447.52 genbc

scalasca • Performance analysis & tuning case studies

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











PRODUCTIVITY



Additional Live-DVD example experiments



- 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)
 - marenostrum_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

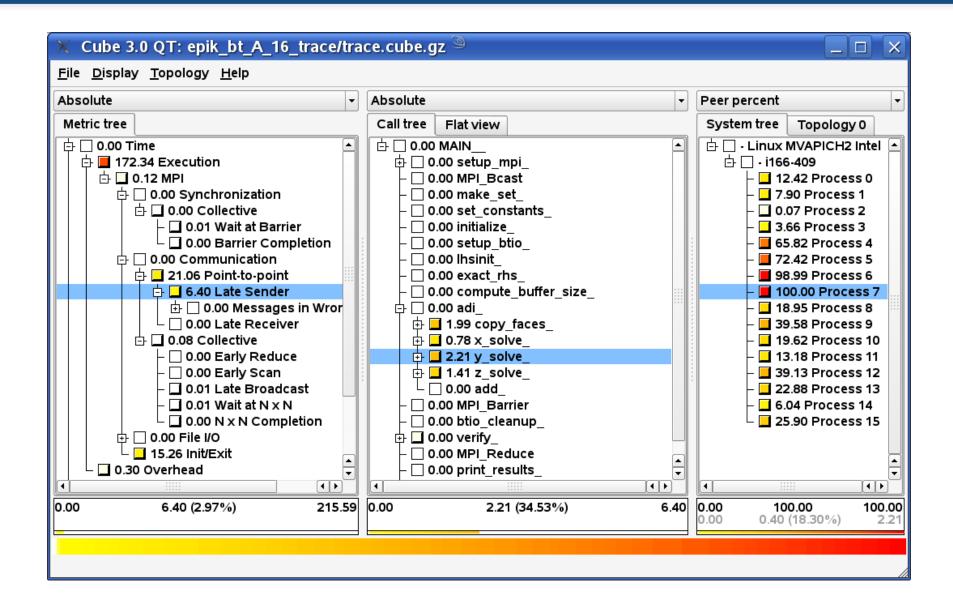
Scalasca NPB-BT experiments



- 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

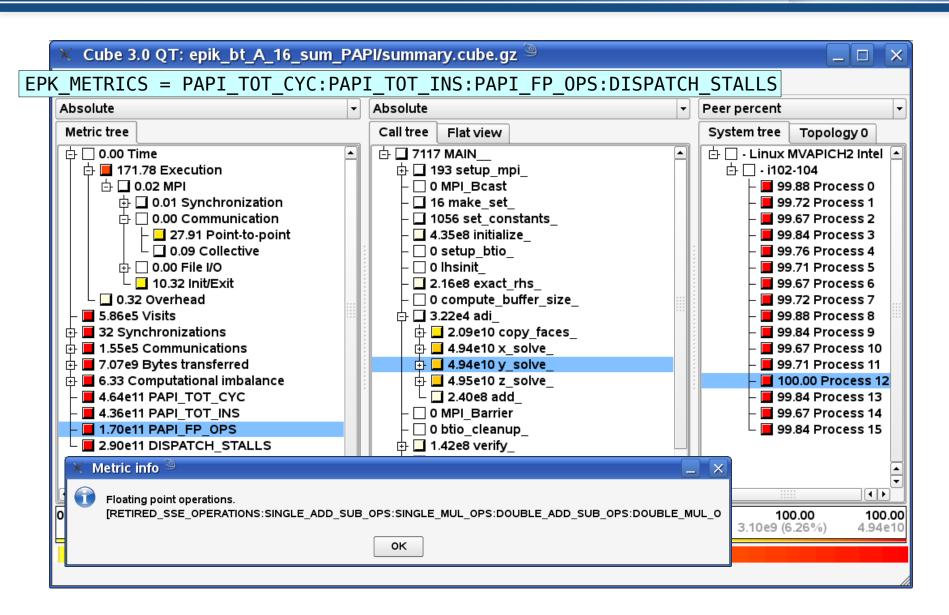
16-process trace analysis





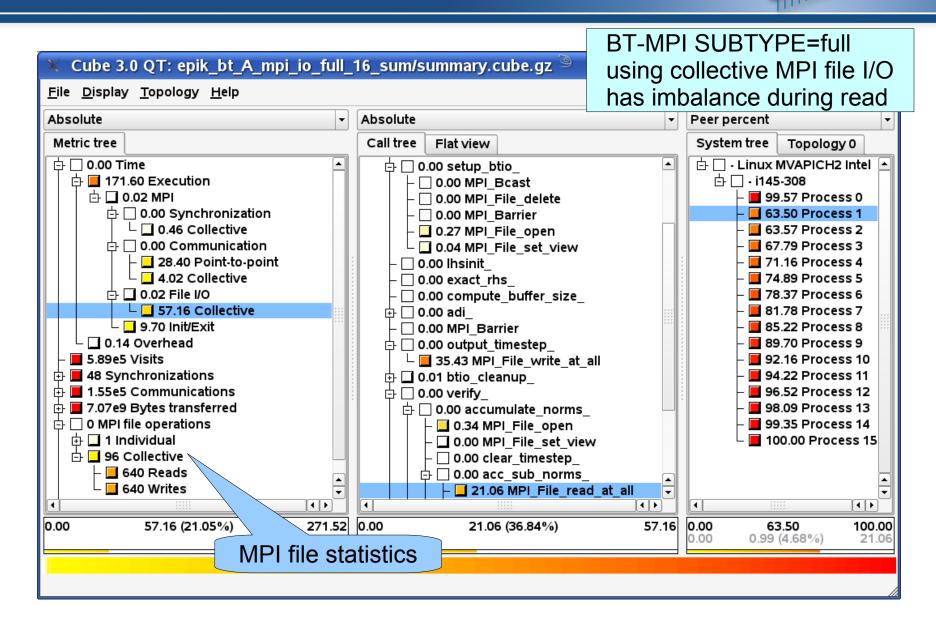
16-process summary analysis with HWC metrics VI-HPS





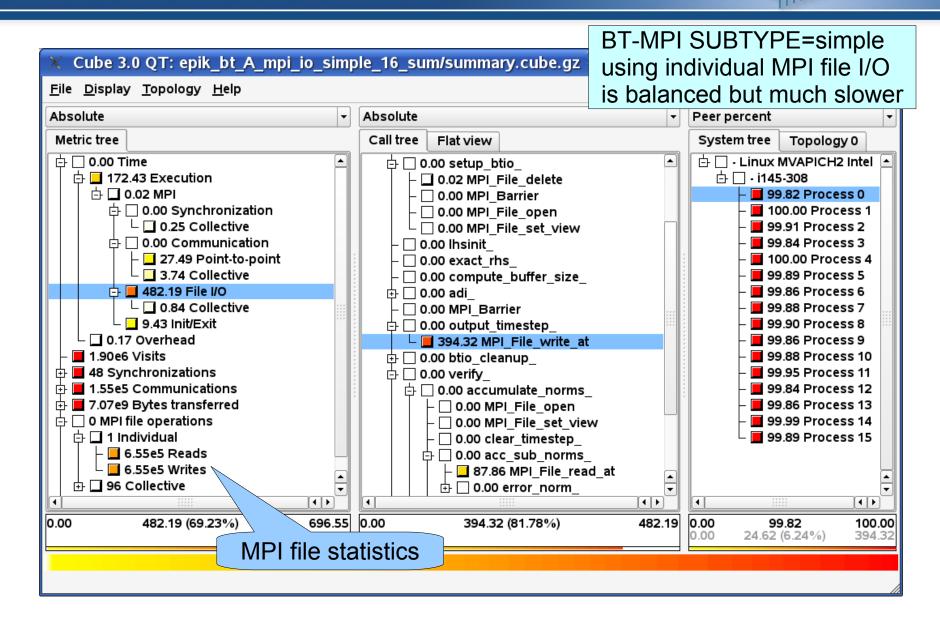
16-process summary analysis: MPI File I/O time





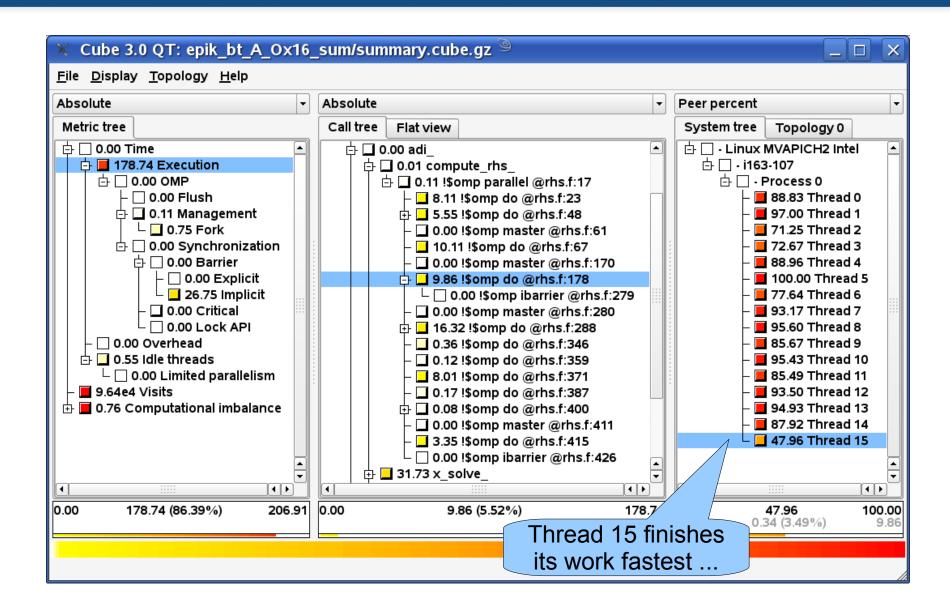
16-process summary analysis: MPI File I/O time





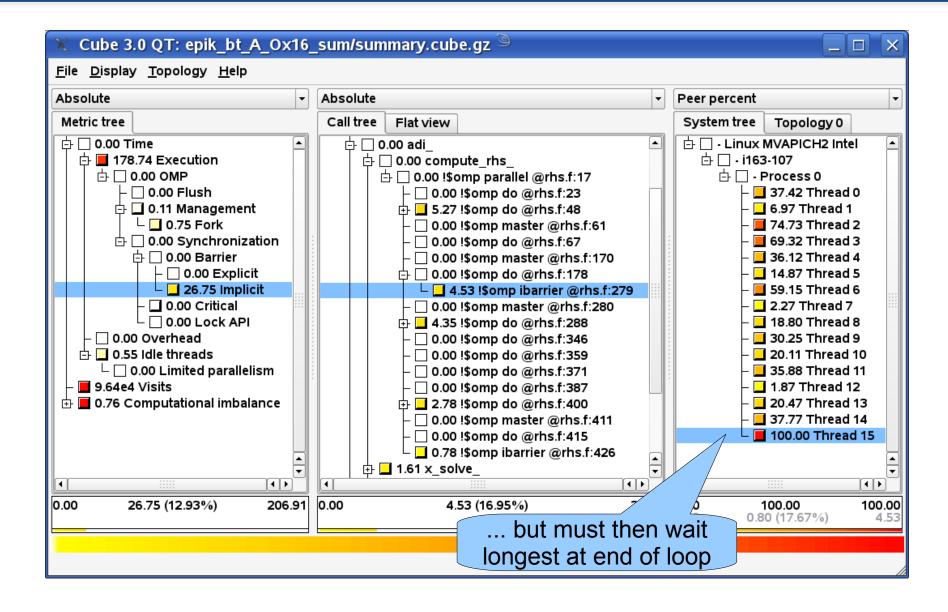
16-thread summary analysis: Execution time





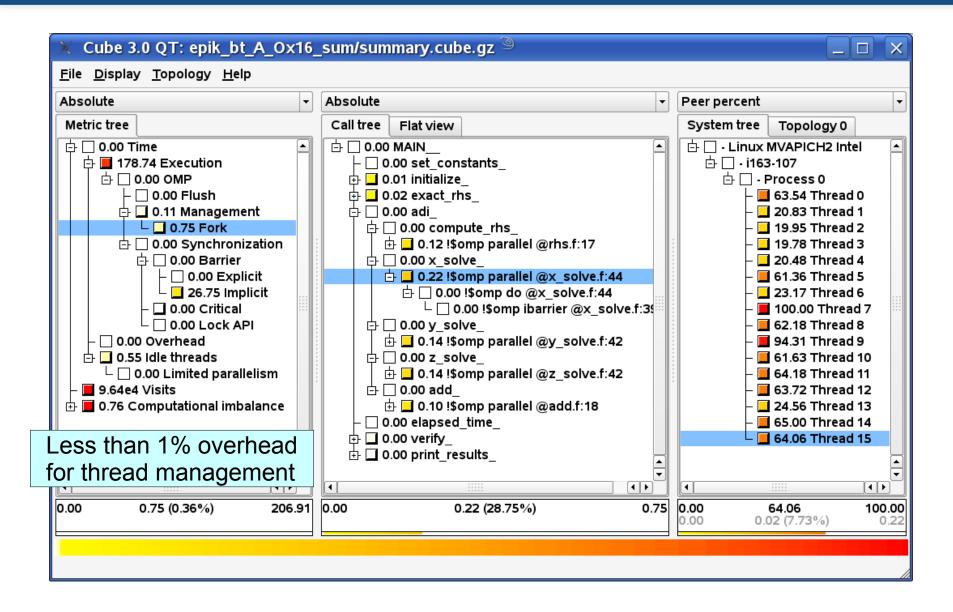
16-thread summary analysis: Implicit barrier time VI-HPS





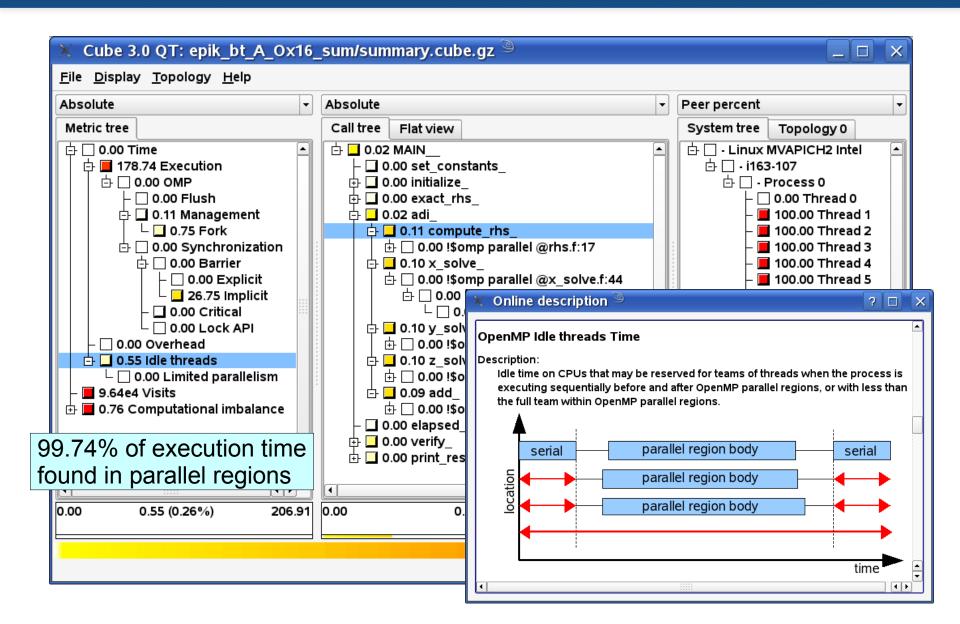
16-thread summary analysis: Thread fork time





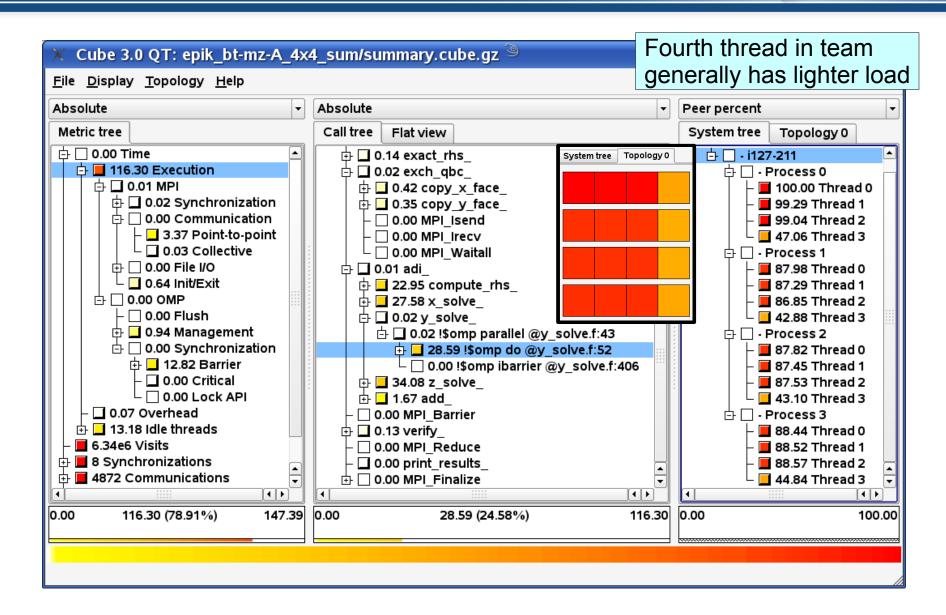
16-thread summary analysis: Idle threads time





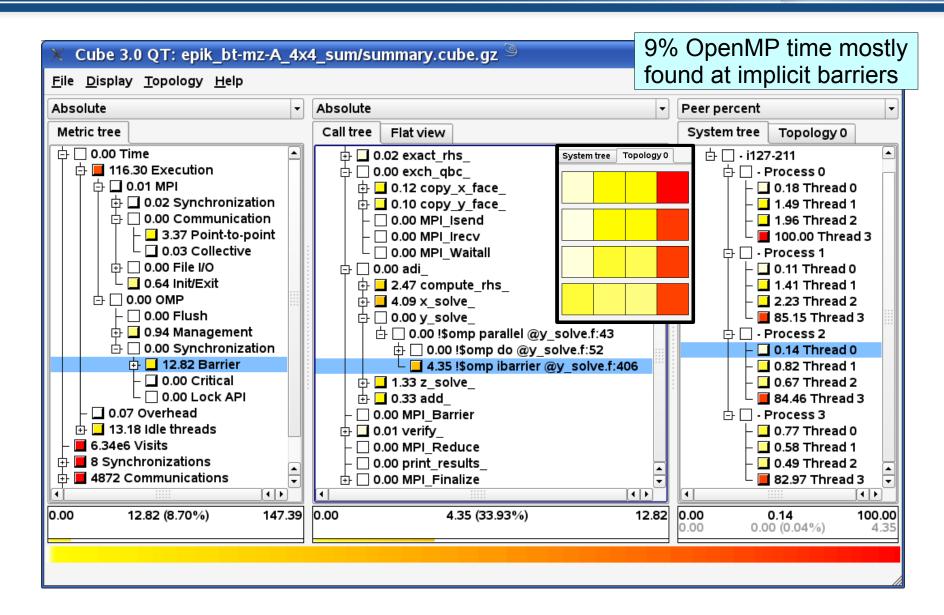
4x4 summary analysis: Execution time





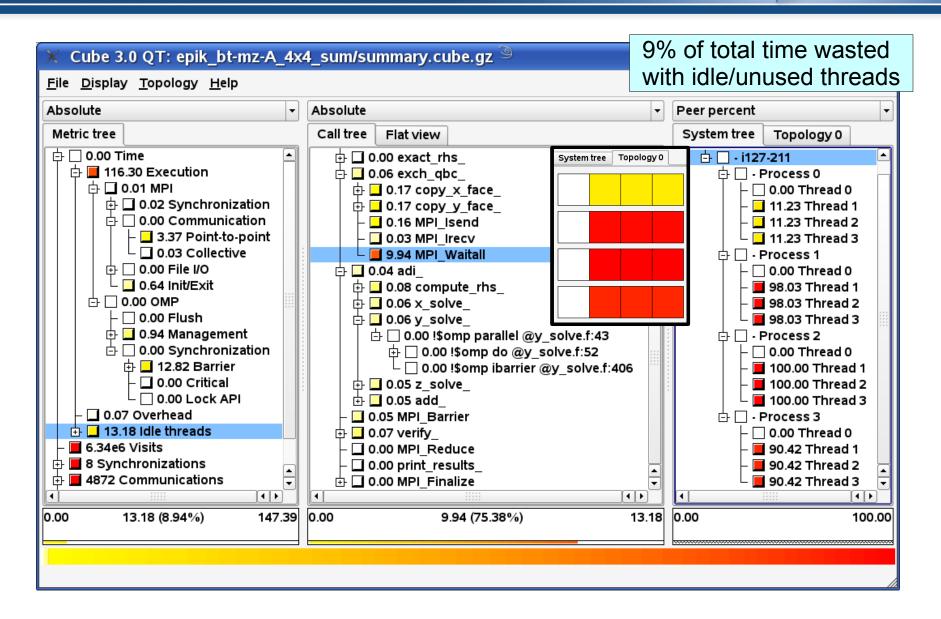
4x4 summary analysis: OpenMP time





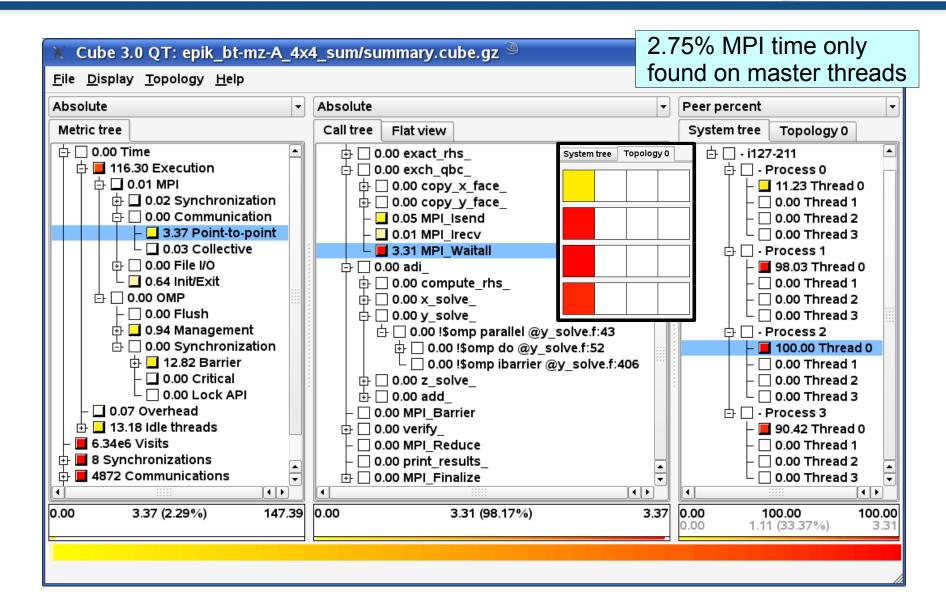
4x4 summary analysis: Idle threads time





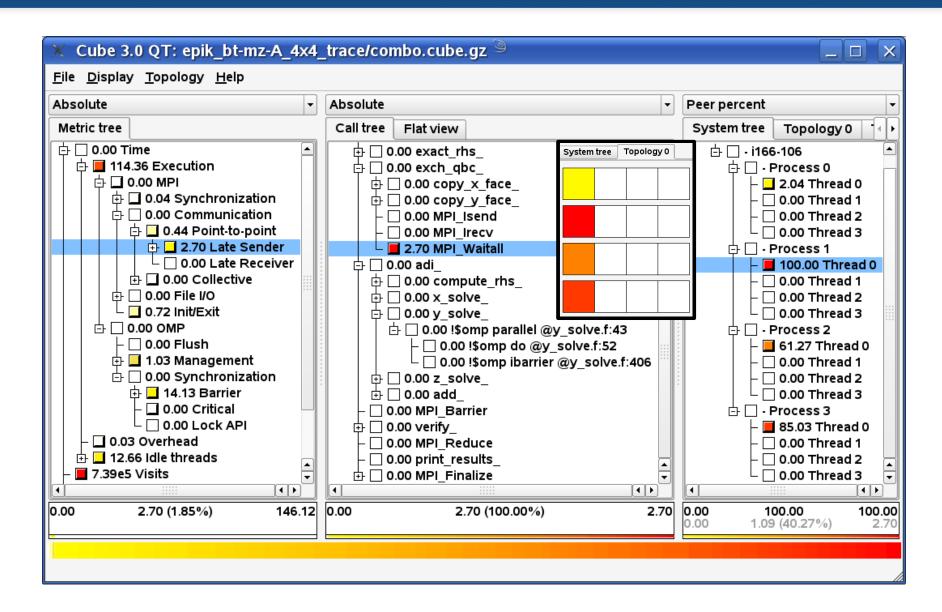
4x4 summary analysis: MPI time





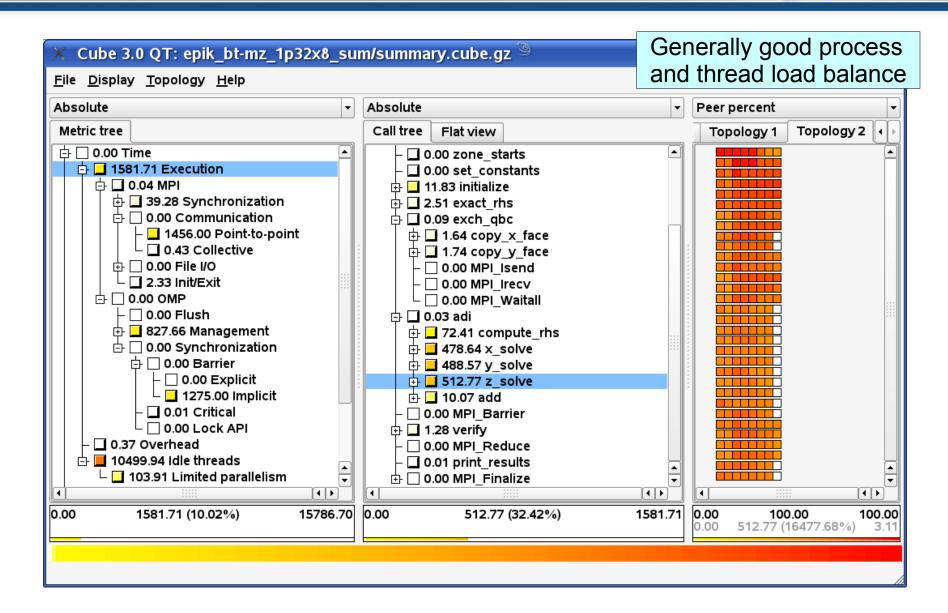
4x4 combined summary & trace analysis





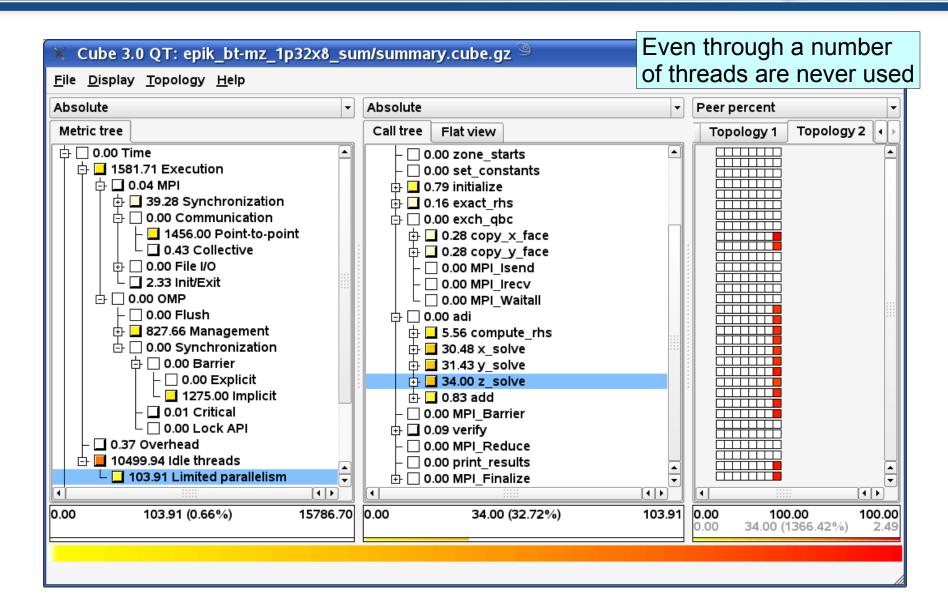
32x8 summary analysis: Excl. execution time





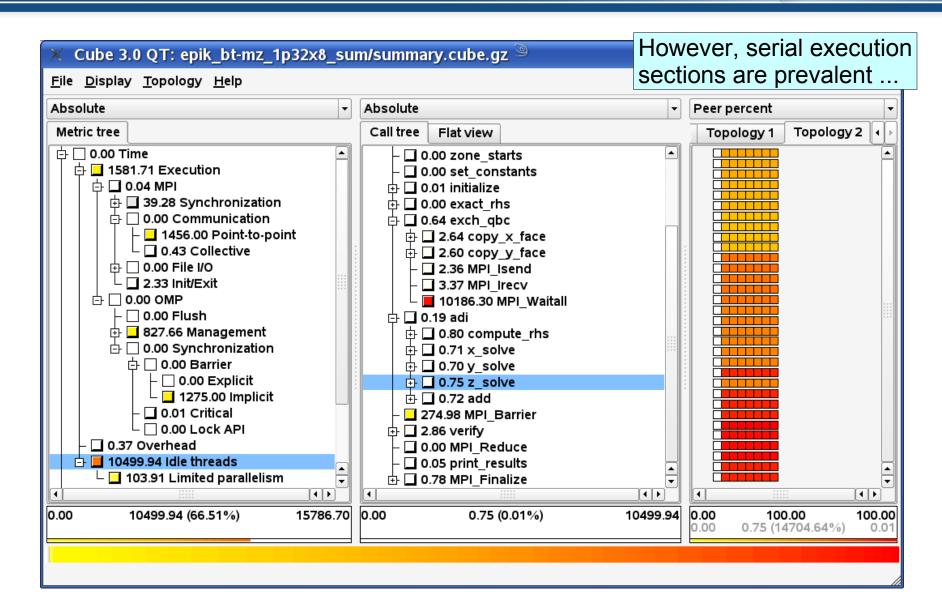
32x8 summary analysis: Limited parallelism



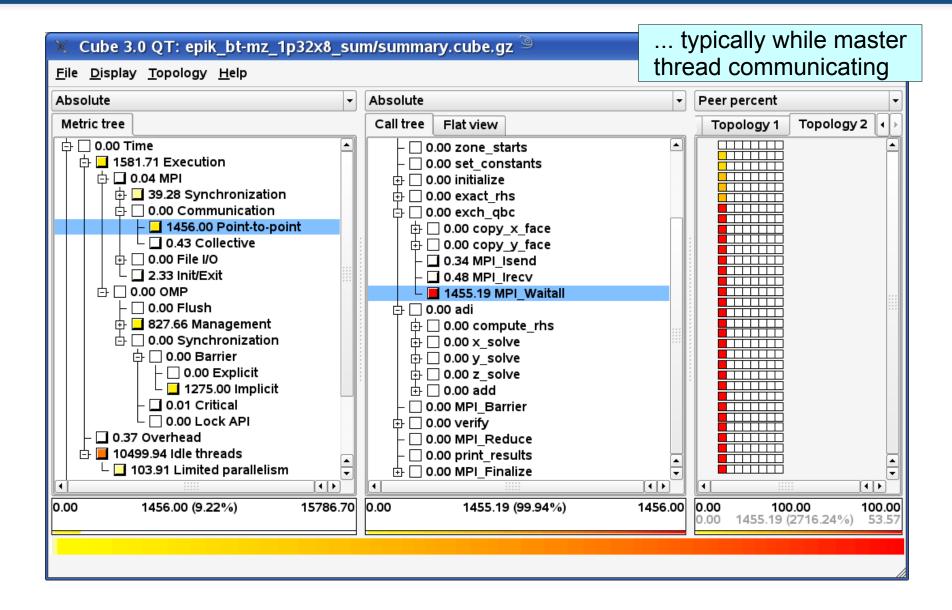


32x8 summary analysis: Idle threads time



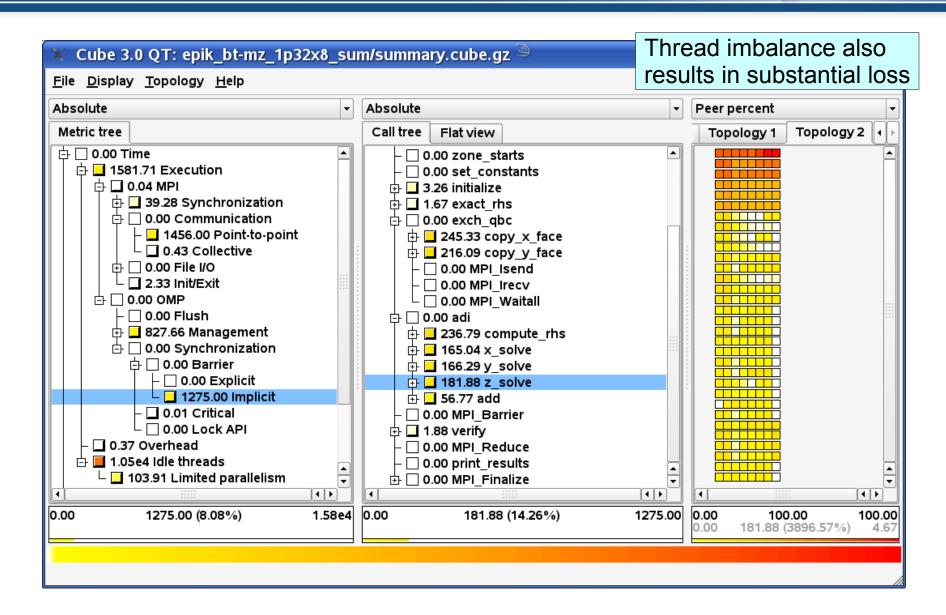


32x8 summary analysis: MPI communication time VI-HPS



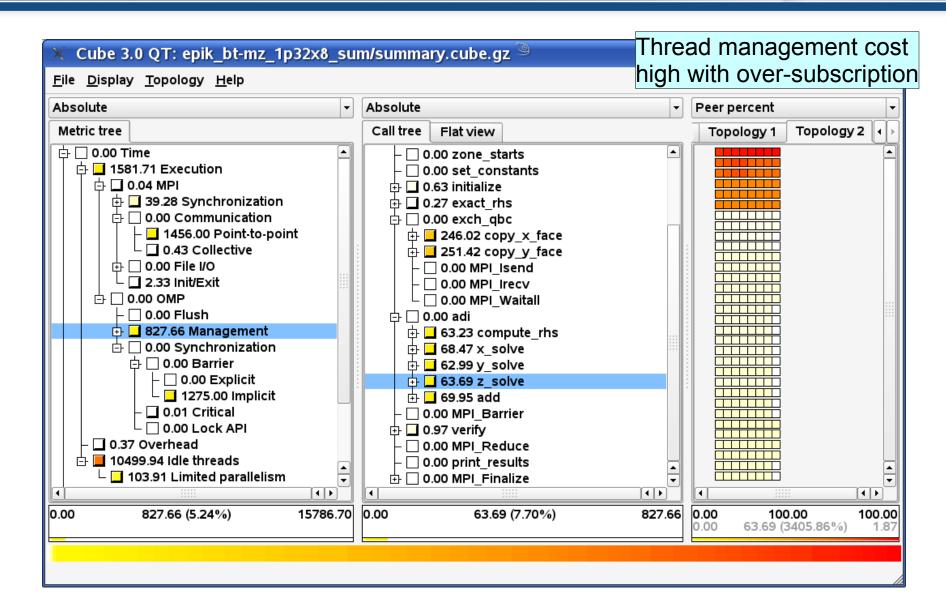
32x8 summary analysis: Implicit barrier time





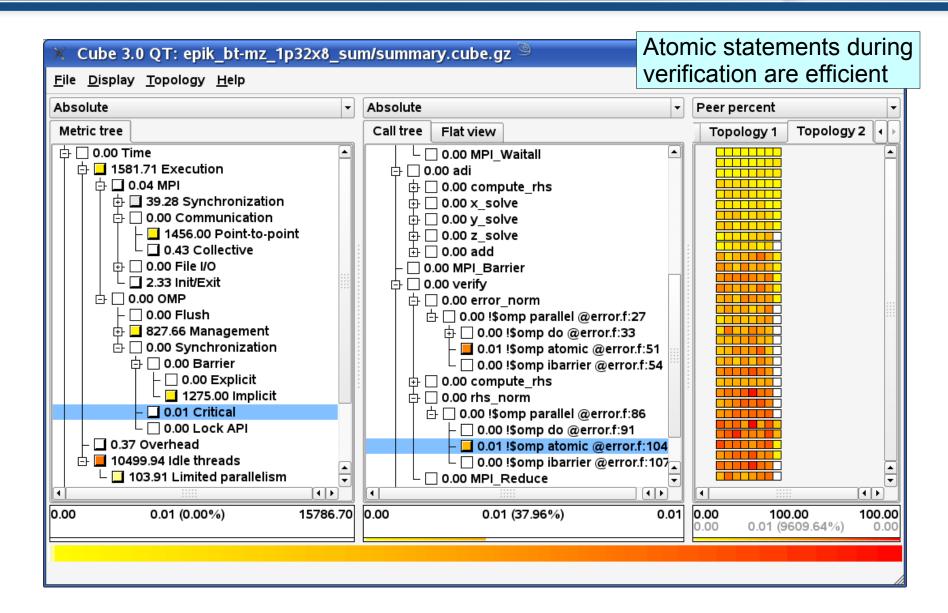
32x8 summary analysis: Thread management





32x8 summary analysis: Critical section time





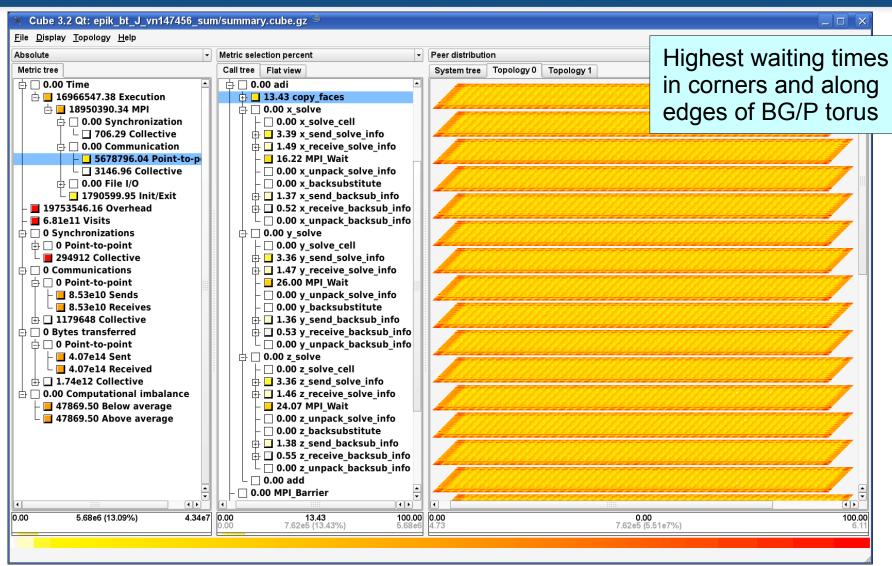
NPB-MPI-BT on BlueGene/P case study



- 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!

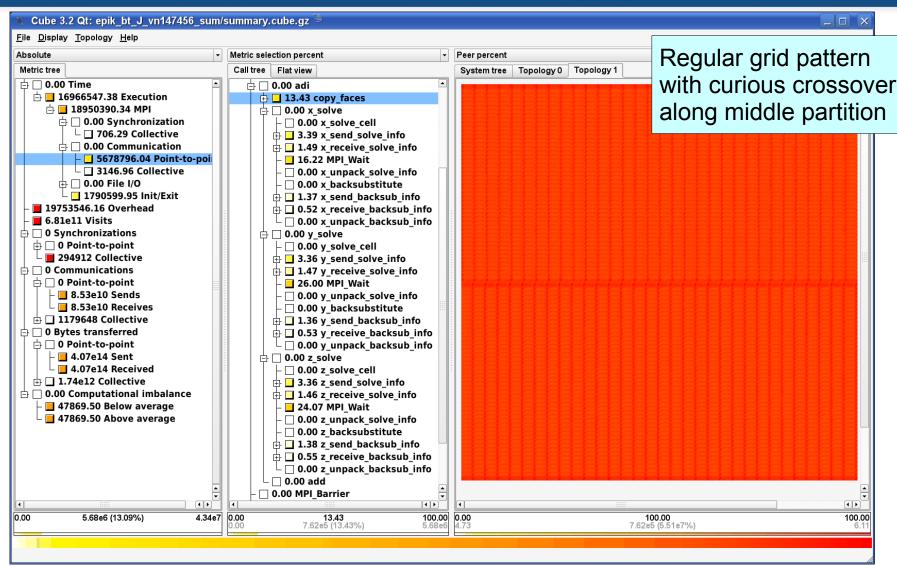
NPB-MPI-BT on jugene@144k summary analysis VI-HPS





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





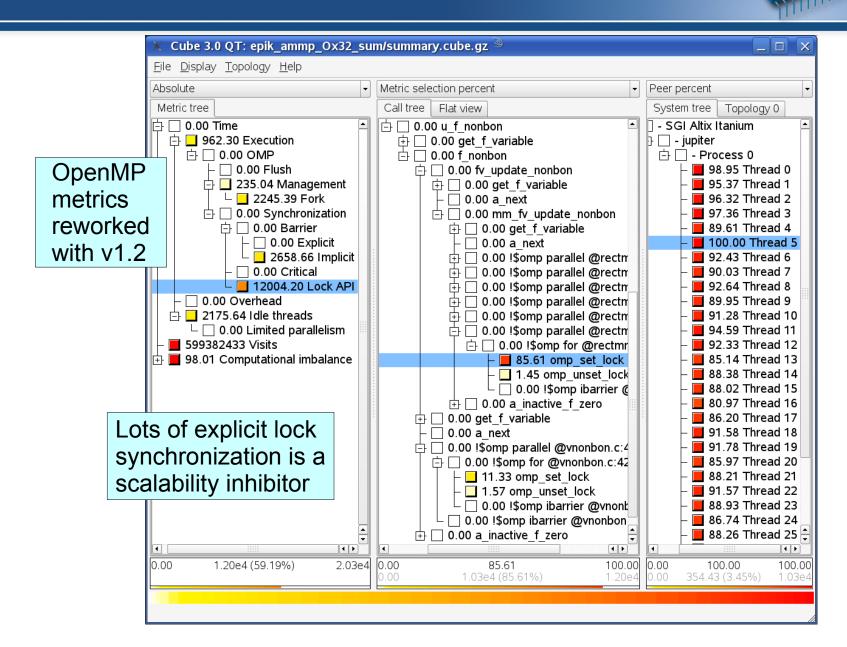
AMMP on Altix case study



- 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

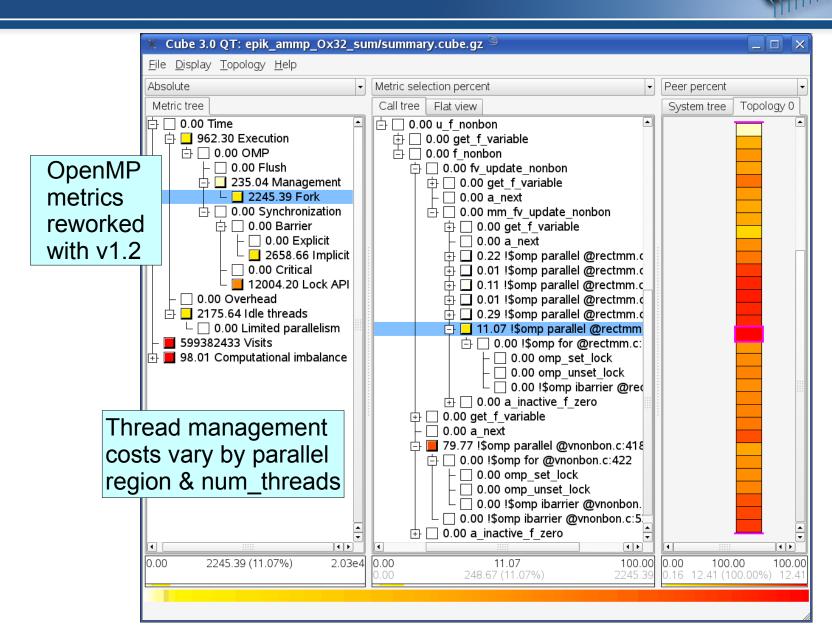
ammp on jupiter@32 OpenMP lock analysis





ammp on jupiter@32 OpenMP fork analysis





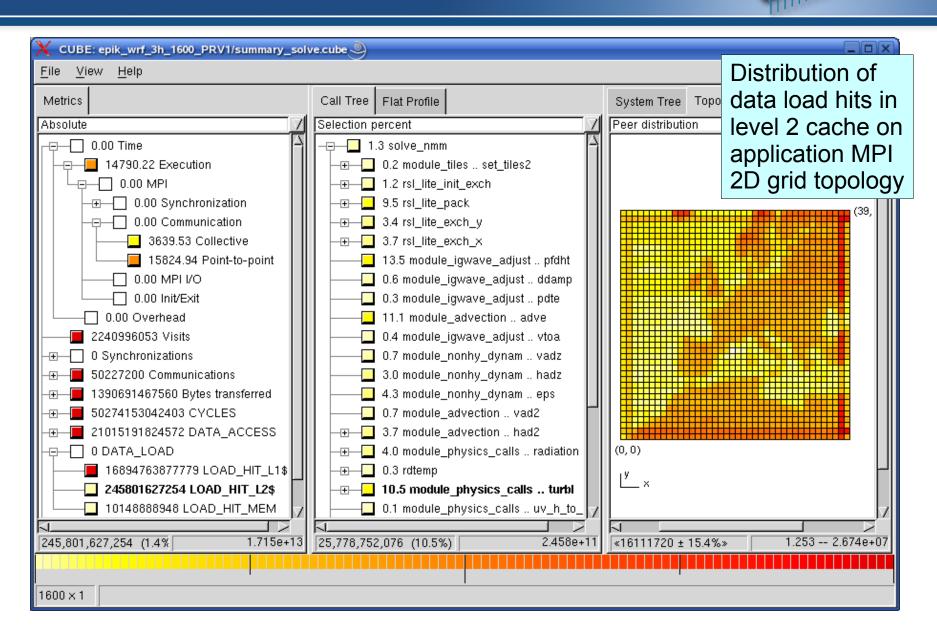
WRF/MareNostrum case study



- 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

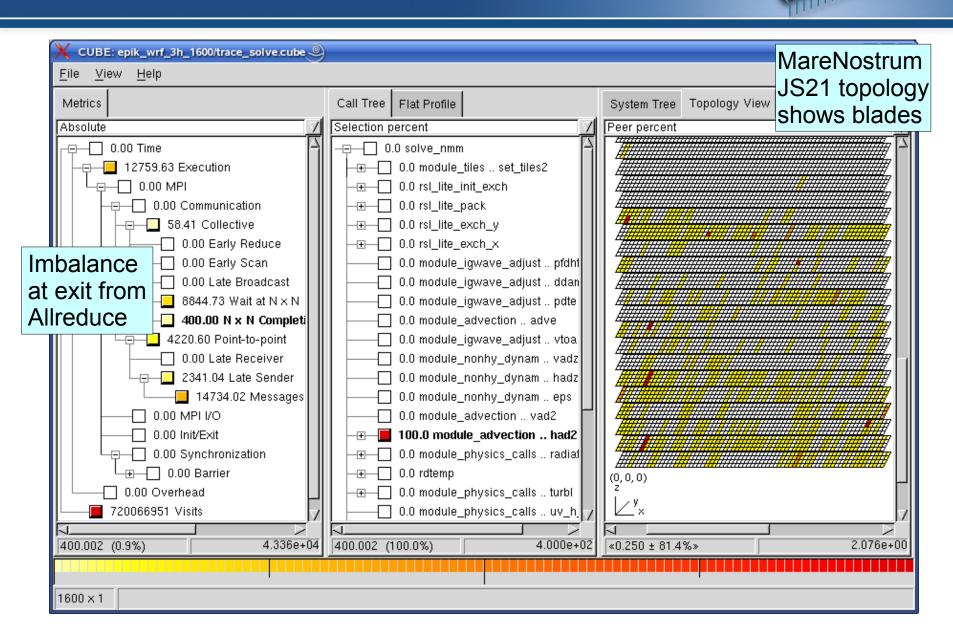
WRF on MareNostrum@1600 with HWC metrics





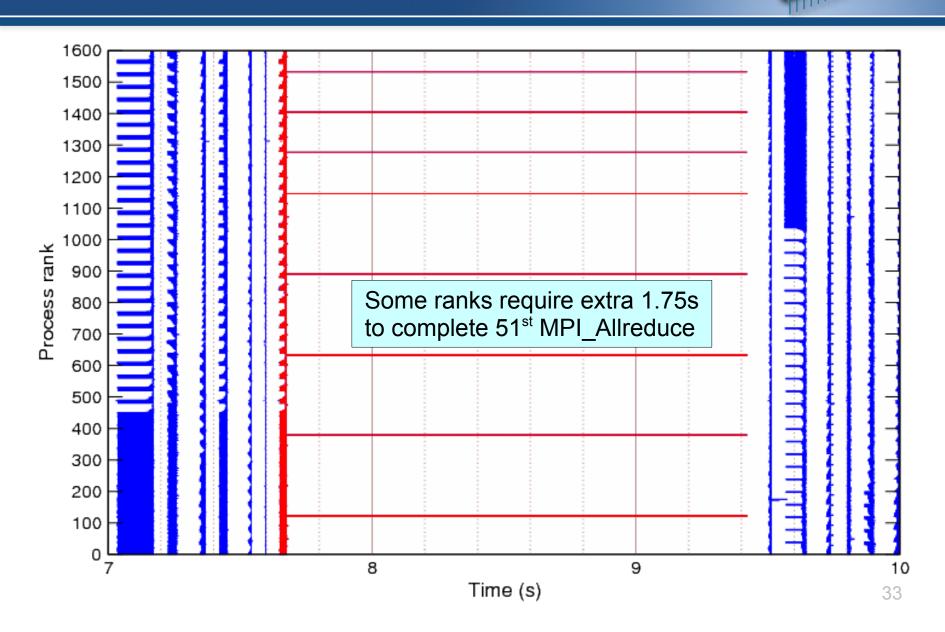
WRF on MareNostrum@1600 trace analysis





WRF on MareNostrum@1600 time-line extract





WRF/MareNostrum experience

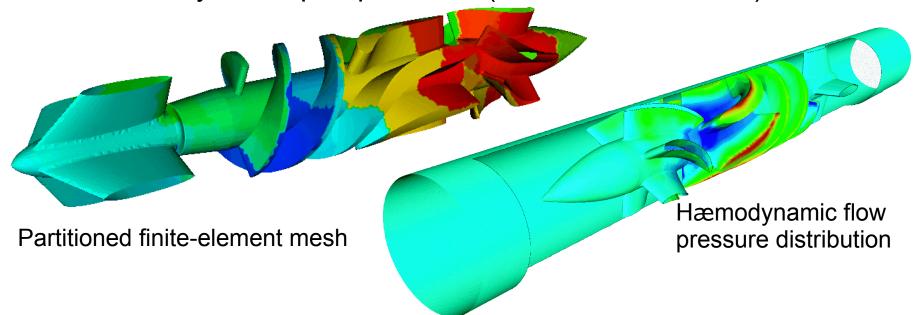


- Limited system I/O requires careful management
 - Selective instrumentation and measurement filtering
- PowerPC hardware counter metrics included in summary
 - Memory/cache data access hierarchy constructed
- Automated trace analysis quantified impact of imbalanced exit from MPI_Allreduce in "NxN completion 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

XNS on BlueGene/L case study

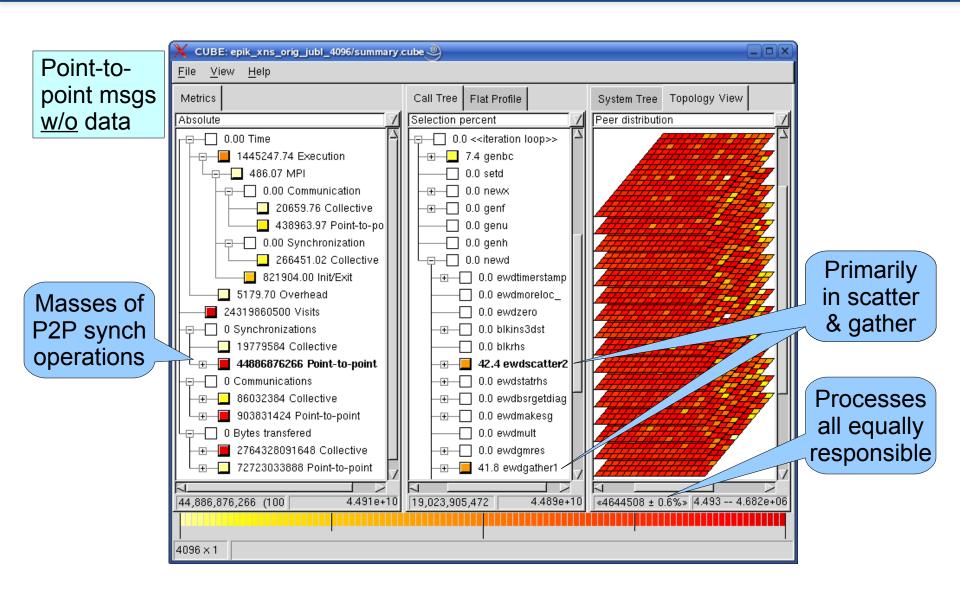


- CFD simulation of unsteady flows
 - developed by RWTH CATS group of Marek Behr
 - exploits finite-element techniques, unstructured 3D meshes, iterative solution strategies
- MPI parallel version (Dec-2006)
 - >40,000 lines of Fortran & C
 - DeBakey blood-pump dataset (3,714,611 elements)



XNS-DeBakey on jubl@4096 summary analysis





XNS-DeBakey scalability on BlueGene/L





XNS on BlueGene/L experience



- 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

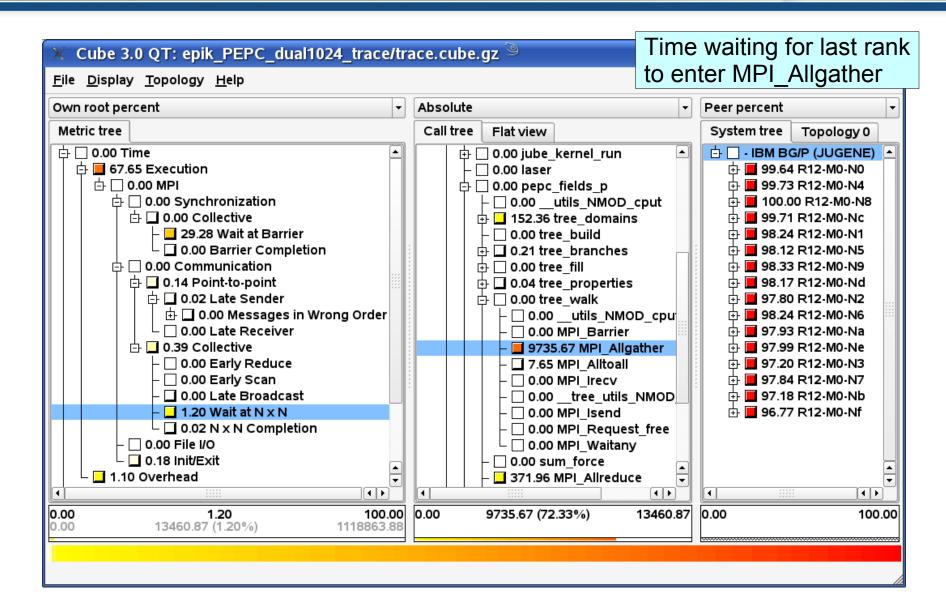
PEPC-B on BG/P & Cray XT case study



- 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

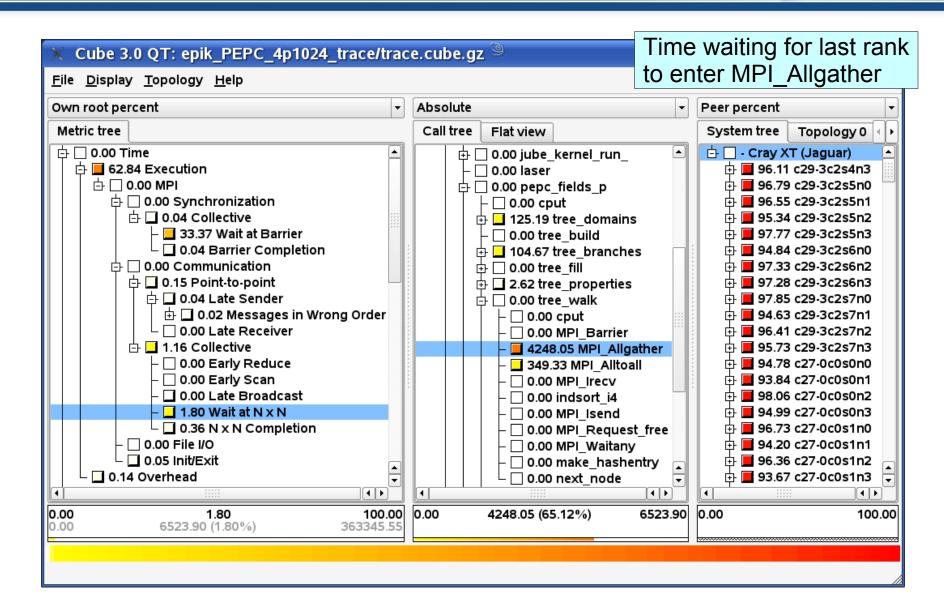
PEPC@1024 on BlueGene/P: Wait at NxN time





PEPC@1024 on Cray XT4: Wait at NxN time





PEPC-B on BG/P & Cray XT experience



- 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

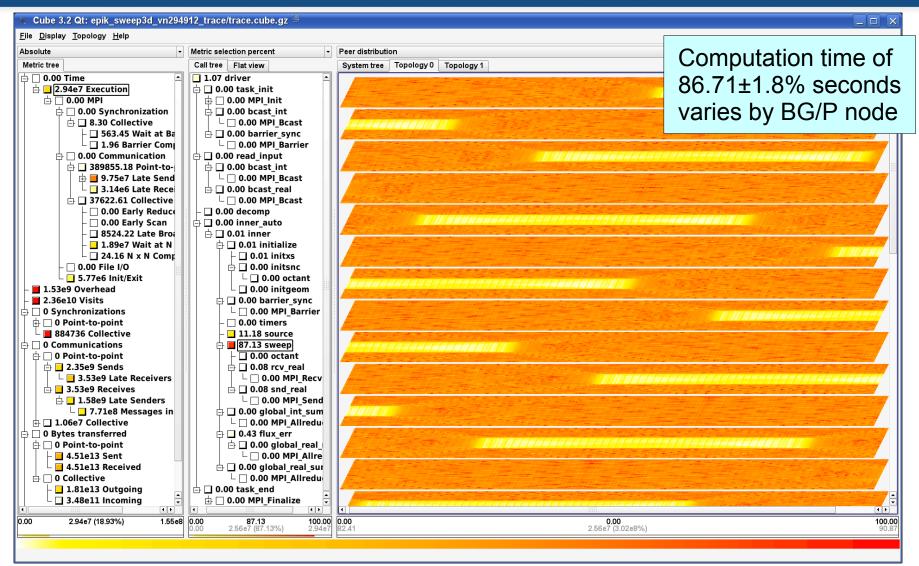
Sweep3d on BlueGene/P case study



- 3D neutron transport simulation
 - ASC benchmark
 - direct order solve uses diagonal sweeps through grid cells
- 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
 - 790GB trace written in 47 minutes, analyzed in 7 minutes
 - ▶ plus 86 minutes just to create 294,912 files (one per MPI rank)
 - ► SIONlib being developed to address management of sets of files
 - Mapping of 576x512 grid of processes onto 3D physical torus results in regular pattern of performance artifacts

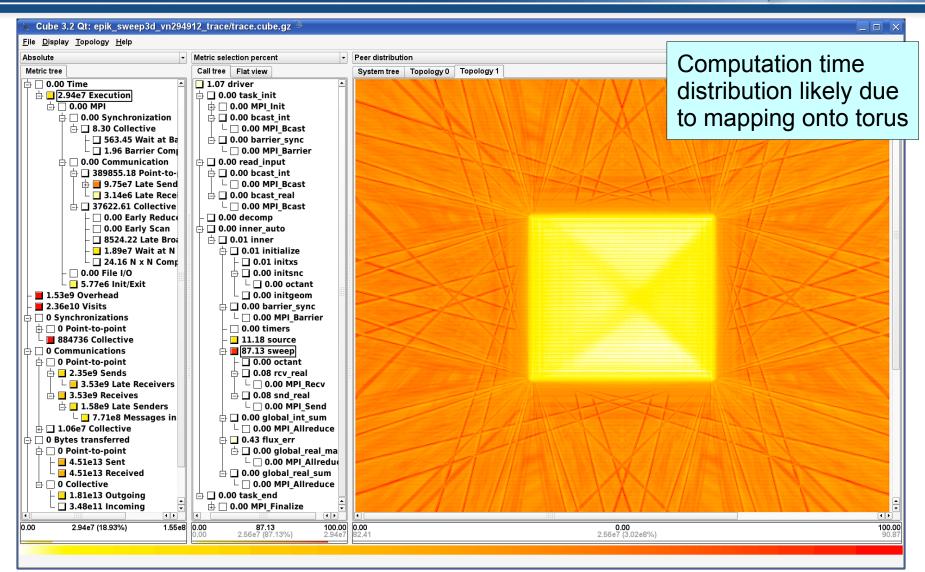
sweep3d on jugene@288k trace analysis





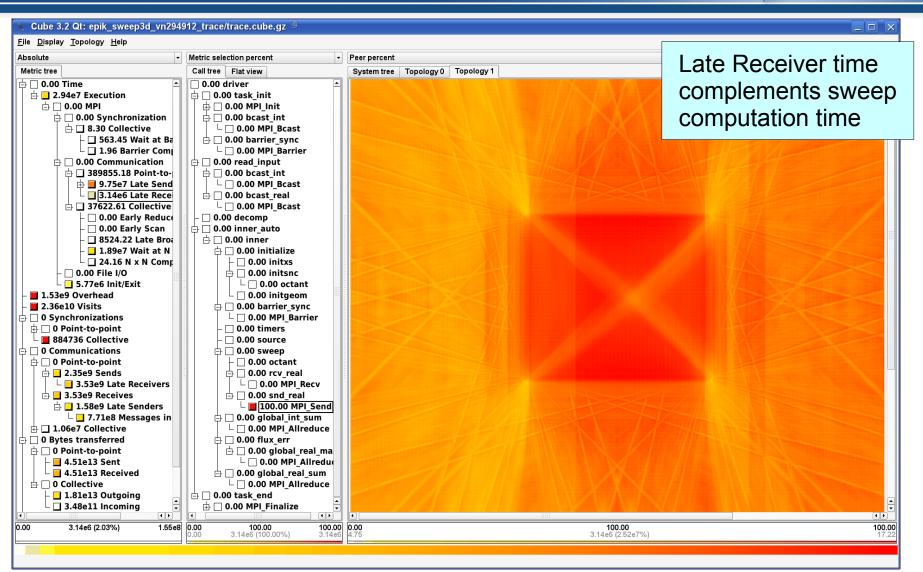
sweep3d on jugene@288k trace analysis





sweep3d on jugene@288k trace (wait) analysis





Acknowledgements



- 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
- 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

