results in extended communication times (11% of execution)



PFLOTRAN jugene@smp8192 trace analysis

VI-HPS

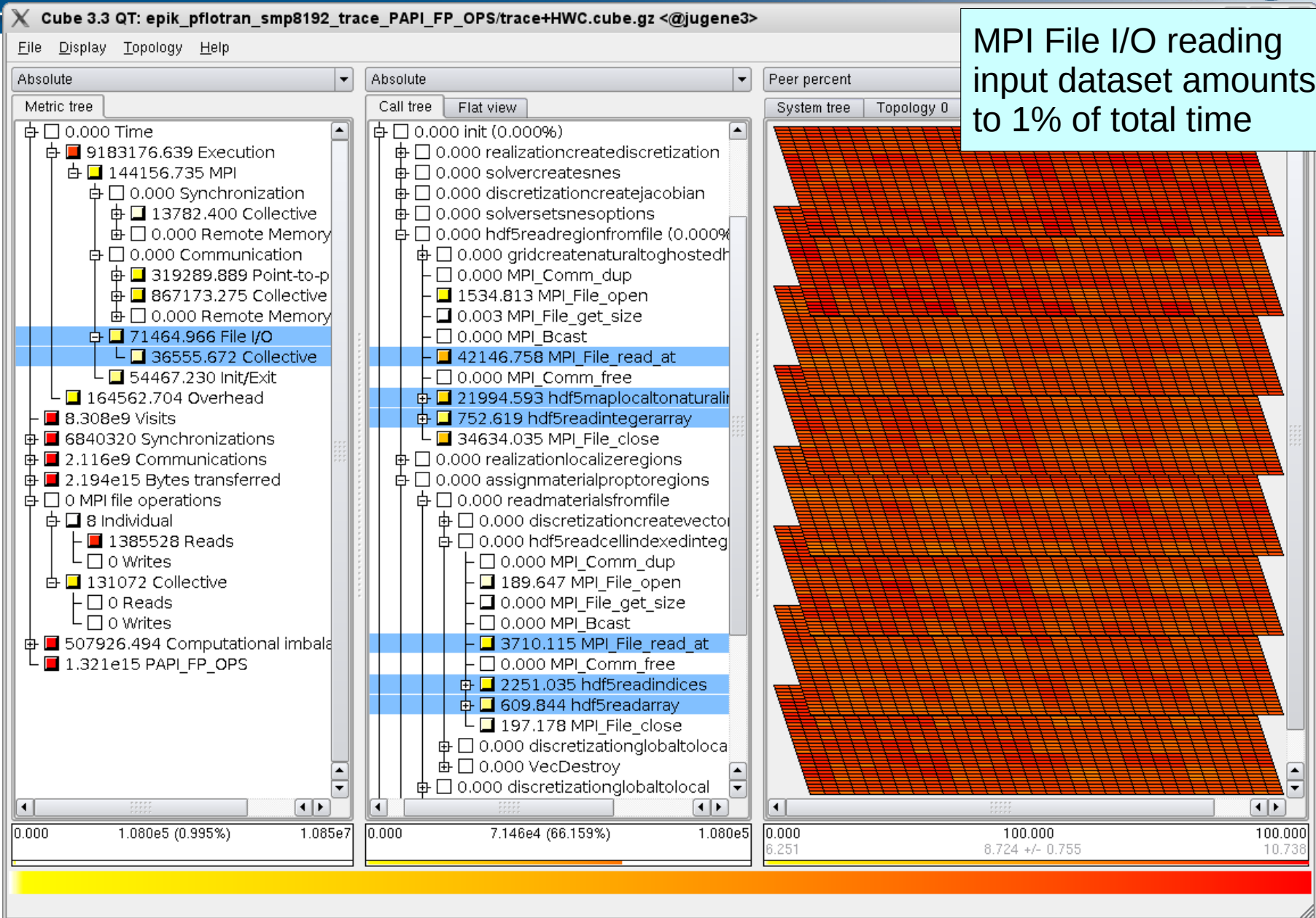




PFLOTRAN jugene@smp8192 trace analysis



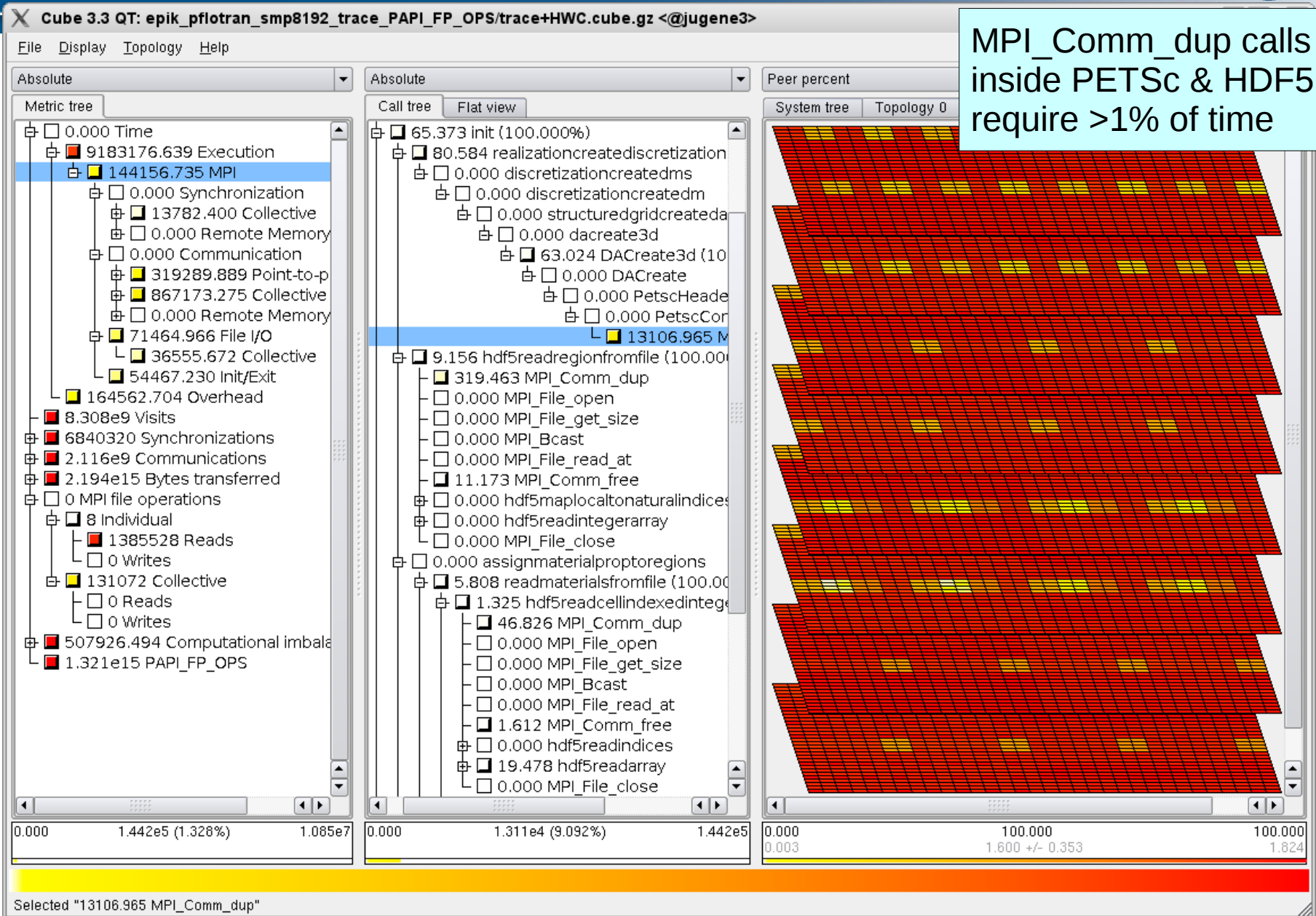
MPI File I/O reading
input dataset amounts
to 1% of total time



PFLOTRAN jugene@smp8192 trace analysis

VI-HPS

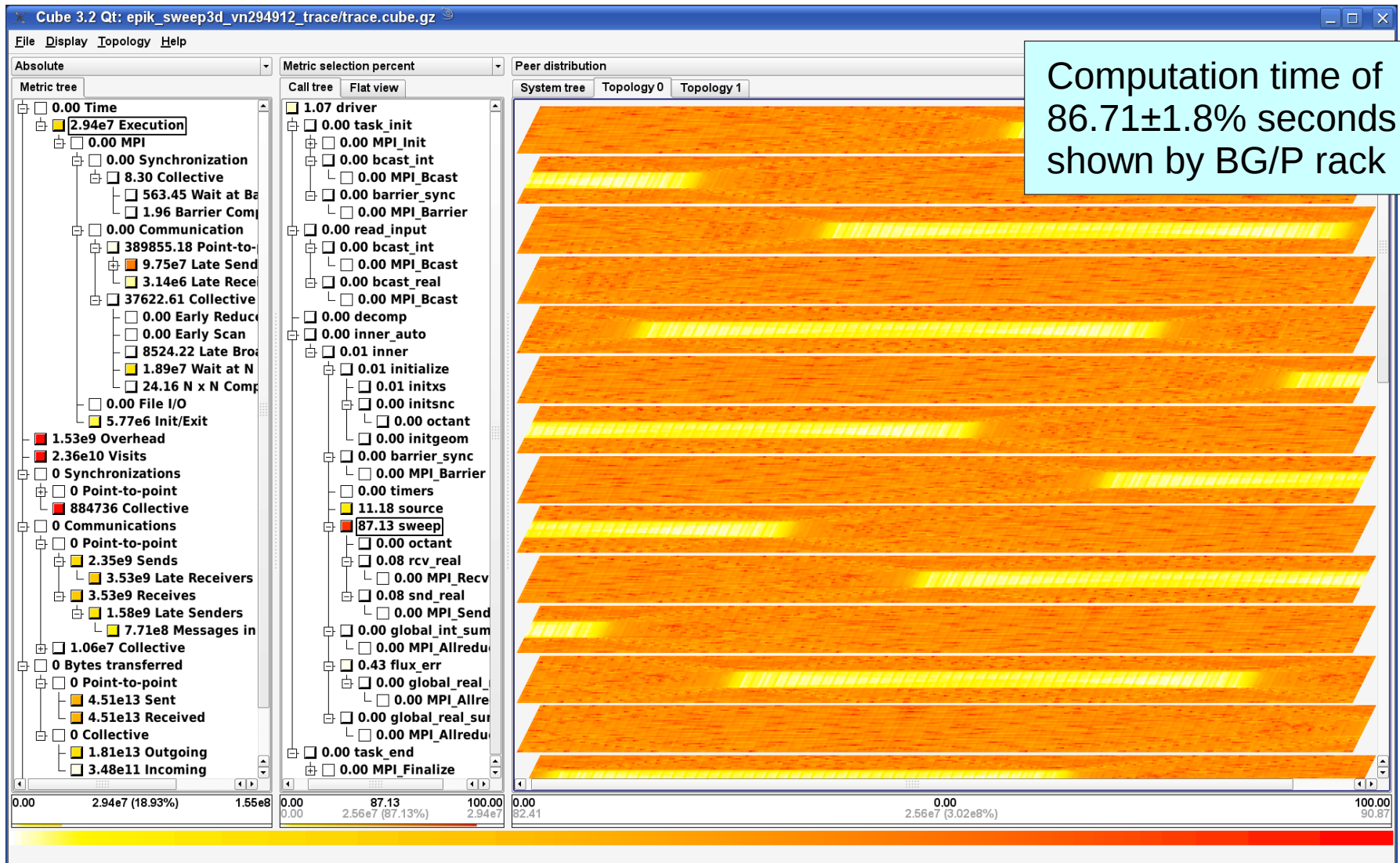
MPI_Comm_dup calls
inside PETSc & HDF5
require >1% of time



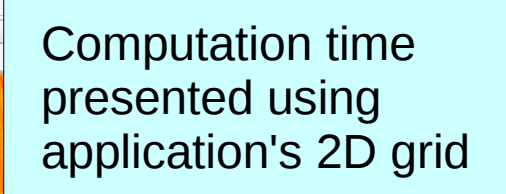
- Initialization phase dominates at larger scales
 - 10% of total execution time spent duplicating communicators with 128k processes on Cray XT5
 - otherwise collective MPI File I/O relatively efficient
 - typically amortized in long simulation runs
- Solver scaled well to 64k processes before degrading
 - similar computation/communication patterns in FLOW & TRAN
 - ▶ callpath profiles distinguish costs
 - ▶ MPI_Allreduce collective communication becomes a bottleneck
 - ▶ communication overhead explodes for smaller FLOW problem
 - TRAN problem is 15x larger due to 15 chemical species
 - inactive processes induce clear computational imbalance
 - ▶ and are associated with large amounts of MPI waiting time
 - ▶ however, they constitute a relatively small minority

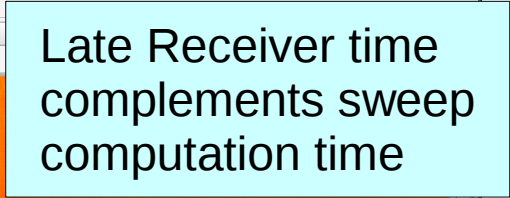
- 3D neutron transport simulation
 - ASC benchmark
 - direct order solve uses diagonal sweeps through grid cells
 - 'fixups' applied to correct unphysical (negative) fluxes
- MPI parallel version 2.2b using 2D domain decomposition
 - ~2,000 lines (12 source modules), mostly Fortran77
- Run on IBM BlueGene/P in VN mode with 288k processes
 - 7.6TB trace written in 17 minutes, analyzed in 10 minutes
 - ▶ of which 10 minutes for SIONlib open/create of 576 physical files
 - ▶ (compared to 86 minutes just to open/create a file per MPI rank)
 - Mapping of metrics onto application's 576x512 process grid reveals regular pattern of performance artifacts
 - ▶ computational imbalance originates from 'fixup' calculations
 - ▶ combined with diagonal wavefront sweeps amplifies waiting times

sweep3d on jugene@288k trace analysis



VI-HPS

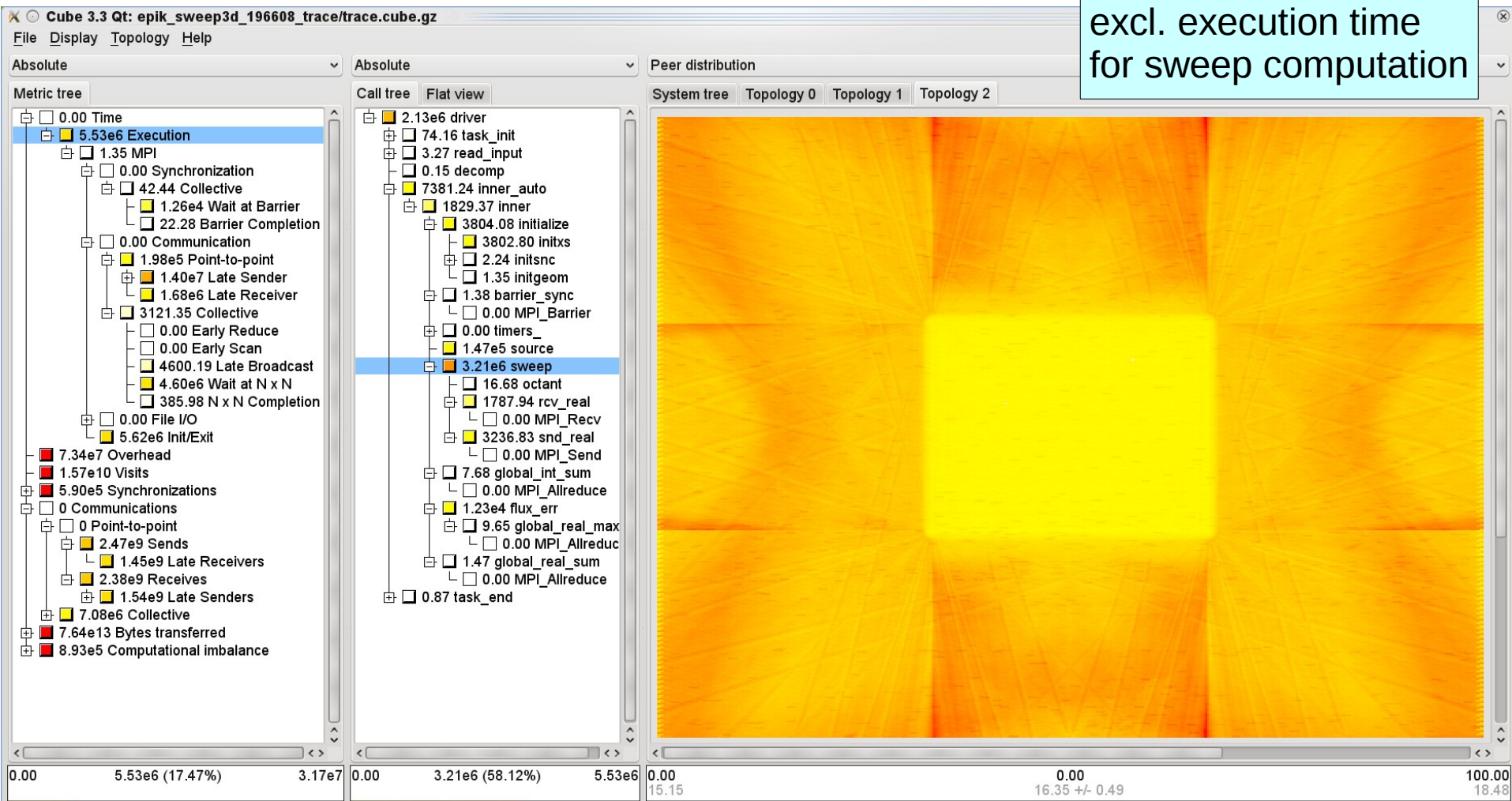




- 3D neutron transport simulation
 - ASC benchmark
 - direct order solve uses diagonal sweeps through grid cells
 - 'fixups' applied to correct unphysical (negative) fluxes
- MPI parallel version 2.2b using 2D domain decomposition
 - ~2,000 lines (12 source modules), mostly Fortran77
- Run on Cray XT5 with 192k processes
 - 0.5TB trace written in 10 minutes, analyzed in 4 minutes
 - ▶ 6 minutes to open/create trace file for each rank
 - ▶ 25s for timestamp correction, 93s for parallel event replay
 - Mapping of metrics onto application's 512x384 process grid reveals regular pattern of performance artifacts
 - ▶ computational imbalance originates from 'fixup' calculations
 - ▶ combined with diagonal wavefront sweeps amplifies waiting times

sweep3d on jaguar@192k trace analysis

Regular imbalance in
excl. execution time
for sweep computation



- 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, CINECA, CMM, CSC, CSCS, DKRZ, EPCC, HLRN, HLRS, ICM, IMAG, JSC, KSL, KTH, LLNL, LRZ, NCAR, NCCS, NICS, RWTH, RZG, SARA, TACC, TGCC, ZIH
 - ▶ Access & usage supported by European Union, German and other national funding organizations
- The Scalasca users for their comprehensive problem reports and improvement requests
 - as well as sharing reports of their analysis & tuning successes
- The Scalasca development team

Scalable performance analysis of **large-scale** parallel applications

- toolset for scalable performance measurement & analysis of MPI, OpenMP & hybrid parallel applications
- supporting most popular HPC computer systems
- available under New BSD open-source license
- sources, documentation & publications:
 - ▶ <http://www.scalasca.org>
 - ▶ [mailto: scalasca@fz-juelich.de](mailto:scalasca@fz-juelich.de)