

scalasca

Performance analysis & tuning case studies

Brian Wylie & Markus Geimer
Jülich Supercomputing Centre
scalasca@fz-juelich.de
October 2010

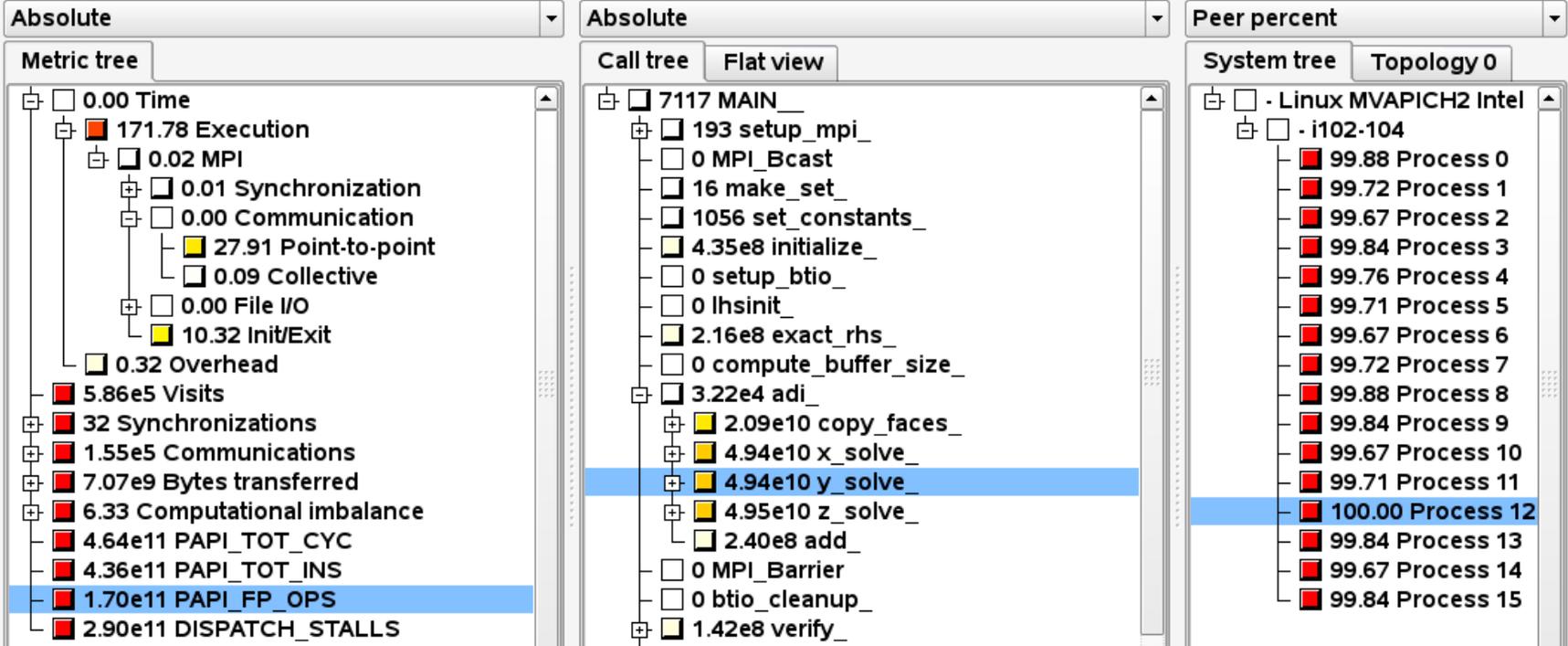
- Example experiment archives provided for examination:
 - jugene_sweep3d
 - ▶ 294,912 & 65,536 MPI processes on BG/P (trace)
 - jump_zeusmp2
 - ▶ 512 MPI processes on p690 cluster (summary & trace)
 - marenostrom_wrf-nmm
 - ▶ 1600 MPI processes on JS21 blade cluster, solver extract
 - ▶ summary analysis 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

16-process summary analysis with HWC metrics

Cube 3.0 QT: epik_bt_A_16_sum_PAPI/summary.cube.gz

EPK_METRICS = PAPI_TOT_CYC:PAPI_TOT_INS:PAPI_FP_OPS:DISPATCH_STALLS



Metric info

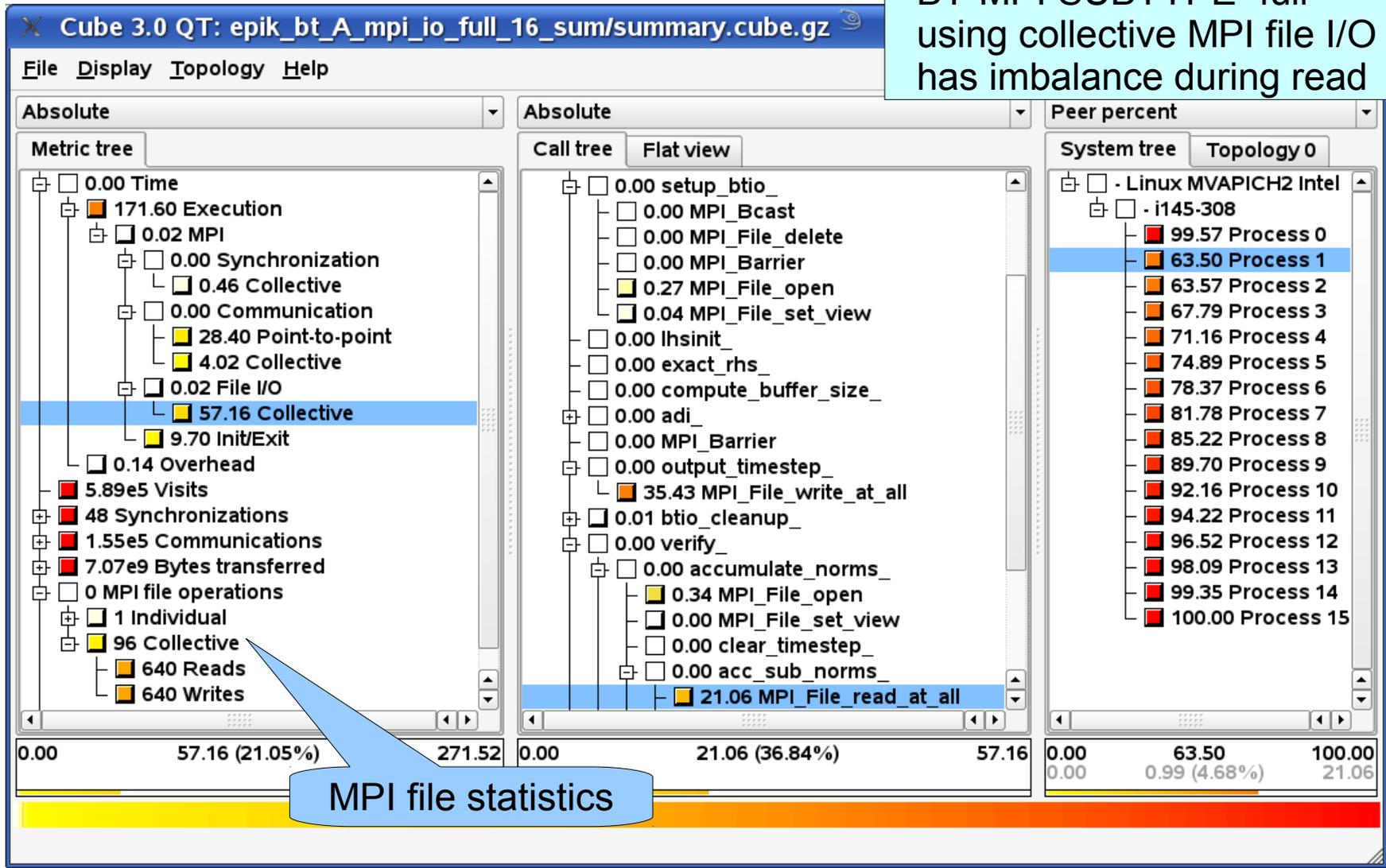
Floating point operations.
 [RETIRED_SSE_OPERATIONS:SINGLE_ADD_SUB_OPS:SINGLE_MUL_OPS:DOUBLE_ADD_SUB_OPS:DOUBLE_MUL_O

OK

100.00	100.00
3.10e9 (6.26%)	4.94e10

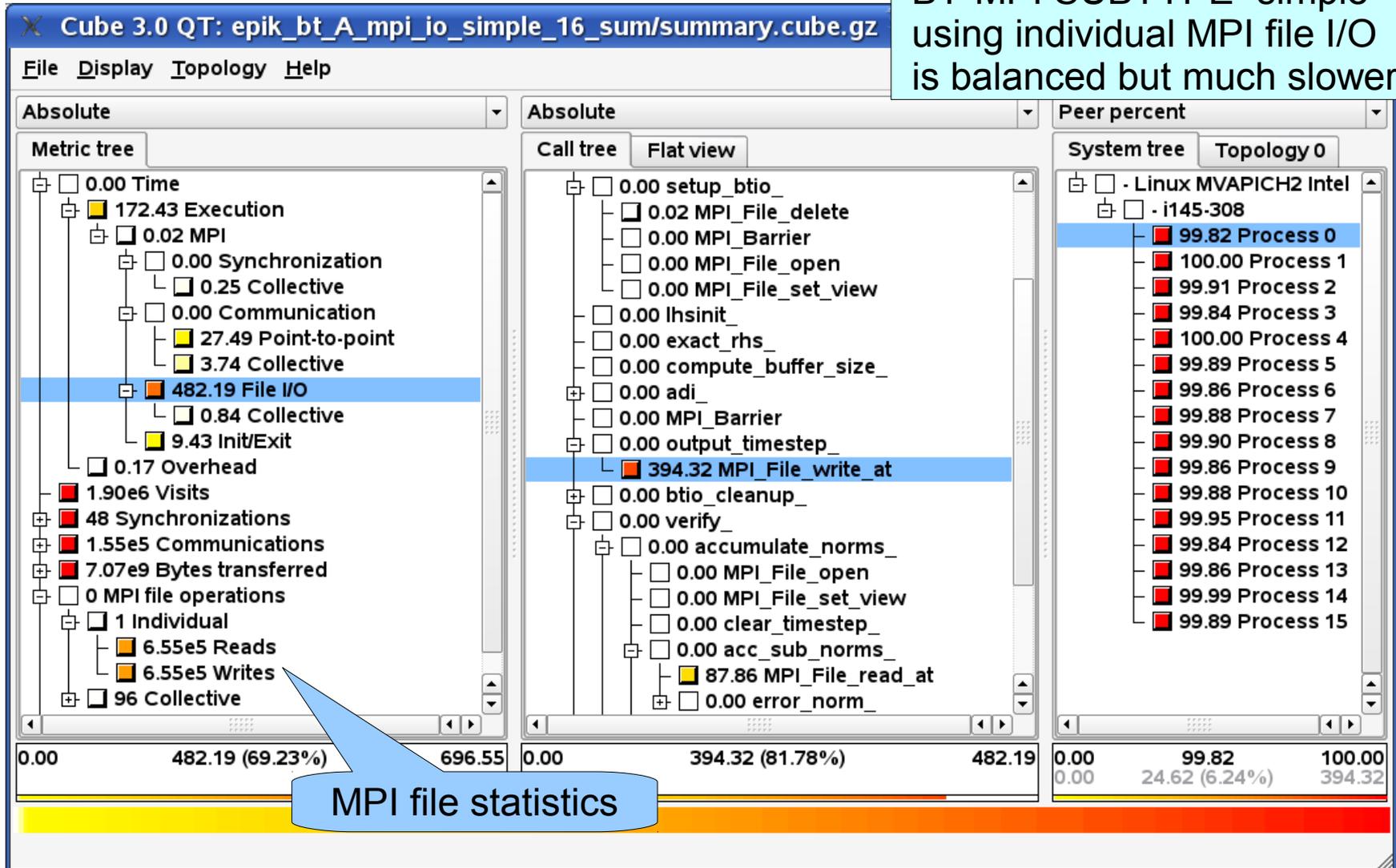
16-process summary analysis: MPI File I/O time VI-HPS

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

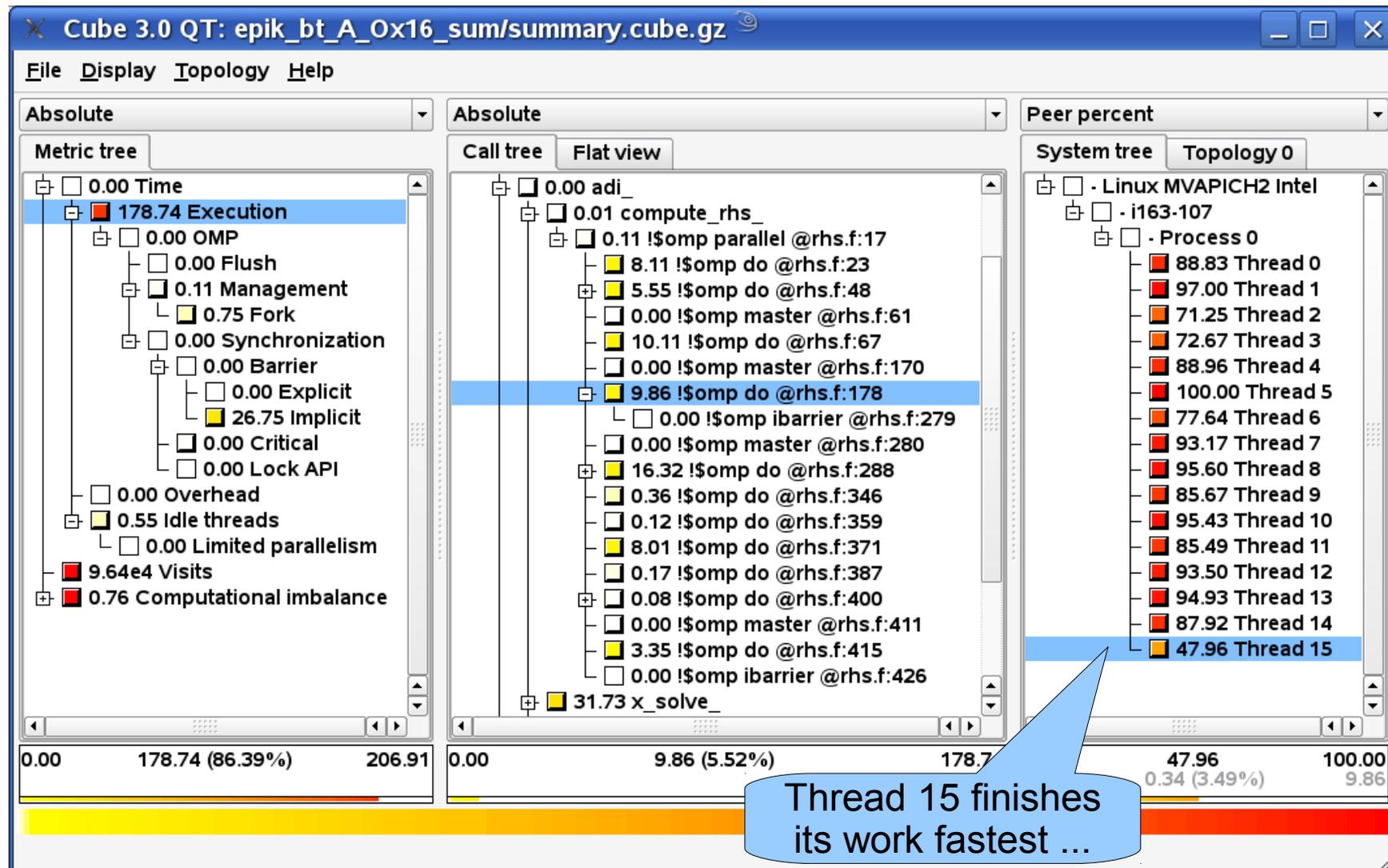


16-process summary analysis: MPI File I/O time

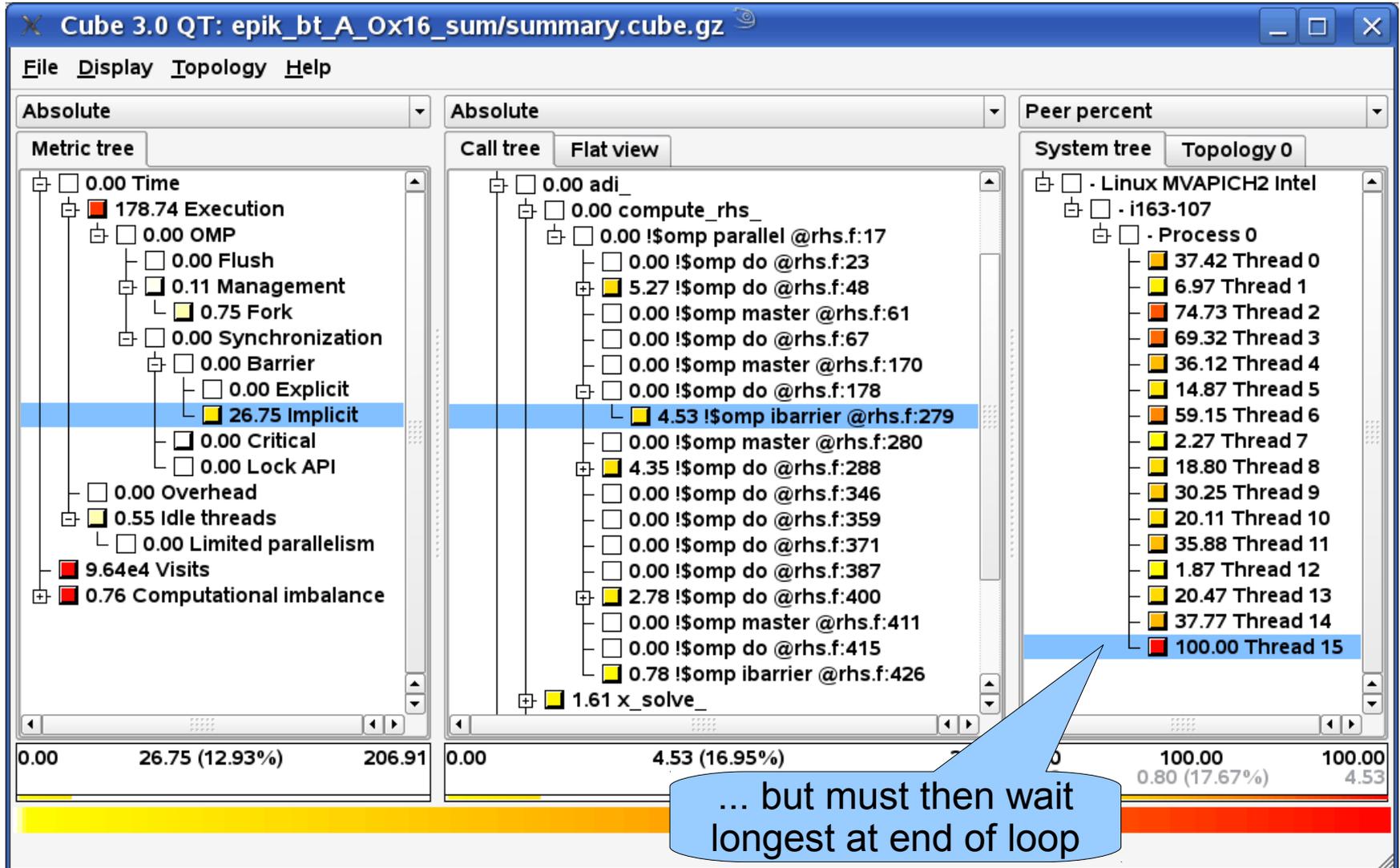
BT-MPI SUBTYPE=simple using individual MPI file I/O is balanced but much slower



16-thread summary analysis: Execution time

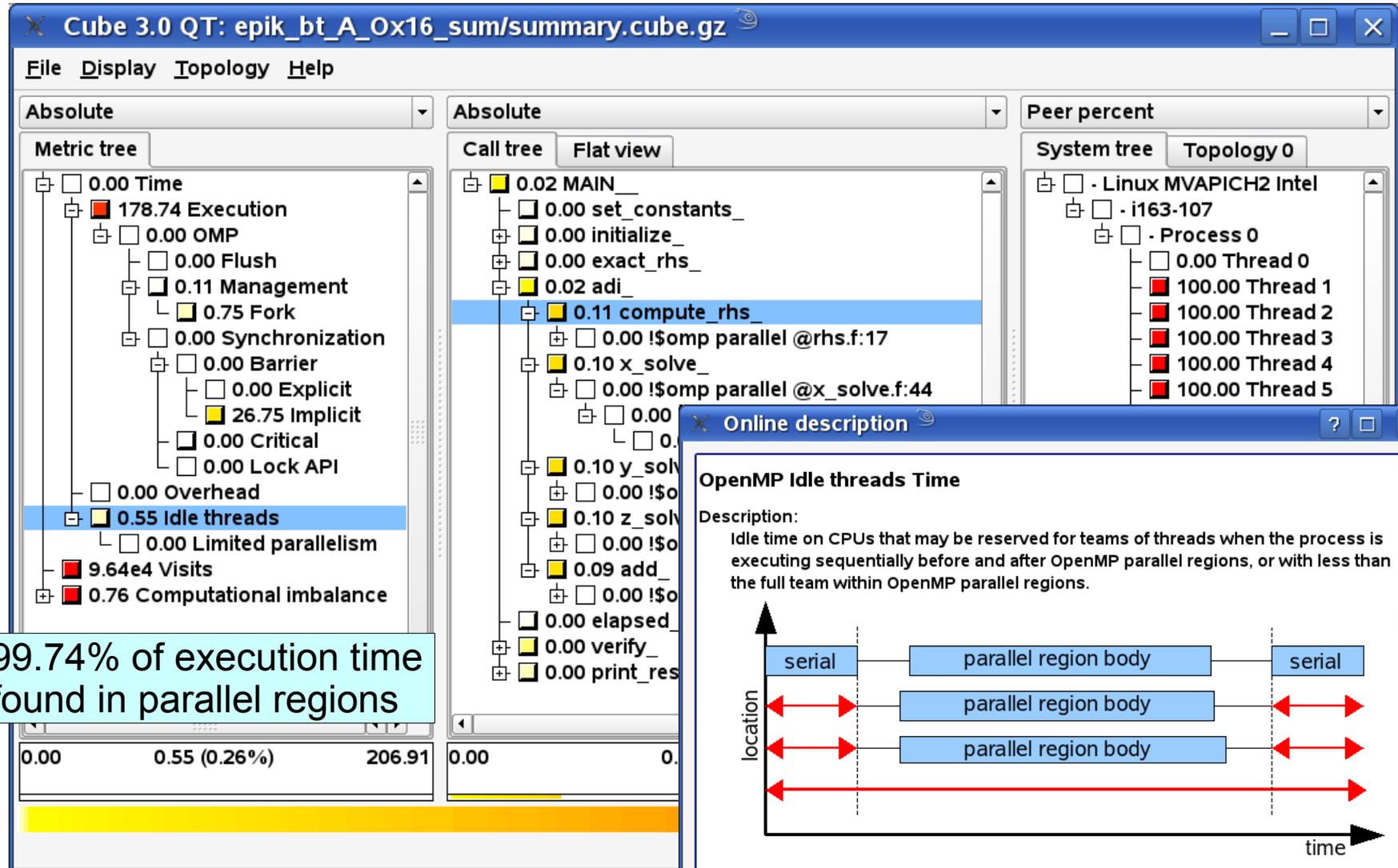


16-thread summary analysis: Implicit barrier time



... but must then wait longest at end of loop

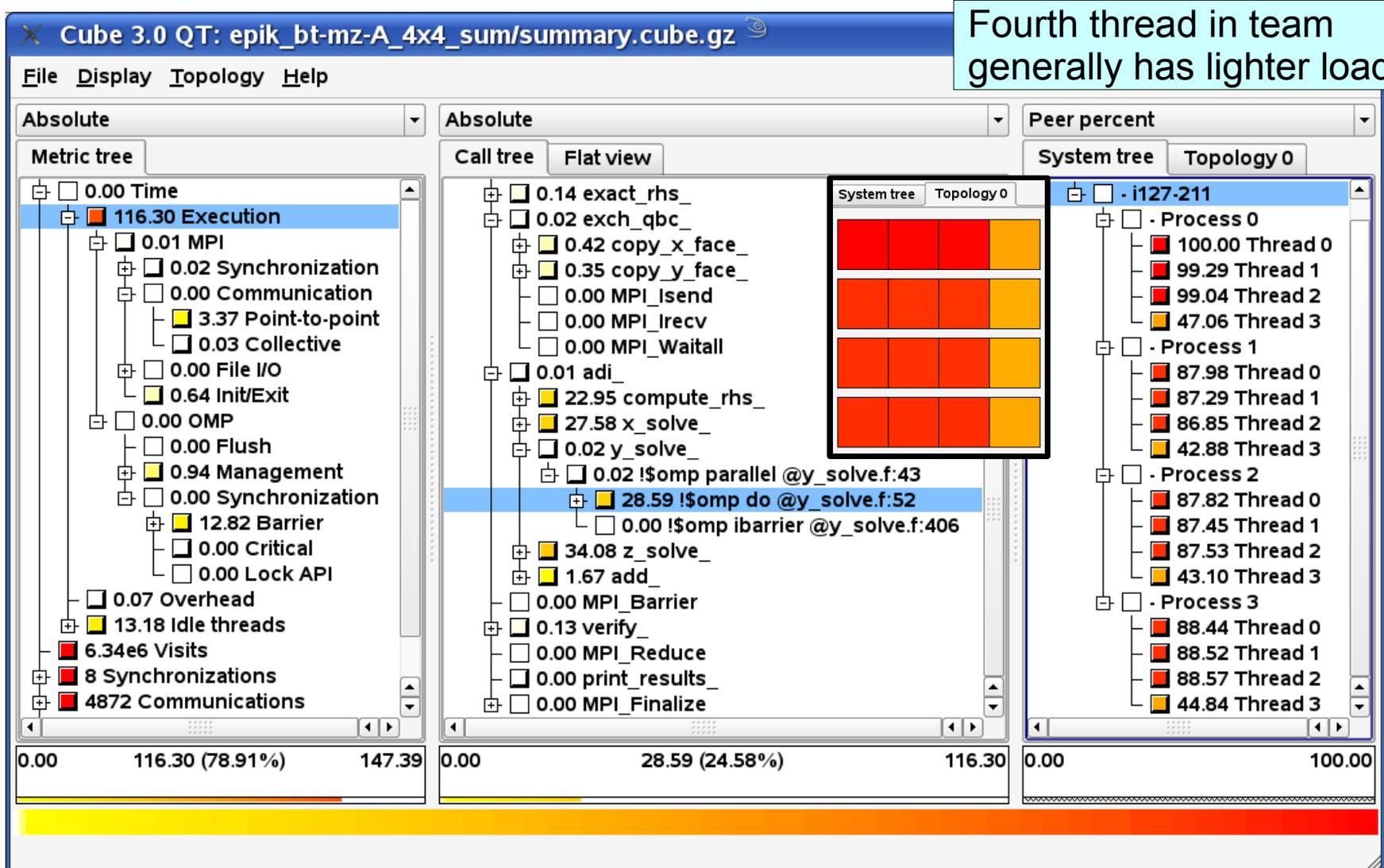
16-thread summary analysis: Idle threads time



99.74% of execution time found in parallel regions

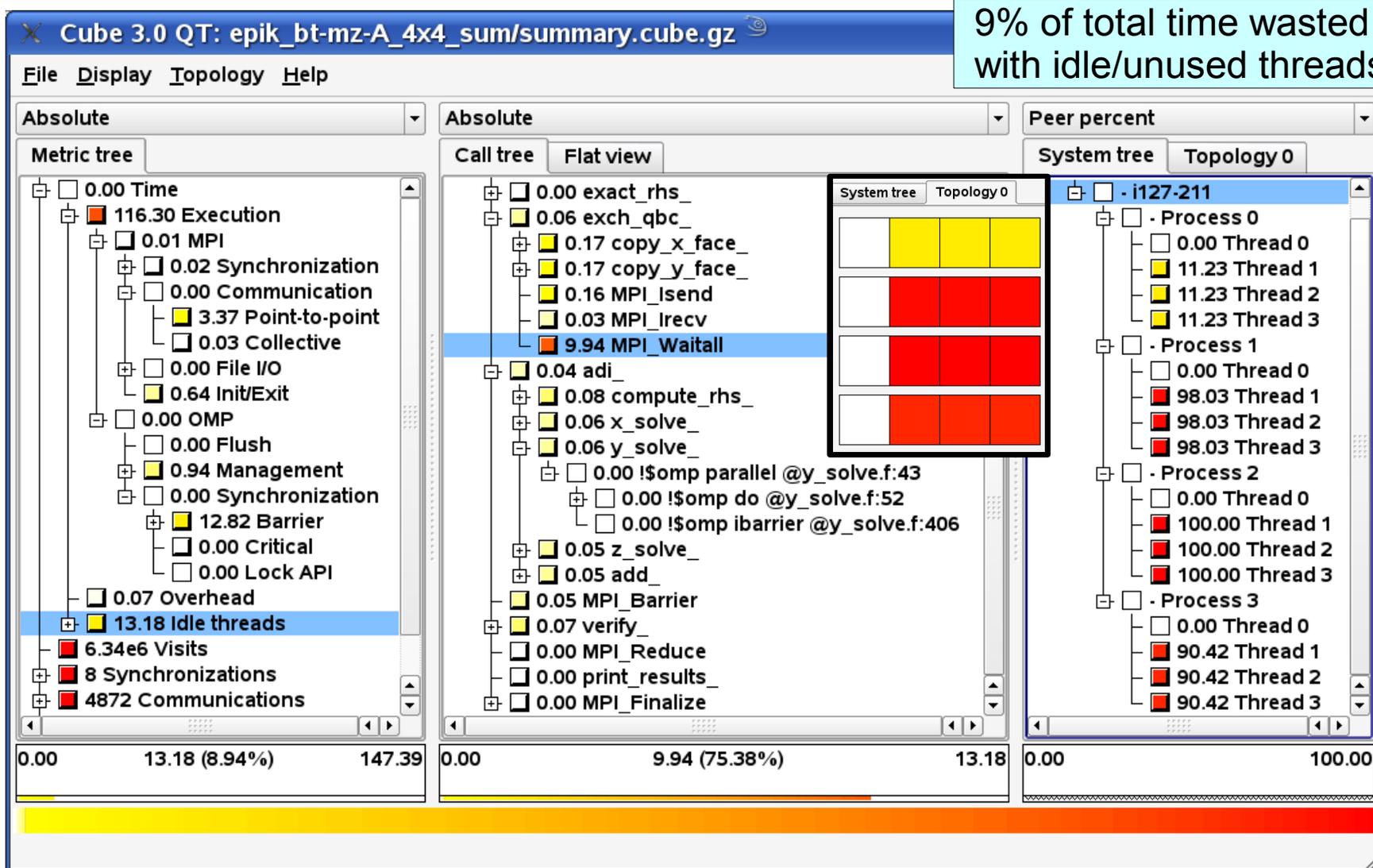
4x4 summary analysis: Execution time

Fourth thread in team generally has lighter load



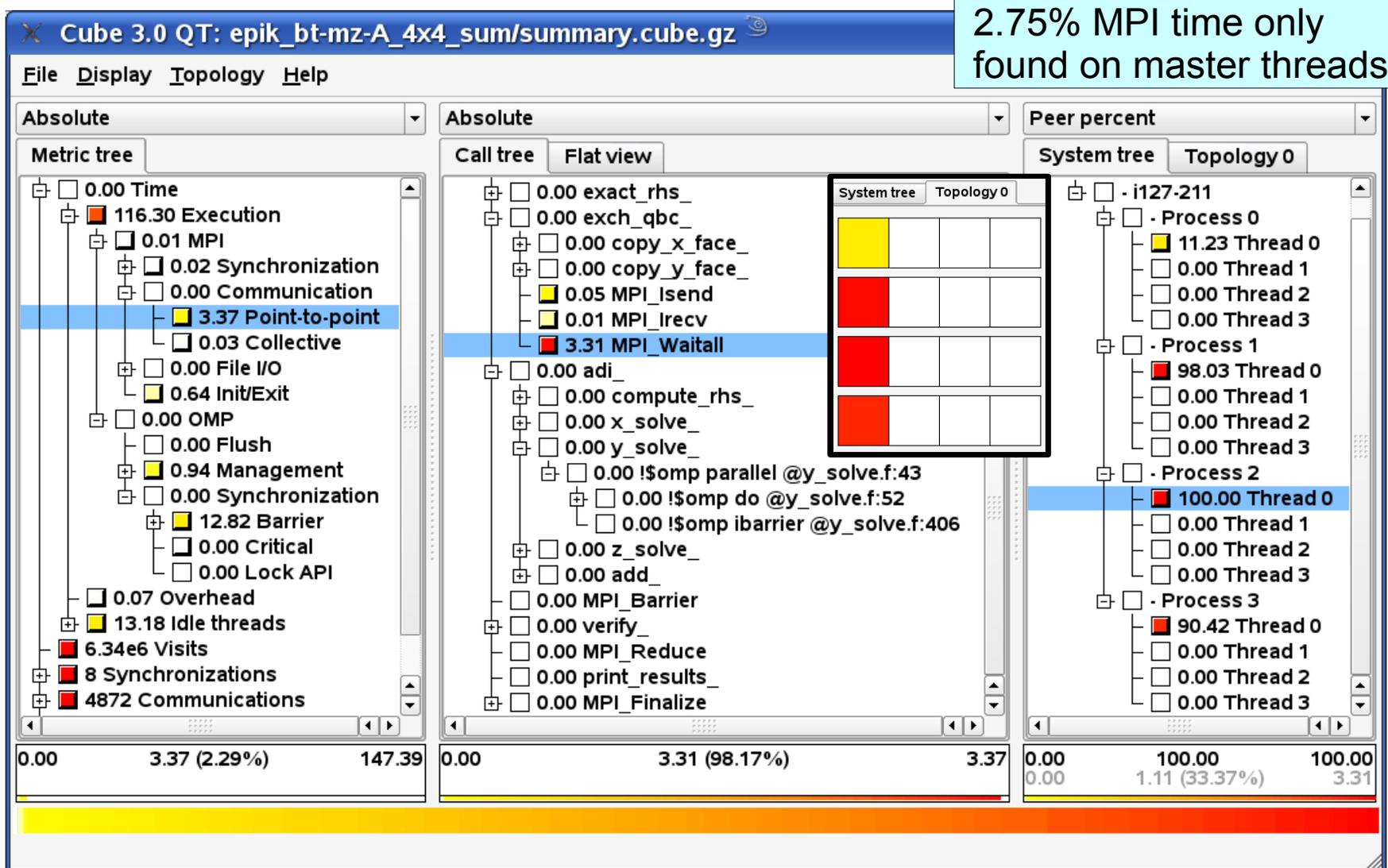
4x4 summary analysis: Idle threads time

9% of total time wasted with idle/unused threads

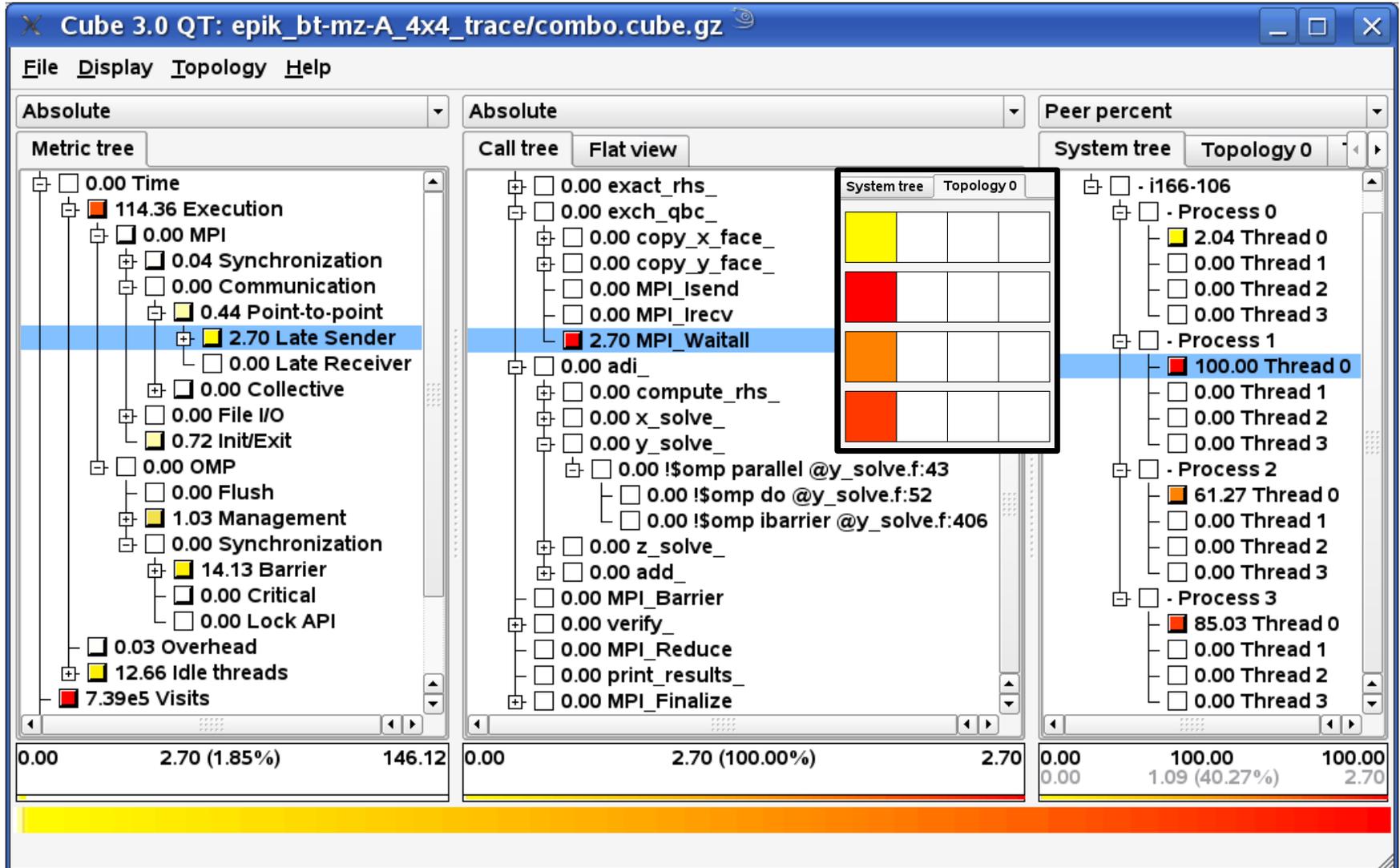


4x4 summary analysis: MPI time

2.75% MPI time only found on master threads



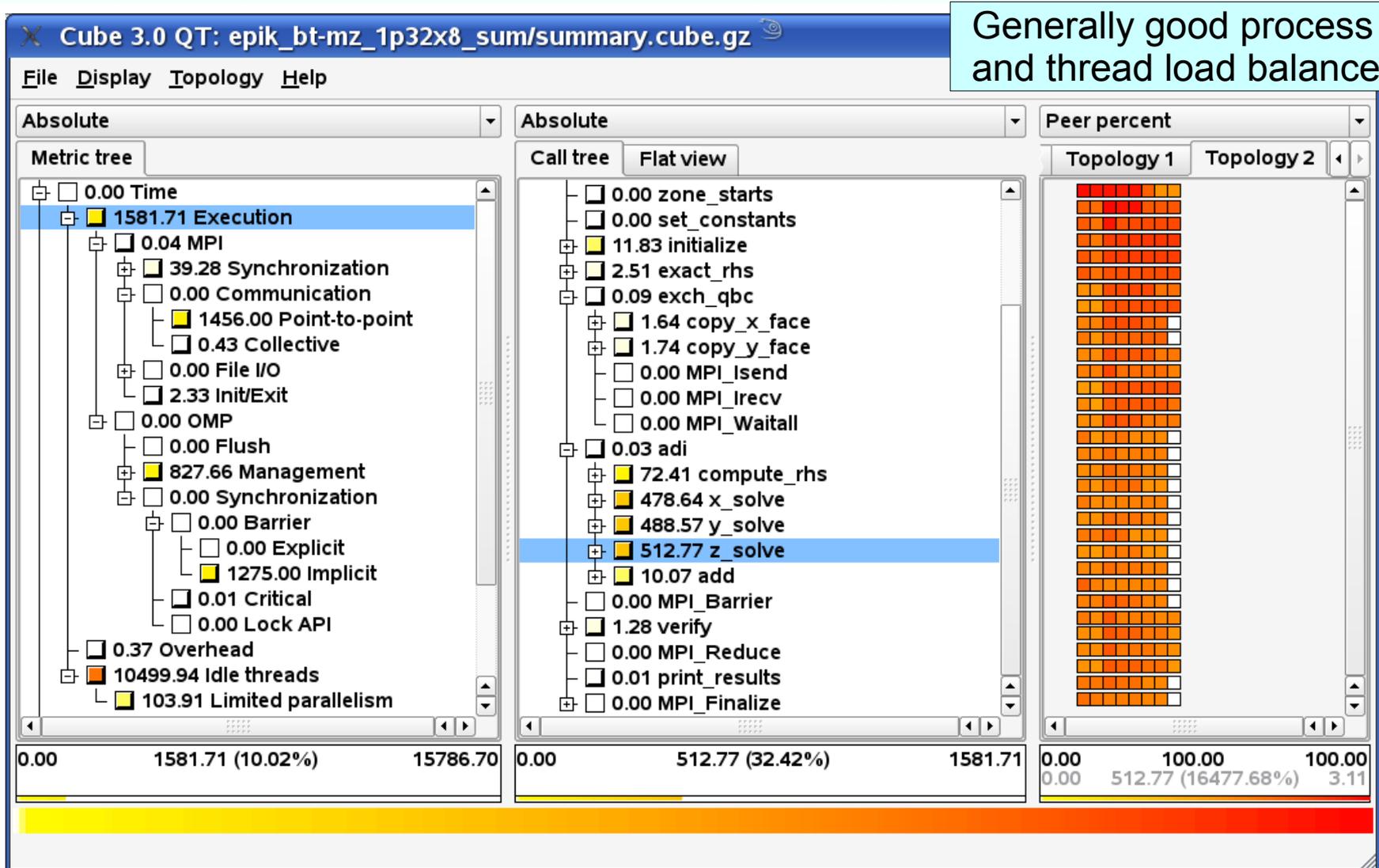
4x4 combined summary & trace analysis



32x8 summary analysis: Excl. execution time



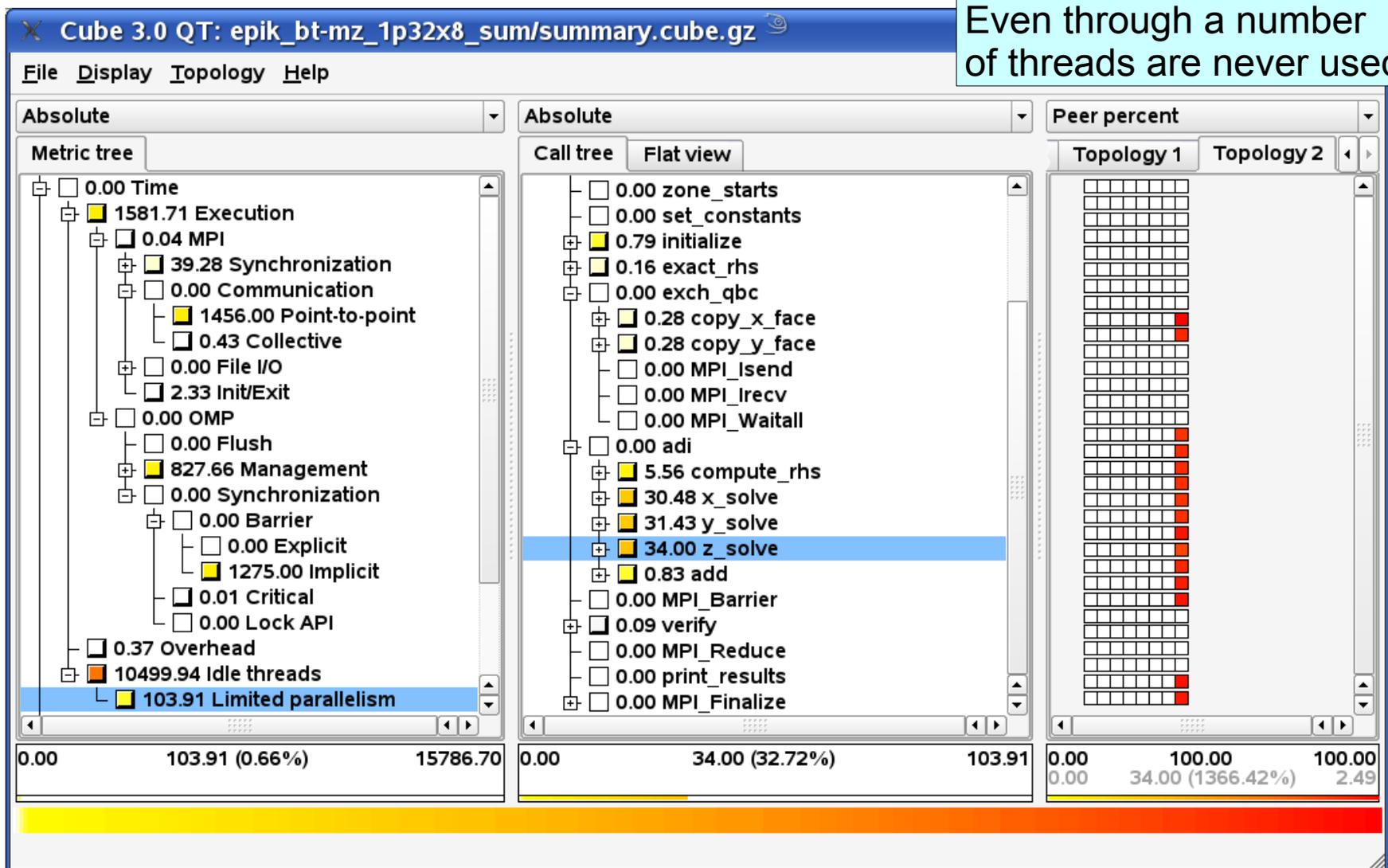
Generally good process and thread load balance



32x8 summary analysis: Limited parallelism

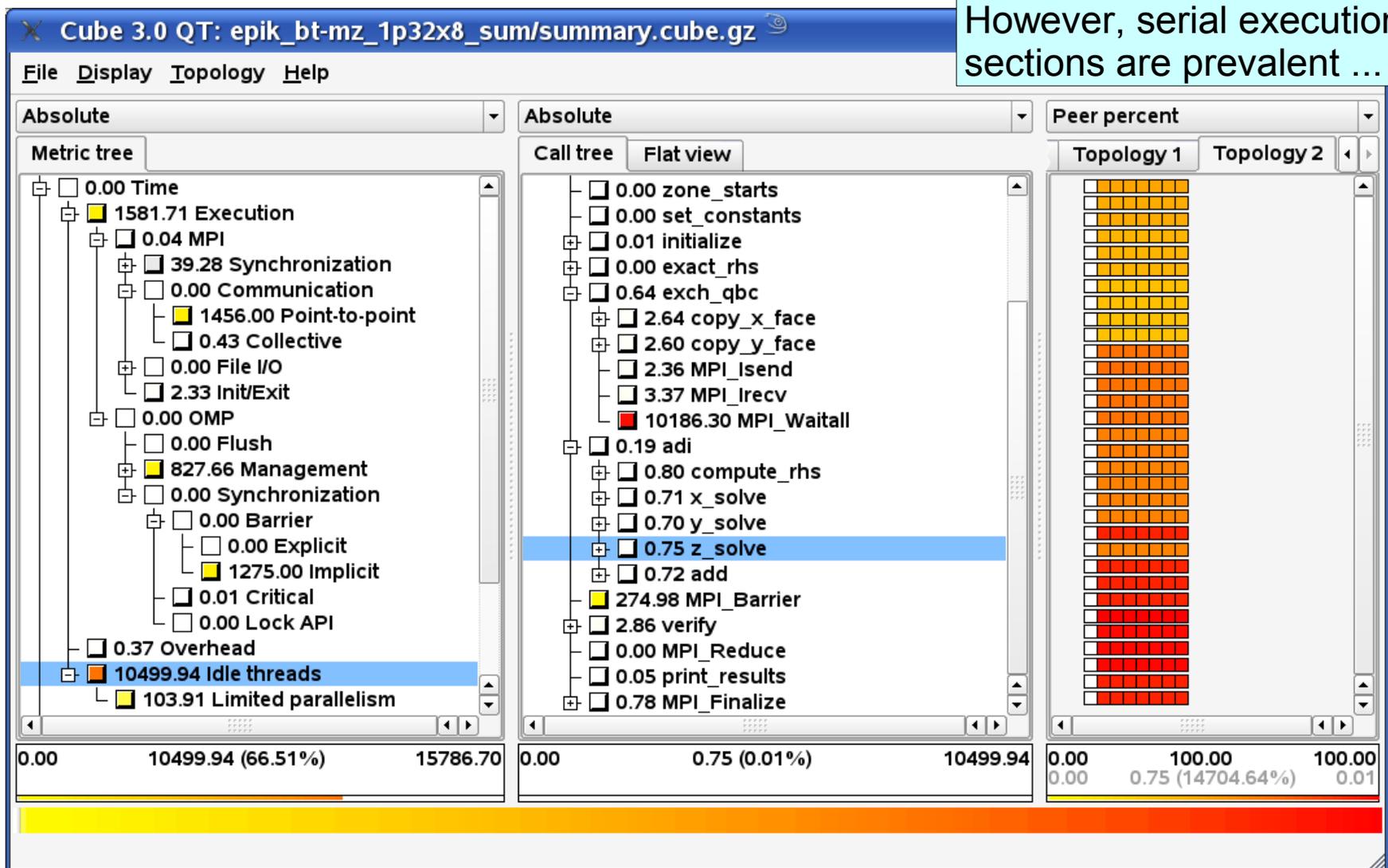


Even through a number of threads are never used



32x8 summary analysis: Idle threads time

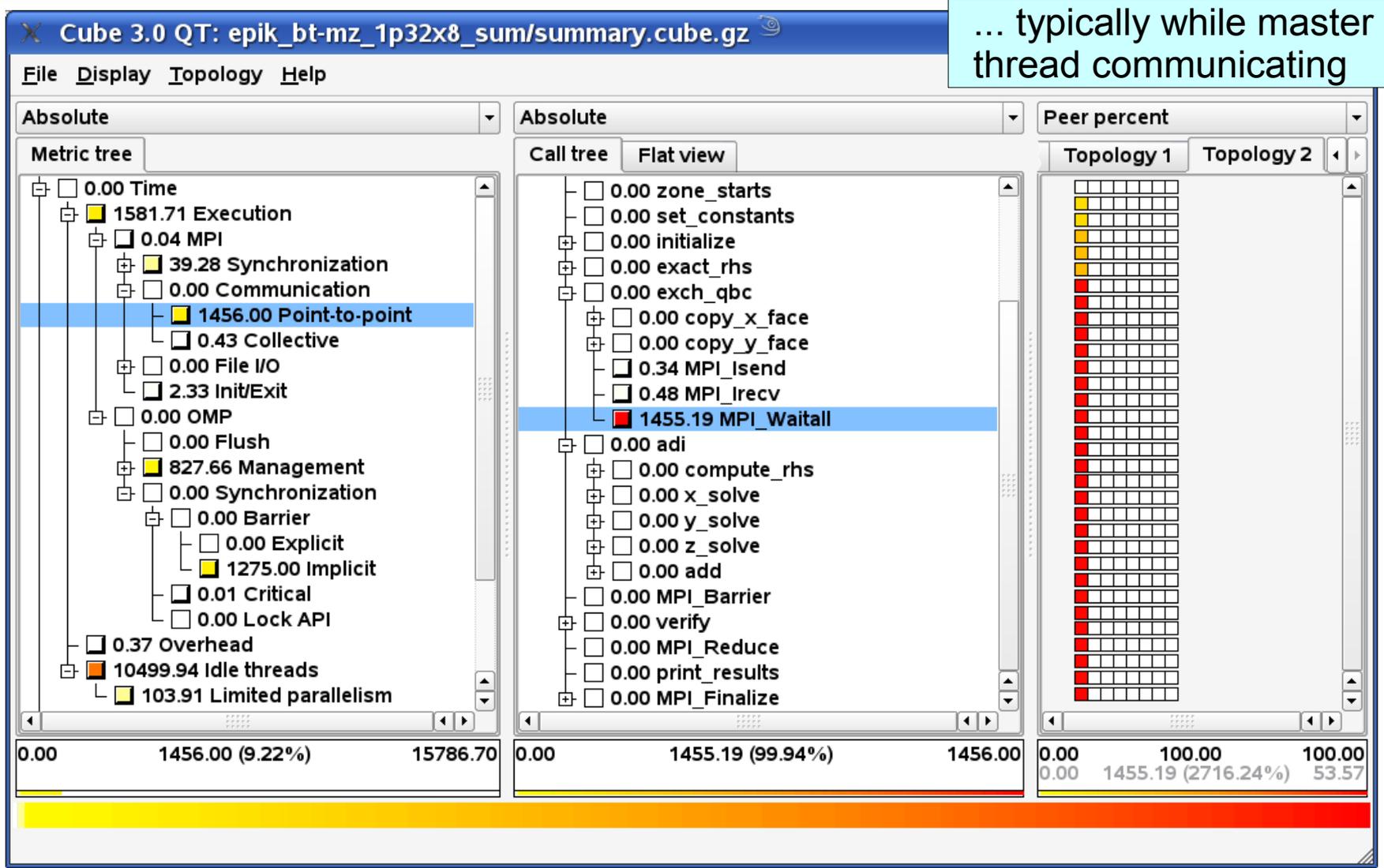
However, serial execution sections are prevalent ...



32x8 summary analysis: MPI communication time



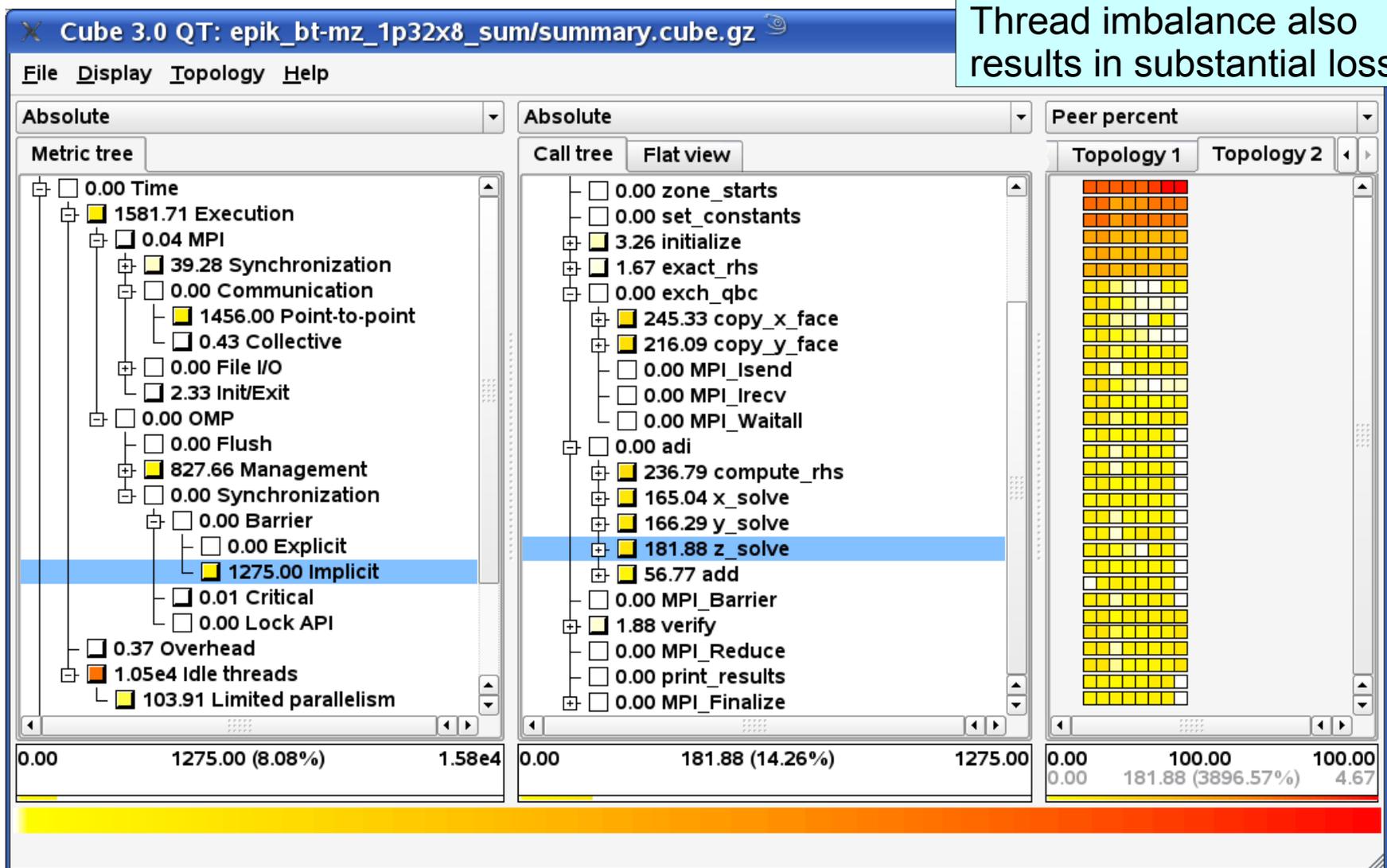
... typically while master thread communicating



32x8 summary analysis: Implicit barrier time



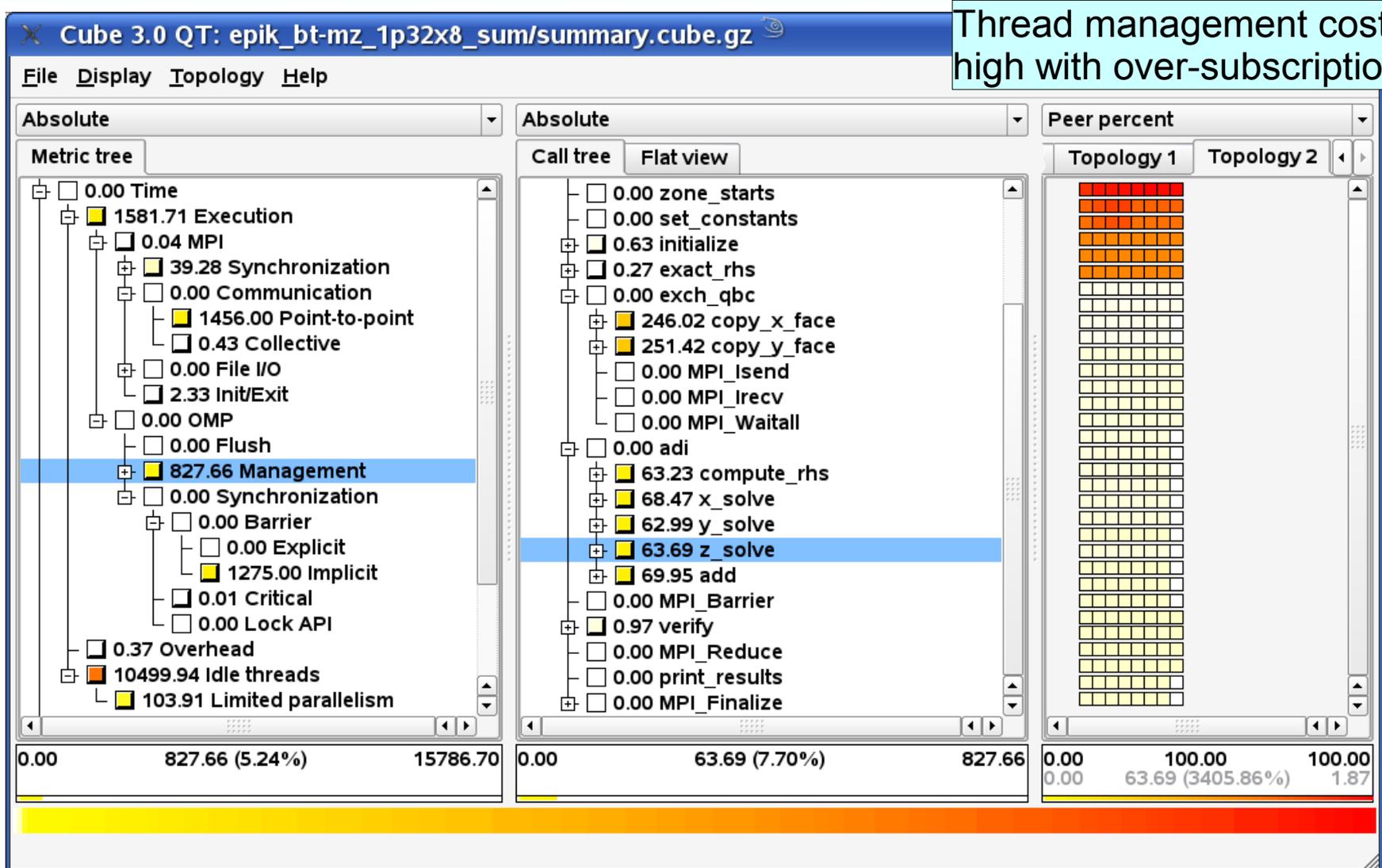
Thread imbalance also results in substantial loss



32x8 summary analysis: Thread management



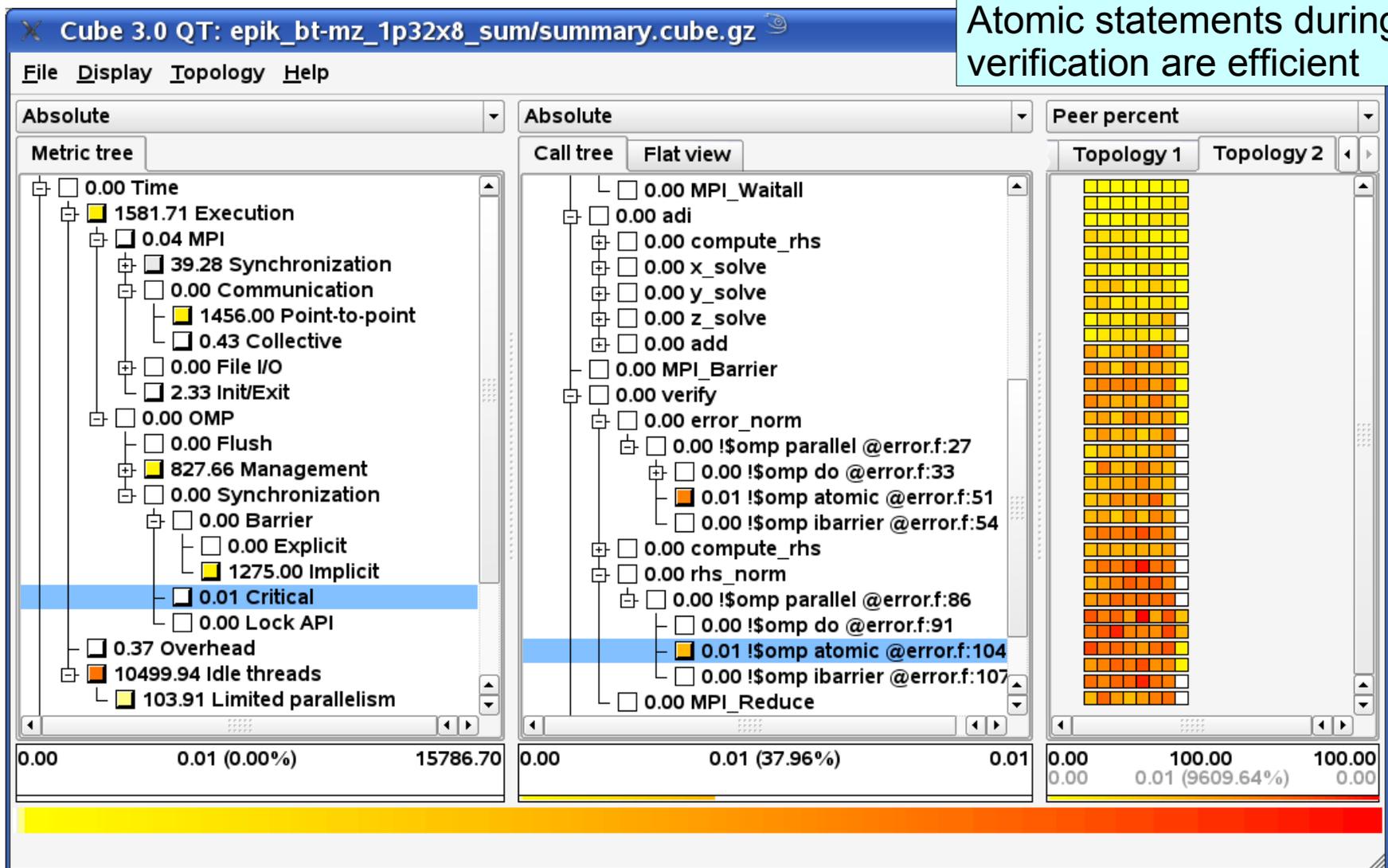
Thread management cost high with over-subscription



32x8 summary analysis: Critical section time

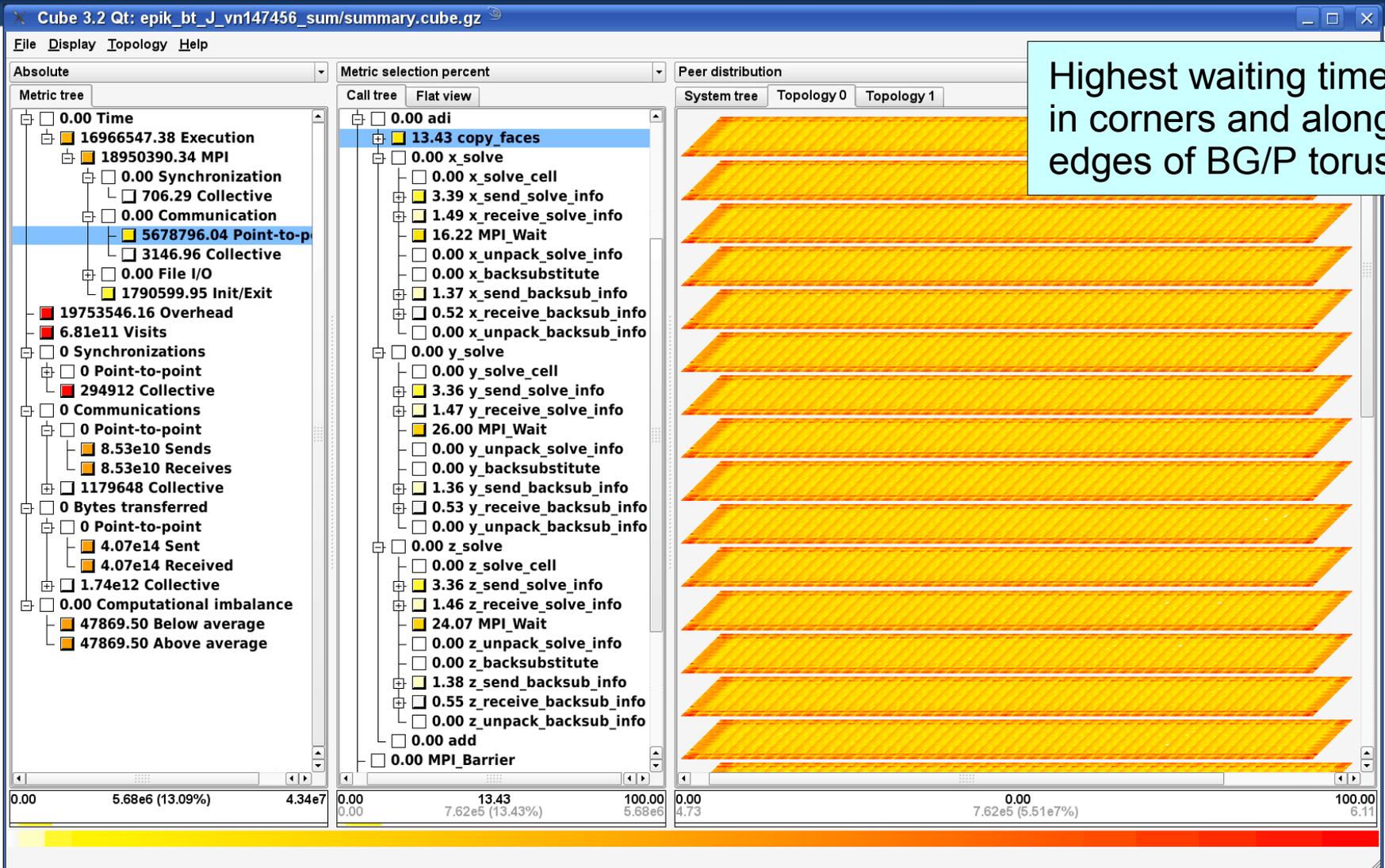


Atomic statements during verification are efficient



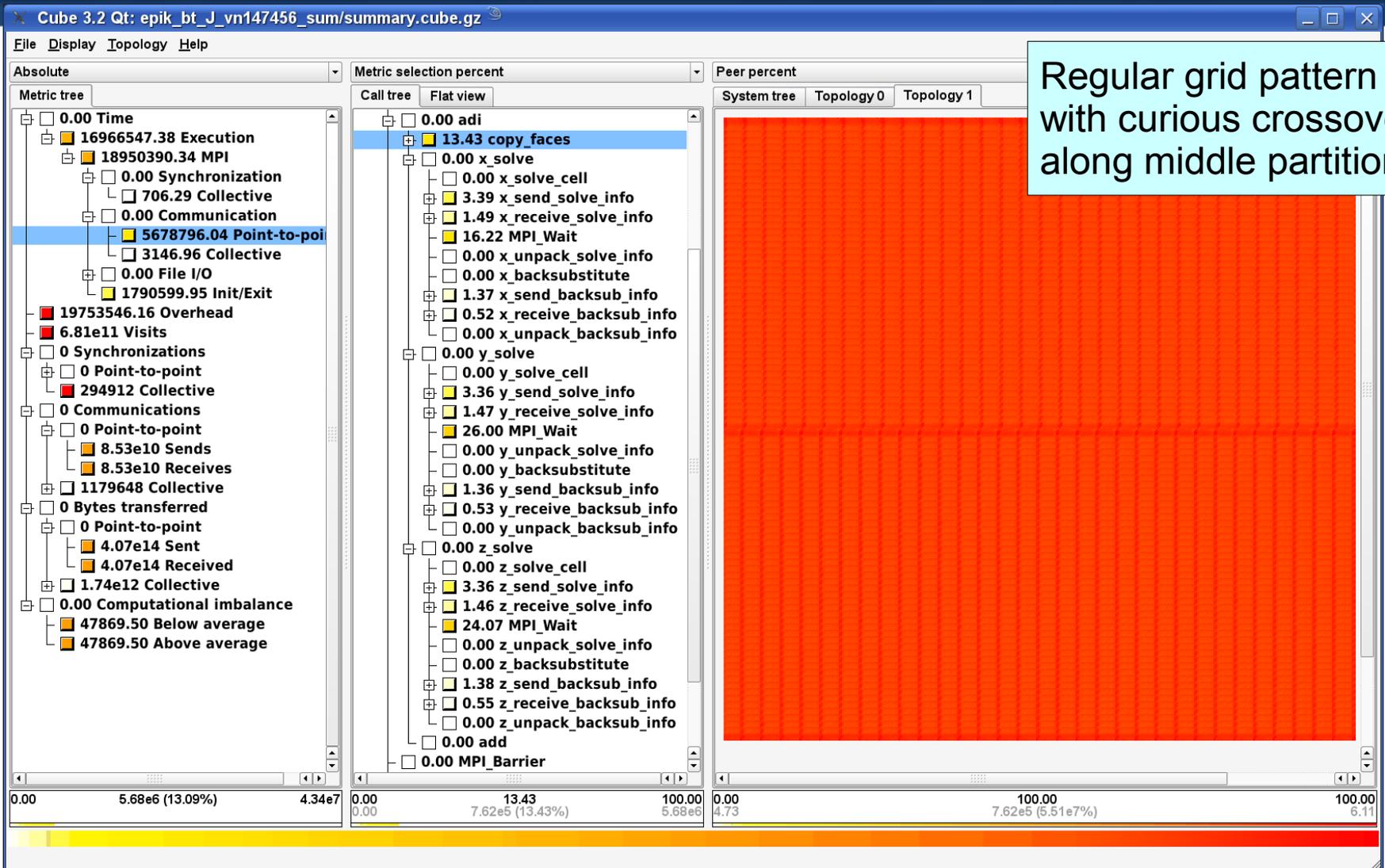
- 3D solution of unsteady, compressible Navier-Stokes eqs
 - NASA NAS parallel benchmark suite Block-Tridiagonal solver
 - series of ADI solve steps in X, Y & Z dimensions
 - ~9,500 lines (20 source modules), mostly Fortran77
- Run on IBM BlueGene/P in VN mode with 144k processes
 - Good scaling when 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



Highest waiting times in corners and along edges of BG/P torus

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

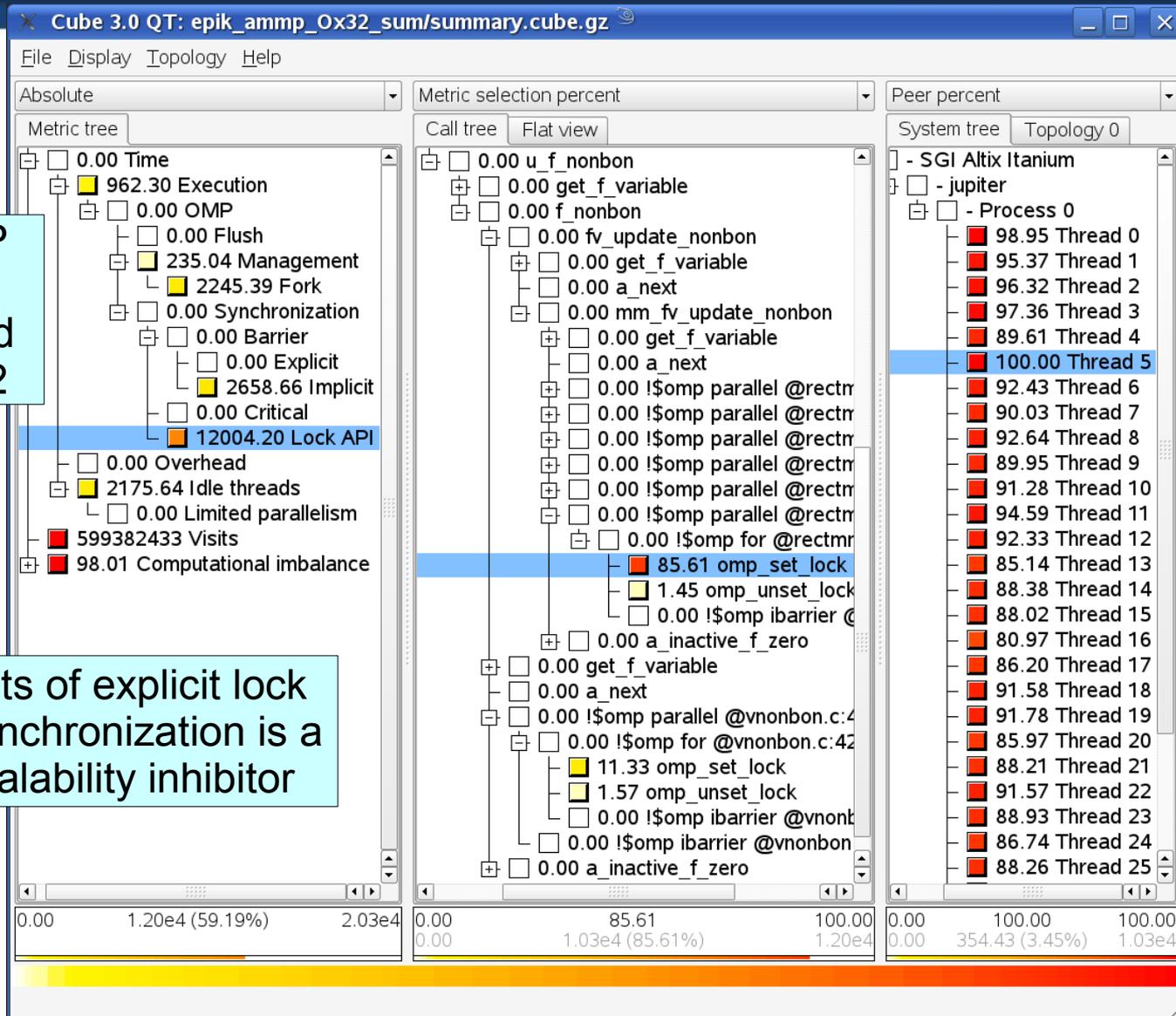


Regular grid pattern with curious crossover along middle partition

- 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

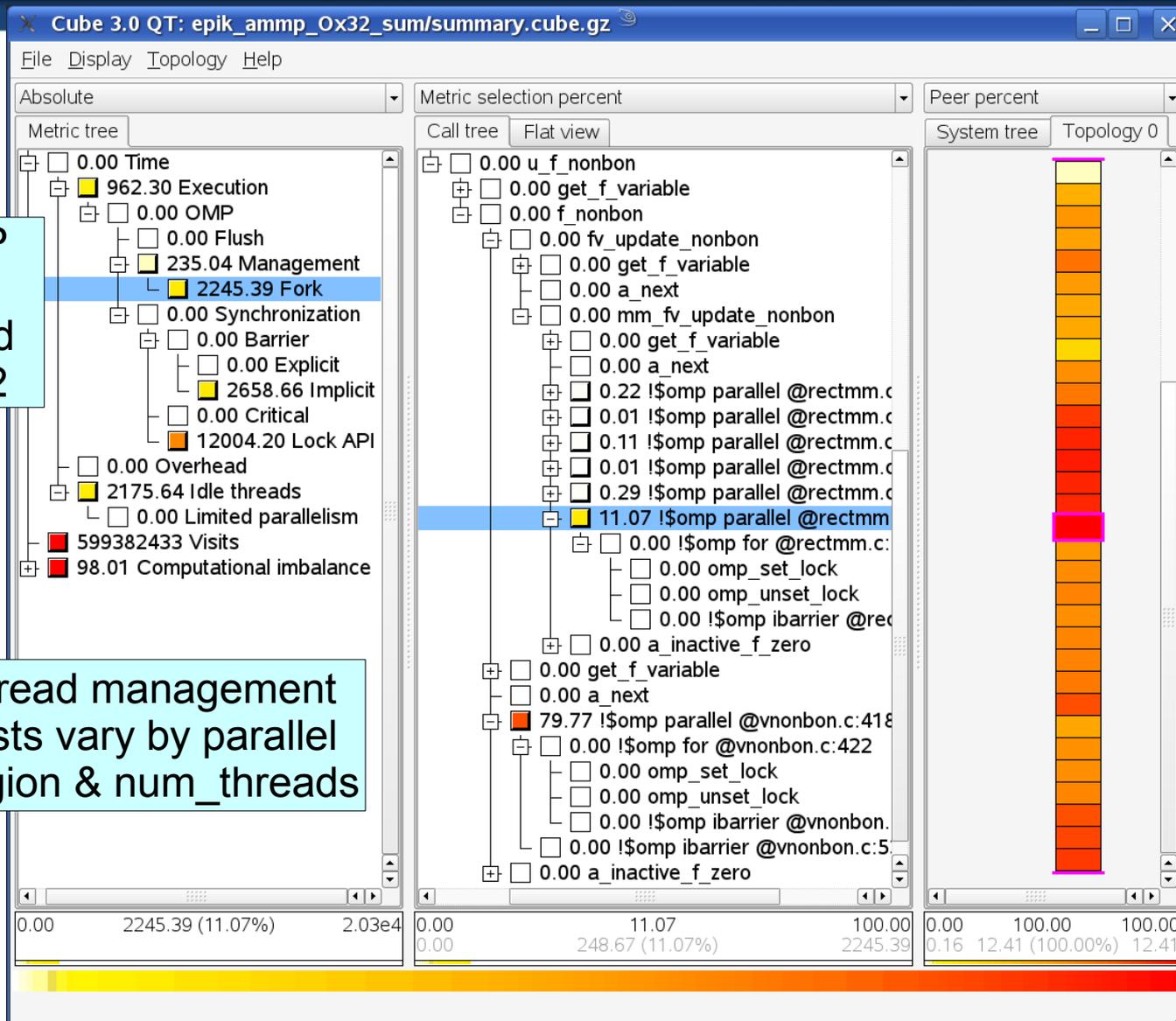
OpenMP metrics reworked with v1.2

Lots of explicit lock synchronization is a scalability inhibitor



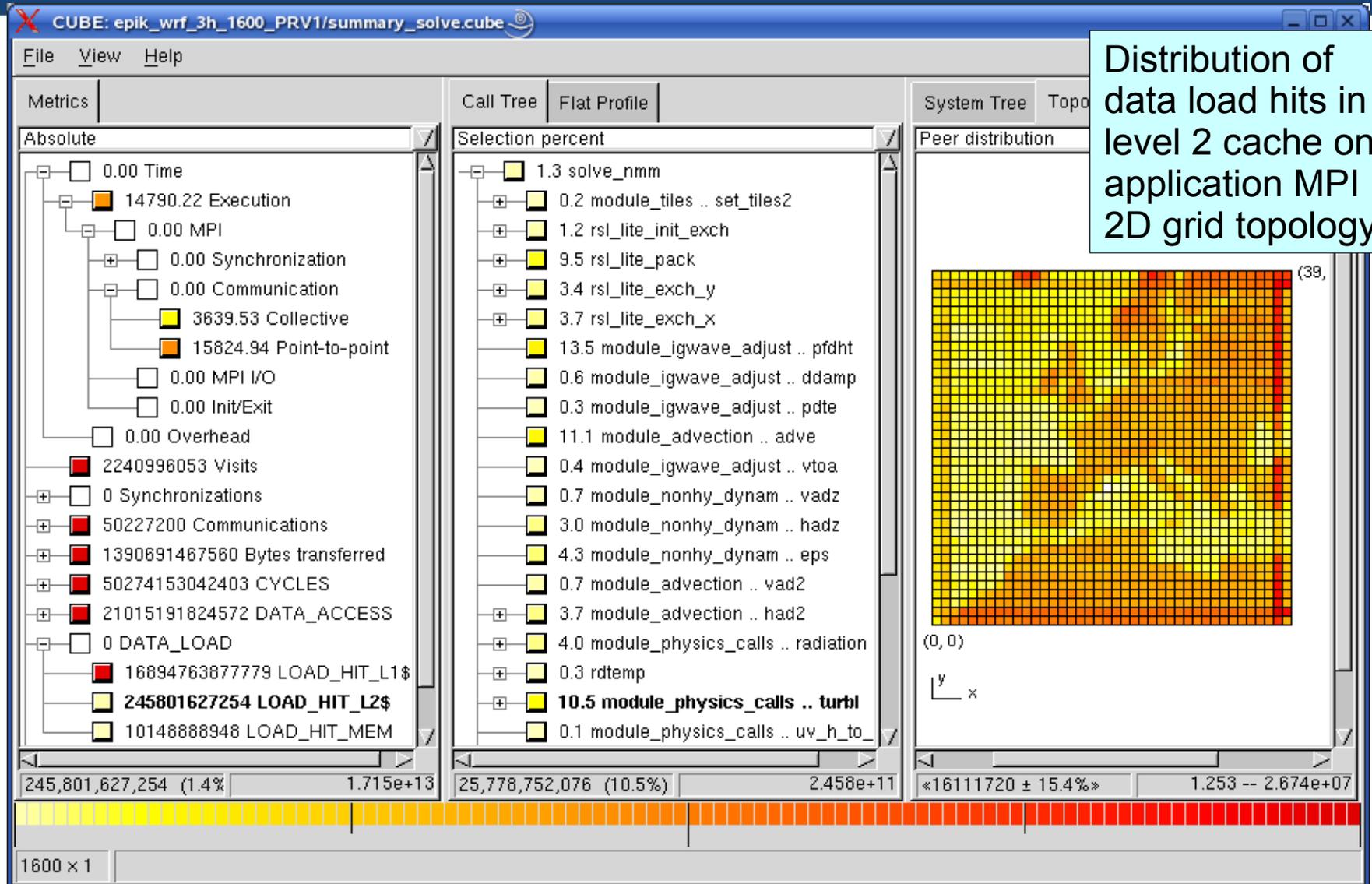
OpenMP metrics reworked with v1.2

Thread management costs vary by parallel region & num_threads



- 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 with HWC metrics

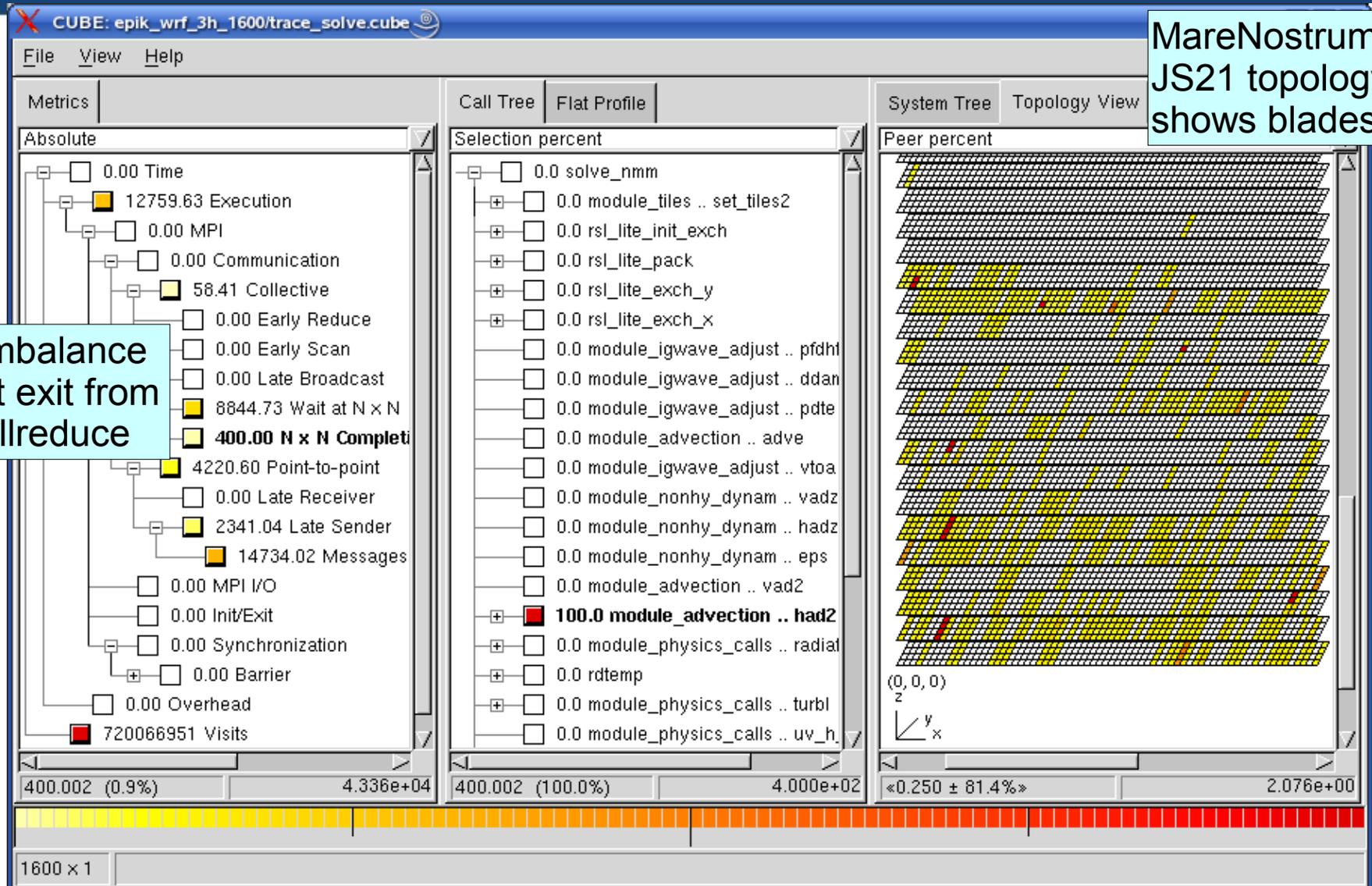


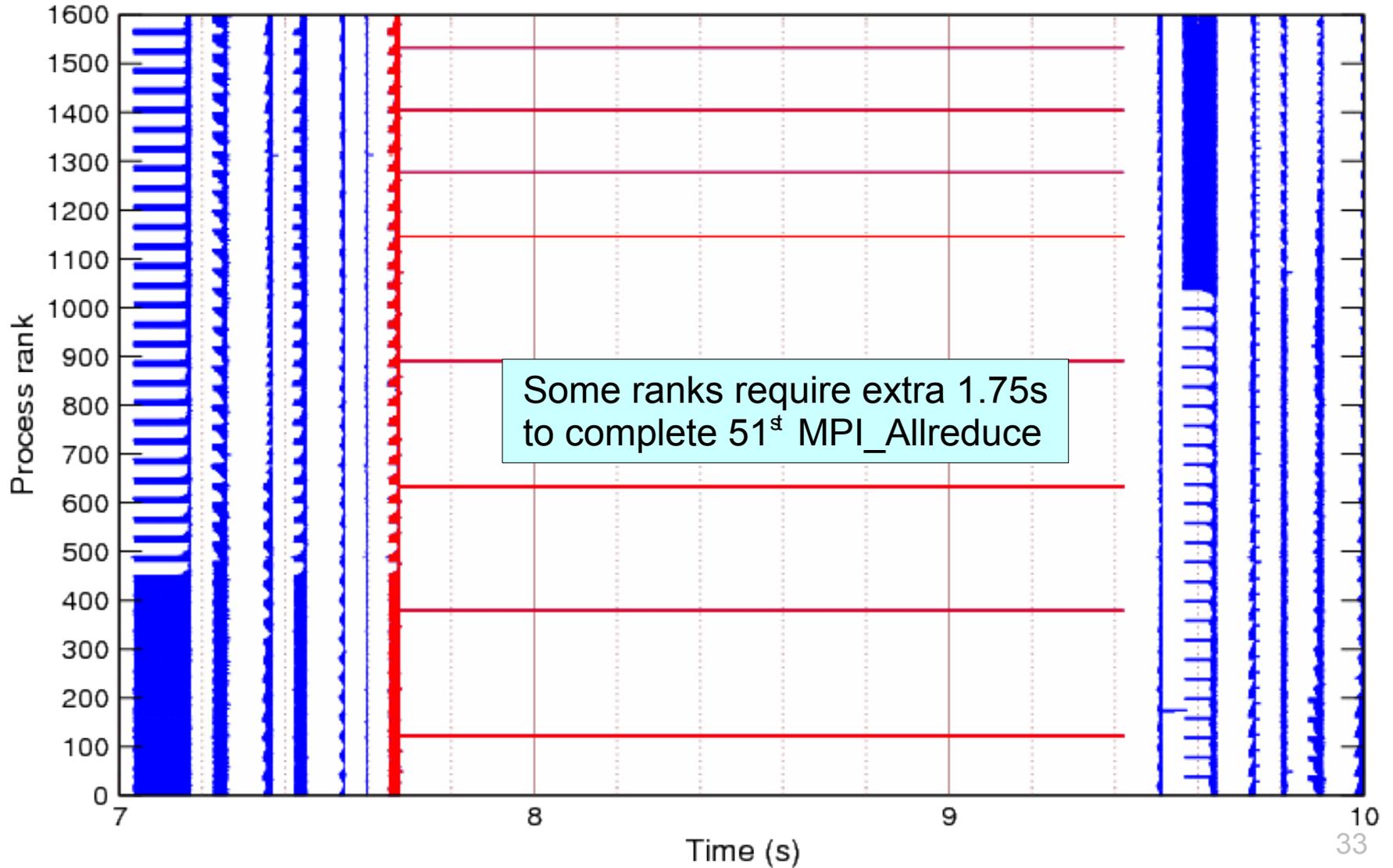
Distribution of data load hits in level 2 cache on application MPI 2D grid topology

WRF on MareNostrum@1600 trace analysis

MareNostrum JS21 topology shows blades

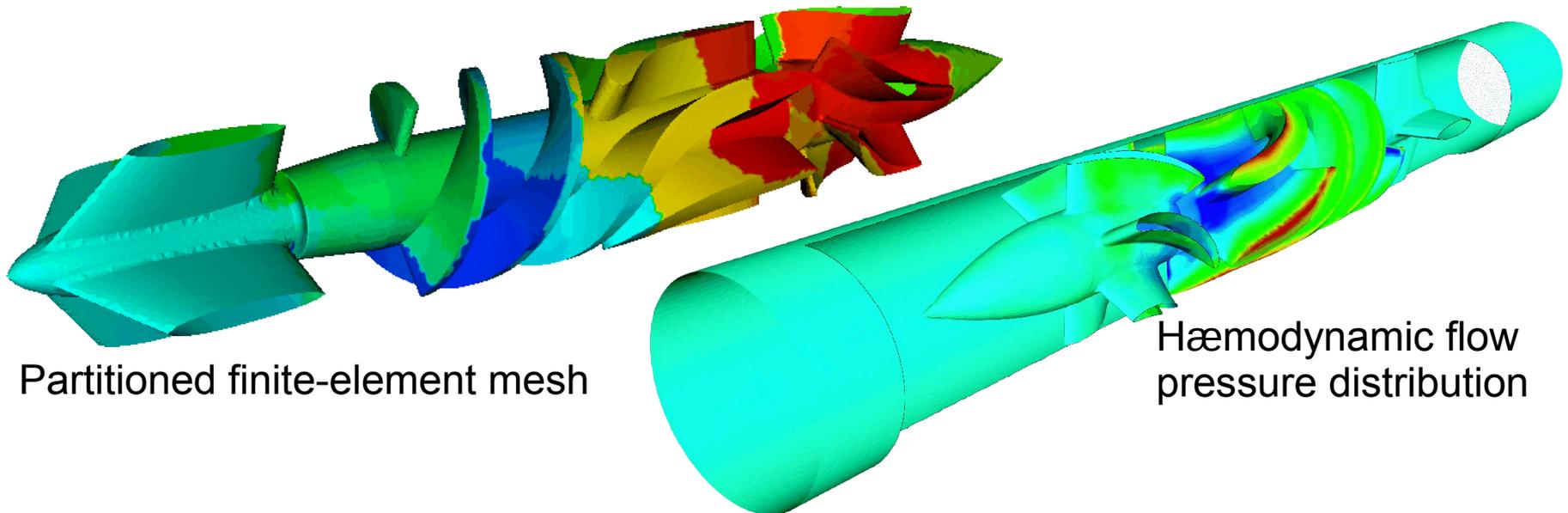
Imbalance at exit from Allreduce





- 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)

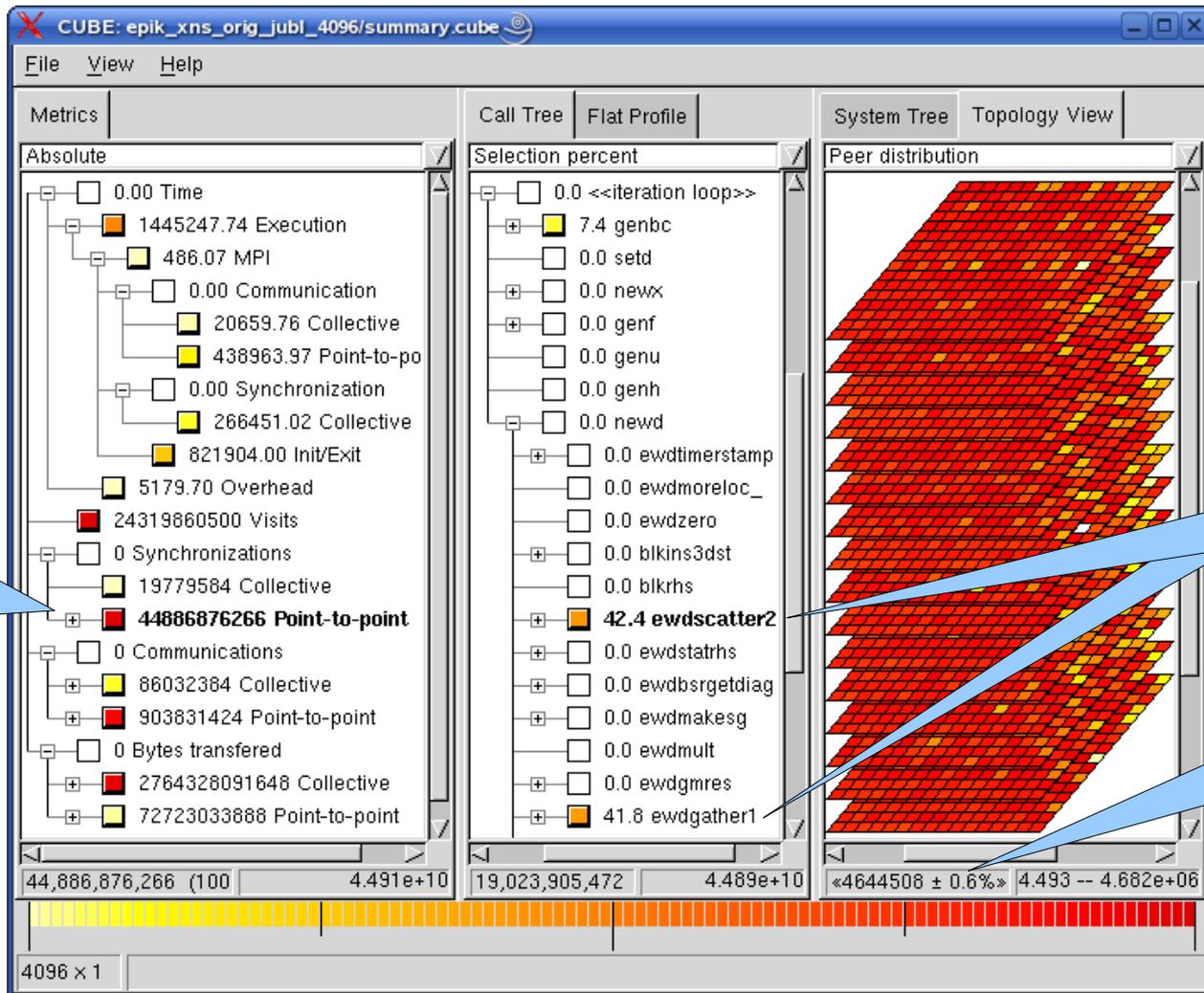


Partitioned finite-element mesh

Haemodynamic flow pressure distribution

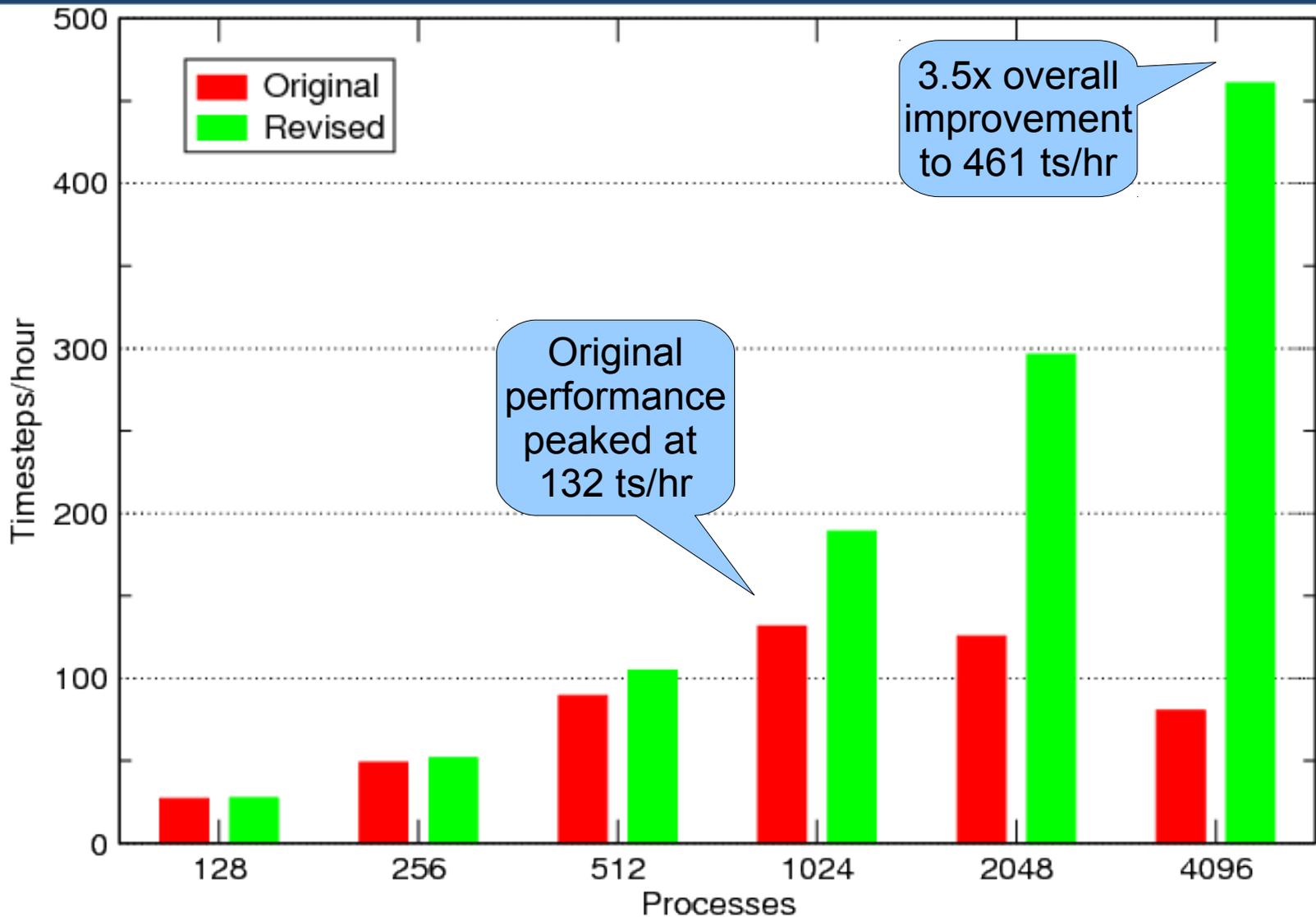
Point-to-point msgs
w/o data

Masses of
P2P synchronizations



Primarily
in scatter
& gather

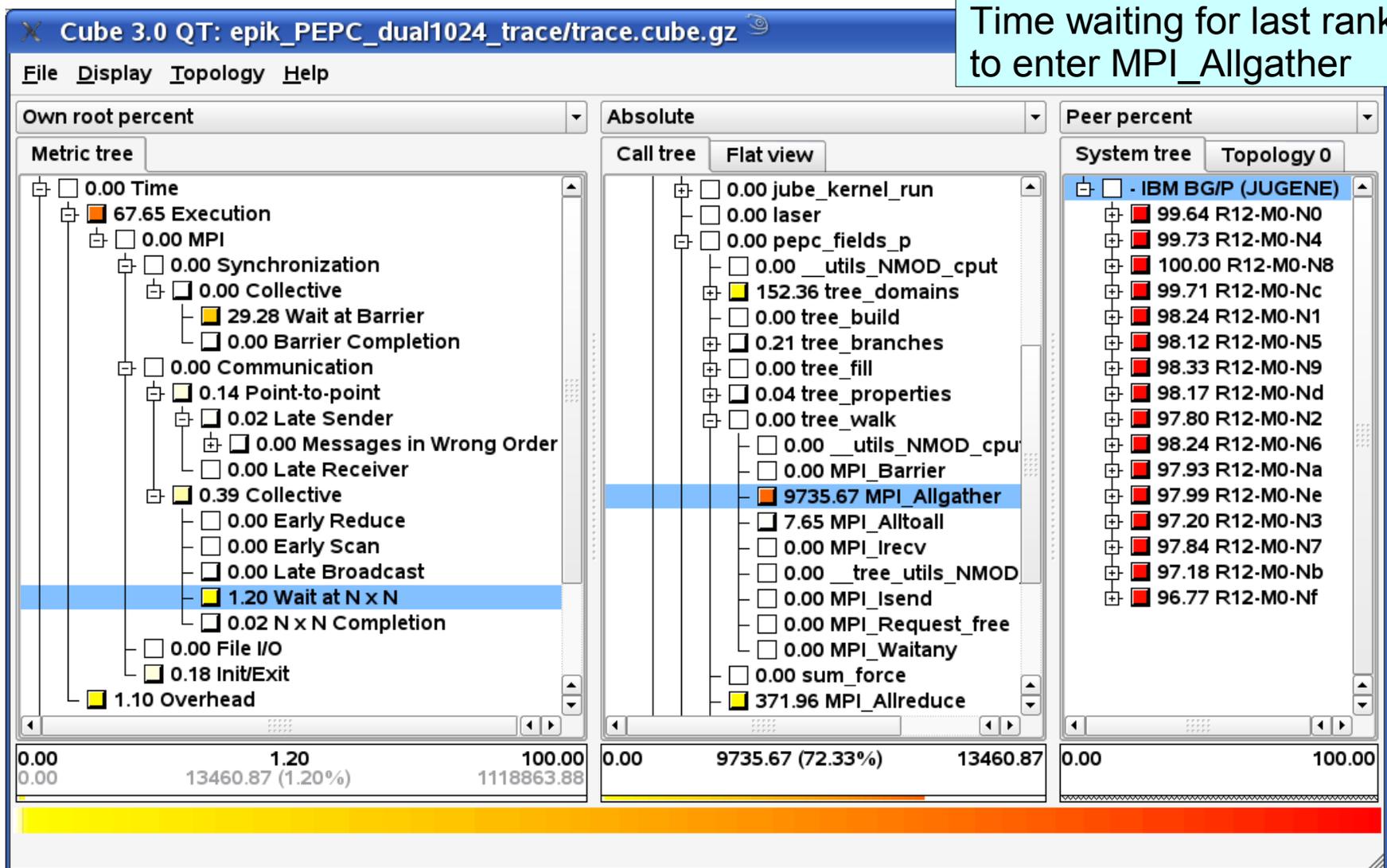
Processes
all equally
responsible



- 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

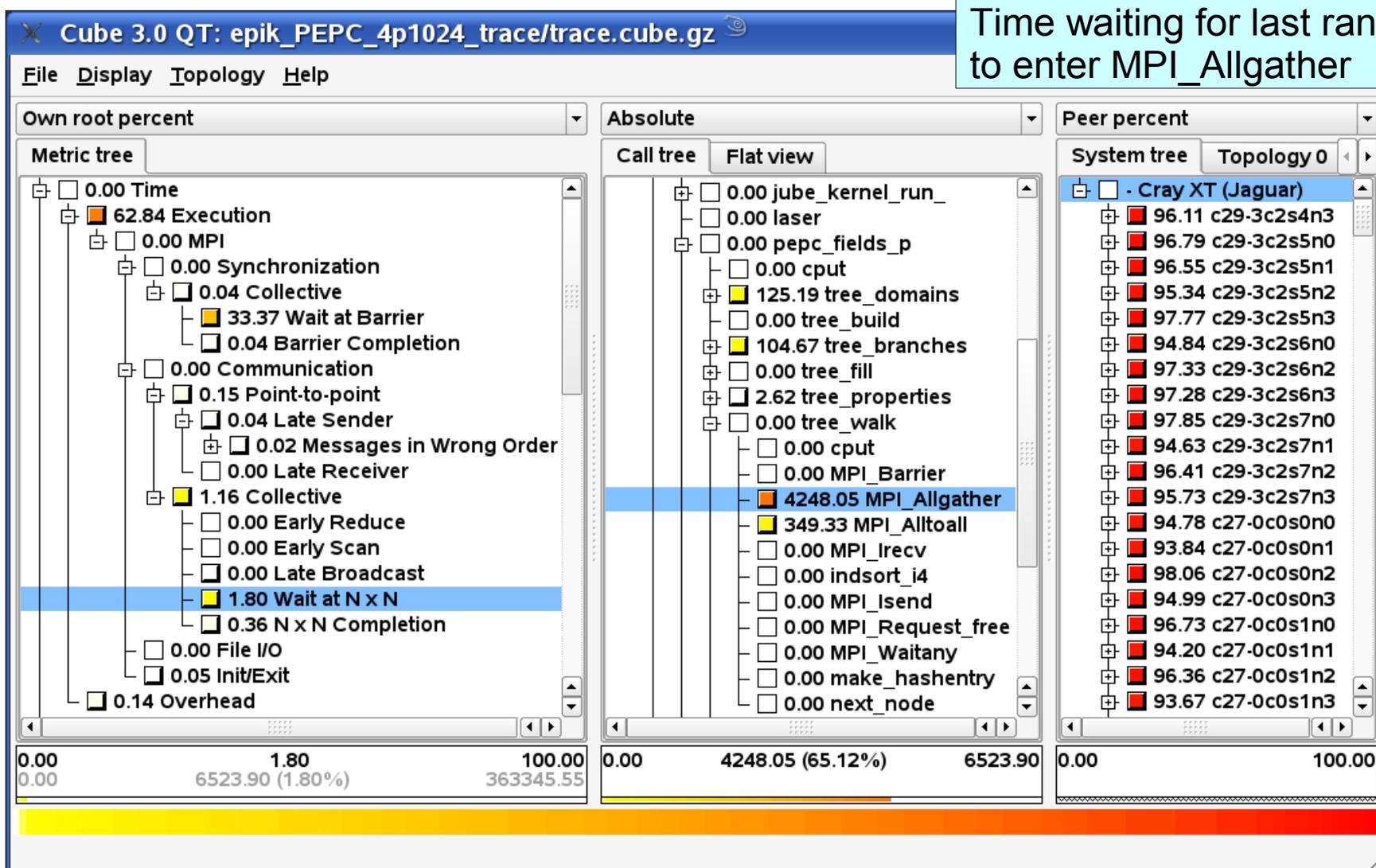
Time waiting for last rank to enter MPI_Allgather



PEPC@1024 on Cray XT4: Wait at NxN time



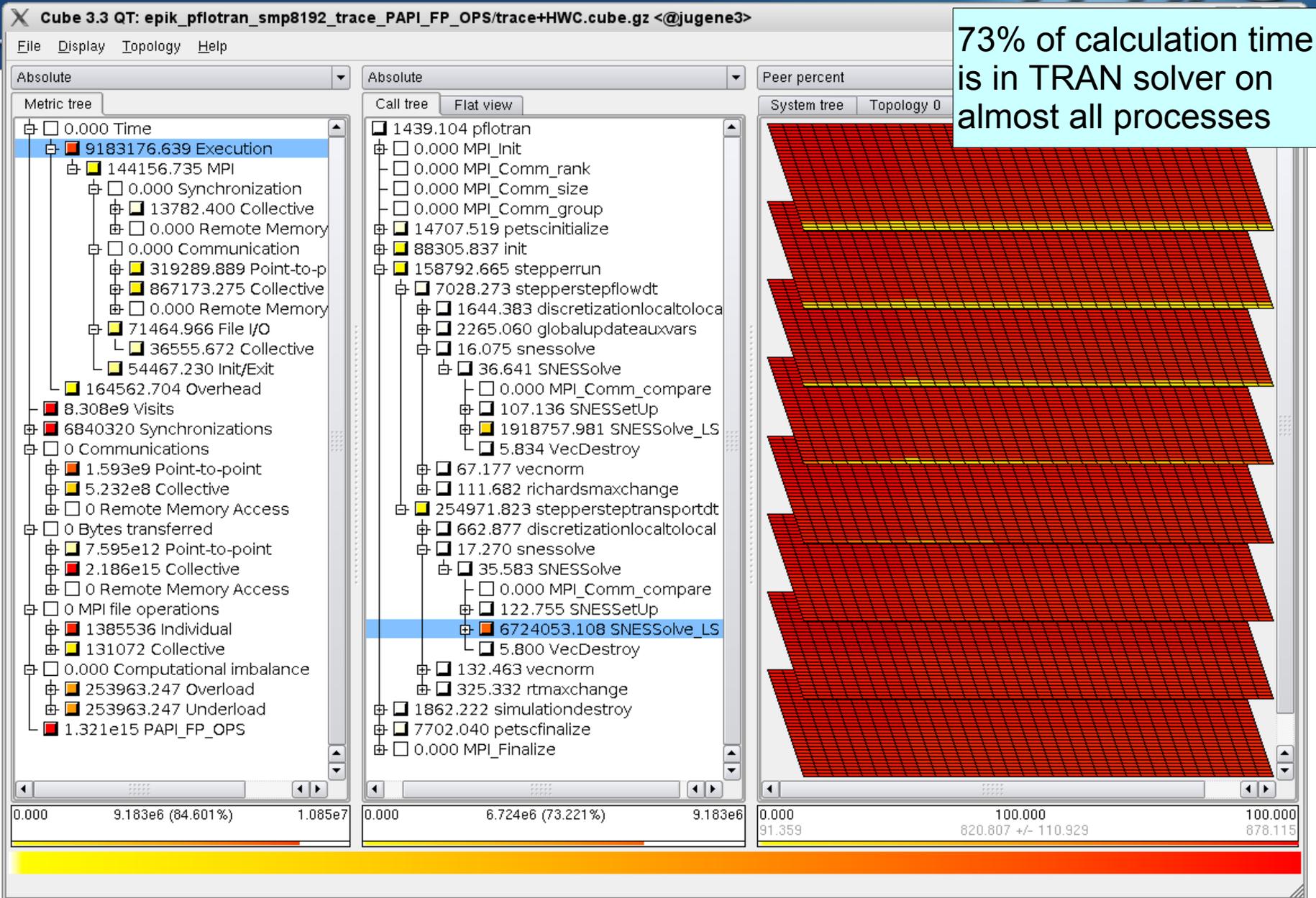
Time waiting for last rank to enter MPI_Allgather



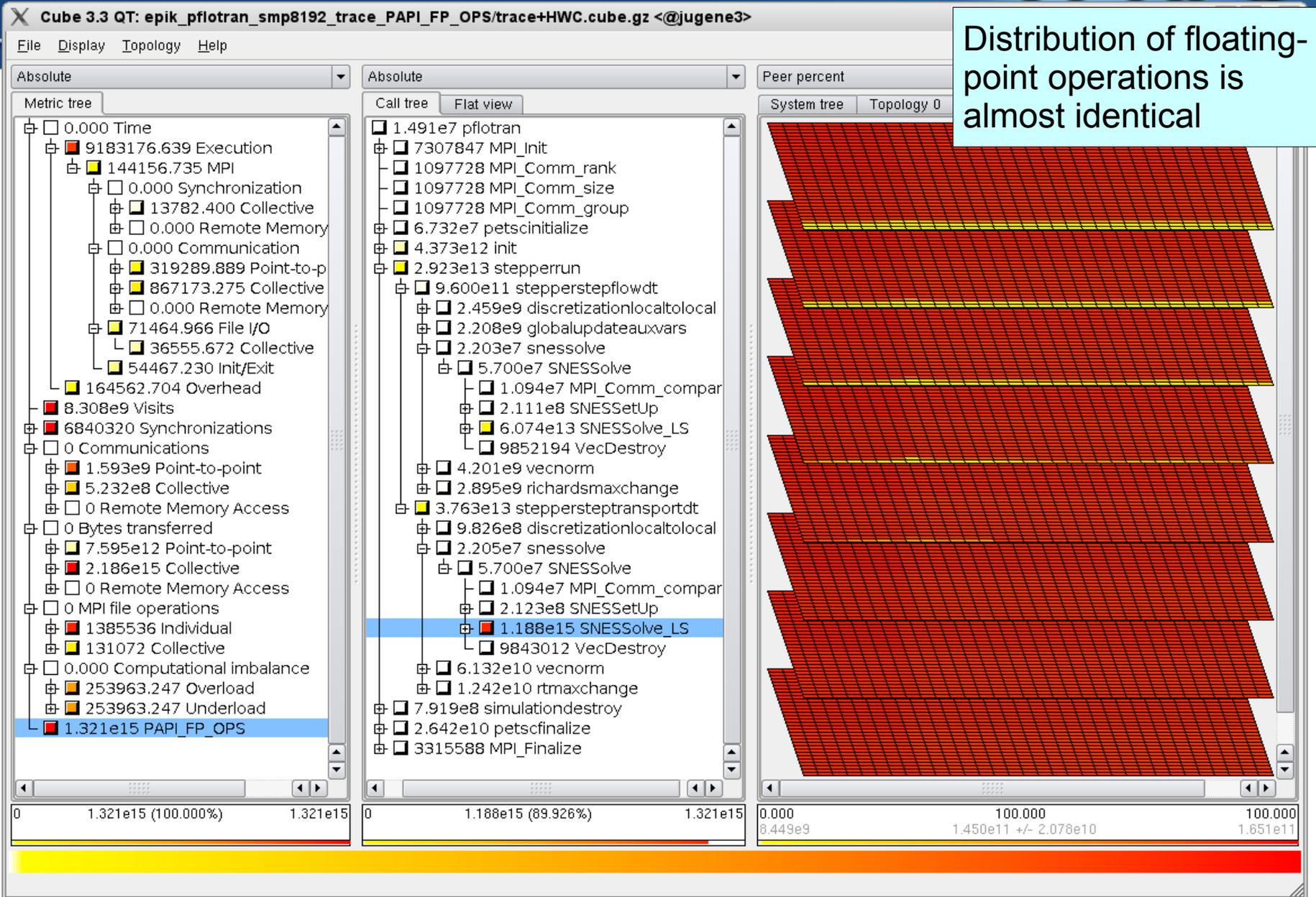
- Despite very different processor and network performance, measurements and analyses can be easily compared
 - different compilers affect function naming & in-lining
- Both spend roughly two-thirds of time in computation
 - tree_walk has expensive computation & communication
- Both waste 30% of time 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]

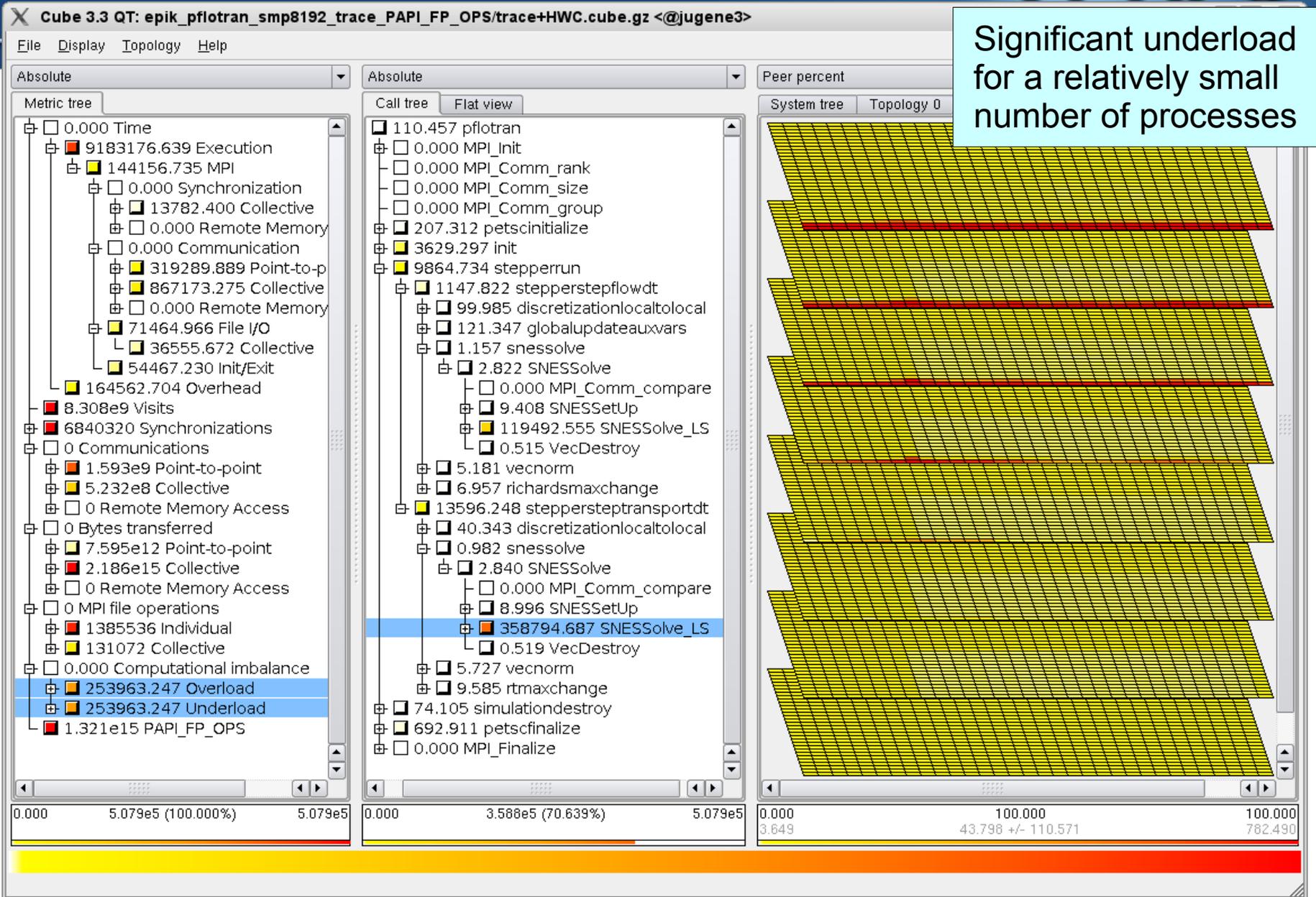
PFLOTRAN jugene@smp8192 trace analysis



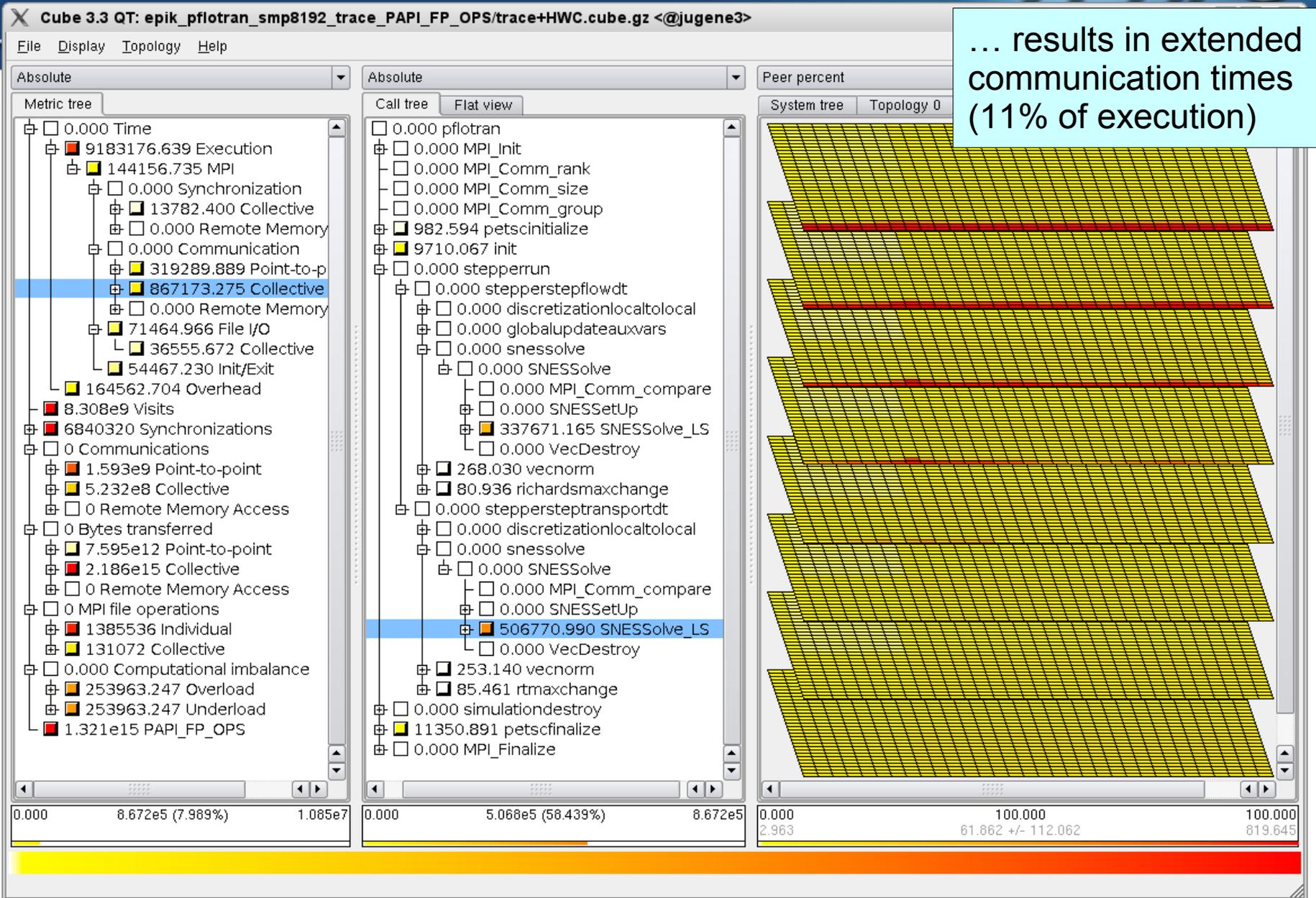
PFLOTRAN jugene@smp8192 trace analysis



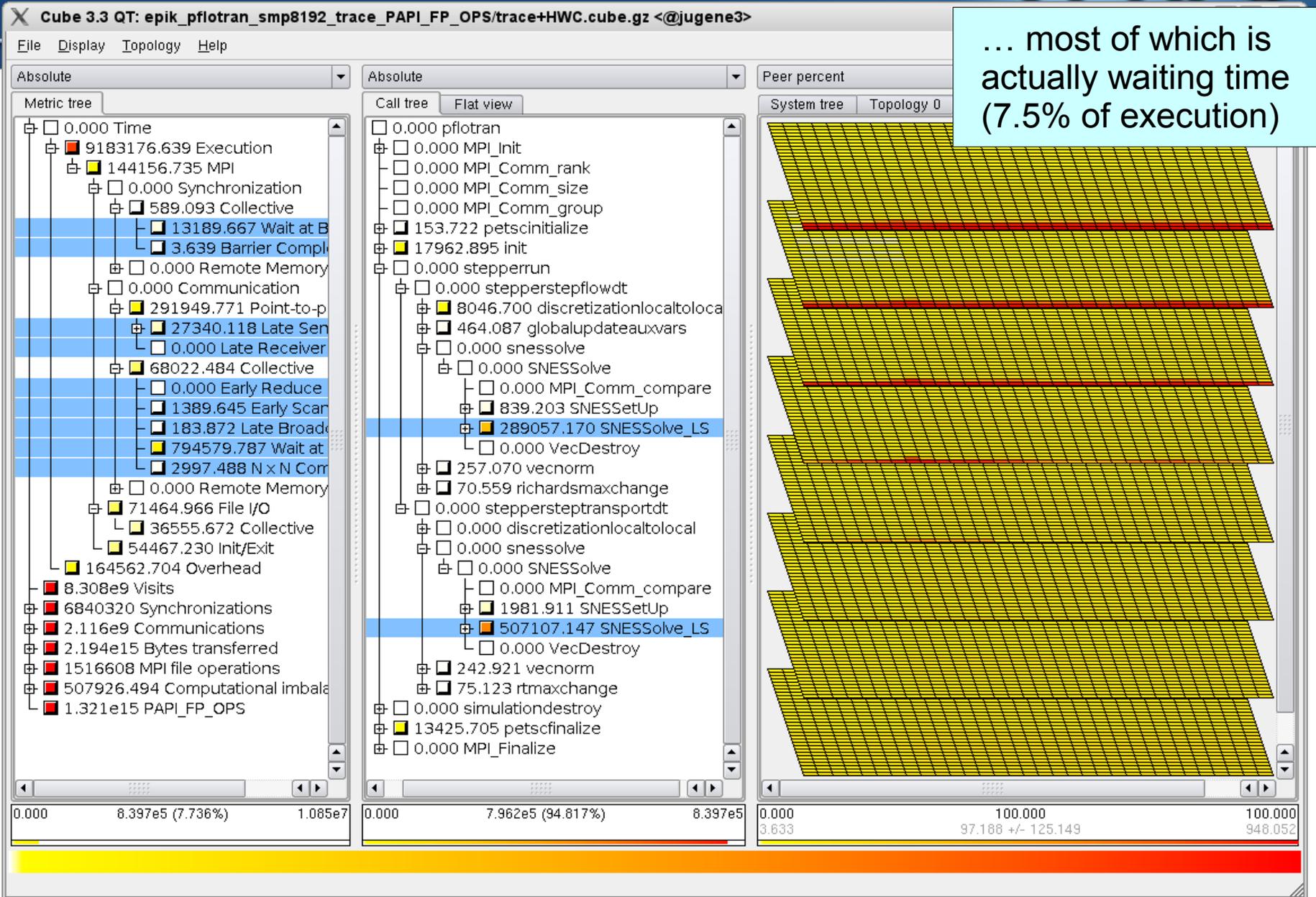
PFLOTRAN jugene@smp8192 trace analysis

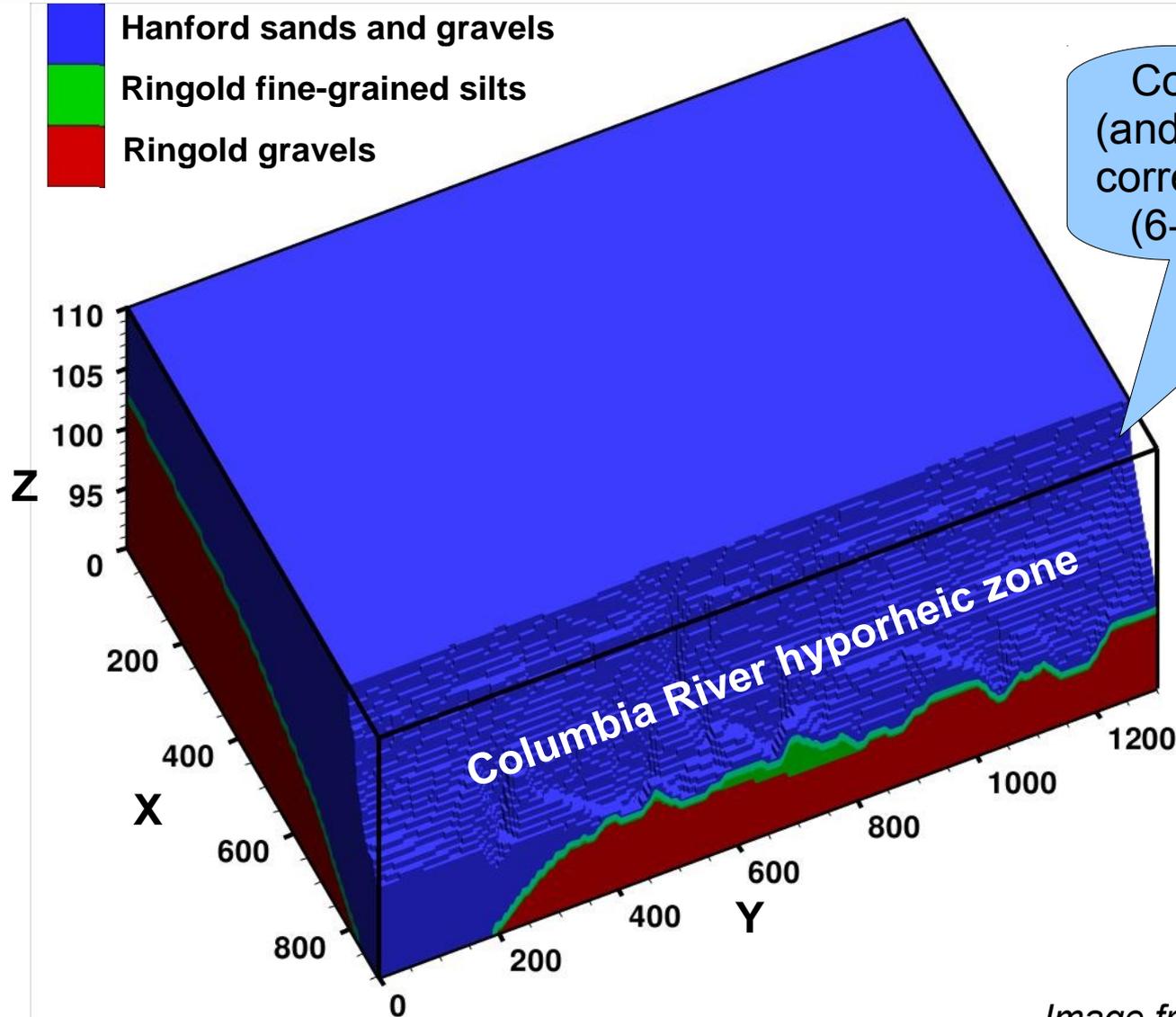


PFLOTRAN jugene@smp8192 trace analysis



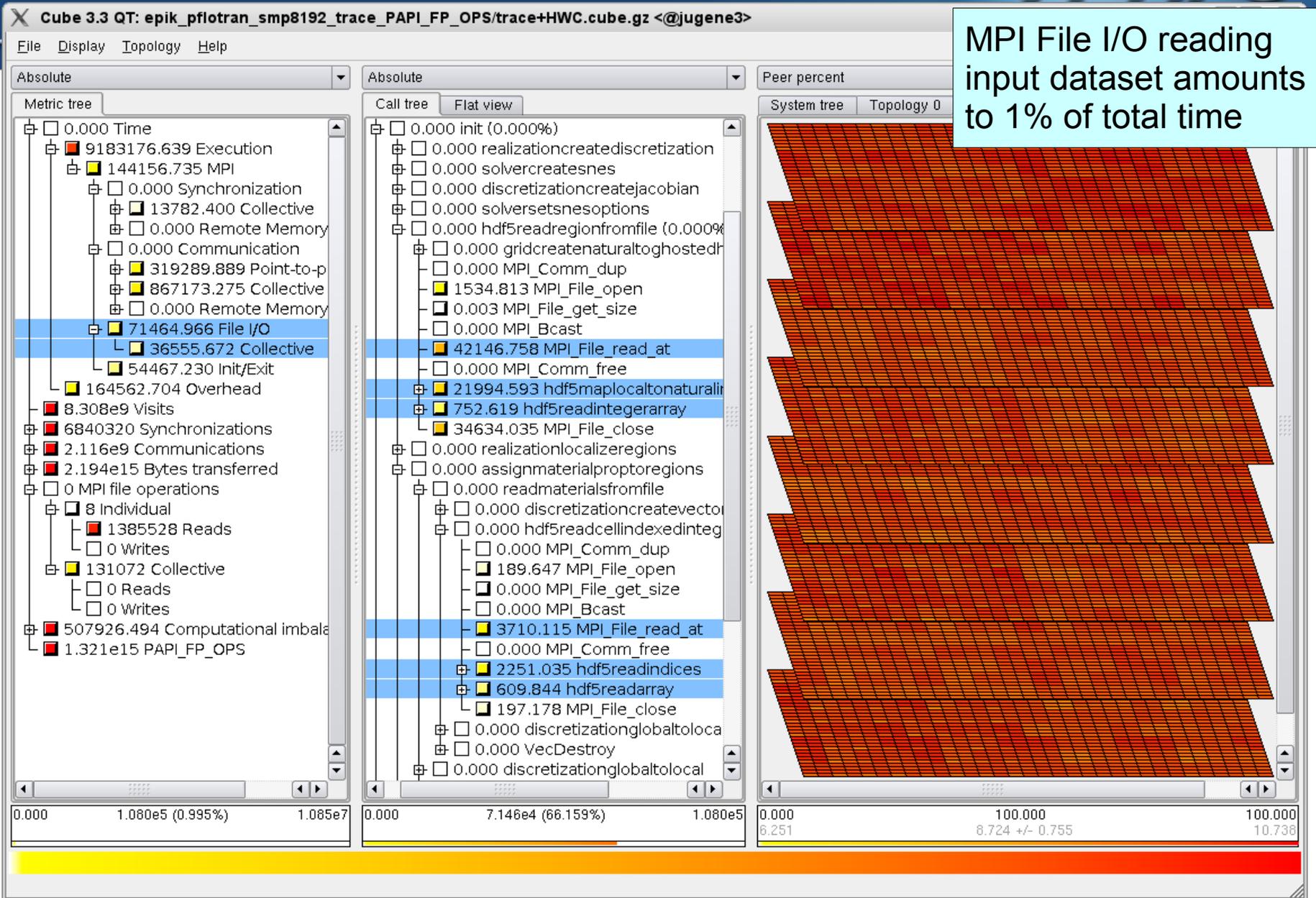
PFLOTRAN jugene@smp8192 trace analysis





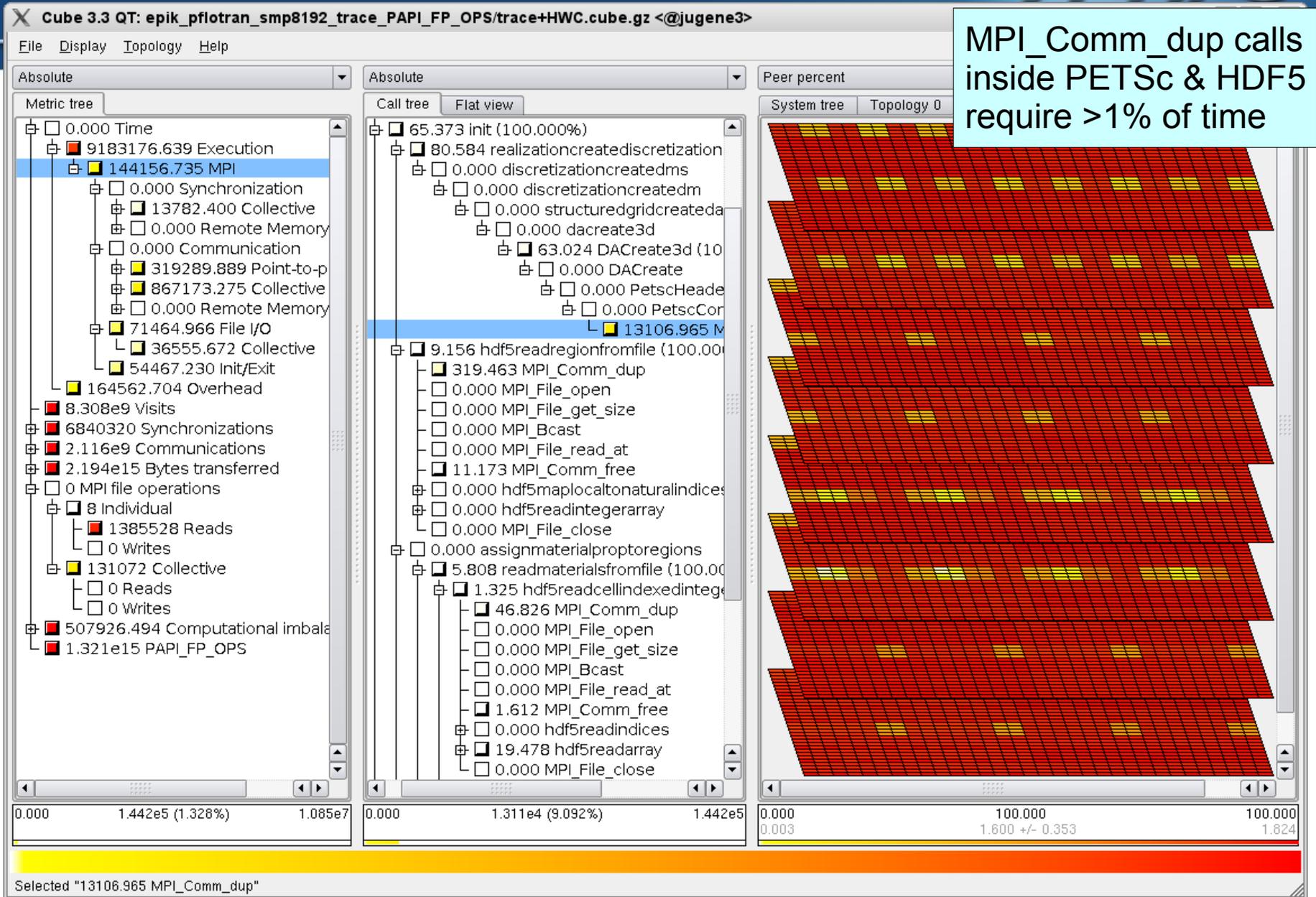
Computational imbalance
(and associated waiting time)
correlate to inactive grid cells
(6-7% of total) within river

PFLOTRAN jugene@smp8192 trace analysis



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

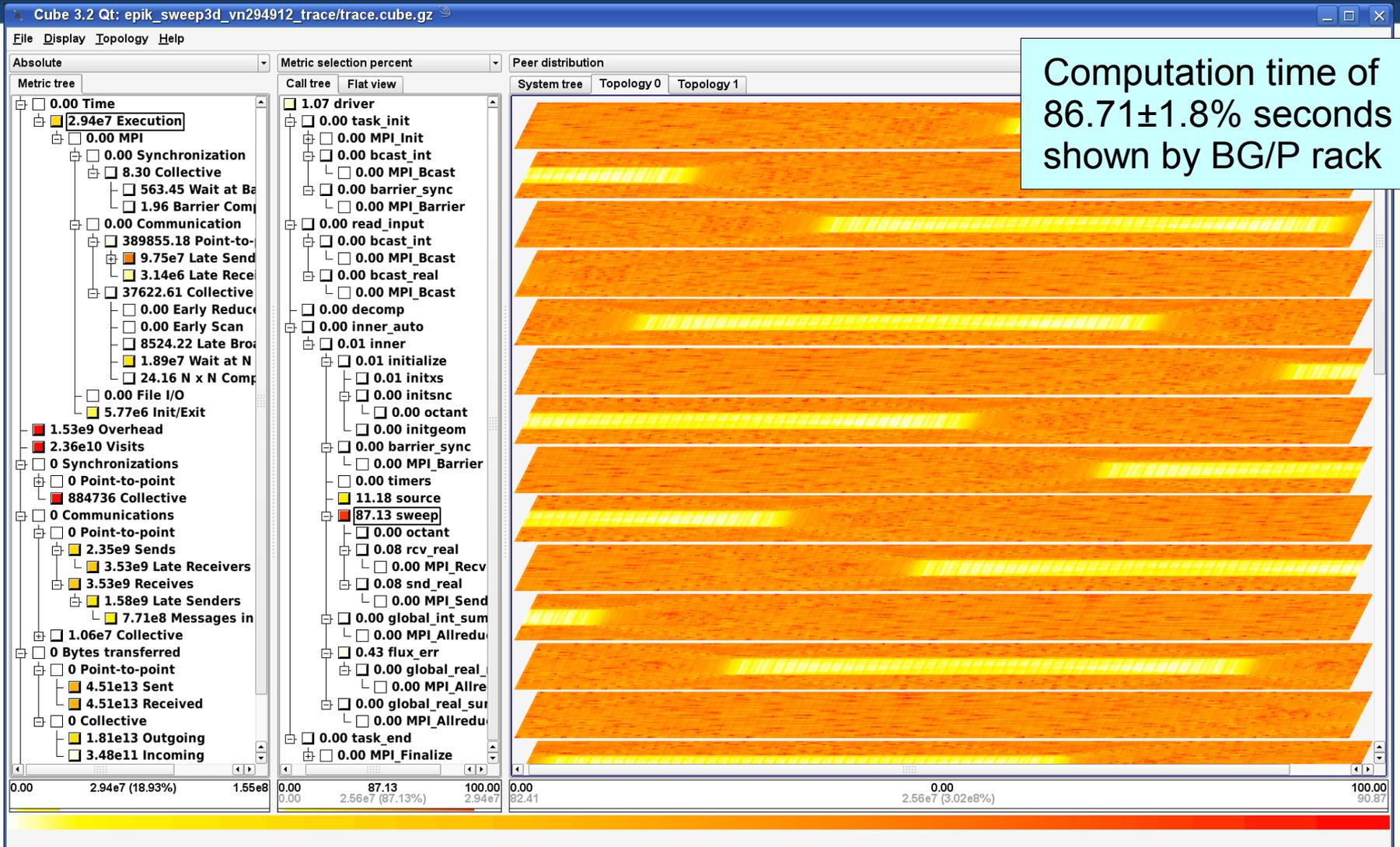
PFLOTRAN jugene@smp8192 trace analysis



- 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

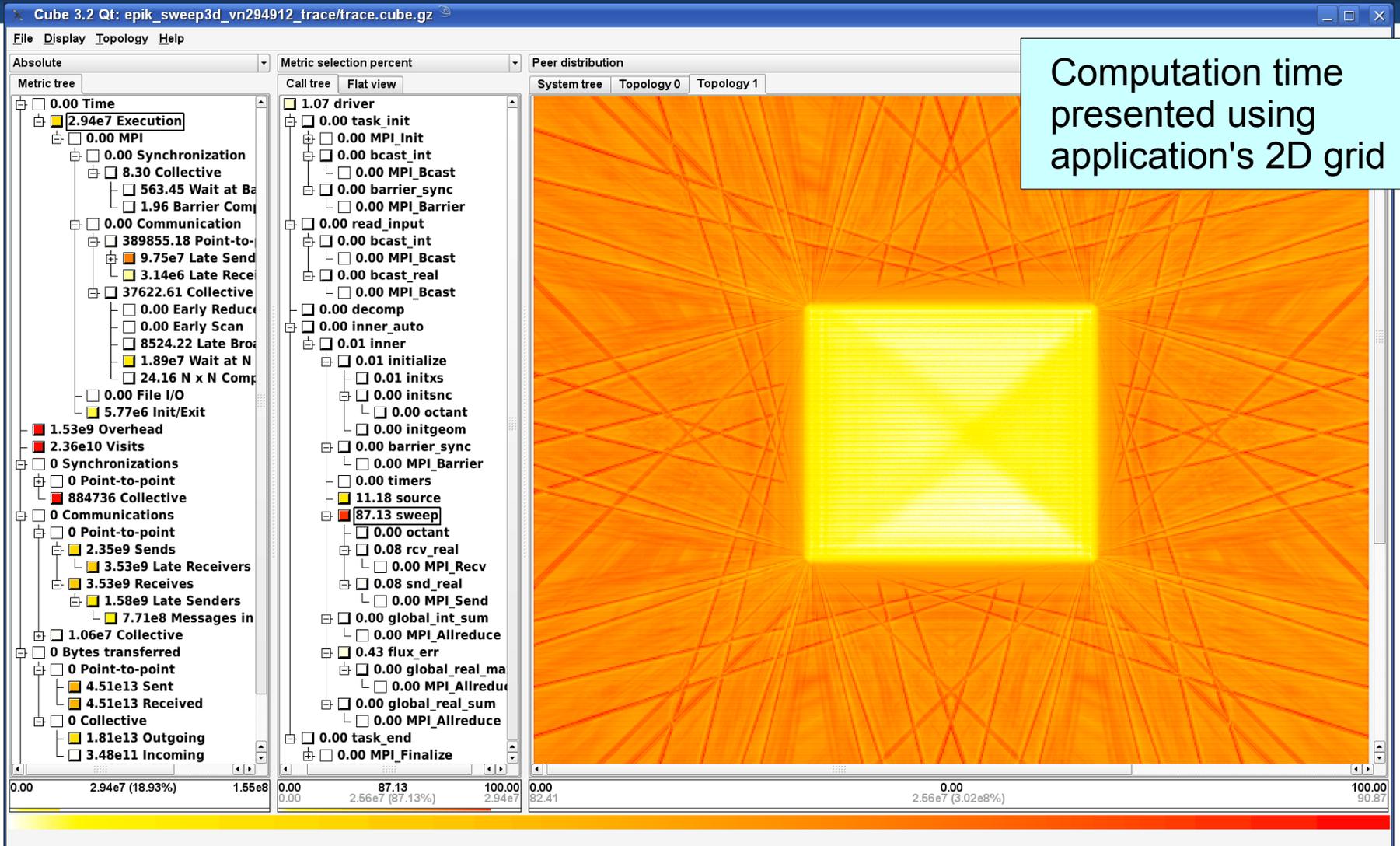
- 3D neutron transport simulation
 - ASC benchmark
 - direct order solve uses diagonal sweeps through grid cells
 - 'fixups' applied to correct unphysical (negative) fluxes
- MPI parallel version 2.2b using 2D domain decomposition
 - ~2,000 lines (12 source modules), mostly Fortran77
- Run on IBM BlueGene/P 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

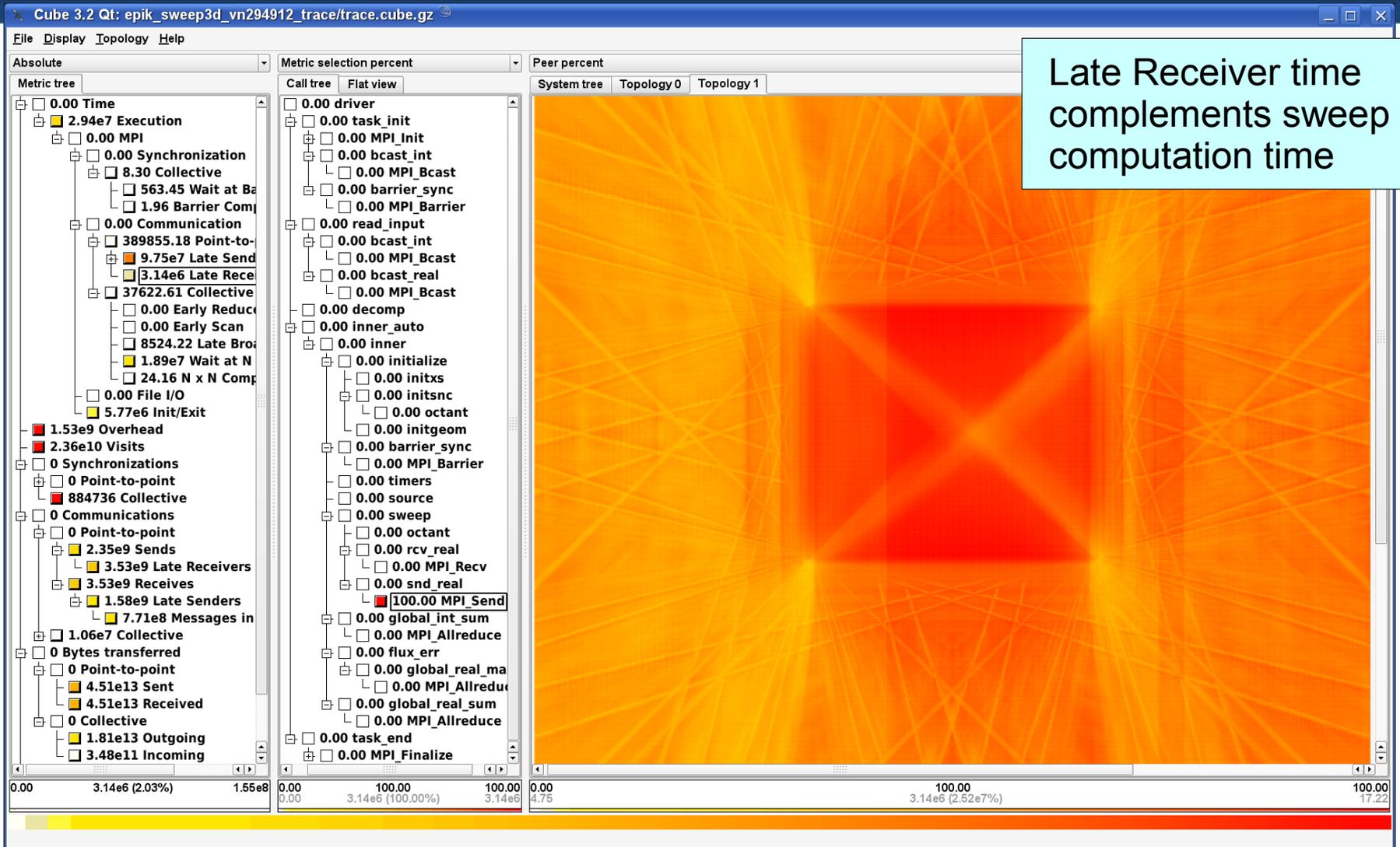
sweep3d on jugene@288k trace analysis



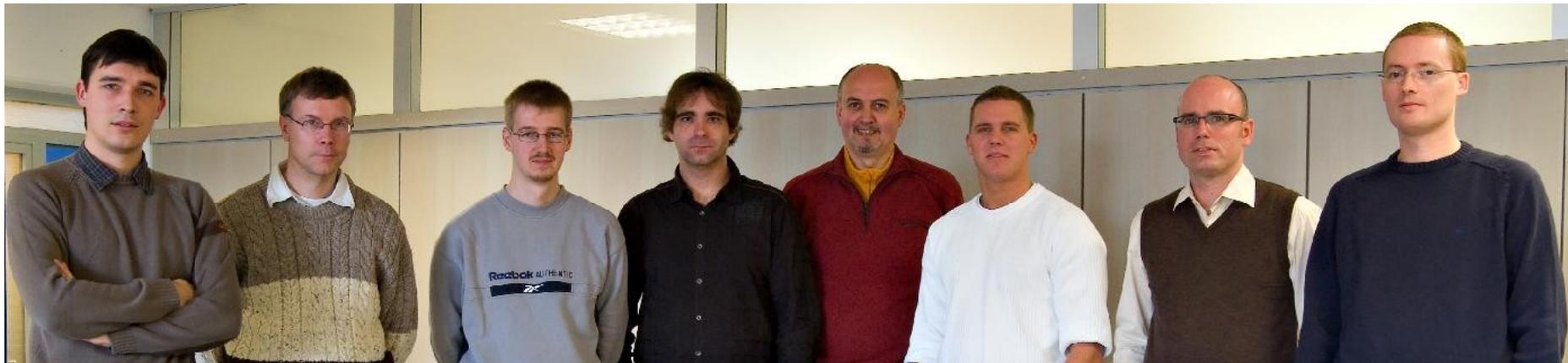
Computation time of $86.71 \pm 1.8\%$ seconds shown by BG/P rack

sweep3d on jugene@288k trace analysis





- 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, KSL, LRZ, NCAR, NCCS, NICS, RWTH, RZG, SARA, TACC, ZIH
 - ▶ Access & usage supported by European Union, German and other national funding organizations
- 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)