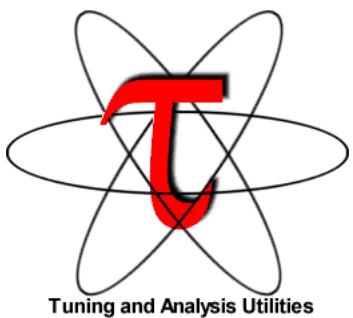


TAU Performance System®

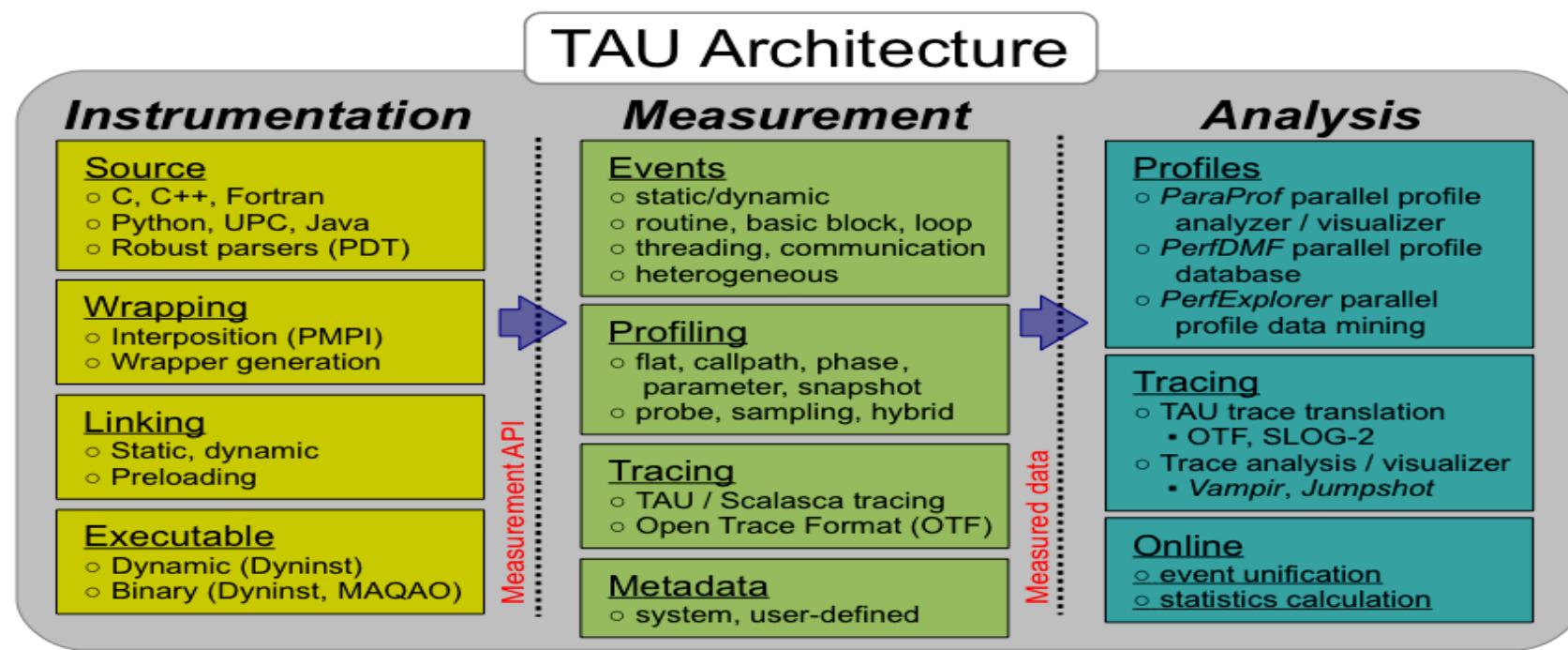
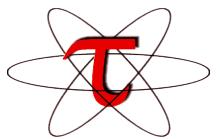


Sameer Shende
sameer@cs.uoregon.edu
Performance Research Laboratory
University of Oregon
<http://tau.uoregon.edu>



TAU Performance System®

- Parallel performance framework and toolkit
 - Supports all HPC platforms, compilers, runtime system
 - Provides portable instrumentation, measurement, analysis



TAU Performance System

- Instrumentation
 - Fortran, C++, C, UPC, Java, Python, Chapel
 - Automatic instrumentation
- Measurement and analysis support
 - MPI, OpenSHMEM, ARMCI, PGAS, DMAPP
 - pthreads, OpenMP, OMPT interface, hybrid, other thread models
 - GPU, CUDA, OpenCL, OpenACC
 - Parallel profiling and tracing
 - Use of Score-P for native OTF2 and CUBEX generation
 - Efficient callpath profiles and trace generation using Score-P
- Analysis
 - Parallel profile analysis (ParaProf), data mining (PerfExplorer)
 - Performance database technology (TAUdb)
 - 3D profile browser

TAU Performance System

- Instrumentation
 - Fortran, C++, C, UPC, Java, Python, Chapel
 - Automatic instrumentation
- Measurement and analysis support
 - MPI, OpenSHMEM, ARMCI, PGAS, DMAPP
 - pthreads, OpenMP, OMPT interface, hybrid, other thread models
 - GPU, CUDA, OpenCL, OpenACC
 - Parallel profiling and tracing
 - Use of Score-P for native OTF2 and CUBEX generation
 - Efficient callpath profiles and trace generation using Score-P
- Analysis
 - Parallel profile analysis (ParaProf), data mining (PerfExplorer)
 - Performance database technology (TAUdb)
 - 3D profile browser

TAU Performance System

- TAU supports both sampling and direct instrumentation
- Memory debugging as well as I/O performance evaluation
- Profiling as well as tracing
- Interfaces with Score-P for more efficient measurements
- TAU's instrumentation covers:
 - Runtime library interposition (`tau_exec`)
 - Compiler-based instrumentation
 - PDT based Source level instrumentation: routine & loop
 - Event based sampling (`TAU_SAMPLING=1` or `tau_exec -ebs`)
 - Callstack unwinding with sampling (`TAU_EBS_UNWIND=1`)
 - OpenMP Tools Interface (OMPT, `tau_exec -T ompt`)
 - CUDA CUPTI, OpenCL (`tau_exec -T cupti -cupti`)

Application Performance Engineering using TAU

- How much time is spent in each application routine and outer *loops*? Within loops, what is the contribution of each *statement*? What is the time spent in OpenMP loops?
- How many instructions are executed in these code regions?
Floating point, Level 1 and 2 *data cache misses*, hits, branches taken? What is the extent of vectorization for loops on Intel MIC?
- What is the memory usage of the code? When and where is memory allocated/de-allocated? Are there any memory leaks? What is the memory footprint of the application? What is the memory high water mark?
- How much energy does the application use in Joules? What is the peak power usage?
- What are the I/O characteristics of the code? What is the peak read and write *bandwidth* of individual calls, total volume?
- What is the contribution of each *phase* of the program? What is the time wasted/spent waiting for collectives, and I/O operations in Initialization, Computation, I/O phases?
- How does the application *scale*? What is the efficiency, runtime breakdown of performance across different core counts?

Using TAU's Source Code Instrumentation

- TAU supports several compilers, measurement, and thread options
 - Intel compilers, profiling with hardware counters using PAPI, MPI library, CUDA...
 - Each measurement configuration of TAU corresponds to a unique stub makefile (configuration file) and library that is generated when you configure it
- To instrument source code automatically using Program Database Toolkit (PDT) static tools

Choose an appropriate TAU stub makefile in <arch>/lib:

```
% source /home/S11505/shared/tools/setup.sh  
% export TAU_MAKEFILE=$TAU/Makefile.tau-mpi-pdt-fujitsu  
% export TAU_OPTIONS=' -optVerbose ...' (see tau_compiler.sh )
```

Use tau_f90.sh, tau_cxx.sh, tau_upc.sh, or tau_cc.sh as F90, C++, UPC, or C compilers respectively:

```
% mpifrpx foo.f90      changes to  
% tau_f90.sh foo.f90
```

- Set runtime environment variables, execute application and analyze performance data:

```
% pprof (for text based profile display)  
% paraprof (for GUI)
```

Automatic Source Instrumentation using PDT

```
% source /home/S11505/shared/tools/setup.sh  
% cat /home/S11505/shared/tools/UO/README.TAU  
% ls $TAU/Makefile.*  
ls $TAU/Makefile.*  
  
Makefile.tau-mpi-pdt-fujitsu  
Makefile.tau-mpi-pdt-scorep-fujitsu  
Makefile.tau-mpi-pdt-openmp-opari-fujitsu  
Makefile.tau-mpi-pdt-openmp-opari-scorep-fujitsu
```

- For an MPI+OpenMP+F90 application with MPI, you may choose

Makefile.tau-mpi-pdt-openmp-opari-fujitsu for native TAU measurement library

Makefile.tau-mpi-pdt-openmp-opari-scorep-fujitsu for Score-P library

- Supports MPI instrumentation & PDT for automatic source instrumentation

```
% export TAU_MAKEFILE=$TAU/Makefile.tau-mpi-pdt-openmp-opari-fujitsu
```

```
% tau_f90.sh matmult.f90 -o matmult  
% mpiexec -n 256 ./matmult  
% paraprof
```

Compile-Time Options

Optional parameters for the TAU_OPTIONS environment variable:

% tau_compiler.sh	
-optVerbose	Turn on verbose debugging messages
-optComInst	Use compiler based instrumentation
-optNoComInst	Do not revert to compiler instrumentation if source instrumentation fails.
-optTrackIO	Wrap POSIX I/O call and calculates vol/bw of I/O operations (configure TAU with <code>-iowrapper</code>)
-optTrackGOMP	Enable tracking GNU OpenMP runtime layer (used without <code>-opari</code>)
-optMemDbg	Enable runtime bounds checking (see <code>TAU_MEMDBG_*</code> env vars)
-optKeepFiles	Does not remove intermediate .pdb and .inst.* files
-optPreProcess	Preprocess sources (OpenMP, Fortran) before instrumentation
-optTauSelectFile=" <code><file></code> "	Specify selective instrumentation file for <code>tau_instrumentor</code>
-optTauWrapFile=" <code><file></code> "	Specify path to <code>link_options.tau</code> generated by <code>tau_gen_wrapper</code>
-optHeaderInst	Enable Instrumentation of headers
-optTrackUPCR	Track UPC runtime layer routines (used with <code>tau_upc.sh</code>)
-optLinking=""	Options passed to the linker. Typically <code>\$(TAU_MPI_FLIBS) \$(TAU_LIBS) \$(TAU_CXXLIBS)</code>
-optCompile=""	Options passed to the compiler. Typically <code>\$(TAU_MPI_INCLUDE) \$(TAU_INCLUDE) \$(TAU_DEFS)</code>
-optPdtF95Opts=""	Add options for Fortran parser in PDT (f95parse/gfparse) ...

Compile-Time Options (contd.)

- Optional parameters for the TAU_OPTIONS environment variable:

% tau_compiler.sh

-optMICOffload	Links code for Intel MIC offloading, requires both host and MIC TAU libraries
-optShared	Use TAU's shared library (libTAU.so) instead of static library (default)
-optPdtCxxOpts=""	Options for C++ parser in PDT (cxxparse).
-optPdtF90Parser=""	Specify a different Fortran parser (e.g., -optPdtF90Parser=/home/S11505/shared/tools/JSC/pdt/REL-3.21/x86_64/bin/gfparse48)
-optPdtCleanscapeParser	Specify the Cleanscape Fortran parser instead of GNU gfparser
-optTau=""	Specify options to the tau_instrumentor
-optTrackDMAPP	Enable instrumentation of low-level DMAPP API calls on Cray
-optTrackPthread	Enable instrumentation of pthread calls
-optPdtCParser	Specify a different C++ parser (-optPdtCParser=cparse4101)

See tau_compiler.sh for a full list of TAU_OPTIONS.

...

Runtime Environment Variables

Environment Variable	Default	Description
TAU_TRACE	0	Setting to 1 turns on tracing
TAU_CALLPATH	0	Setting to 1 turns on callpath profiling
TAU_TRACK_MEMORY_FOOTPRINT	0	Setting to 1 turns on tracking memory usage by sampling periodically the resident set size and high water mark of memory usage
TAU_TRACK_POWER	0	Tracks power usage by sampling periodically.
TAU_CALLPATH_DEPTH	2	Specifies depth of callpath. Setting to 0 generates no callpath or routine information, setting to 1 generates flat profile and context events have just parent information (e.g., Heap Entry: foo)
TAU_SAMPLING	1	Setting to 1 enables event-based sampling.
TAU_TRACK_SIGNALS	0	Setting to 1 generate debugging callstack info when a program crashes
TAU_COMM_MATRIX	0	Setting to 1 generates communication matrix display using context events
TAU_THROTTLE	1	Setting to 0 turns off throttling. Throttles instrumentation in lightweight routines that are called frequently
TAU_THROTTLE_NUMCALLS	100000	Specifies the number of calls before testing for throttling
TAU_THROTTLE_PERCALL	10	Specifies value in microseconds. Throttle a routine if it is called over 100000 times and takes less than 10 usec of inclusive time per call
TAU_COMPENSATE	0	Setting to 1 enables runtime compensation of instrumentation overhead
TAU_PROFILE_FORMAT	Profile	Setting to “merged” generates a single file. “snapshot” generates xml format
TAU_METRICS	TIME	Setting to a comma separated list generates other metrics. (e.g., ENERGY,TIME,P_VIRTUAL_TIME,PAPI_FP_INS,PAPI_NATIVE_<event>:<subevent>)

Runtime Environment Variables (contd.)

Environment Variable	Default	Description
TAU_TRACK_MEMORY_LEAKS	0	Tracks allocates that were not de-allocated (needs –optMemDbg or tau_exec –memory)
TAU_EBS_SOURCE	TIME	Allows using PAPI hardware counters for periodic interrupts for EBS (e.g., TAU_EBS_SOURCE=PAPI_TOT_INS when TAU_SAMPLING=1)
TAU_EBS_PERIOD	100000	Specifies the overflow count for interrupts
TAU_MEMDBG_ALLOC_MIN/MAX	0	Byte size minimum and maximum subject to bounds checking (used with TAU_MEMDBG_PROTECT_*)
TAU_MEMDBG_OVERHEAD	0	Specifies the number of bytes for TAU's memory overhead for memory debugging.
TAU_MEMDBG_PROTECT_BELOW/ABOVE	0	Setting to 1 enables tracking runtime bounds checking below or above the array bounds (requires –optMemDbg while building or tau_exec –memory)
TAU_MEMDBG_ZERO_MALLOC	0	Setting to 1 enables tracking zero byte allocations as invalid memory allocations.
TAU_MEMDBG_PROTECT_FREE	0	Setting to 1 detects invalid accesses to deallocated memory that should not be referenced until it is reallocated (requires –optMemDbg or tau_exec –memory)
TAU_MEMDBG_ATTEMPT_CONTINUE	0	Setting to 1 allows TAU to record and continue execution when a memory error occurs at runtime.
TAU_MEMDBG_FILL_GAP	Undefined	Initial value for gap bytes
TAU_MEMDBG_ALIGNMENT	Sizeof(int)	Byte alignment for memory allocations
TAU_EVENT_THRESHOLD	0.5	Define a threshold value (e.g., .25 is 25%) to trigger marker events for min/max

Simplifying TAU's usage (`tau_exec`)

- Uninstrumented execution
 - % mpiexec -n 4 ./a.out

- Track MPI performance
 - % mpiexec -n 4 **tau_exec** ./a.out

- Using TAU with Score-P and OTF2 with Vampir
 - % export SCOREP_ENABLE_TRACING=1
 - % mpiexec -n 4 **tau_exec -T scorep,papi,mpi,pdt** ./a.out

- Track memory operations
 - % export TAU_TRACK_MEMORY_LEAKS=1
 - % mpiexec -n 8 **tau_exec -memory_debug** ./a.out (bounds check)

- Use event based sampling (compile with -g)
 - % mpiexec -n 8 **tau_exec -ebs** ./a.out
 - Also -ebs_source=<PAPI_COUNTER> -ebs_period=<overflow_count>

- Load wrapper interposition library
 - % mpiexec -n 8 **tau_exec -loadlib=<path/libwrapper.so>** ./a.out

- Track GPGPU operations
 - % mpiexec -n 8 **tau_exec -cuhti** ./a.out
 - % mpiexec -n 8 **tau_exec -opencl** ./a.out
 - % mpiexec -n 8 **tau_exec -openacc** ./a.out

Binary Rewriting Instrumentation

- Support for both static and dynamic executables
- Specify a list of routines to instrument
- Specify the TAU measurement library to be injected
- MAQAO [UVSQ, Intel Exascale Labs]:

```
% tau_rewrite -T [tags] a.out -o a.inst
```

- DyninstAPI [U. Maryland and U. Wisconsin, Madison]:

```
% tau_run -T [tags] a.out -o a.inst
```

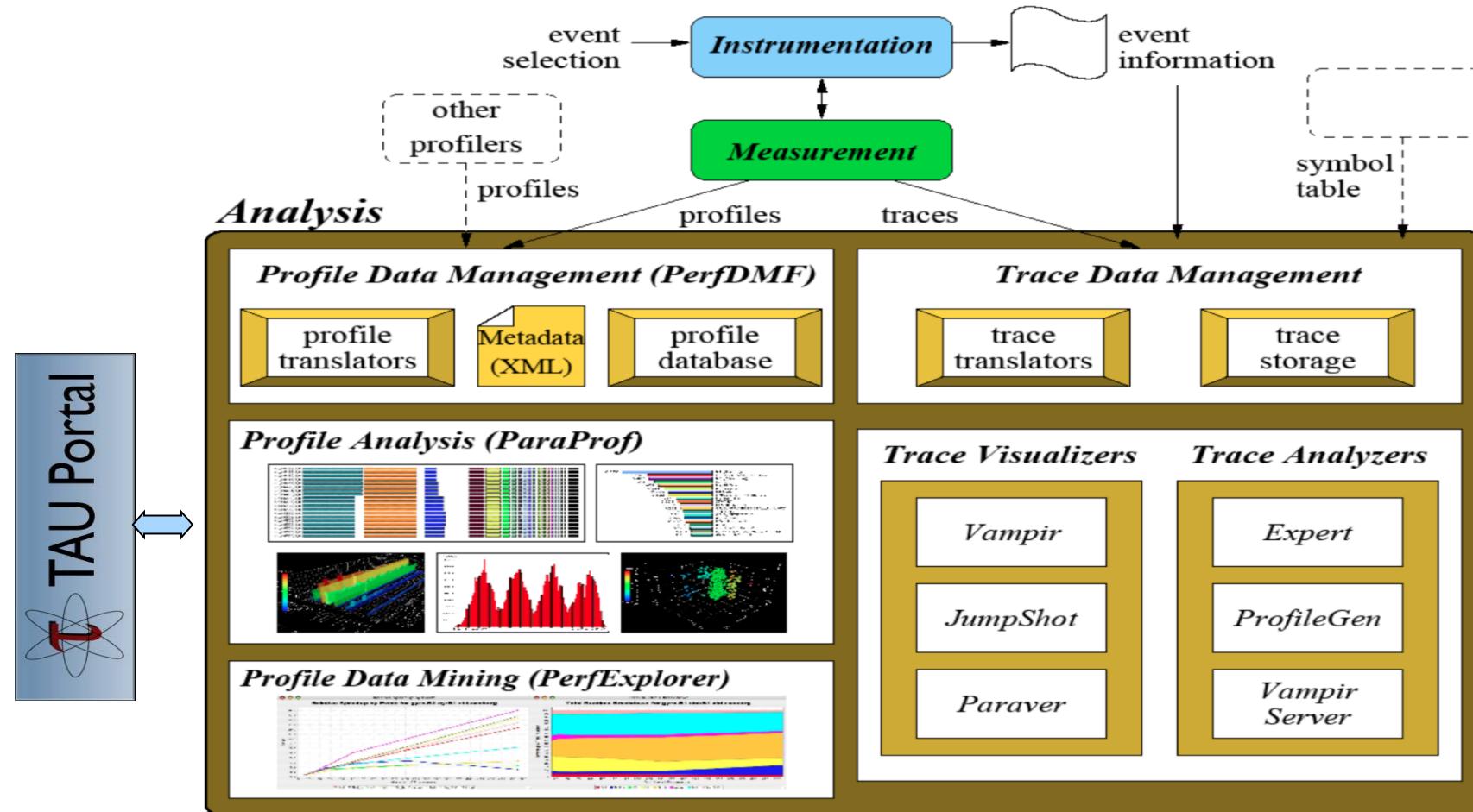
- Pebil [UC San Diego]:

```
% tau_pebil_rewrite -T [tags] a.out -o a.inst
```

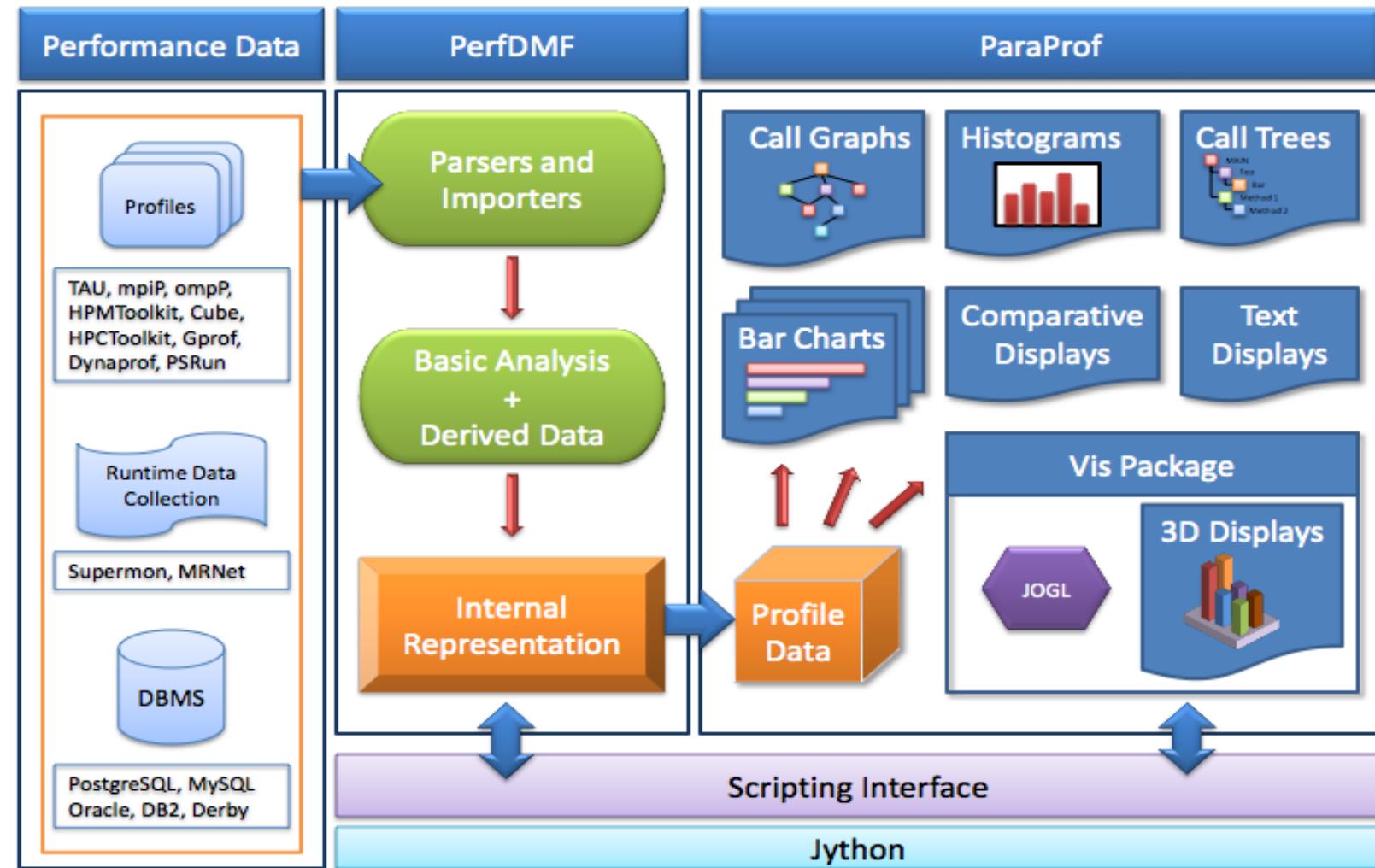
- Execute the application to get measurement data:

```
% mpieexec -n 256 ./a.inst
```

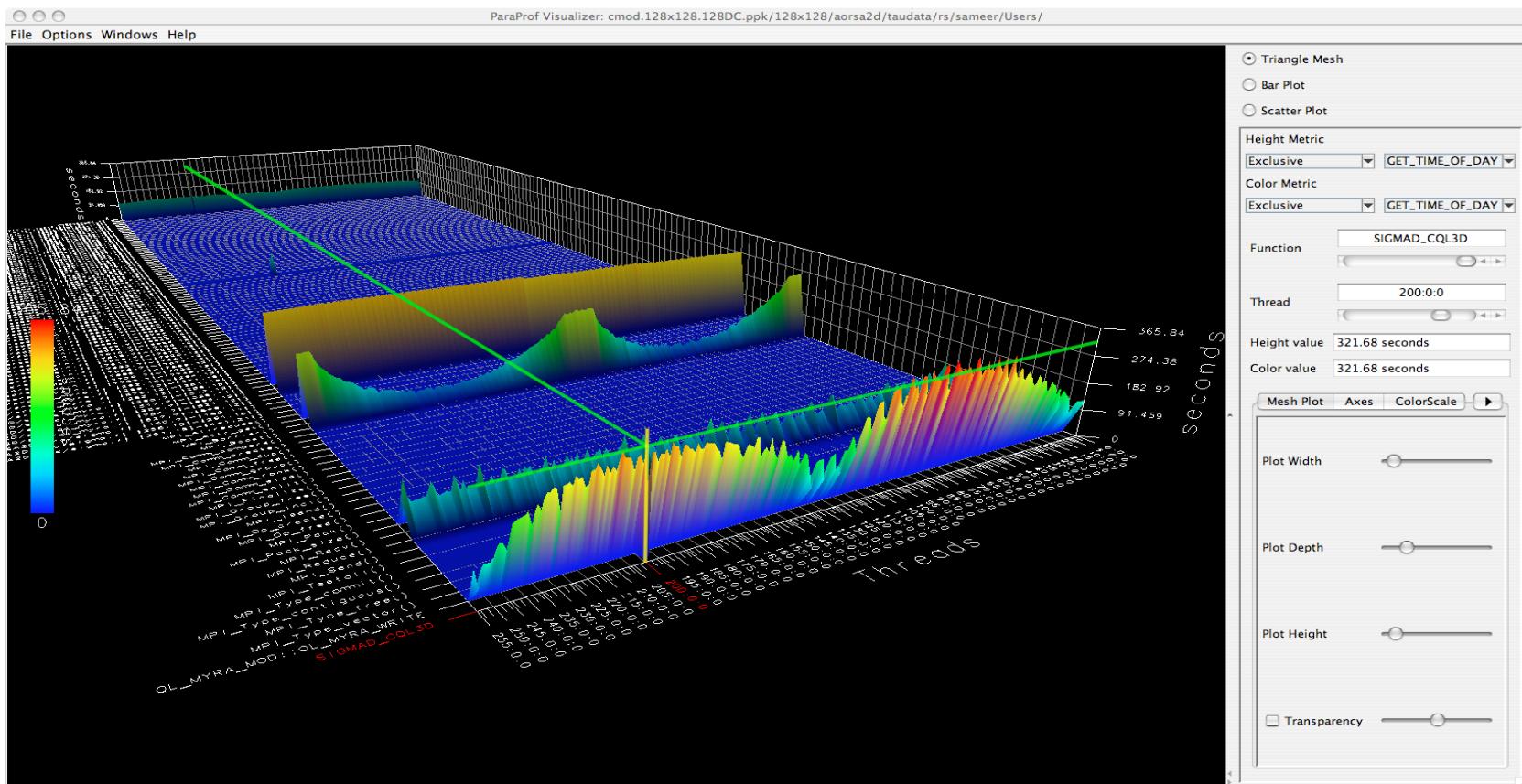
TAU Analysis



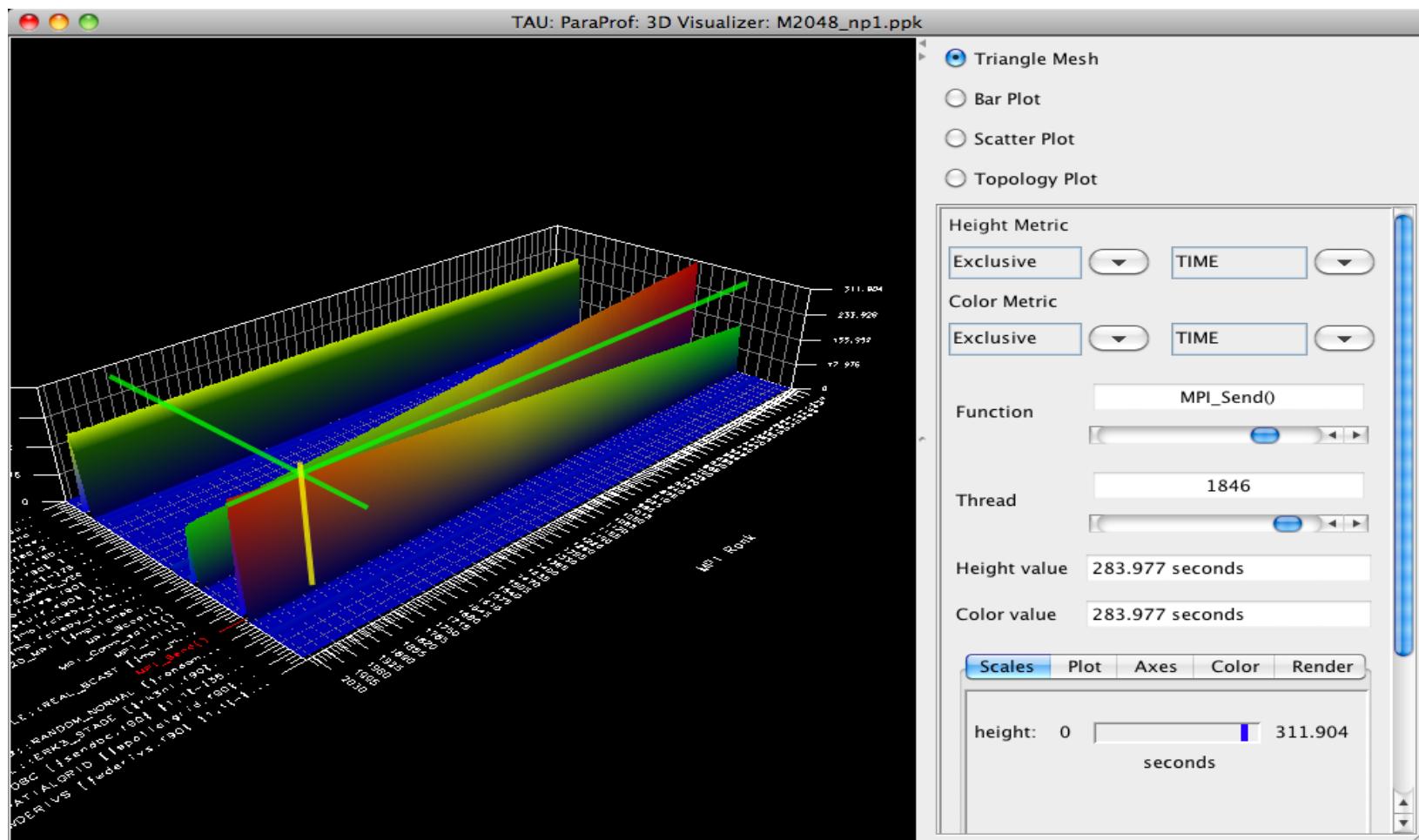
ParaProf Profile Analysis Framework



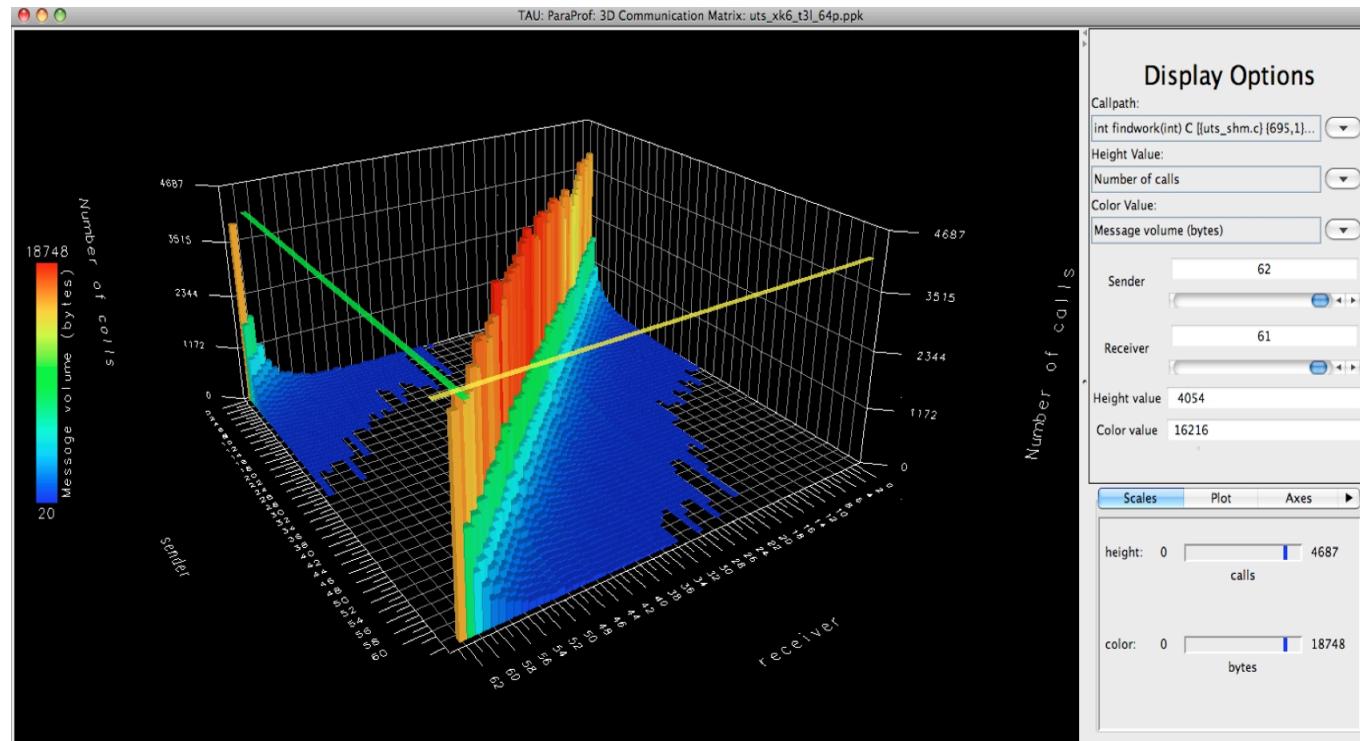
Parallel Profile Visualization: ParaProf



Parallel Profile Visualization: ParaProf

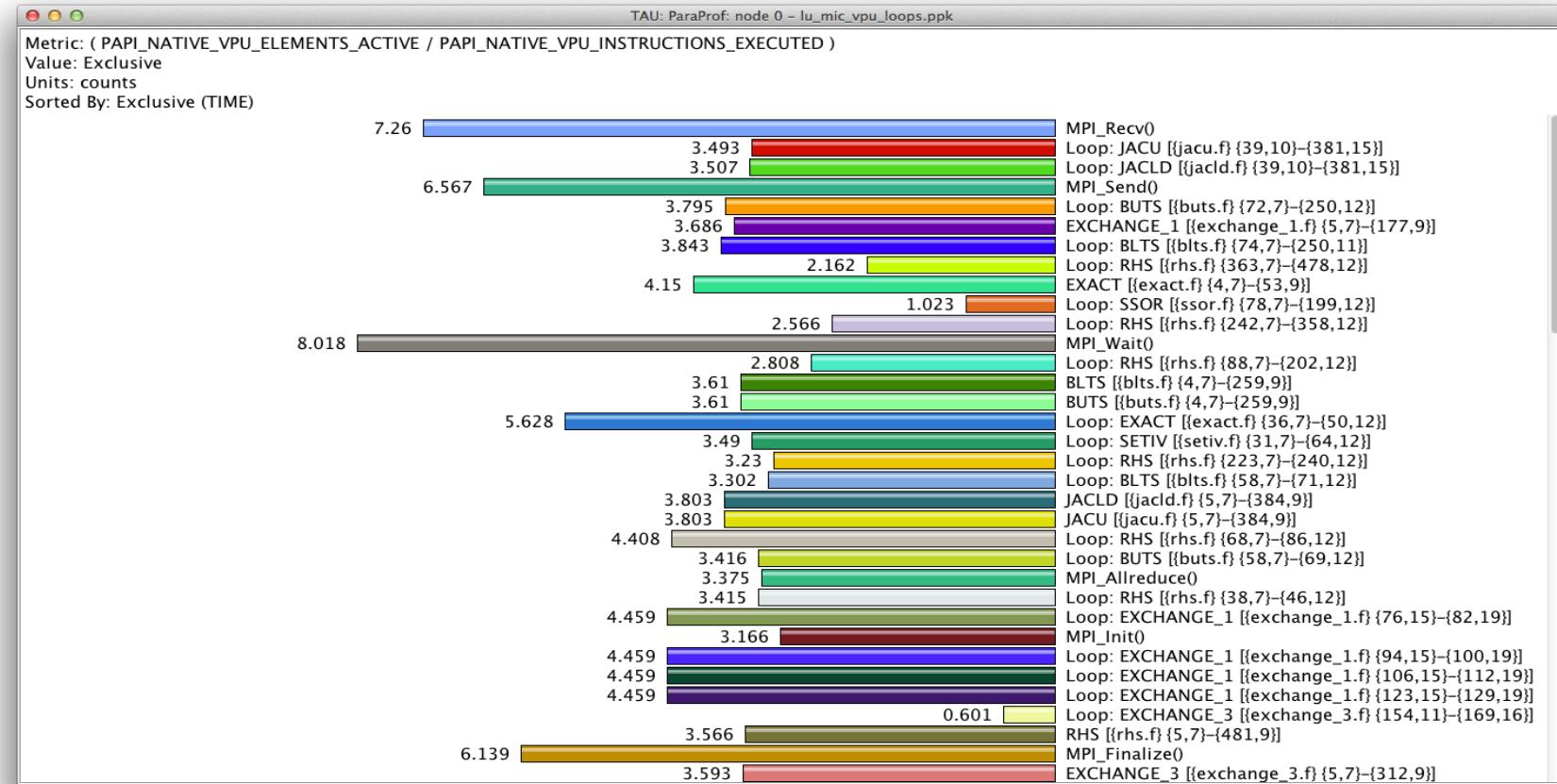


ParaProf 3D Communication Matrix



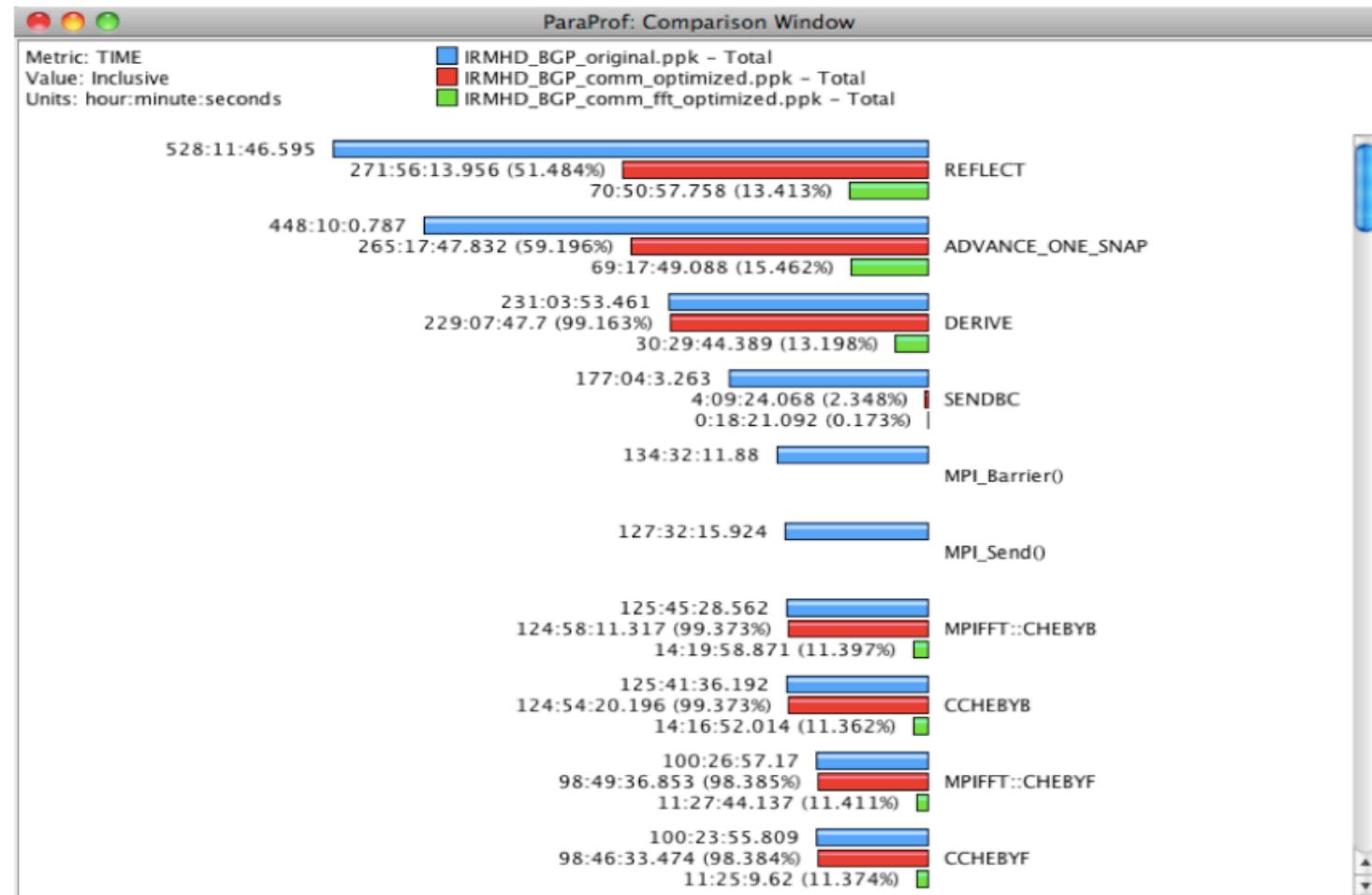
% export TAU_COMM_MATRIX=1

ParaProf Derived Metric Window: Intel MIC

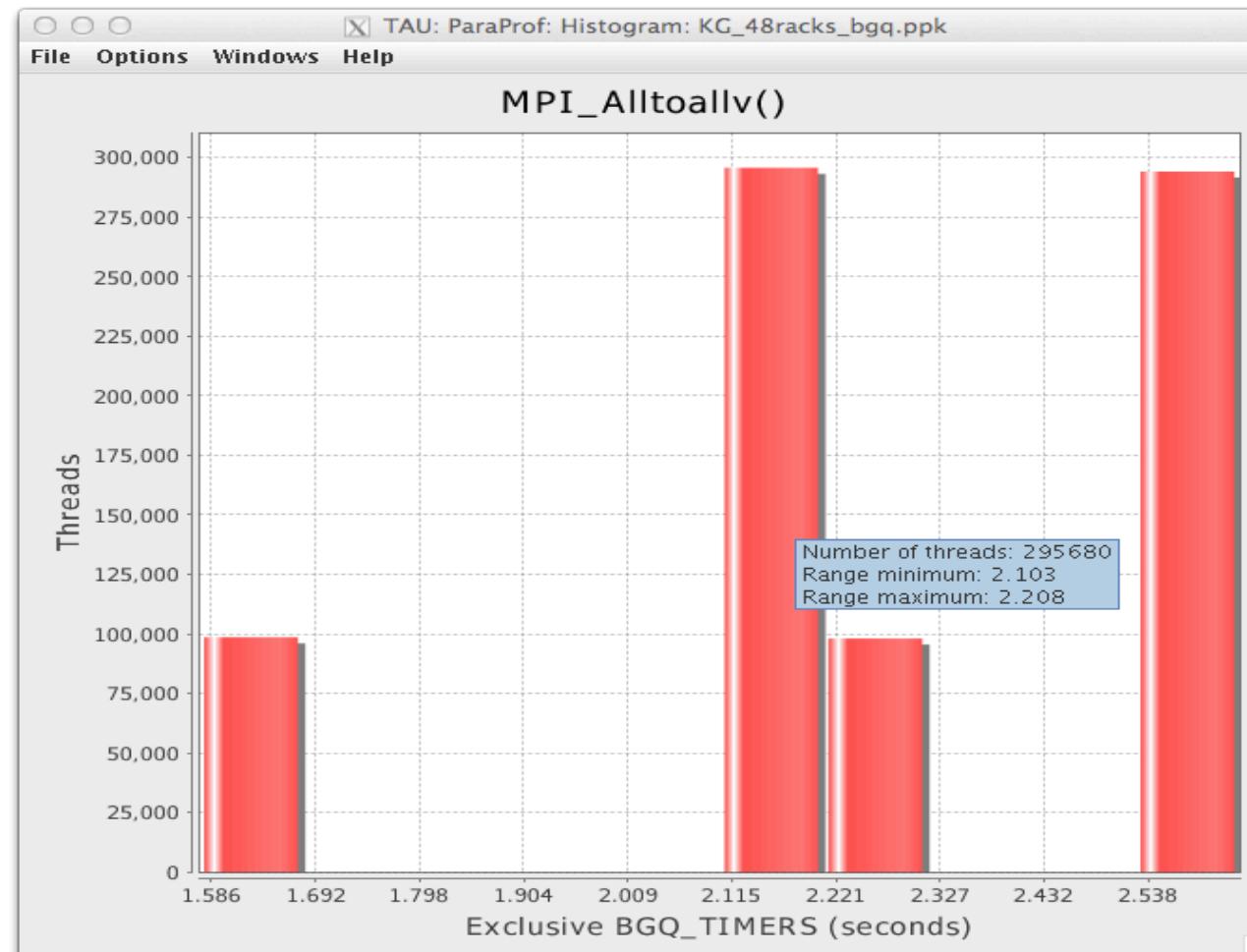


```
% export TAU_MAKEFILE=$TAUROOT/mic_linux/lib/Makefile.tau-icpc-papi-mpi-pdt
% export TAU_METRICS=TIME,PAPI_NATIVE_VPU_ELEMENTS_ACTIVE,PAPI_NATIVE_VPU_INSTRUCTIONS_EXECUTED
```

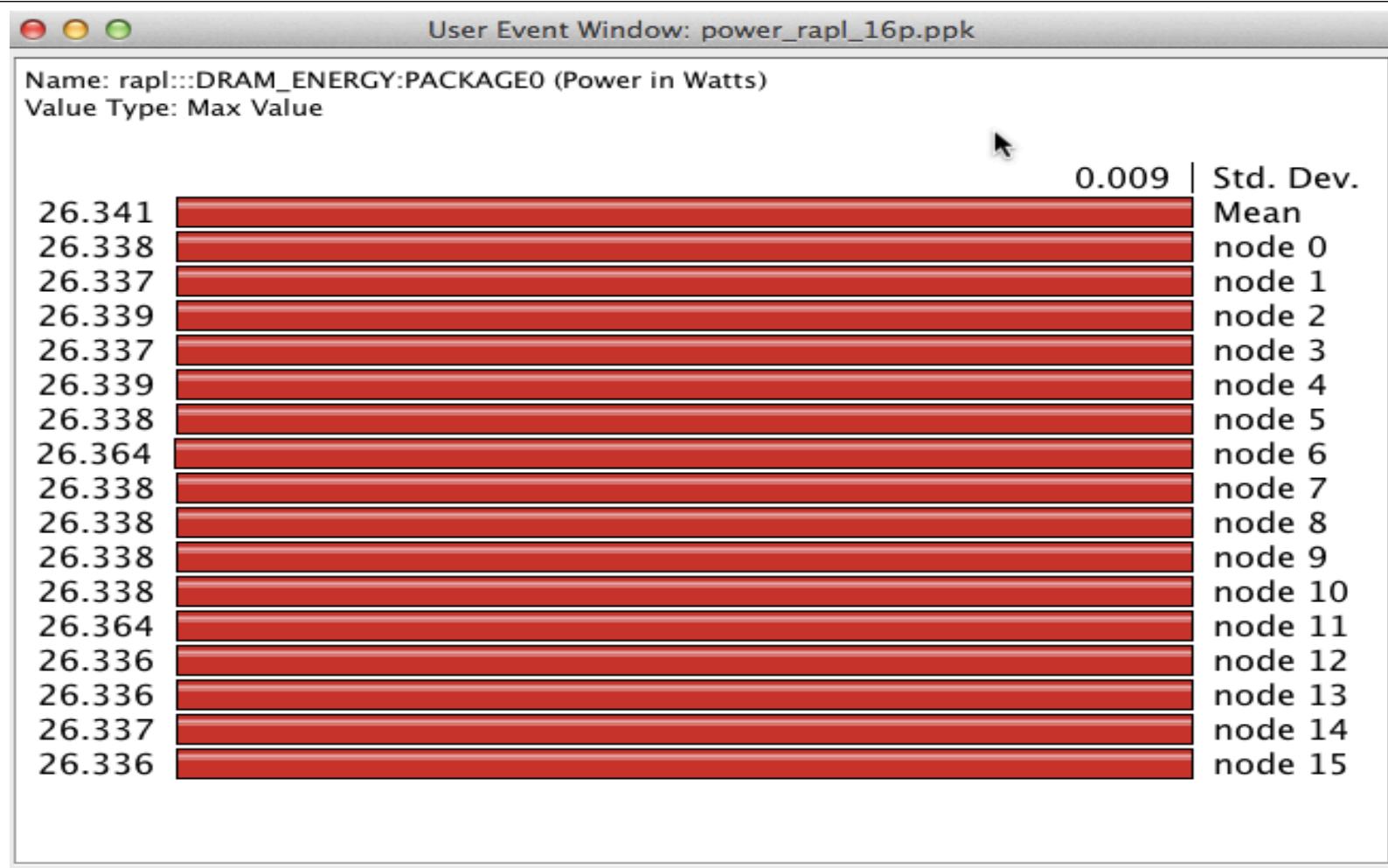
ParaProf Comparison Window



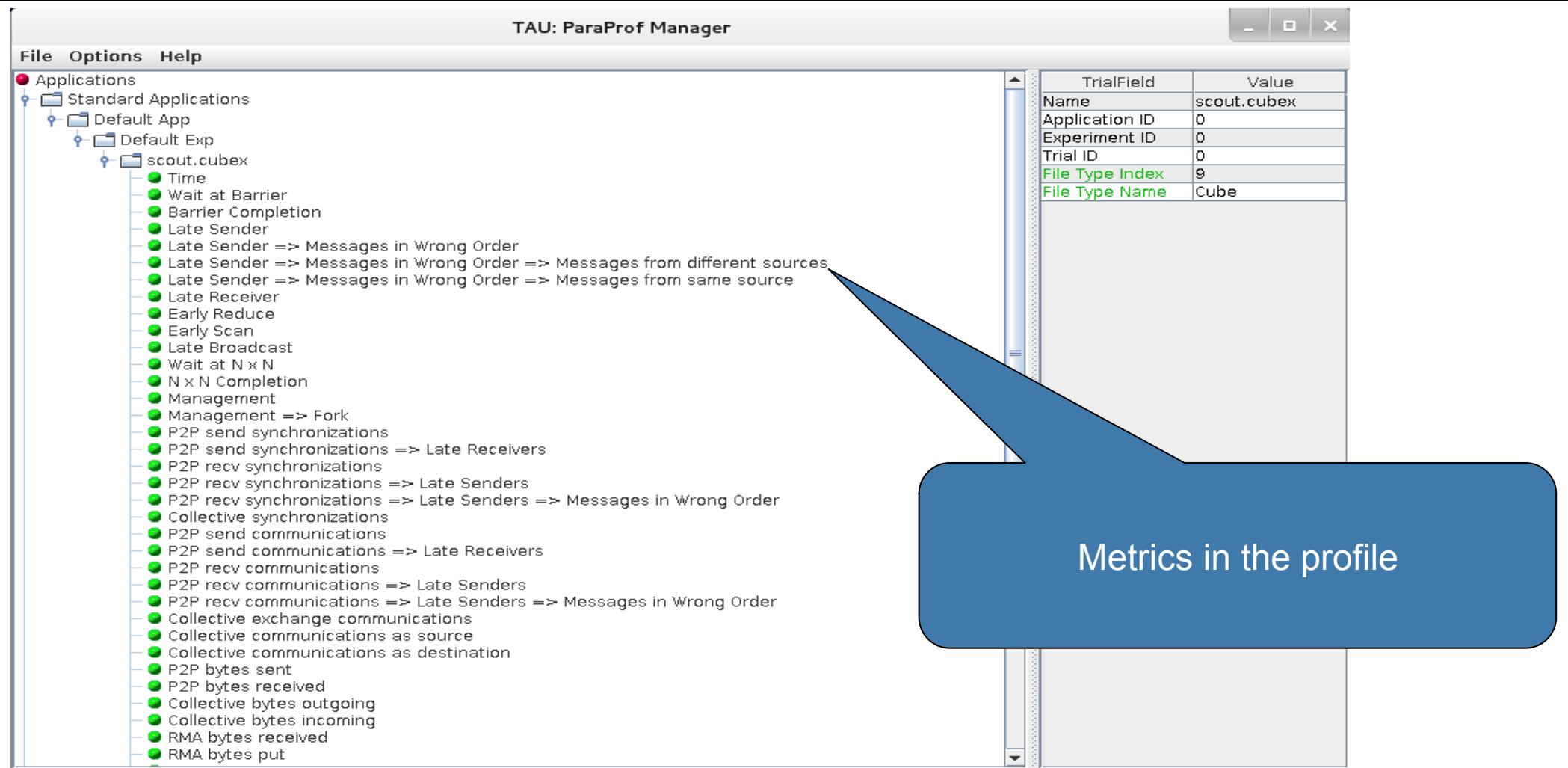
ParaProf Histogram Window



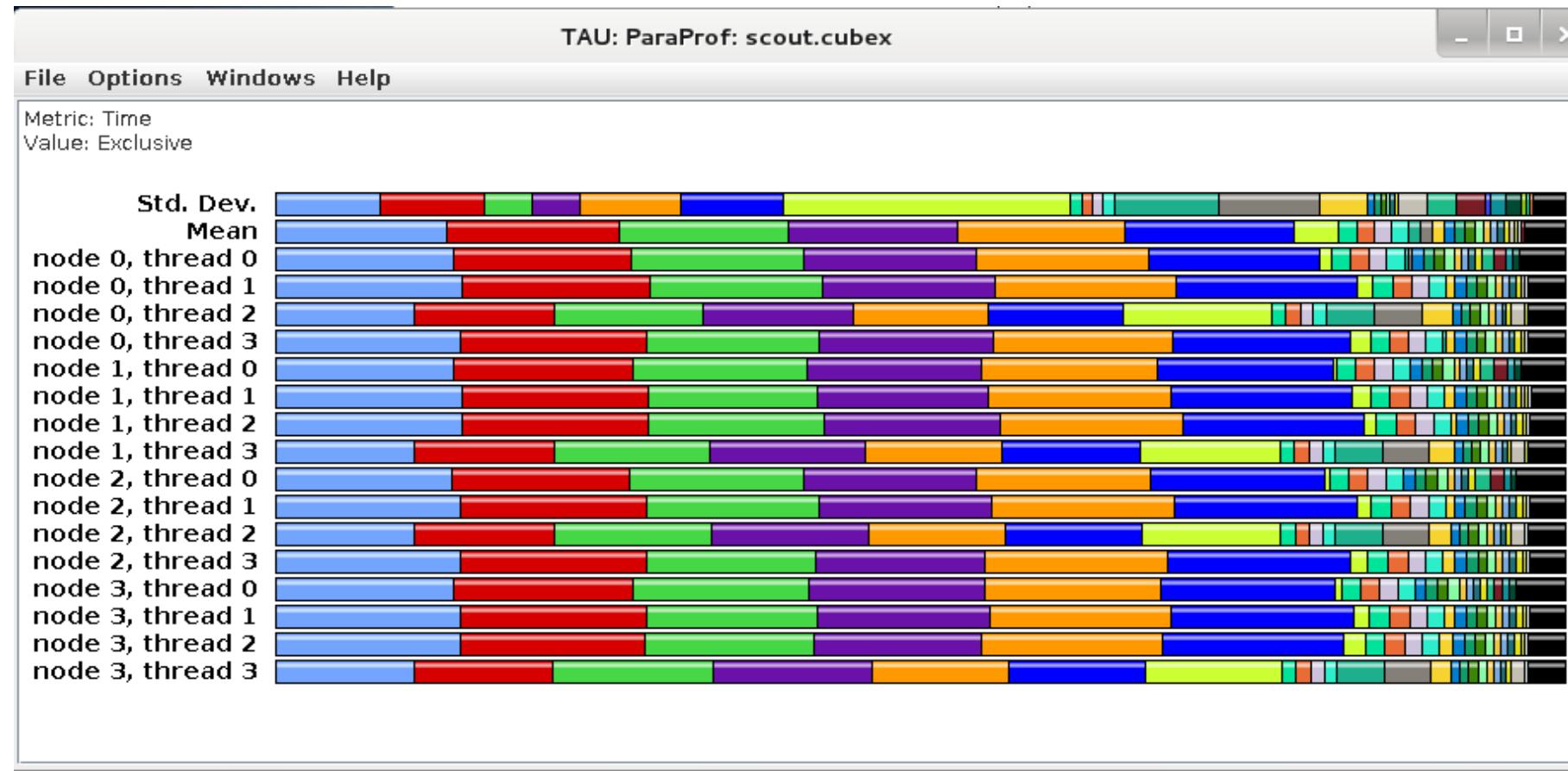
Energy Profiling in TAU



ParaProf Manager Window: scout.cubex



ParaProf: Main Window



ParaProf: Thread Statistics Table

TAU: ParaProf: Statistics for: node 0, thread 0 – scout.cubex

File Options Windows Help

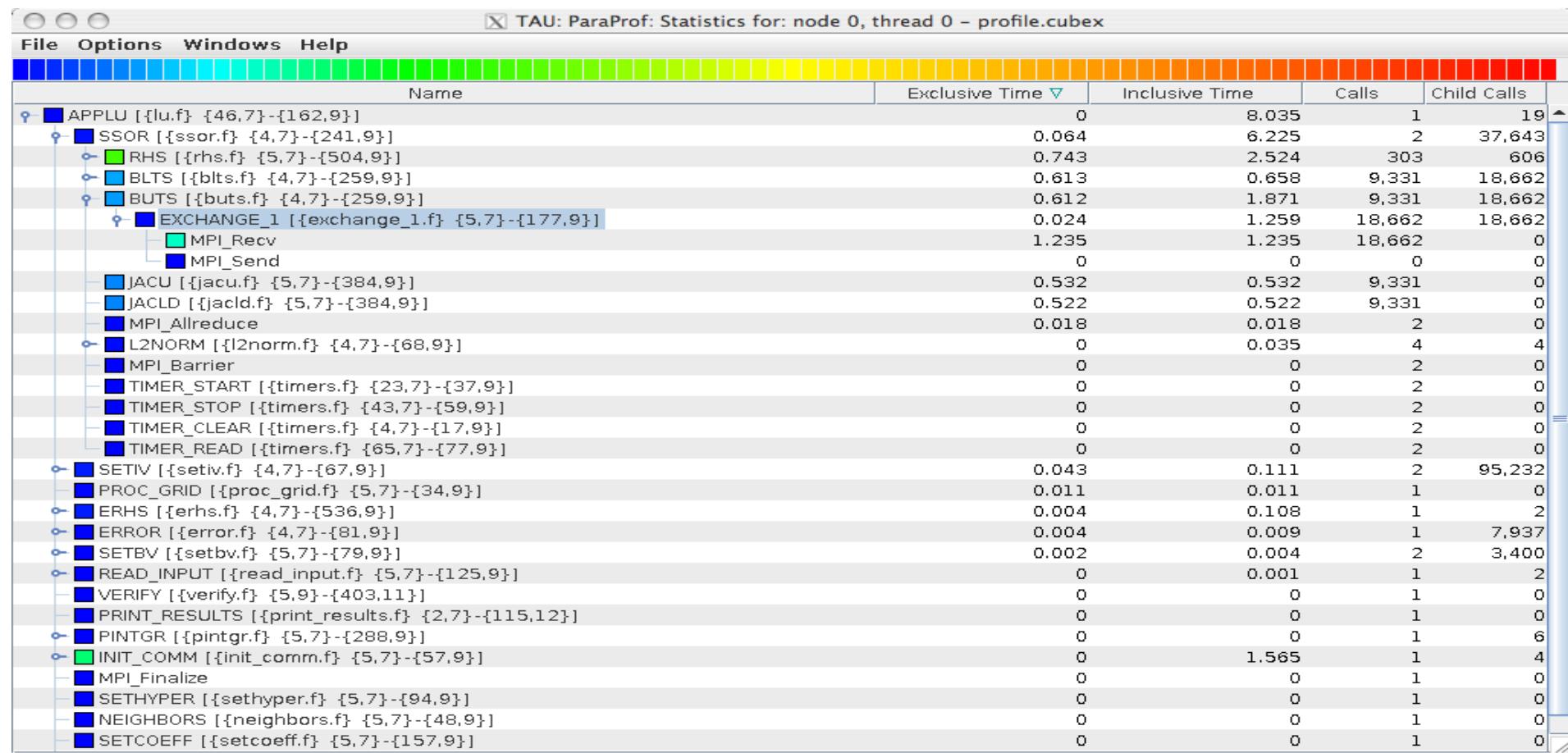
Time

Name	Exclusive Time	Inclusive Time	Calls	Child Calls
!\$omp do @y_solve.f:52	5.817	5.817	3,216	0
!\$omp do @z_solve.f:52	5.657	5.657	3,216	0
!\$omp do @x_solve.f:54	5.609	5.609	3,216	0
!\$omp do @rhs.f:191	0.609	0.609	3,232	0
!\$omp do @rhs.f:80	0.583	0.583	3,232	0
MPI_Waitall	0.402	0.402	603	0
!\$omp implicit barrier	0.402	0.402	0	0
!\$omp do @rhs.f:301	0.36	0.36	0	0
!\$omp implicit barrier	0.026	0.026	0	0
!\$omp implicit barrier	0	0	0	0
!\$omp do @rhs.f:37	0.343	0.343	0	0
!\$omp do @rhs.f:62	0.225	0.225	0	0
!\$omp implicit barrier	0.004	0.004	3,216	0
!\$omp implicit barrier	0	0	16	0
MPI_Init_thread	0.218	0.218	1	0
!\$omp do @rhs.f:384	0.199	0.199	3,232	0
!\$omp parallel do @add.f:22	0.099	0.111	3,216	3,216
!\$omp do @rhs.f:428	0.069	0.069	3,232	0
MPI_Isend	0.043	0.043	603	0
!\$omp do @initialize.f:50	0.04	0.04	32	0
!\$omp parallel @rhs.f:28	0.03	2.536	3,232	51,712
!\$omp parallel do @exch_qbc.f:215	0.021	0.029	6,432	6,432
!\$omp parallel do @exch_qbc.f:255	0.02	0.033	6,432	6,432
!\$omp parallel @exch_qbc.f:255	0.02	0.053	6,432	6,432
!\$omp parallel @exch_qbc.f:244	0	0	0	0

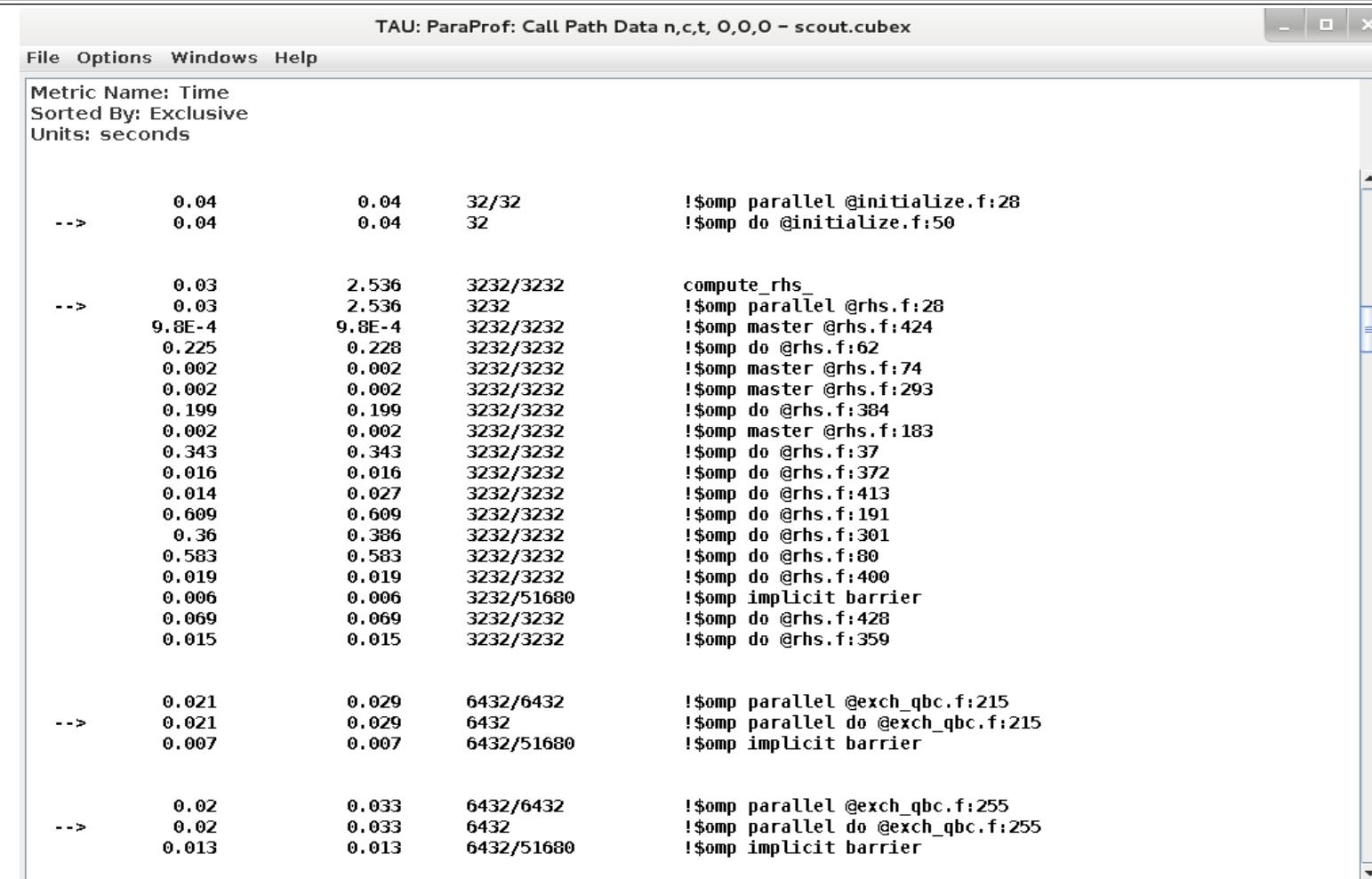
Click to sort by a given metric, drag and move to rearrange columns

FinderScreenSnapz003.png

Example: Score-P with TAU (NPB LU)



ParaProf: Callpath Thread Relations Window



TAU demo objectives

- Familiarise with usage of TAU tools
 - complementary tools' capabilities & interoperability
- Prepare to apply tools productively to *your* applications(s)
- Exercise is based on a small portable benchmark code
 - unlikely to have significant optimisation opportunities

- Optional (recommended) exercise extensions
 - analyse performance of alternative configurations
 - investigate effectiveness of system-specific compiler/MPI optimisations and/or placement/binding/affinity capabilities
 - investigate scalability and analyse scalability limiters
 - compare performance on different HPC platforms
 - ...

Local Installation (Fujitsu HPC Prime FX10 *Pi, U. Kobe*)

- Setup preferred program environment compilers
 - Default set Intel Compilers with Intel MPI
 - Generate profile files using Score-P

```
% source /home/S11505/shared/tools/setup.sh  
% paraprof profile.cubex &
```

For PerfExplorer:

```
% wget http://tau.uoregon.edu/data.tgz; tar zxf data.tgz; cd data  
% cat README  
And follow the steps
```

NPB-MZ-MPI Suite

- The NAS Parallel Benchmark suite (MPI+OpenMP version)
 - Available from:

<http://www.nas.nasa.gov/Software/NPB>

- 3 benchmarks in Fortran77
- Configurable for various sizes & classes
- Move into the NPB3.3-MZ-MPI root directory

```
% ls
bin/    common/   jobscript/  Makefile  README.install  SP-MZ/
BT-MZ/   config/   LU-MZ/     README     README.tutorial  sys/
```

- Subdirectories contain source code for each benchmark
 - plus additional configuration and common code
- The provided distribution has already been configured for the tutorial, such that it's ready to "make" one or more of the benchmarks and install them into a (tool-specific) "bin" subdirectory

NPB-MZ-MPI / BT: config/make.def

```
#           SITE- AND/OR PLATFORM-SPECIFIC DEFINITIONS.  
#  
#-----  
# Configured for generic MPI with GCC compiler  
#-----  
#OPENMP = -fopenmp      # GCC compiler  
OPENMP = -openmp        # Intel compiler  
  
...  
#-----  
# The Fortran compiler used for MPI programs  
#-----  
# MPIF77 = mpifrtpx # Fujitsu compiler  
  
# Alternative variant to perform instrumentation  
MPIF77 = tau_f90.sh -tau_makefile=<path>/Makefile.tau-[options]  
  
# PREP is a generic preposition macro for instrumentation preparation  
#MPIF77 = $(PREP) mpif77 -f77=ifort  
#MPIF77 = scorep ...  
...
```

Default (no instrumentation)

Uncomment TAU's compiler
wrapper to do source
instrumentation with TAU
Comment out Score-P wrapper

Building an NPB-MZ-MPI Benchmark

```
% make  
=====  
= NAS PARALLEL BENCHMARKS 3.3 =  
= MPI+OpenMP Multi-Zone Versions =  
= F77 =  
=====
```

To make a NAS multi-zone benchmark type

```
make <benchmark-name> CLASS=<class> NPROCS=<nprocs>
```

where <benchmark-name> is "bt-mz", "lu-mz", or "sp-mz"
<class> is "S", "W", "A" through "F"
<nprocs> is number of processes

[...]

```
*****  
* Custom build configuration is specified in config/make.def *  
* Suggested tutorial exercise configuration for HPC systems: *  
*   make bt-mz CLASS=B NPROCS=8 *  
*****
```

- Type “make” for instructions

Building an NPB-MZ-MPI Benchmark

```
% make suite
make[1]: Entering directory `BT-MZ'
make[2]: Entering directory `sys'
cc -o setparams setparams.c -lm
make[2]: Leaving directory `sys'
../sys/setparams bt-mz 8 B
make[2]: Entering directory `../BT-MZ'
tau_f90.sh -c -O3 -g -openmp          bt.f
[...]
tau_f90.sh -c -O3 -g -openmp          mpi_setup.f
cd ..;/common; mpiifort -c -O3 -g -openmp      print_results.f
cd ..;/common; mpiifort -c -O3 -g -openmp      timers.f
tau_f90.sh -O3 -g -openmp -o ../bin.tau/bt-mz_C.8 bt.o
  initialize.o exact_solution.o exact_rhs.o set_constants.o adi.o
  rhs.o zone_setup.o x_solve.o y_solve.o exch_qbc.o solve_subs.o
  z_solve.o add.o error.o verify.o mpi_setup.o ../common/print_results.o
  ../common/timers.o
make[2]: Leaving directory `BT-MZ'
Built executable ../bin/bt-mz_B.8
make[1]: Leaving directory `BT-MZ'
```

- Specify the benchmark configuration
 - benchmark name: **bt-mz**, lu-mz, sp-mz
 - the number of MPI processes: NPROCS=**8**
 - the benchmark class (S, W, A, B, C, D, E): **CLASS=B**

Shortcut: % **make suite**

NPB-MZ-MPI / BT without TAU

```
% make bt-mz CLASS=B NPROCS=4
% cd bin
% cp ..../jobscript/fx10/run.sh .
% pbsub run.sh
% cat mzmpibt.o<job_id>
NAS Parallel Benchmarks (NPB3.3-MZ-MPI) - BT-MZ MPI+OpenMP Benchmark
Number of zones: 8 x 8
Iterations: 200 dt: 0.000300
Number of active processes: 8
Total number of threads: 32 ( 4.0 threads/process)

Time step 1
Time step 20
[...]
Time step 180
Time step 200
Verification Successful

BT-MZ Benchmark Completed.
Time in seconds = 16.03
%
```

- Copy jobscript and launch as a hybrid MPI+OpenMP application

Hint: save the benchmark output (or note the run time) to be able to refer to it later

NPB-MZ-MPI / BT (Block Tridiagonal Solver)

- What does it do?
 - Solves a discretized version of the unsteady, compressible Navier-Stokes equations in three spatial dimensions
 - Performs 200 time-steps on a regular 3-dimensional grid
- Implemented in 20 or so Fortran77 source modules

- Uses MPI & OpenMP in combination
 - 8 processes each with 4 threads should be reasonable
 - bt-mz_B.8 should take around 16 seconds

tau_exec

```
$ tau_exec

Usage: tau_exec [options] [--] <exe> <exe options>

Options:
  -v          Verbose mode
  -s          Show what will be done but don't actually do anything (dryrun)
  -qsub       Use qsub mode (BG/P only, see below)
  -io         Track I/O
  -memory    Track memory allocation/deallocation
  -memory_debug Enable memory debugger
  -cuda       Track GPU events via CUDA
  -cupti     Track GPU events via CUPTI (Also see env. variable TAU_CUPTI_API)
  -opencl    Track GPU events via OpenCL
  -openacc   Track GPU events via OpenACC (currently PGI only)
  -ompt      Track OpenMP events via OMPT interface
  -armci     Track ARMCI events via PARMCI
  -ebs       Enable event-based sampling
  -ebs_period=<count> Sampling period (default 1000)
  -ebs_source=<counter> Counter (default itimer)
  -um        Enable Unified Memory events via CUPTI
  -T <DISABLE,GNU,ICPC,MPI,OMPT,OPENMP,PAPI,PDT,PROFILE,PTHREAD,SCOREP,SERIAL> : Specify TAU tags
  -loadlib=<file.so>  : Specify additional load library
  -XrunTAUsh-<options> : Specify TAU library directly
  -gdb       Run program in the gdb debugger
```

Notes:

Defaults if unspecified: -T MPI
MPI is assumed unless SERIAL is specified

- Tau_exec preloads the TAU wrapper libraries and performs measurements.

No need to recompile the application!

tau_exec Example (continued)

Example:

```
mpiexec -n 2 tau_exec -T icpc,ompt,mpi -ompt ./a.out  
mpiexec -n 2 tau_exec -io ./a.out
```

Example - event-based sampling with samples taken every 1,000,000 FP instructions

```
mpiexec -n 8 tau_exec -ebs -ebs_period=1000000 -ebs_source=PAPI_FP_INS ./ring
```

Examples - GPU:

```
tau_exec -T serial,cupti -cupti ./matmult (Preferred for CUDA 4.1 or later)  
tau_exec -openacc ./a.out
```

```
tau_exec -T serial -opencl ./a.out (OPENCL)
```

```
mpiexec -n 2 tau_exec -T mpi,cupti,papi -cupti -um ./a.out (Unified Virtual Memory in CUDA 6.0+)
```

qsub mode (IBM BG/Q only):

Original:

```
qsub -n 1 --mode smp -t 10 ./a.out
```

With TAU:

```
tau_exec -qsub -io -memory -- qsub -n 1 ... -t 10 ./a.out
```

Memory Debugging:

-memory option:

Tracks heap allocation/deallocation and memory leaks.

-memory_debug option:

Detects memory leaks, checks for invalid alignment, and checks for array overflow. This is exactly like setting TAU_TRACK_MEMORY_LEAKS=1 and TAU_MEMDBG_PROTECT_ABOVE=1 and running with -memory

- tau_exec can enable event based sampling while launching the executable using the **-ebs** flag!

Using tau_exec -ebs

- Load TAU tools in your path

```
% source /home/S11505/shared/tools/setup.sh
```

- Make clean and build new tool-specific executable; compile application with -g

```
% make clean  
% make suite  
Built executable .../bin/bt-mz_B.8
```

- Change to the directory containing the new executable before running it with the desired tool configuration

```
% cp ..../jobscript/fx10/run.sh .  
EDIT run.sh and add tau_exec -ebs:  
mpirun -np $NPROCS $EXE           changes to  
mpirun -np $NPROCS tau_exec -ebs $EXE  
% pbsub run.sh  
% paraprof &
```

ParaProf: Thread Statistics Table (right click on node 0)

% paraprof

Name	Exclusive TIME	Inclusive TIME	Calls	Child Calls
.TAU application	16.237	18.326	1	4,245
[CONTEXT].TAU application	0	15.956	3,192	0
[SUMMARY] binvcrhs_ [{/home/S11505/sameer/NPB3.3-MZ-MPI/BT-MZ/solve_subsf}]	5.728	5.728	1,129	0
[SAMPLE] binvcrhs_ [{/home/S11505/sameer/NPB3.3-MZ-MPI/BT-MZ/solve_subsf} {225}]	0.191	0.191	40	0
[SAMPLE] binvcrhs_ [{/home/S11505/sameer/NPB3.3-MZ-MPI/BT-MZ/solve_subsf} {393}]	0.173	0.173	38	0
[SAMPLE] binvcrhs_ [{/home/S11505/sameer/NPB3.3-MZ-MPI/BT-MZ/solve_subsf} {363}]	0.168	0.168	30	0
[SAMPLE] binvcrhs_ [{/home/S11505/sameer/NPB3.3-MZ-MPI/BT-MZ/solve_subsf} {312}]	0.119	0.119	22	0
[SAMPLE] binvcrhs_ [{/home/S11505/sameer/NPB3.3-MZ-MPI/BT-MZ/solve_subsf} {439}]	0.115	0.115	16	0
[SAMPLE] binvcrhs_ [{/home/S11505/sameer/NPB3.3-MZ-MPI/BT-MZ/solve_subsf} {460}]	0.108	0.108	19	0
[SAMPLE] binvcrhs_ [{/home/S11505/sameer/NPB3.3-MZ-MPI/BT-MZ/solve_subsf} {334}]	0.107	0.107	17	0
[SAMPLE] binvcrhs_ [{/home/S11505/sameer/NPB3.3-MZ-MPI/BT-MZ/solve_subsf} {328}]	0.089	0.089	17	0
[SAMPLE] binvcrhs_ [{/home/S11505/sameer/NPB3.3-MZ-MPI/BT-MZ/solve_subsf} {476}]	0.087	0.087	14	0
[SAMPLE] binvcrhs_ [{/home/S11505/sameer/NPB3.3-MZ-MPI/BT-MZ/solve_subsf} {357}]	0.083	0.083	21	0
[SAMPLE] binvcrhs_ [{/home/S11505/sameer/NPB3.3-MZ-MPI/BT-MZ/solve_subsf} {433}]	0.081	0.081	17	0
[SAMPLE] binvcrhs_ [{/home/S11505/sameer/NPB3.3-MZ-MPI/BT-MZ/solve_subsf} {226}]	0.08	0.08	16	0
[SAMPLE] binvcrhs_ [{/home/S11505/sameer/NPB3.3-MZ-MPI/BT-MZ/solve_subsf} {474}]	0.079	0.079	14	0
[SAMPLE] binvcrhs_ [{/home/S11505/sameer/NPB3.3-MZ-MPI/BT-MZ/solve_subsf} {335}]	0.077	0.077	12	0
[SAMPLE] binvcrhs_ [{/home/S11505/sameer/NPB3.3-MZ-MPI/BT-MZ/solve_subsf} {299}]	0.074	0.074	20	0
[SAMPLE] binvcrhs_ [{/home/S11505/sameer/NPB3.3-MZ-MPI/BT-MZ/solve_subsf} {326}]	0.07	0.07	14	0
[SAMPLE] binvcrhs_ [{/home/S11505/sameer/NPB3.3-MZ-MPI/BT-MZ/solve_subsf} {244}]	0.07	0.07	14	0
[SAMPLE] binvcrhs_ [{/home/S11505/sameer/NPB3.3-MZ-MPI/BT-MZ/solve_subsf} {472}]	0.069	0.069	13	0
[SAMPLE] binvcrhs_ [{/home/S11505/sameer/NPB3.3-MZ-MPI/BT-MZ/solve_subsf} {458}]	0.068	0.068	11	0
[SAMPLE] binvcrhs_ [{/home/S11505/sameer/NPB3.3-MZ-MPI/BT-MZ/solve_subsf} {290}]	0.068	0.068	10	0
[SAMPLE] binvcrhs_ [{/home/S11505/sameer/NPB3.3-MZ-MPI/BT-MZ/solve_subsf} {286}]	0.068	0.068	10	0
[SAMPLE] binvcrhs_ [{/home/S11505/sameer/NPB3.3-MZ-MPI/BT-MZ/solve_subsf} {465}]	0.067	0.067	9	0
[SAMPLE] binvcrhs_ [{/home/S11505/sameer/NPB3.3-MZ-MPI/BT-MZ/solve_subsf} {342}]	0.062	0.062	16	0
[SAMPLE] binvcrhs_ [{/home/S11505/sameer/NPB3.3-MZ-MPI/BT-MZ/solve_subsf} {449}]	0.061	0.061	13	0
[SAMPLE] binvcrhs_ [{/home/S11505/sameer/NPB3.3-MZ-MPI/BT-MZ/solve_subsf} {413}]	0.06	0.06	12	0
[SAMPLE] binvcrhs_ [{/home/S11505/sameer/NPB3.3-MZ-MPI/BT-MZ/solve_subsf} {254}]	0.06	0.06	12	0
[SAMPLE] binvcrhs_ [{/home/S11505/sameer/NPB3.3-MZ-MPI/BT-MZ/solve_subsf} {111}]	0.06	0.06	12	0

Right click on node 0 -> Show Thread Statistics Table. Click on column to sort, open keys on left

TAU Source Instrumentation

- See TAU commands

```
% cat /home/S11505/shared/tools/UO/README.TAU
```

- Make clean and build new tool-specific executable

```
% make clean  
% export TAU_MAKEFILE=$TAU/Makefile.tau-mpi-pdt-openmp-opari-scorep-fujitsu  
% make MPIF77=tau_f77.sh suite  
Built executable ../bin.tau/bt-mz_B.8
```

- Change to the directory containing the new executable before running it with the desired tool configuration

```
% cd bin.tau  
% cp ../jobscript/fx10/run.sh .  
% pbsub run.sh  
% paraprof &
```

TAU Analysis Tools: paraprof

- Launch paraprof

```
% paraprof
```

Metric

TAU: ParaProf Manager

The screenshot shows the TAU ParaProf Manager application window. On the left, there is a tree view under the 'Applications' section. It includes 'Standard Applications', 'Default App' (which is expanded), 'Default Exp' (also expanded), and a file named 'bt_ompt.ppk'. A green dot icon labeled 'TIME' is placed next to the 'bt_ompt.ppk' entry. Below this tree is a connection string: 'Default (jdbc:h2:/Users/sameer/.ParaProf/perfdmf/perfdmf;AUTO_SERVER=TRUE)'. On the right side of the window, there is a large table titled 'TAU: ParaProf Manager' with two columns: 'TrialField' and 'Value'. The table lists various parameters and their values, such as Name (bt_ompt.ppk), Application ID (0), Experiment ID (0), Trial ID (0), CPU Cores (8), CPU MHz (2600.000), CPU Type (Intel(R) Xeon(R) CPU E5-2670 0 @ 2.60GHz), CPU Vendor (GenuineIntel), CWD (/scratch/sameer/NPB3.3-MZ-MPI/bin), Cache Size (20480 KB), Command Line (./bt-mz_C.8), Executable (/scratch/sameer/NPB3.3-MZ-MPI/bin/bt-mz_C.8), File Type Index (0), File Type Name (ParaProf Packed Profile), Hostname (frog9), Local Time (2015-05-18T00:37:38+02:00), MPI Processor Name (frog9), Memory Size (65944056 kB), Node Name (frog9), OMP_CHUNK_SIZE (1), OMP_DYNAMIC (off), OMP_MAX_THREADS (4), OMP_NESTED (off), OMP_NUM_PROCS (4), OMP_SCHEDULE (UNKNOWN), OS Machine (x86_64), OS Name (Linux), OS Release (2.6.32-279.5.2.b16.Bull.33.x86_64), and OS Version (#1 SMP Sat Nov 10 01:48:00 CET 2012).

TrialField	Value
Name	bt_ompt.ppk
Application ID	0
Experiment ID	0
Trial ID	0
CPU Cores	8
CPU MHz	2600.000
CPU Type	Intel(R) Xeon(R) CPU E5-2670 0 @ 2.60GHz
CPU Vendor	GenuineIntel
CWD	/scratch/sameer/NPB3.3-MZ-MPI/bin
Cache Size	20480 KB
Command Line	./bt-mz_C.8
Executable	/scratch/sameer/NPB3.3-MZ-MPI/bin/bt-mz_C.8
File Type Index	0
File Type Name	ParaProf Packed Profile
Hostname	frog9
Local Time	2015-05-18T00:37:38+02:00
MPI Processor Name	frog9
Memory Size	65944056 kB
Node Name	frog9
OMP_CHUNK_SIZE	1
OMP_DYNAMIC	off
OMP_MAX_THREADS	4
OMP_NESTED	off
OMP_NUM_PROCS	4
OMP_SCHEDULE	UNKNOWN
OS Machine	x86_64
OS Name	Linux
OS Release	2.6.32-279.5.2.b16.Bull.33.x86_64
OS Version	#1 SMP Sat Nov 10 01:48:00 CET 2012

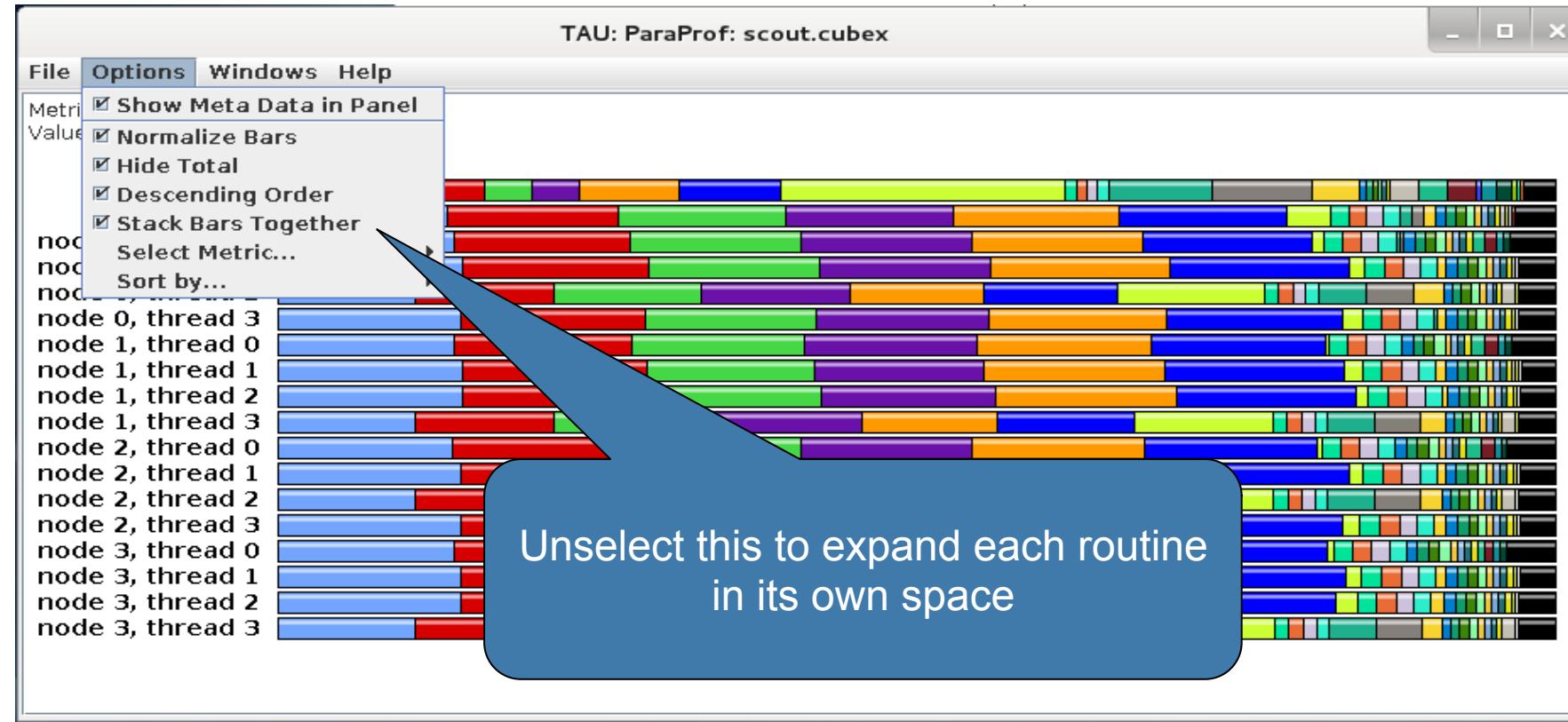
Paraprof main window



Colors represent code regions

Options -> uncheck Stack Bars Together

Paraprof main window



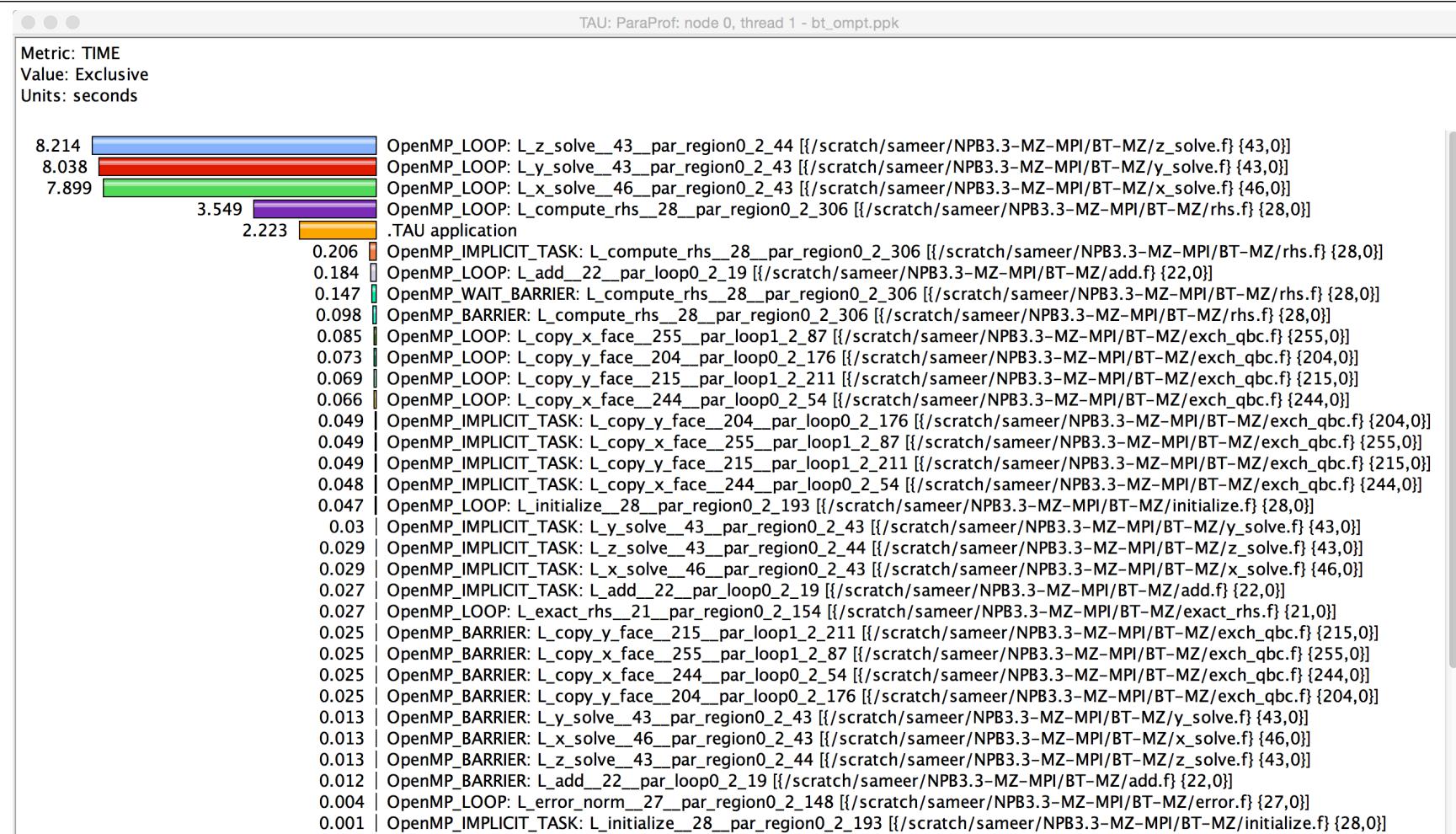
Paraprof main window



Left/right
click here

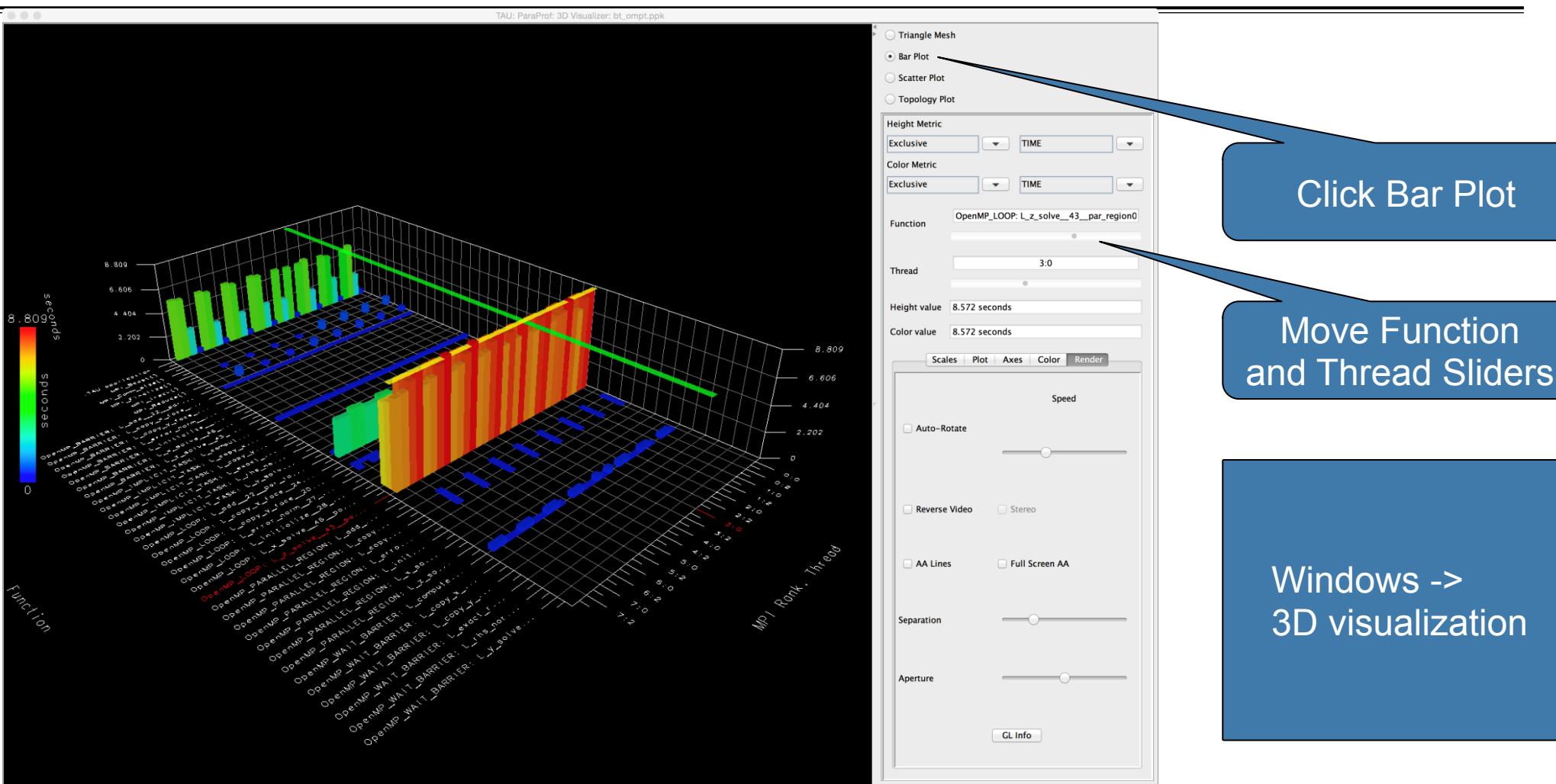
Each routine occupies its own space.
Can see the extent of imbalance
across all threads.

Paraprof node window (function barchart window)

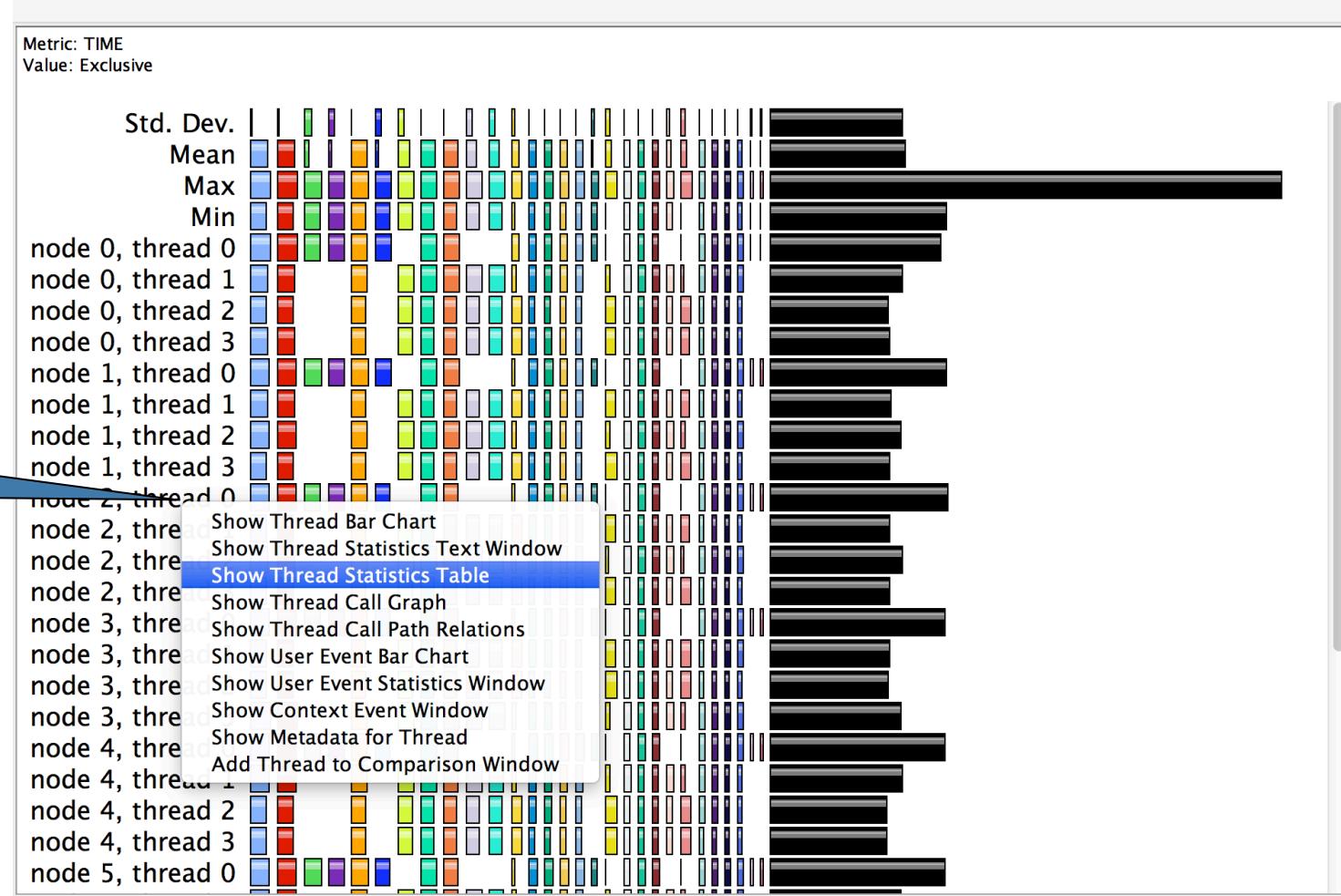


Exclusive time spent in each code region (OpenMP loop) is shown here for MPI rank 0 thread 1

Paraprof 3D visualization window

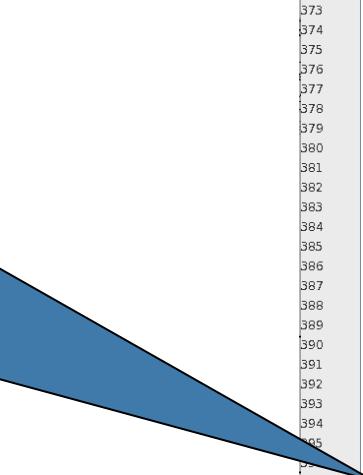


Paraprof Thread Statistics Table with TAU_SAMPLING=1



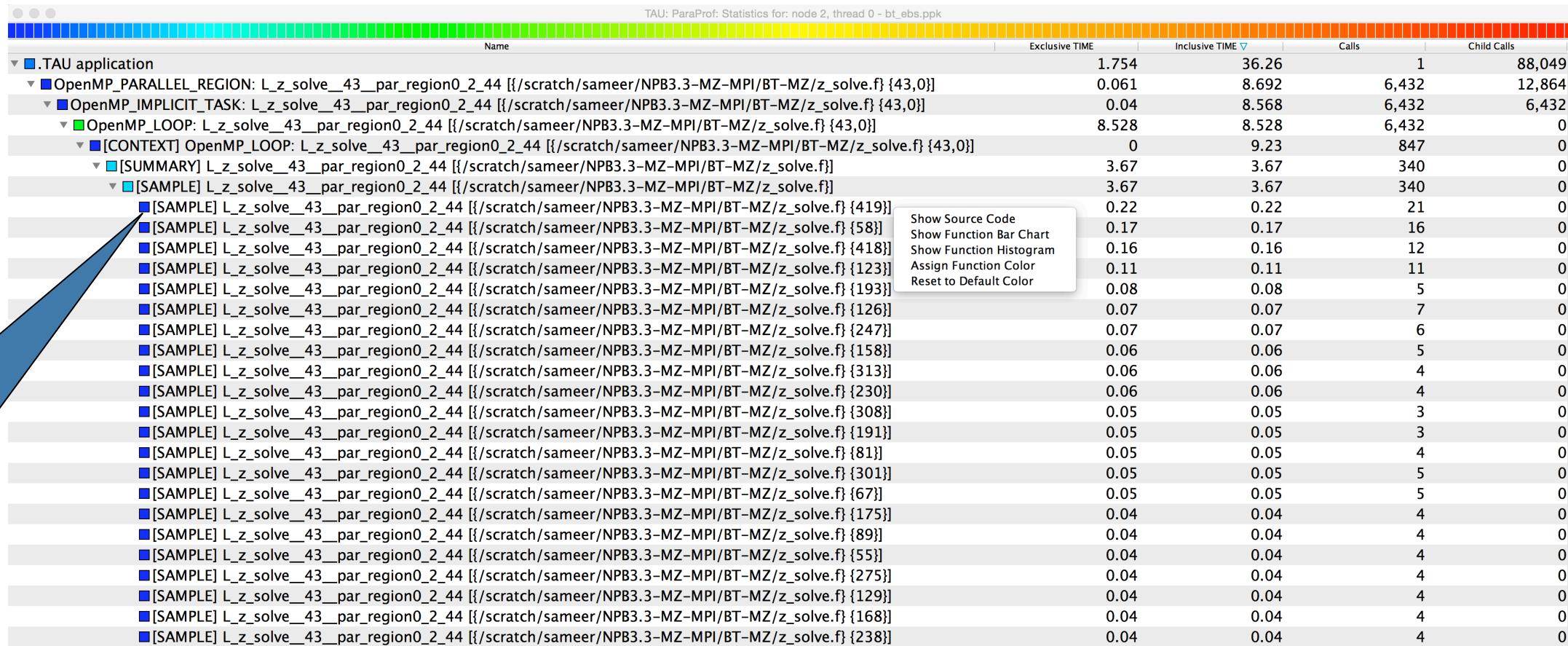
Statement Level Profiling with TAU

Source location where samples are taken.
Compute intensive region.



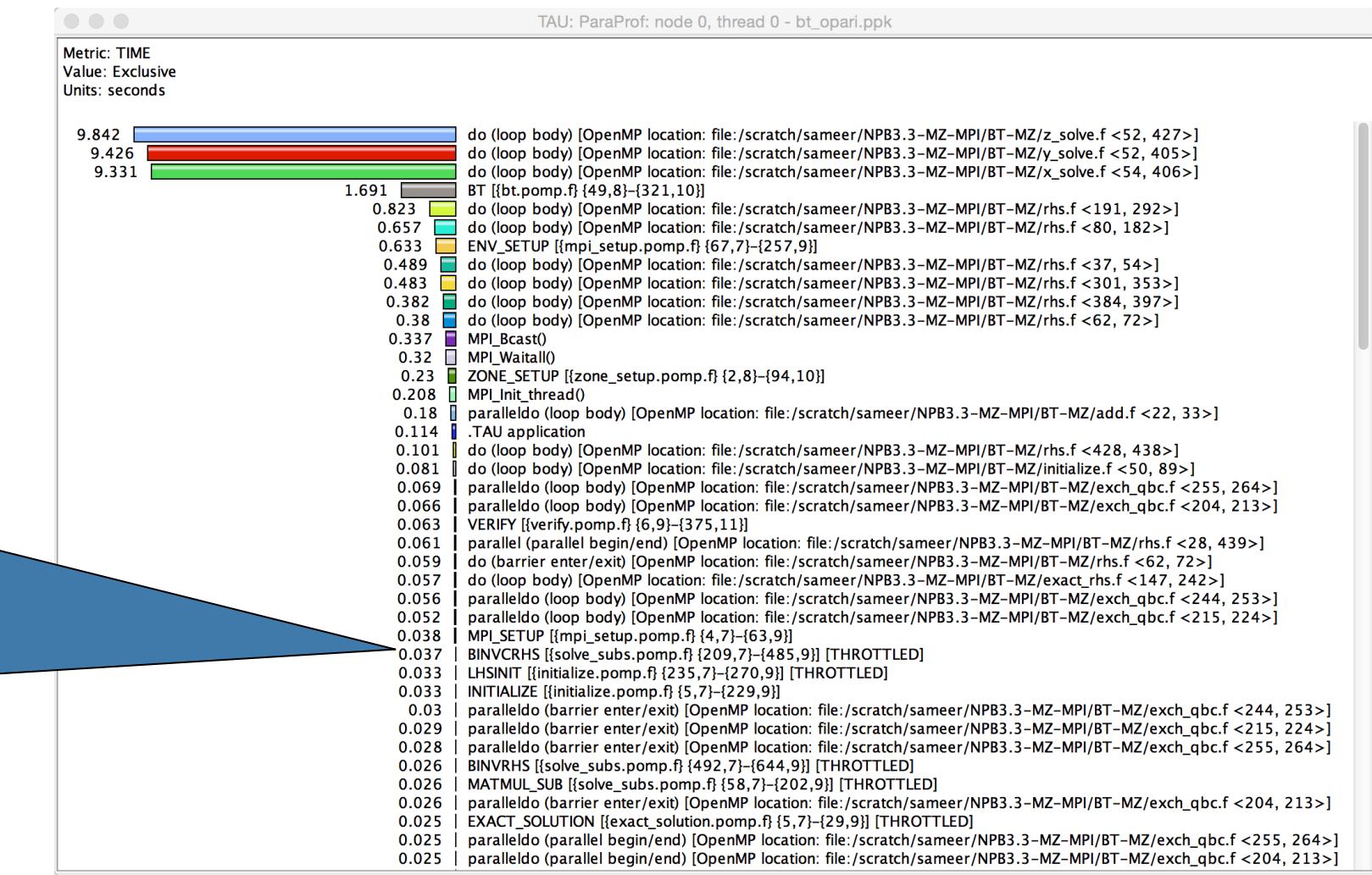
```
File Help
353     call matmul_sub(lhs(1,1,aa,i),
354             >           lhs(1,1,cc,i-1),
355             >           lhs(1,1,bb,i))
356
357
358 C-----c
359 c   multiply c(i,j,k) by b_inverse and copy back to c
360 c   multiply rhs(1,j,k) by b_inverse(1,j,k) and copy to rhs
361 C-----c
362     call binvcrhs( lhs(1,1,bb,i),
363             >           lhs(1,1,cc,i),
364             >           rhs(1,i,j,k) )
365
366     enddo
367
368 C-----c
369 c   rhs(isize) = rhs(isize) - A*rhs(isize-1)
370 C-----c
371     call matvec_sub(lhs(1,1,aa,isize),
372             >           rhs(1,isize-1,j,k),rhs(1,isize,j,k))
373
374 C-----c
375 c   B(isize) = B(isize) - C(isize-1)*A(isize)
376 C-----c
377     call matmul_sub(lhs(1,1,aa,isize),
378             >           lhs(1,1,cc,isize-1),
379             >           lhs(1,1,bb,isize))
380
381 C-----c
382 c   multiply rhs() by b_inverse() and copy to rhs
383 C-----c
384     call binvrhs( lhs(1,1,bb,isize),
385             >           rhs(1,isize,j,k) )
386
387 C-----c
388 c   back solve: if last cell, then generate U(isize)=rhs(isize)
389 c   else assume U(isize) is loaded in un pack backsub_info
390 c   so just use it
391 c   after call u(istart) will be sent to next cell
392 C-----c
393
394 do i=isize-1,0,-1
395     do m=1,BLOCK_SIZE
396         do n=1,BLOCK_SIZE
397             rhs(m,i,j,k) = rhs(m,i,j,k)
398             >           - lhs(m,n,cc,i)*rhs(n,i+1,j,k)
399
400         enddo
401     enddo
402 enddo
```

Paraprof Thread Statistics Table



Instrumenting Source Code with PDT and Opari

Frequently executing lightweight routines are automatically throttled at runtime. Reduces runtime dilation.

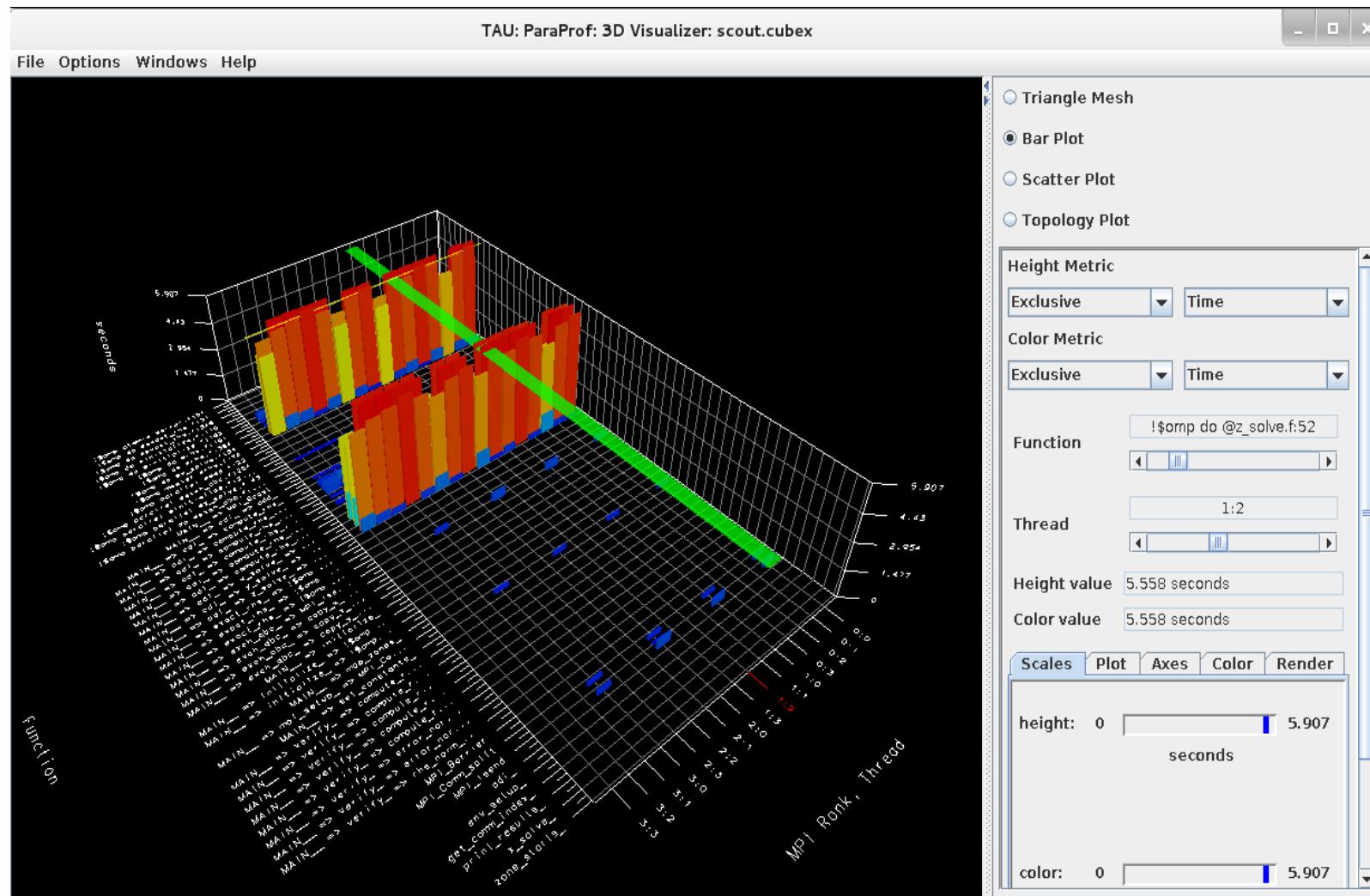


Marker Events in TAU Show Sudden Spikes

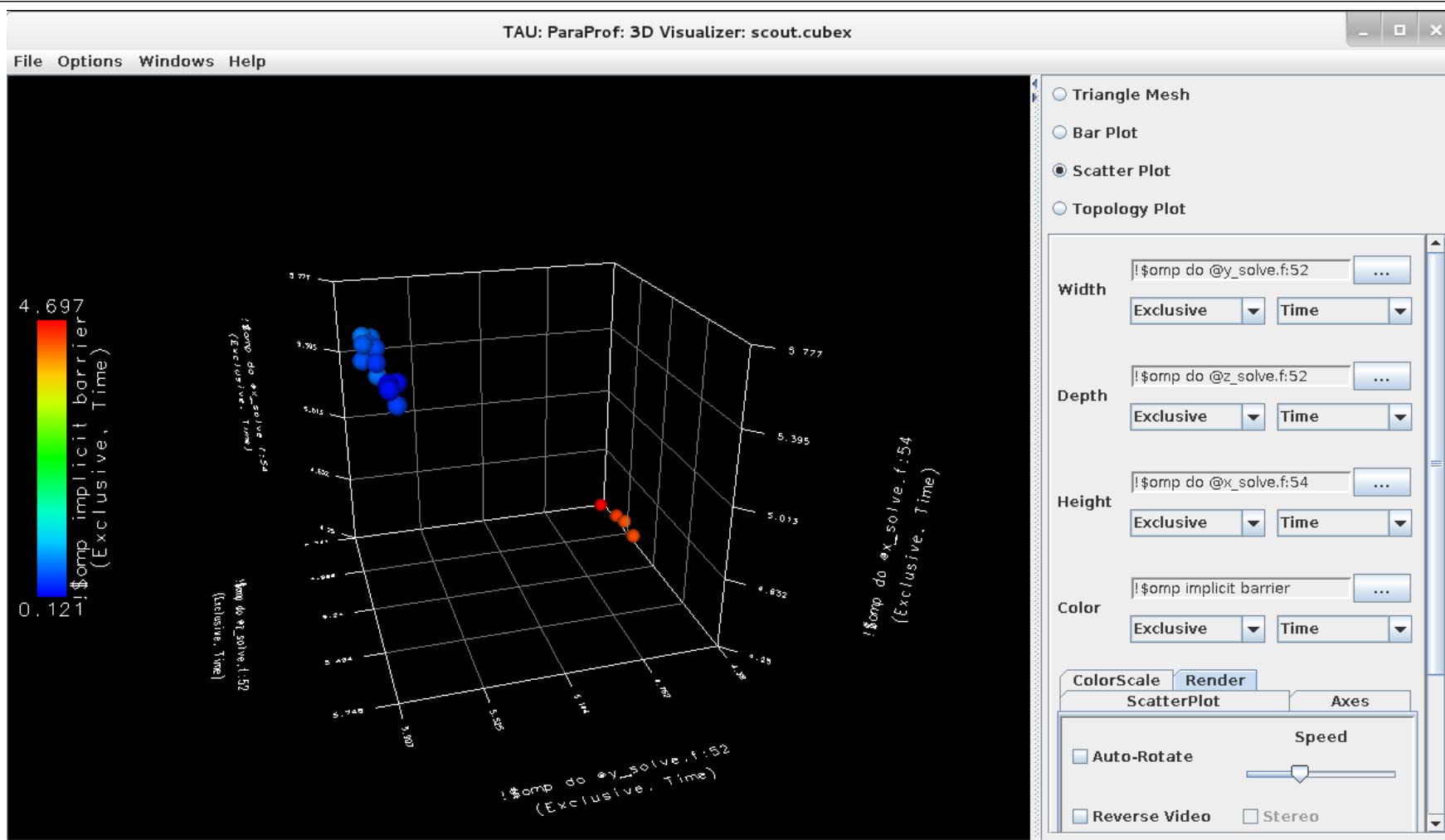
Name	MaxValue	MinValue	NumSamples	MeanValue	Std. Dev.
TAU: ParaProf: Context Events for: node 2 – rapl_marker_16p.ppk					
int main(int, char **) C [{matmult.c}{165,1}–{237,1}]					
double do_work(void) C [{matmult.c}{126,1}–{157,1}]					
void compute(double **, double **, double **, int, int, int) C [{matmult.c}{84,1}–{103,1}]					
[GROUP=MAX_MARKER] rapl:::DRAM_ENERGY:PACKAGE0 (Power in Watts)	17.585	17.469	5	17.521	0.037
[GROUP=MAX_MARKER] rapl:::DRAM_ENERGY:PACKAGE1 (Power in Watts)	15.261	15.218	4	15.237	0.016
[GROUP=MAX_MARKER] rapl:::PACKAGE_ENERGY:PACKAGE0 (Power in Watts)	118.903	114.923	22	116.98	1.201
[GROUP=MAX_MARKER] rapl:::PACKAGE_ENERGY:PACKAGE1 (Power in Watts)	113.466	110.207	22	111.778	0.996
[GROUP=MAX_MARKER] rapl:::PPO_ENERGY:PACKAGE0 (Power in Watts)	100.138	96.266	24	98.206	1.13
[GROUP=MAX_MARKER] rapl:::PPO_ENERGY:PACKAGE1 (Power in Watts)	95.846	92.758	24	94.319	0.937
[GROUP=MIN_MARKER] rapl:::DRAM_ENERGY:PACKAGE0 (Power in Watts)	17.397	17.303	4	17.358	0.035
[GROUP=MIN_MARKER] rapl:::DRAM_ENERGY:PACKAGE1 (Power in Watts)	15.048	15.042	2	15.045	0.003
int mysleep(int) C [{matmult.c}{46,1}–{49,1}]					
[GROUP=MIN_MARKER] rapl:::DRAM_ENERGY:PACKAGE0 (Power in Watts)	15.84	15.84	1	15.84	0
[GROUP=MIN_MARKER] rapl:::DRAM_ENERGY:PACKAGE1 (Power in Watts)	14.275	14.275	1	14.275	0
[GROUP=MIN_MARKER] rapl:::PACKAGE_ENERGY:PACKAGE1 (Power in Watts)	96.853	96.853	1	96.853	0
[GROUP=MIN_MARKER] rapl:::PACKAGE_ENERGY:PACKAGE0 (Power in Watts)	93.125	93.125	1	93.125	0
[GROUP=MIN_MARKER] rapl:::PPO_ENERGY:PACKAGE0 (Power in Watts)	75.096	75.096	1	75.096	0
[GROUP=MIN_MARKER] rapl:::PPO_ENERGY:PACKAGE1 (Power in Watts)	79.646	79.646	1	79.646	0
void compute_interchange(double **, double **, double **, int, int, int) C [{matmult.c}{105,1}–{124,1}]					
[GROUP=MAX_MARKER] rapl:::DRAM_ENERGY:PACKAGE0 (Power in Watts)	26.064	25.711	2	25.887	0.176
[GROUP=MAX_MARKER] rapl:::DRAM_ENERGY:PACKAGE1 (Power in Watts)	24.373	23.965	4	24.232	0.159
[GROUP=MAX_MARKER] rapl:::PACKAGE_ENERGY:PACKAGE0 (Power in Watts)	126.872	125.182	6	125.732	0.557
[GROUP=MAX_MARKER] rapl:::PACKAGE_ENERGY:PACKAGE1 (Power in Watts)	124.377	116.689	5	122.428	2.885
[GROUP=MAX_MARKER] rapl:::PPO_ENERGY:PACKAGE0 (Power in Watts)	103.981	102.21	6	102.769	0.584
[GROUP=MAX_MARKER] rapl:::PPO_ENERGY:PACKAGE1 (Power in Watts)	102.615	101.693	4	102.115	0.33
rapl:::DRAM_ENERGY:PACKAGE0 (Power in Watts)	26.064	15.84	36	19.053	3.39
rapl:::DRAM_ENERGY:PACKAGE1 (Power in Watts)	24.373	14.275	36	16.435	3.155
rapl:::PACKAGE_ENERGY:PACKAGE0 (Power in Watts)	126.872	93.125	36	117.729	5.403
rapl:::PACKAGE_ENERGY:PACKAGE1 (Power in Watts)	124.377	96.853	36	112.961	4.776
rapl:::PPO_ENERGY:PACKAGE0 (Power in Watts)	103.981	75.096	36	98.208	4.466
rapl:::PPO_ENERGY:PACKAGE1 (Power in Watts)	102.615	79.646	36	94.872	3.662

```
% export TAU_EVENT_THRESHOLD 0.5
```

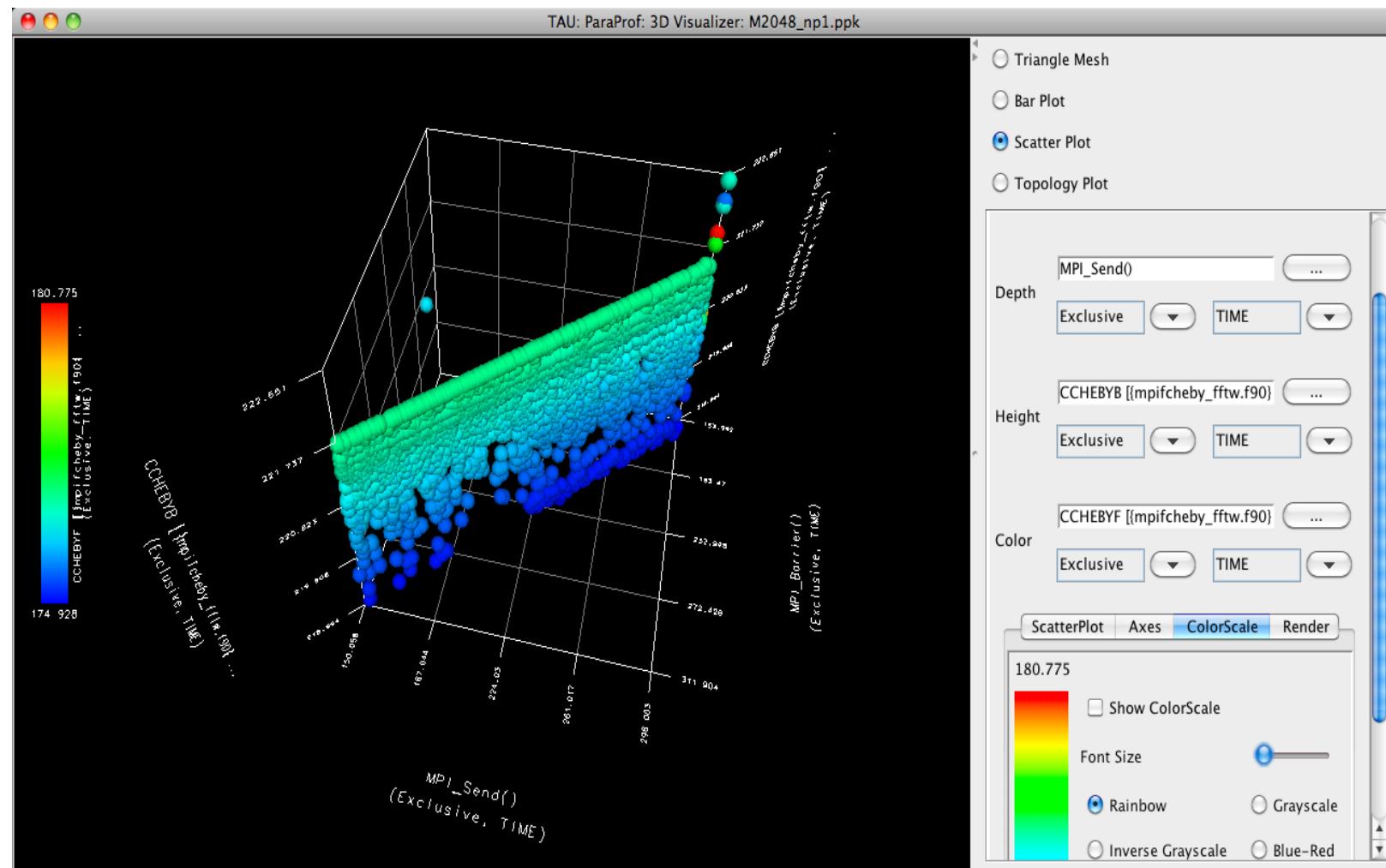
ParaProf: 3D Visualization Window



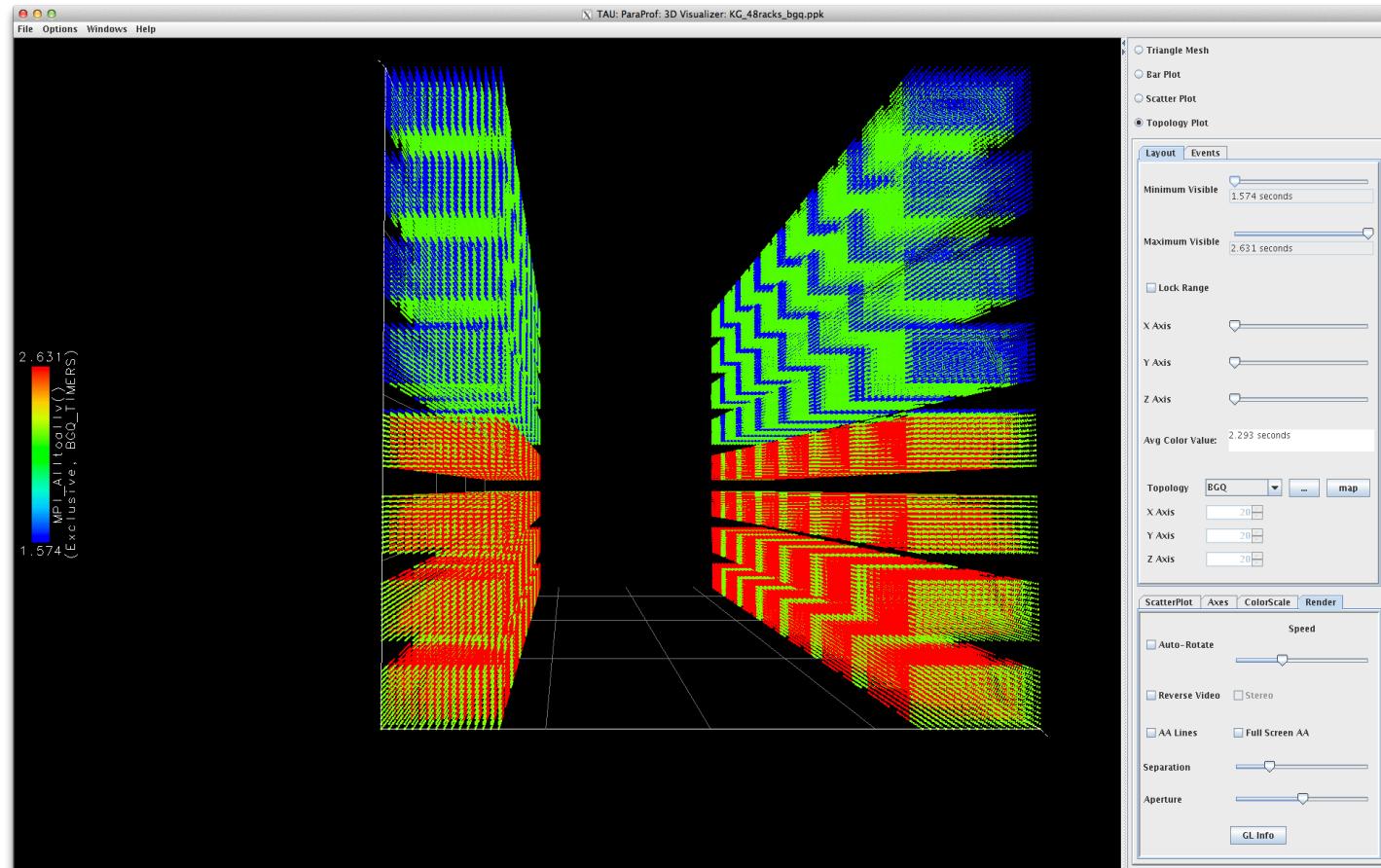
ParaProf: 3D Scatter Plot



ParaProf: Scatter Plot



ParaProf Topology View: Torus

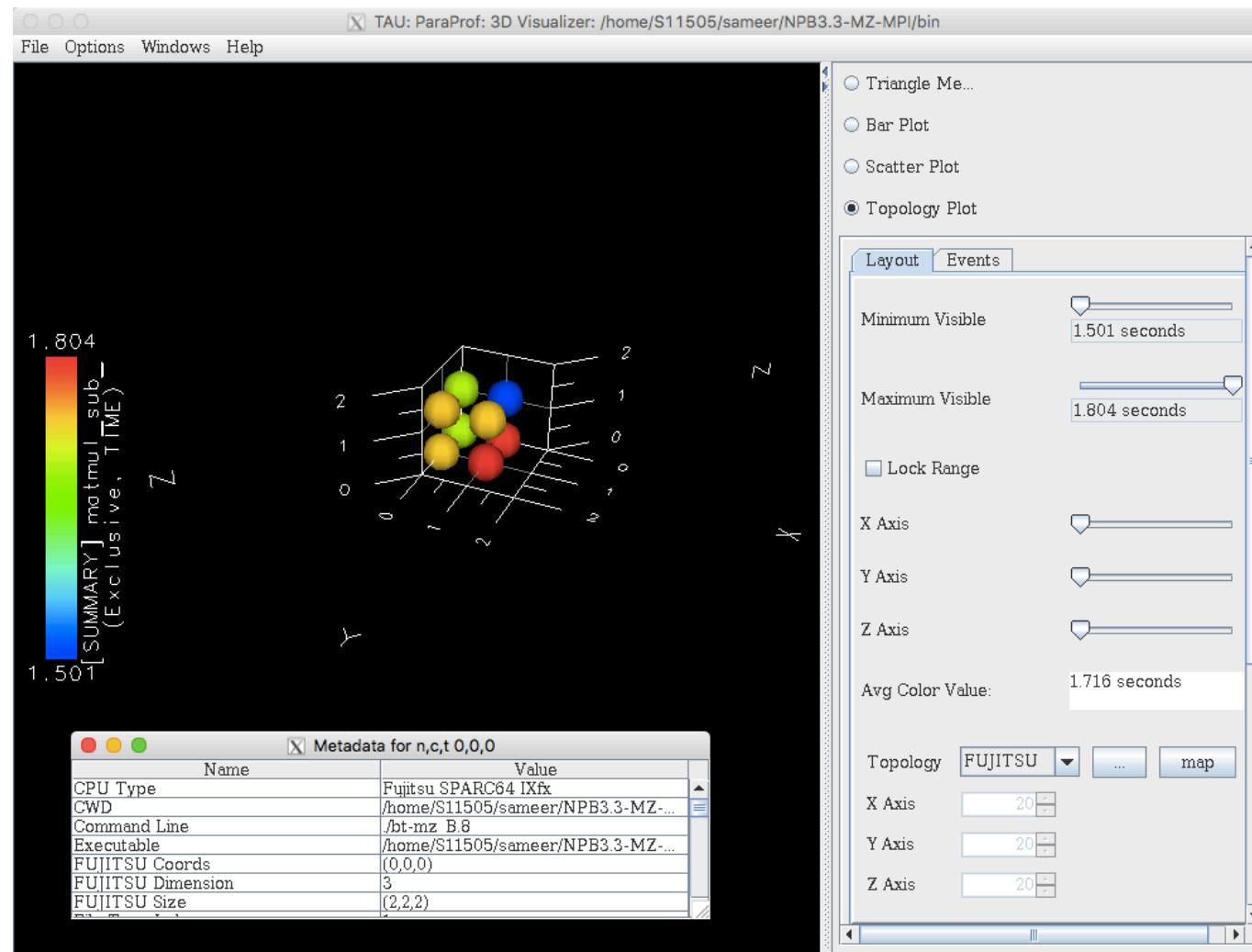


ParaProf MetaData: FX10

Name	Value
FUJITSU Coords	(3,1,0)
FUJITSU Dimension	3
FUJITSU Size	(6,6,6)
File Type Index	0
File Type Name	ParaProf Packed Profile
Hostname	e09t14226
Local Time	2012-11-12T02:14:16+09:00
MPI Processor Name	e09t14226
Memory Size	32836968 kB
Node Name	e09t14226
OS Machine	s64fx
OS Name	Linux
OS Release	2.6.25.8
OS Version	#1 SMP Tue Sep 11 11:04:02 JST 2...
Starting Timestamp	1352654056461761
TAU Architecture	sparc64fx

```
% pbsub –interact –L node=“6x6x6”
% mpiexec –n 216 ./a.out
```

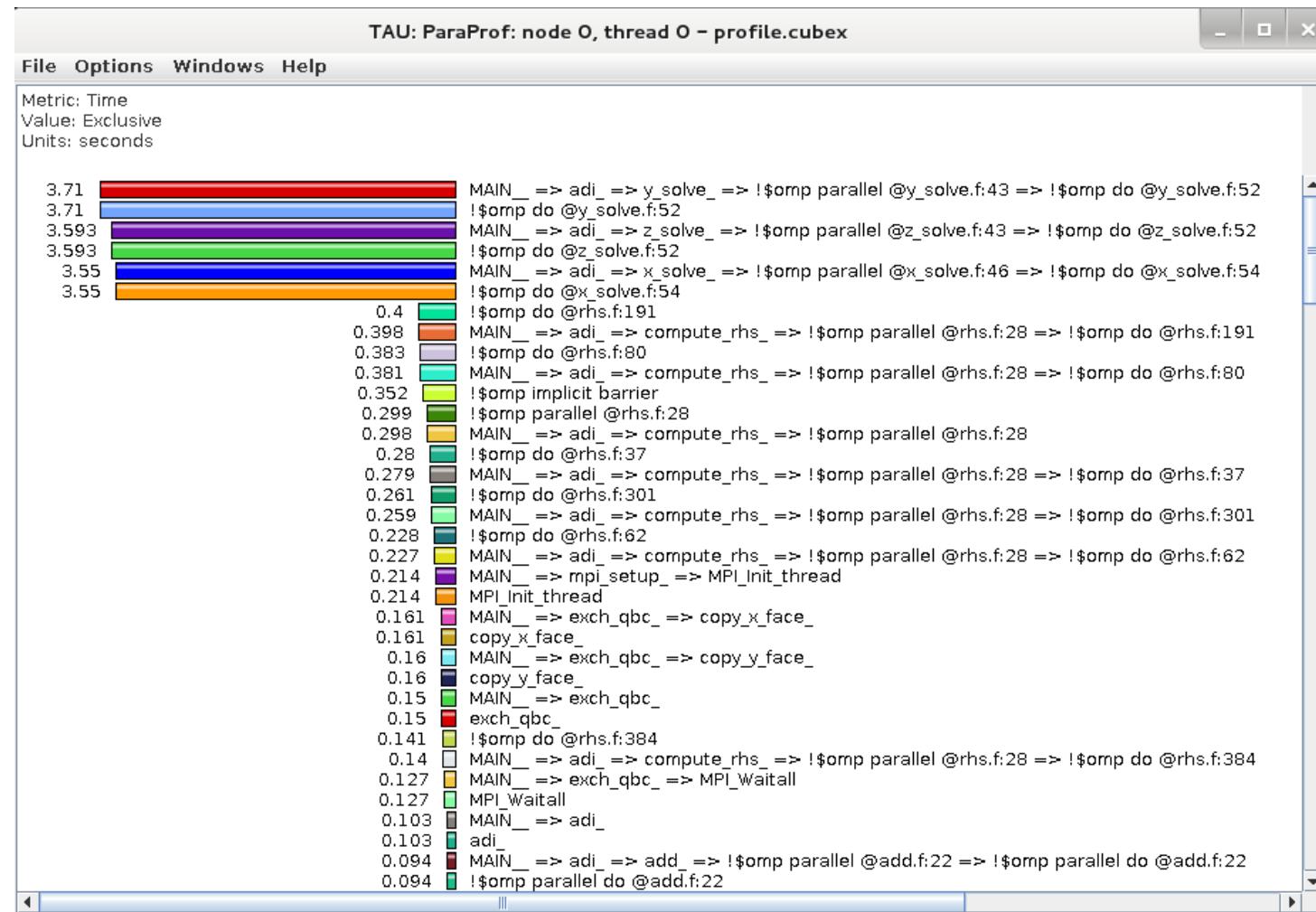
ParaProf 3D Topology Display for Fujitsu FX10



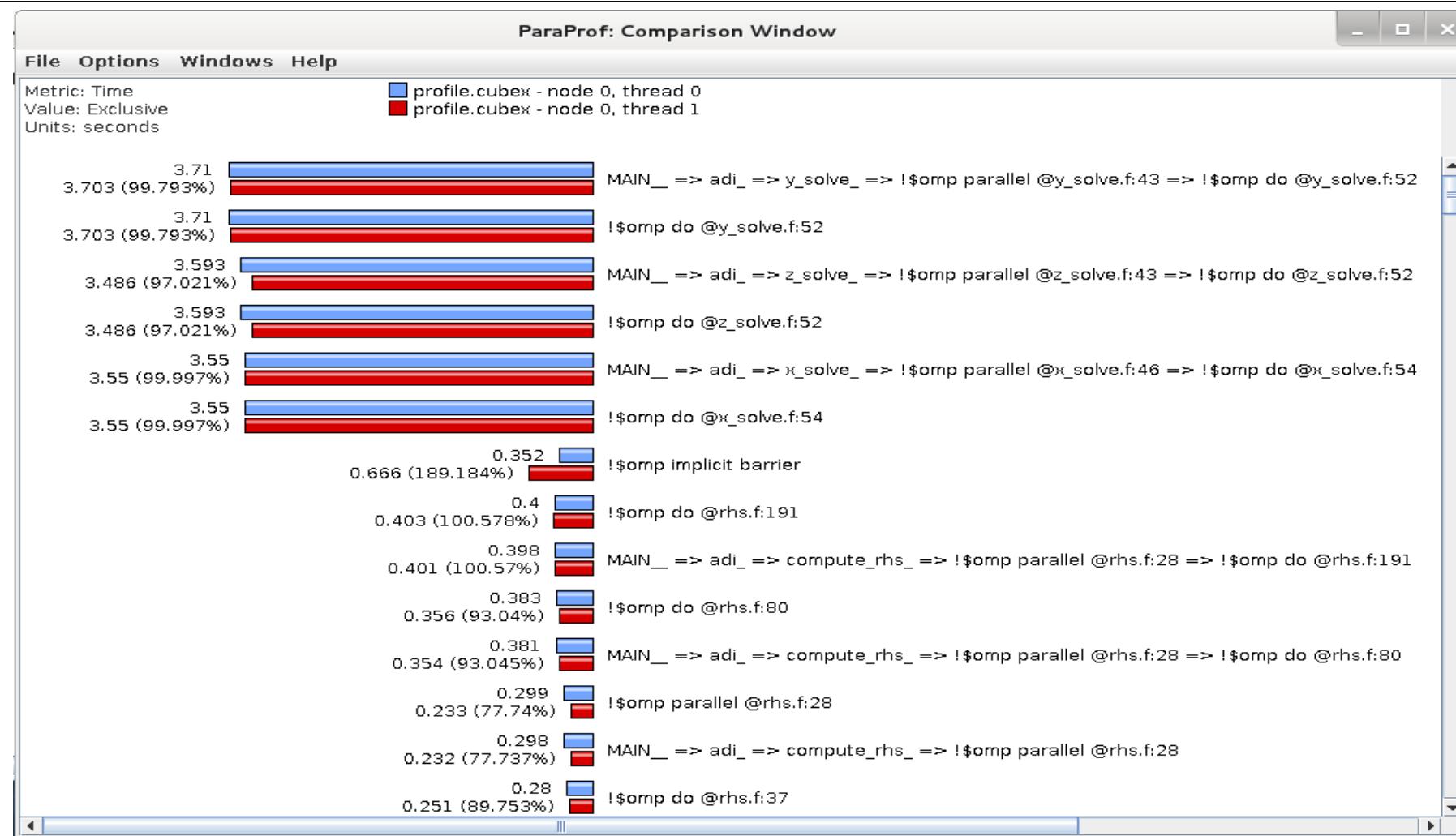
#PJSUB -L node="2x2x2"

mpiexec -n 8 tau_exec -ebs \$EXE

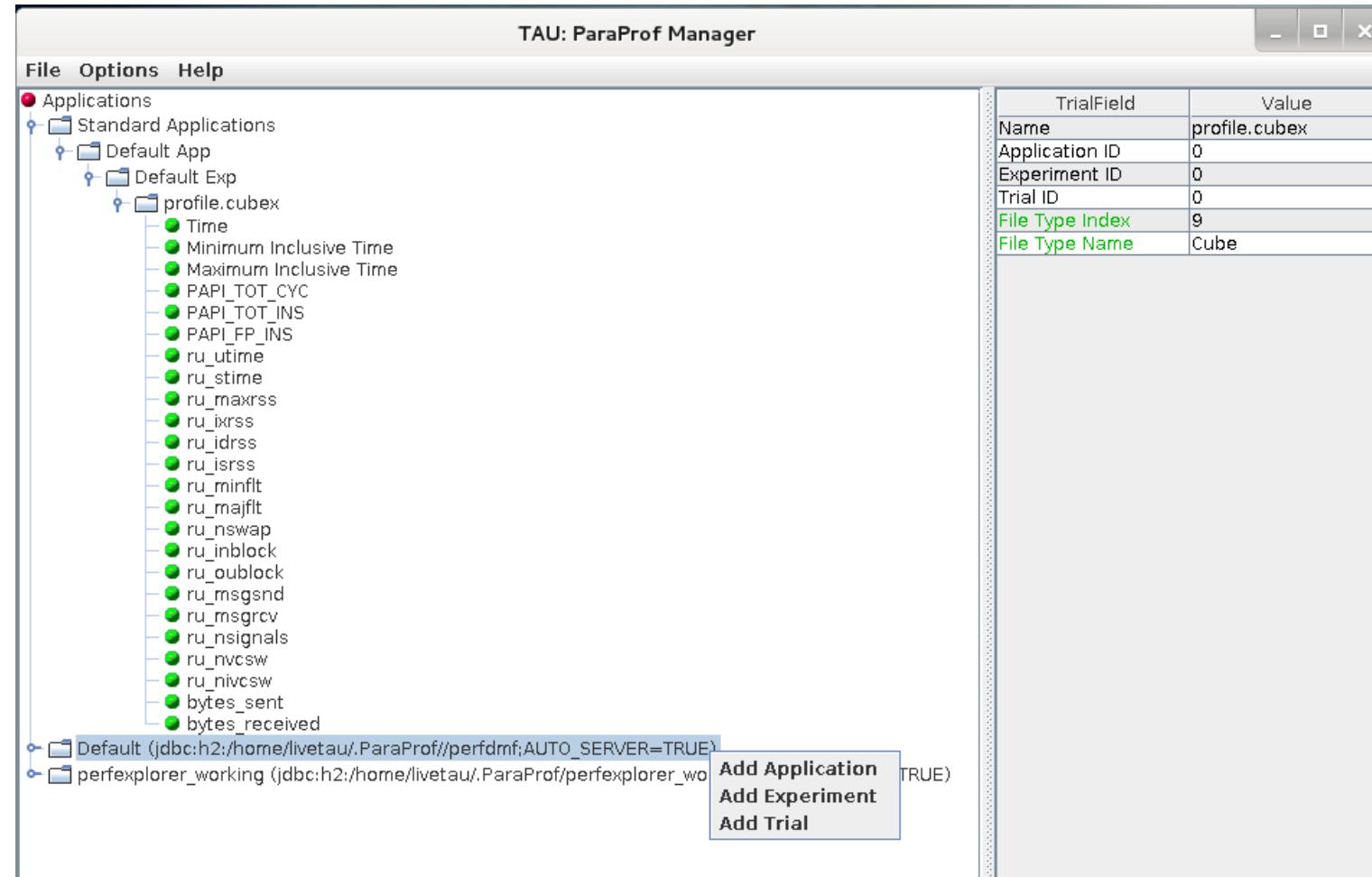
ParaProf: Node View



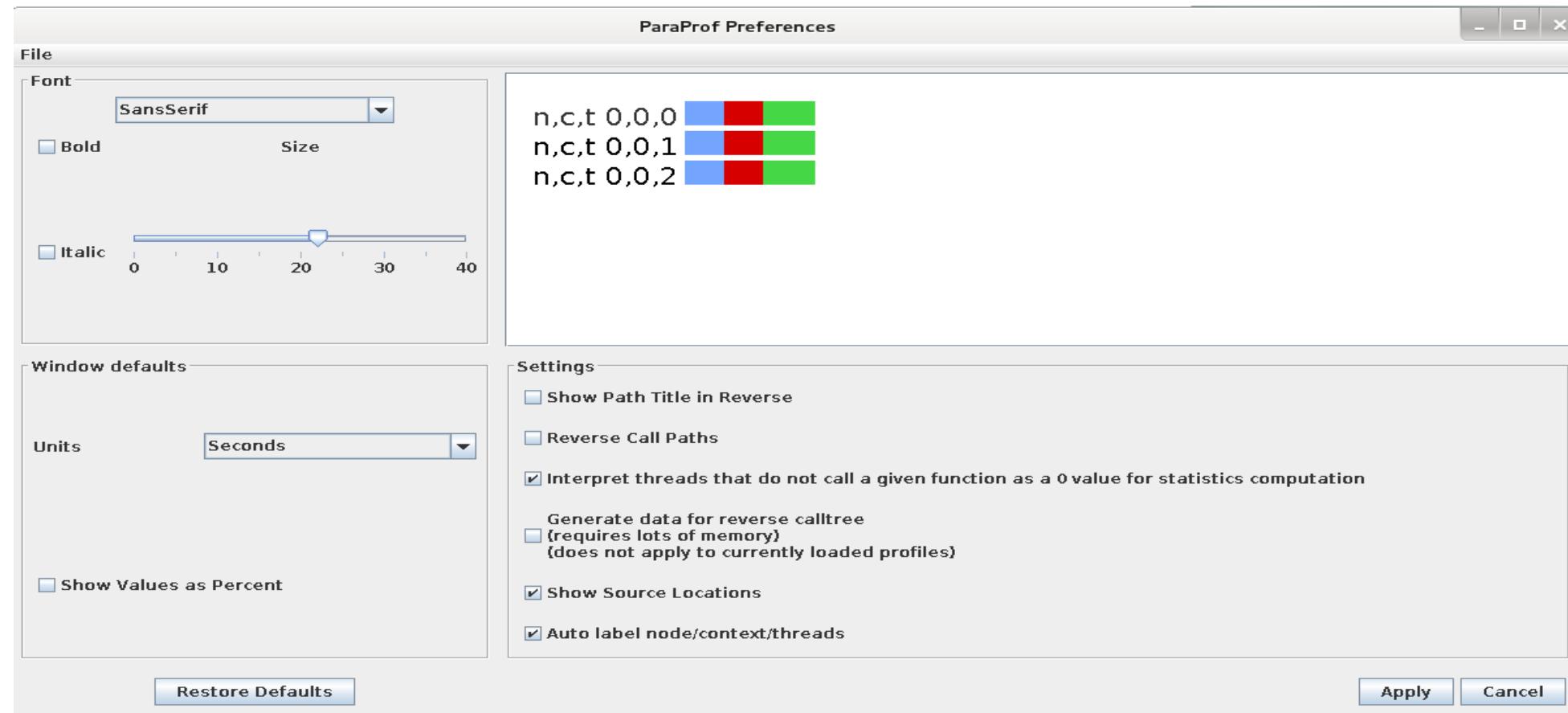
ParaProf: Add Thread to Comparison Window



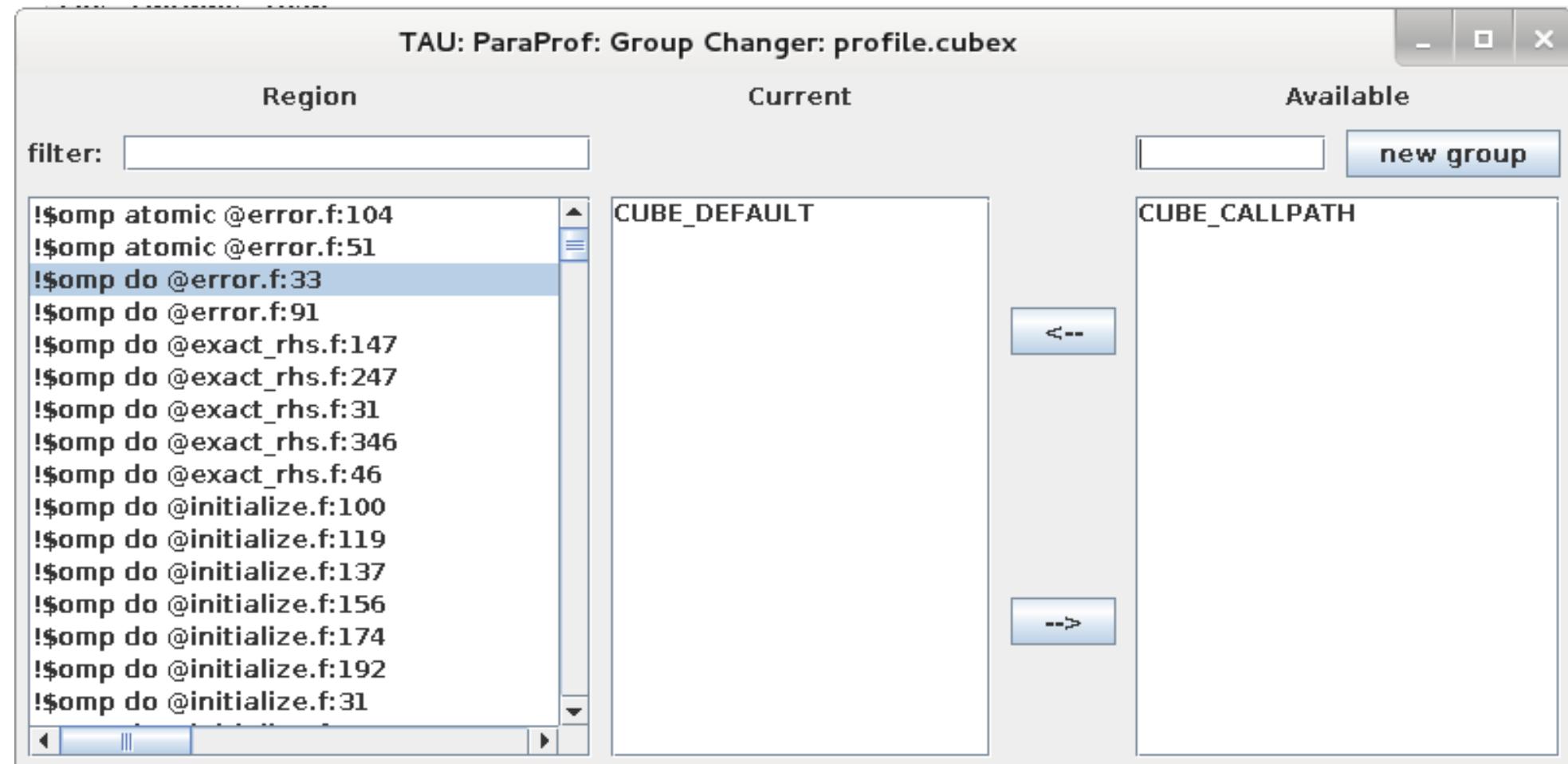
ParaProf: Score-P Profile Files, Database



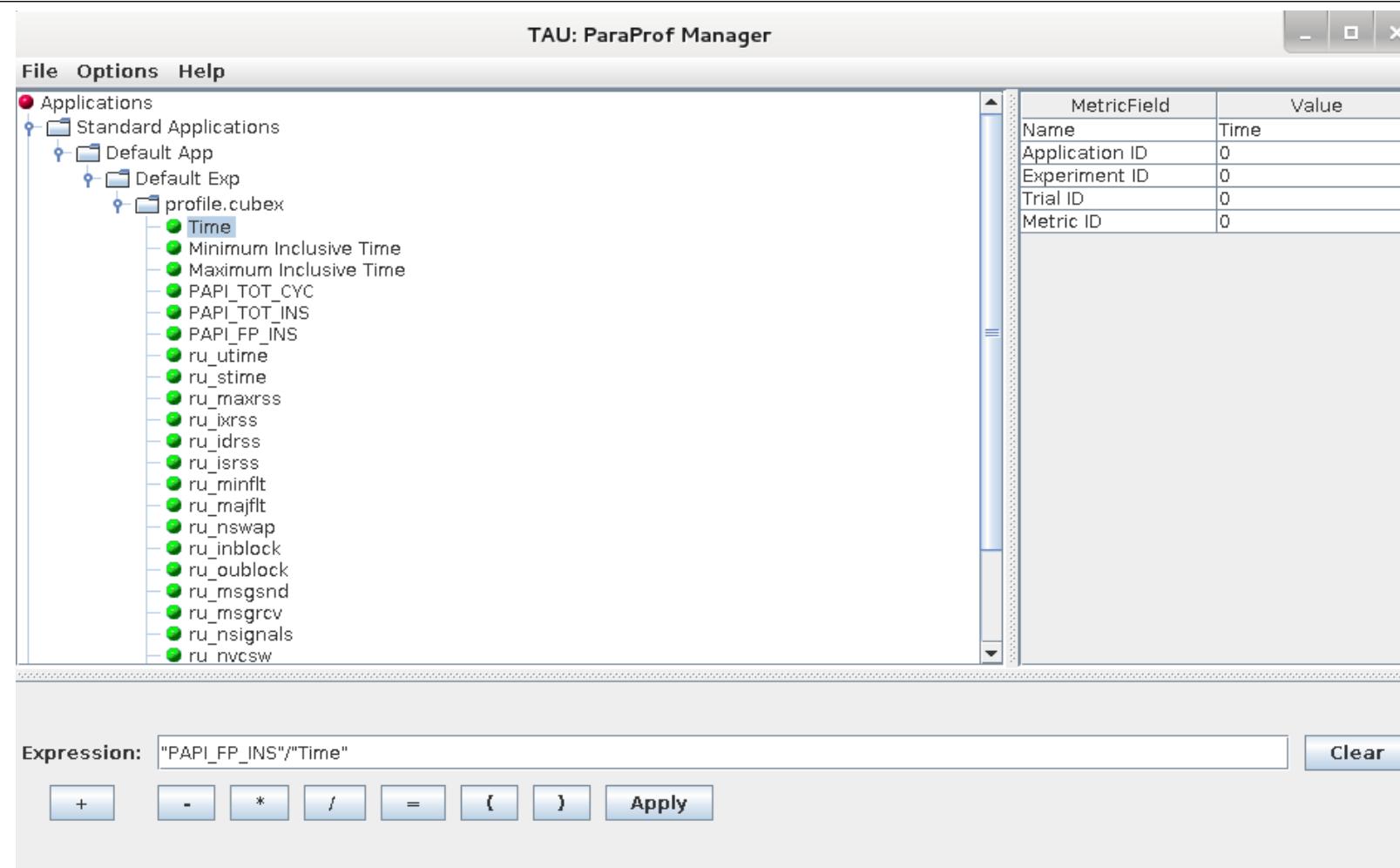
ParaProf: File Preferences Window



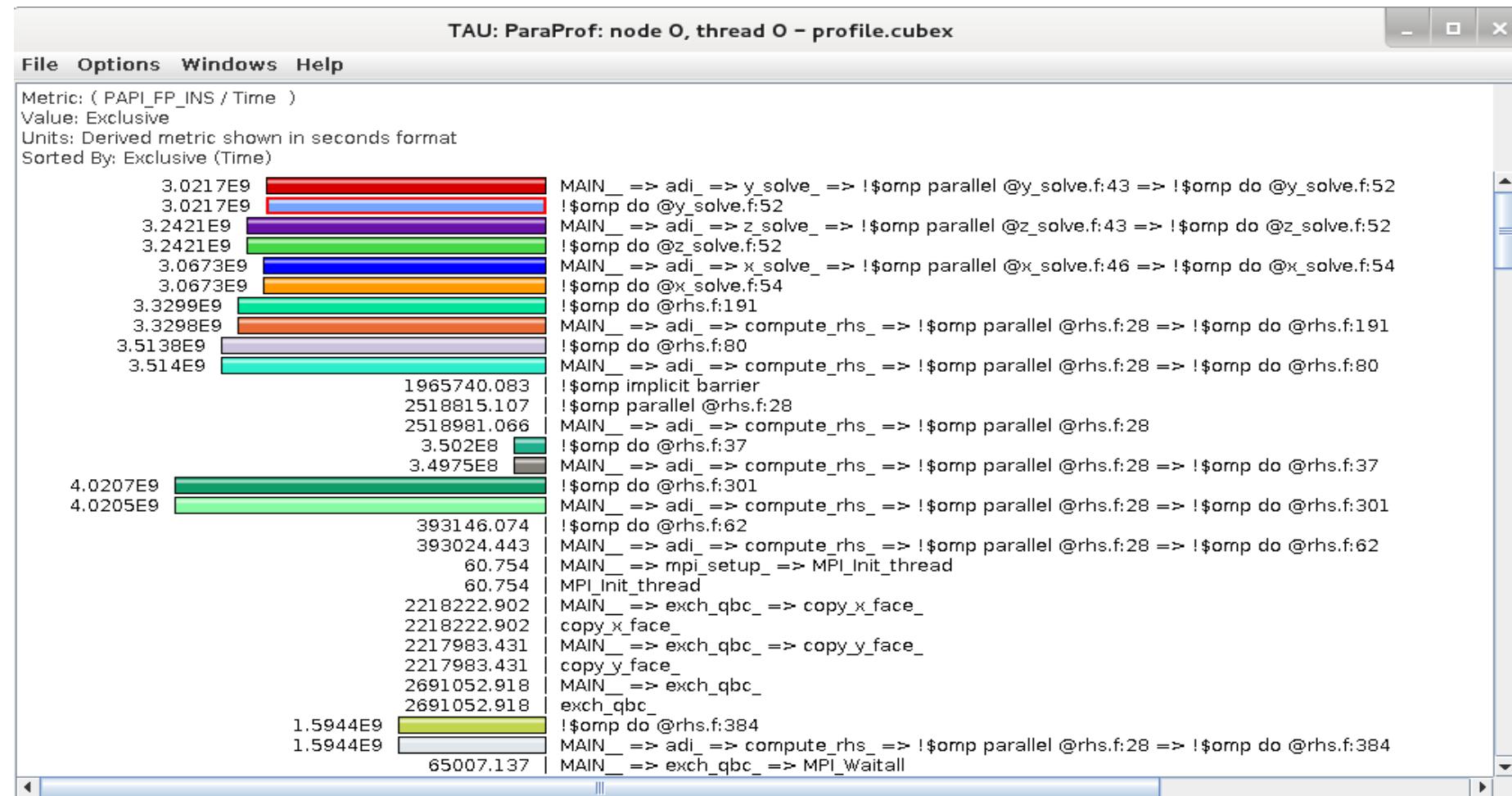
ParaProf: Group Changer Window



ParaProf: Derived Metric Panel in Manager Window

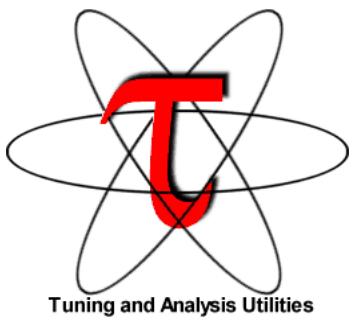


Sorting Derived FLOPS Metric by Exclusive Time





Performance Data Management with TAU PerfExplorer

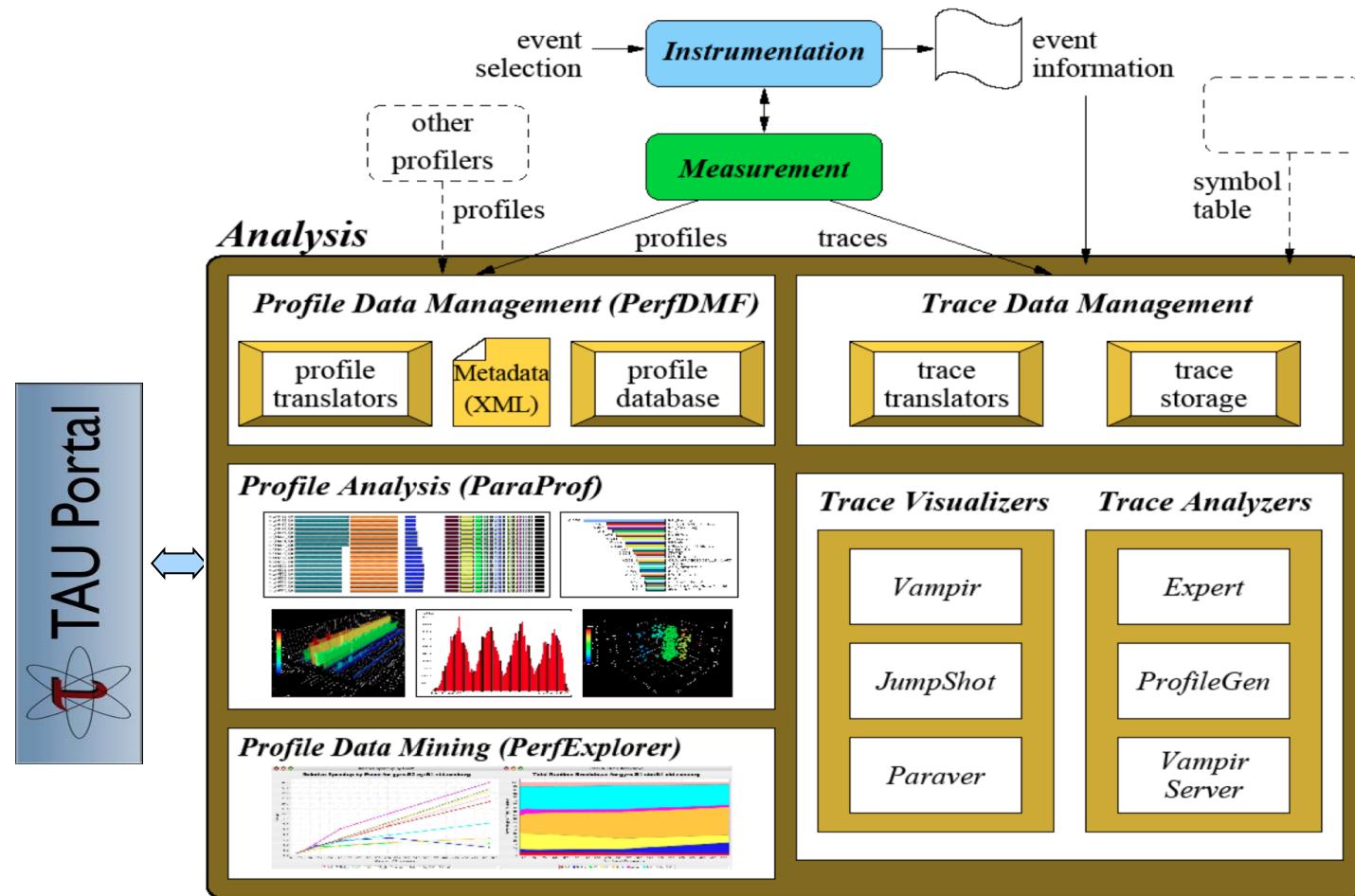


Tuning and Analysis Utilities

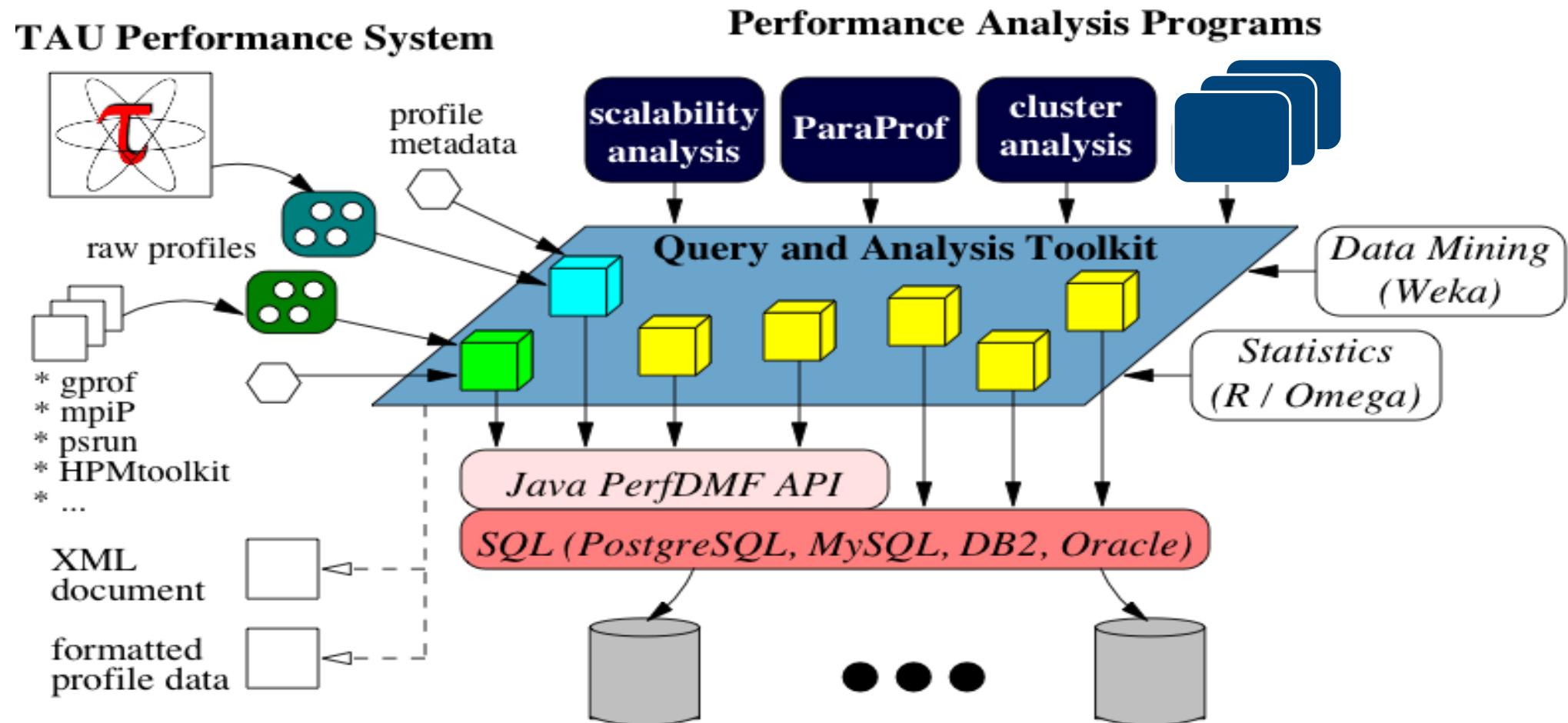
Sameer Shende
sameer@cs.uoregon.edu
Performance Research Laboratory
University of Oregon
<http://tau.uoregon.edu>



TAU Analysis



TAUdb: Performance Data Management Framework



Using TAUdb

- Configure TAUdb (Done by each user)

```
% taudb_configure --create-default
```

- Choose derby, PostgreSQL, MySQL, Oracle or DB2
- Hostname
- Username
- Password
- Say yes to downloading required drivers (we are not allowed to distribute these)
- Stores parameters in your ~/.ParaProf/taudb.cfg file

- Configure PerfExplorer (Done by each user)

```
% perfexplorer_configure
```

- Execute PerfExplorer

```
% perfexplorer
```

Using PerfExplorer

```
% wget http://tau.uoregon.edu/data.tgz (Contains CUBE profiles from Score-P)
% taudb_configure --create-default
(Chooses derby, blank user/passwd, yes to save passwd, defaults)
% perfexplorer_configure
(Yes to load schema, defaults)
% paraprof
(load each trial: DB -> Add Trial -> Type (Paraprof Packed Profile) -> OK)
OR use taudb_loadtrial -a "app" -x "experiment" -n "name" file.ppk
Then,
% tar zxf $TAU/data.tgz; cd data/tau;
% taudb_loadtrial -a BT_MZ -x "Class_B" bt-mz_B.*.ppk
% perfexplorer
(Select experiment, Menu: Charts -> Speedup)
```

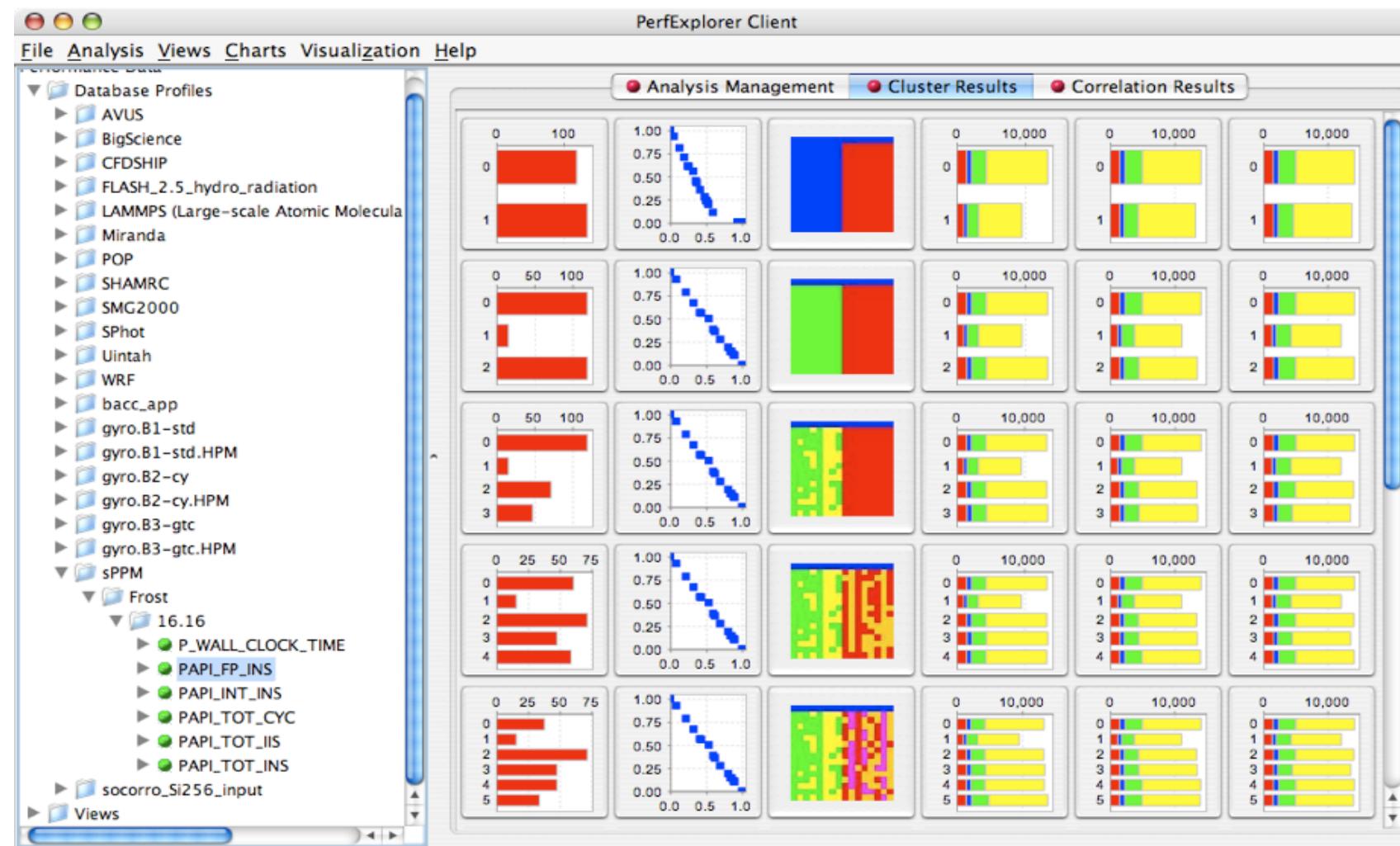
Performance Data Mining (PerfExplorer)

- Performance knowledge discovery framework
 - Data mining analysis applied to parallel performance data
 - comparative, clustering, correlation, dimension reduction, ...
 - Use the existing TAU infrastructure
 - TAU performance profiles, taudb
 - Client-server based system architecture
- Technology integration
 - Java API and toolkit for portability
 - taudb
 - R-project/Omegahat, Octave/Matlab statistical analysis
 - WEKA data mining package
 - JFreeChart for visualization, vector output (EPS, SVG)

PerfExplorer: Using Cluster Analysis

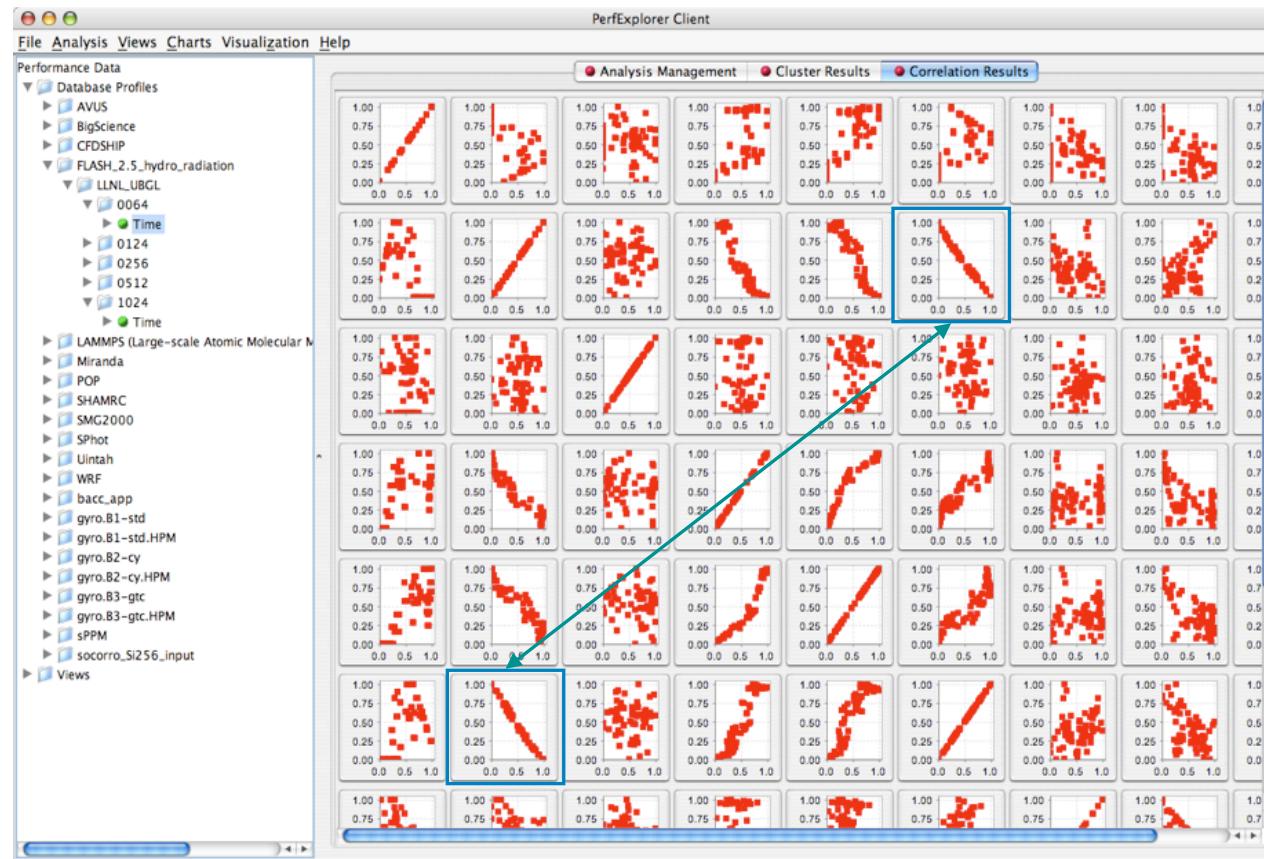
- Performance data represented as vectors - each dimension is the cumulative time for an event
- k -means: k random centers are selected and instances are grouped with the "closest" (Euclidean) center
- New centers are calculated and the process repeated until stabilization or max iterations
- Dimension reduction necessary for meaningful results
- Virtual topology, summaries constructed

PerfExplorer - Cluster Analysis (sPPM)



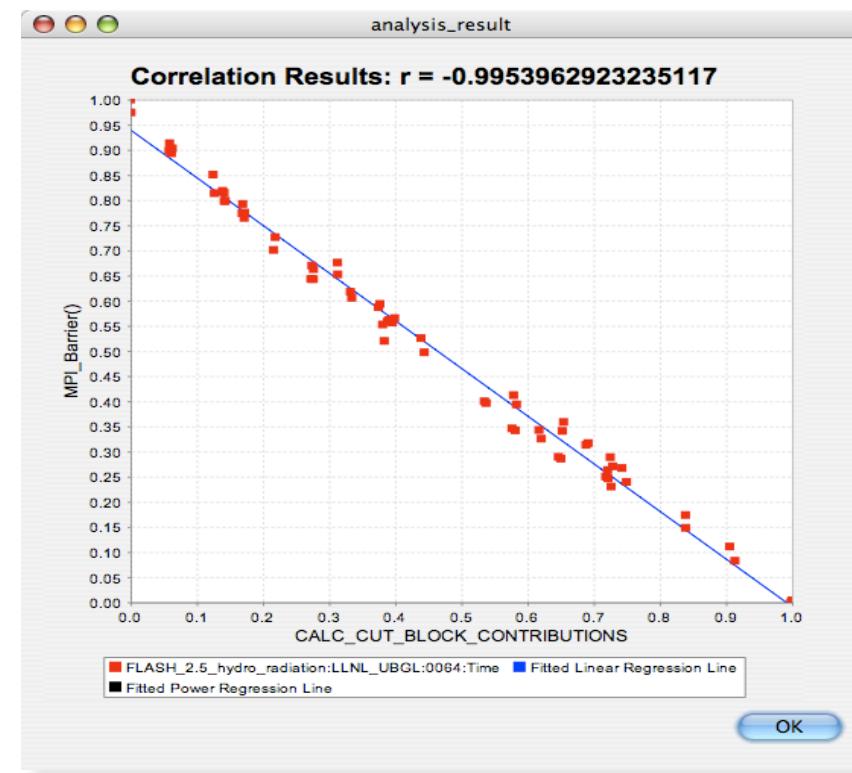
PerfExplorer - Correlation Analysis (Flash)

- Describes strength and direction of a linear relationship between two variables (events) in the data



PerfExplorer - Correlation Analysis (Flash)

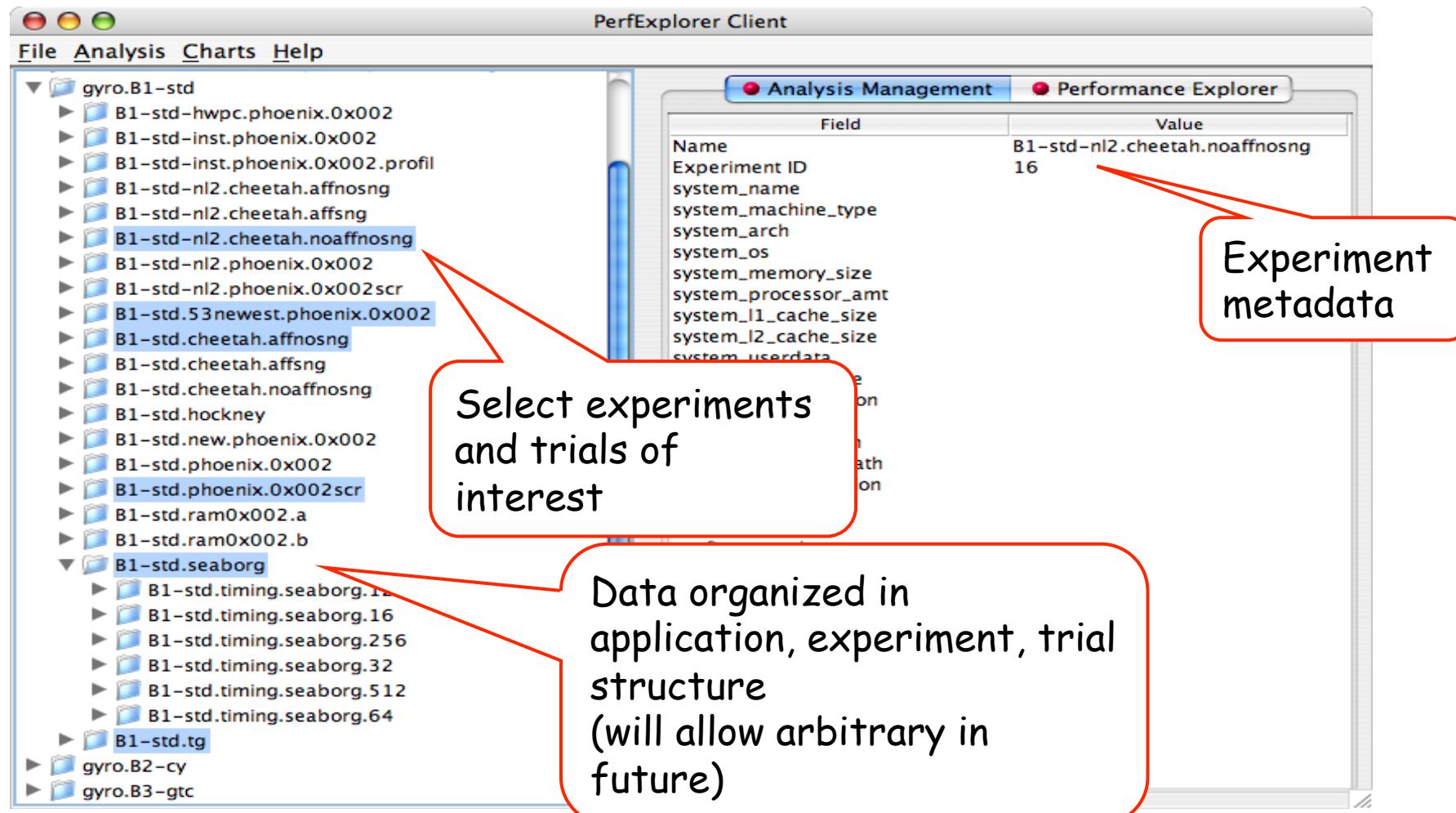
- 0.995 indicates strong, negative relationship
- As CALC_CUT_BLOCK_CONTRIBUTIONS() increases in execution time, MPI_Barrier() decreases



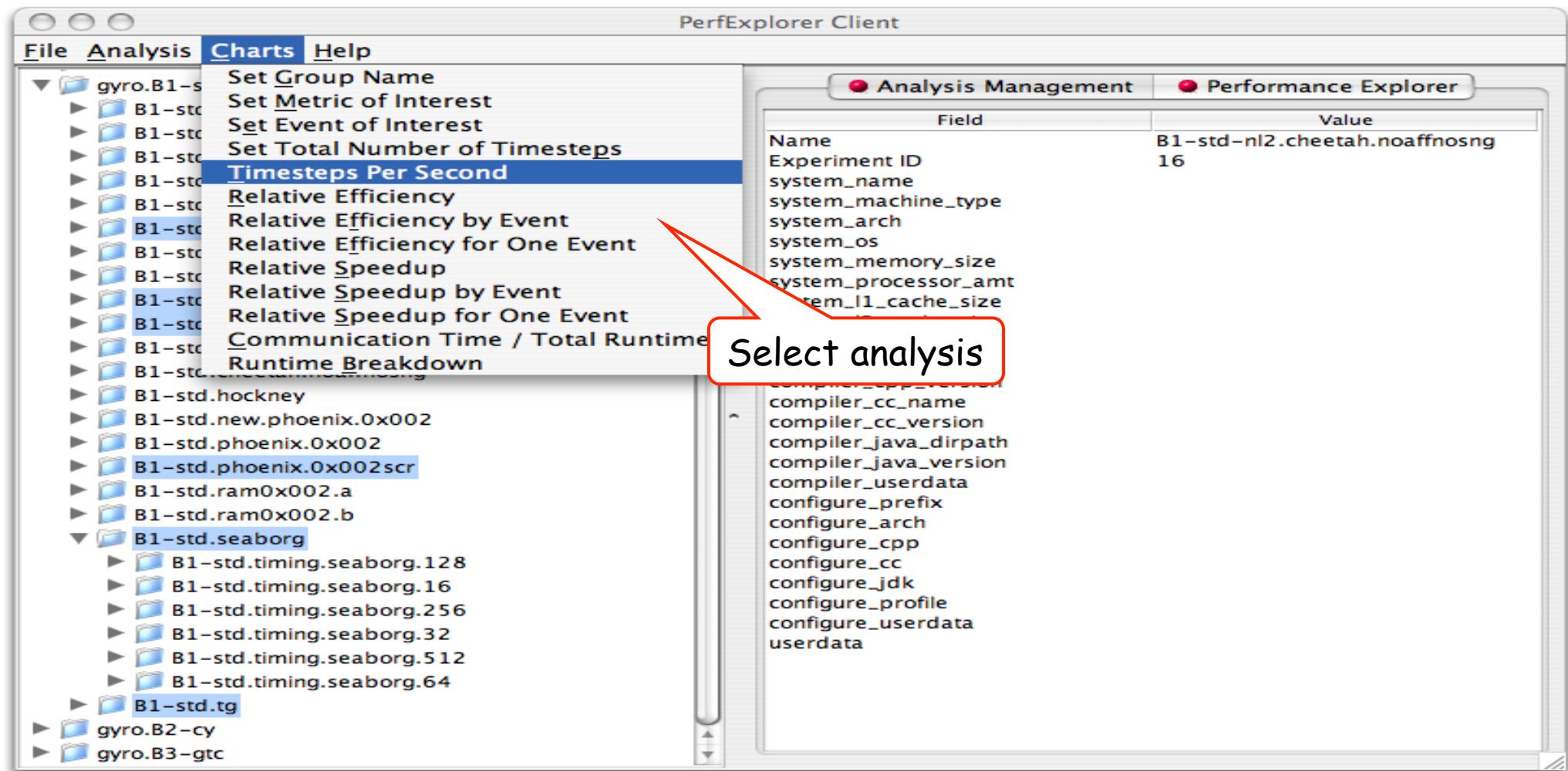
PerfExplorer - Comparative Analysis

- Relative speedup, efficiency
 - total runtime, by event, one event, by phase
- Breakdown of total runtime
- Group fraction of total runtime
- Correlating events to total runtime
- Timesteps per second

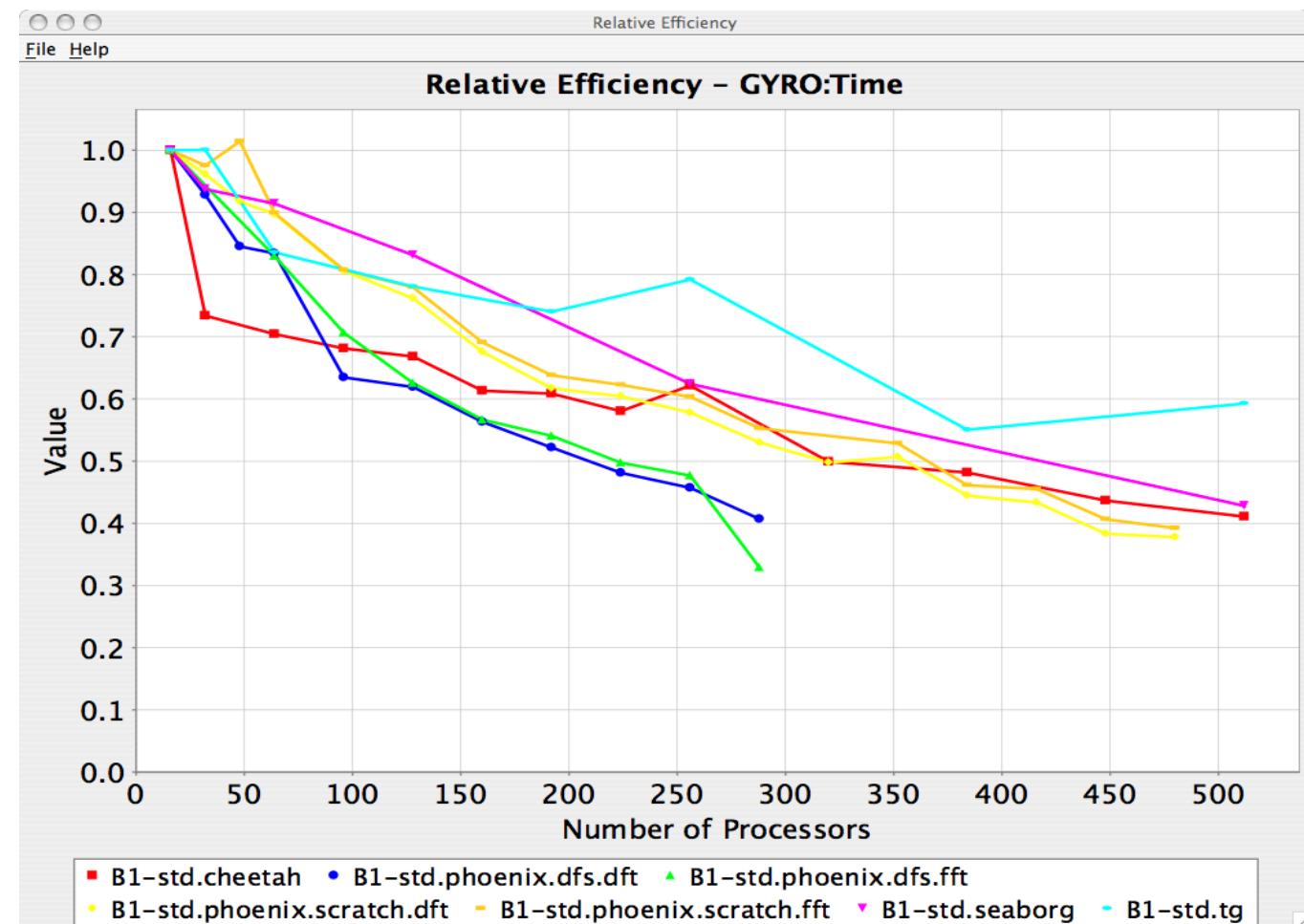
PerfExplorer - Interface



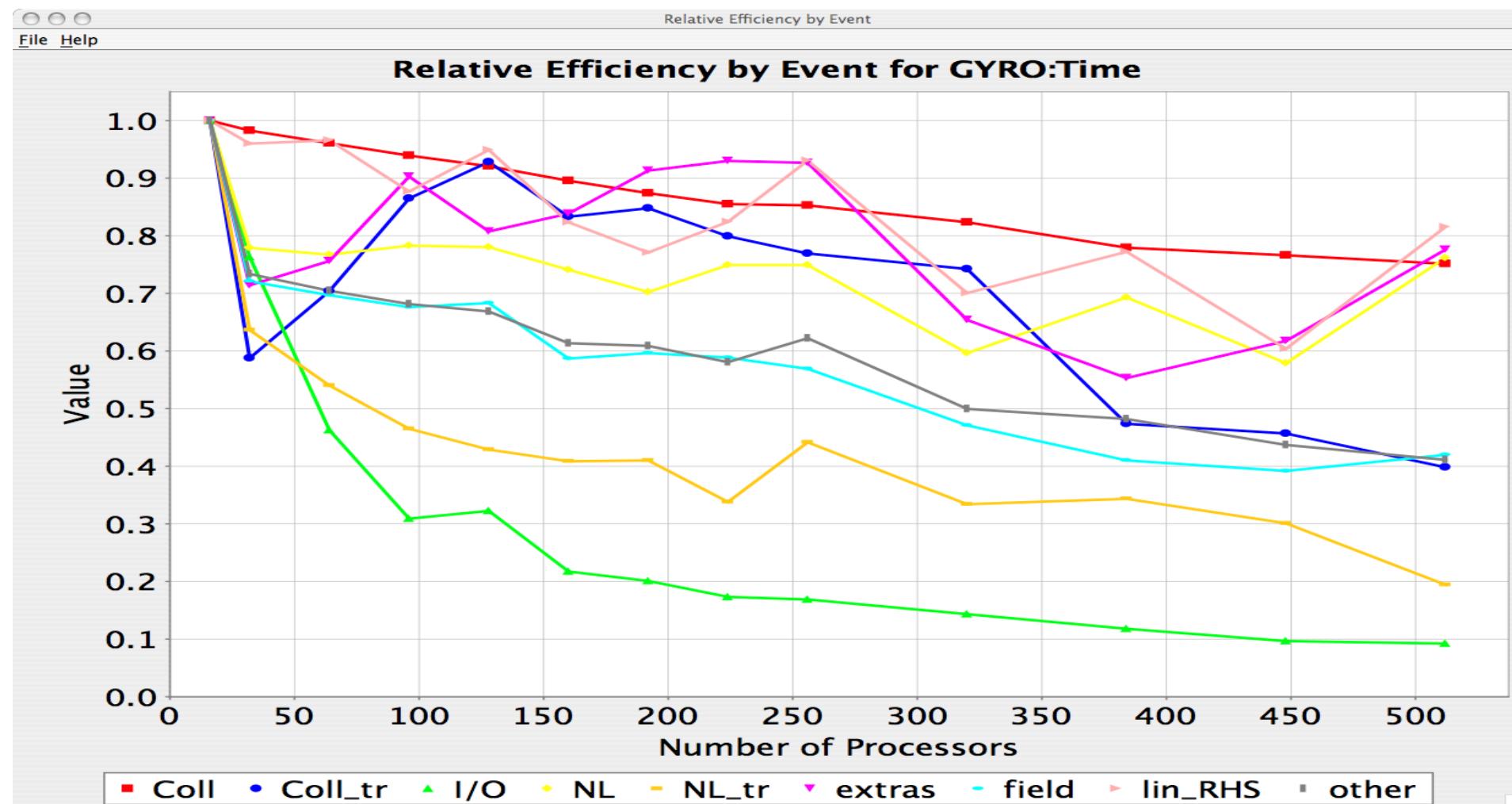
PerfExplorer - Interface



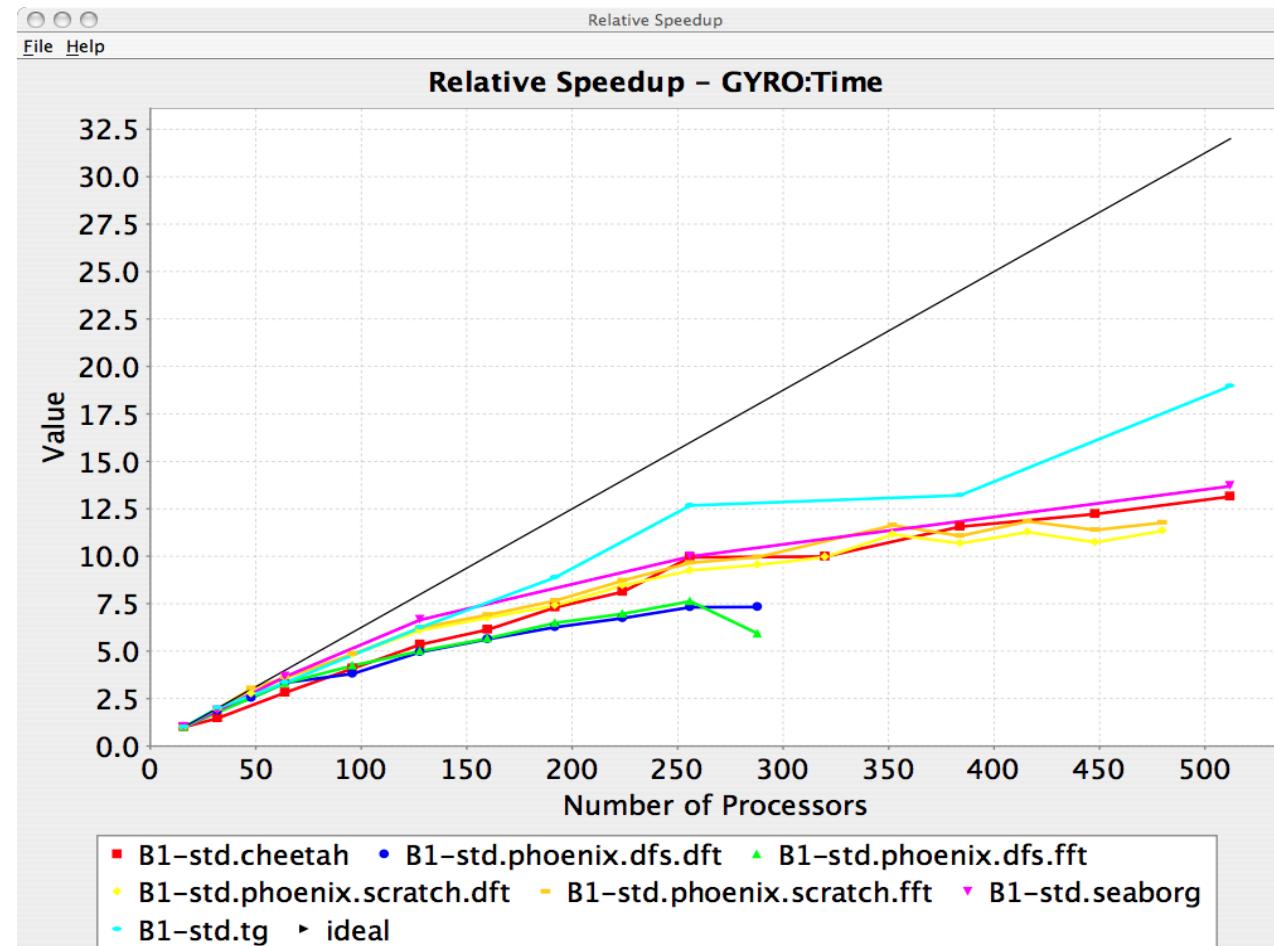
PerfExplorer - Relative Efficiency Plots



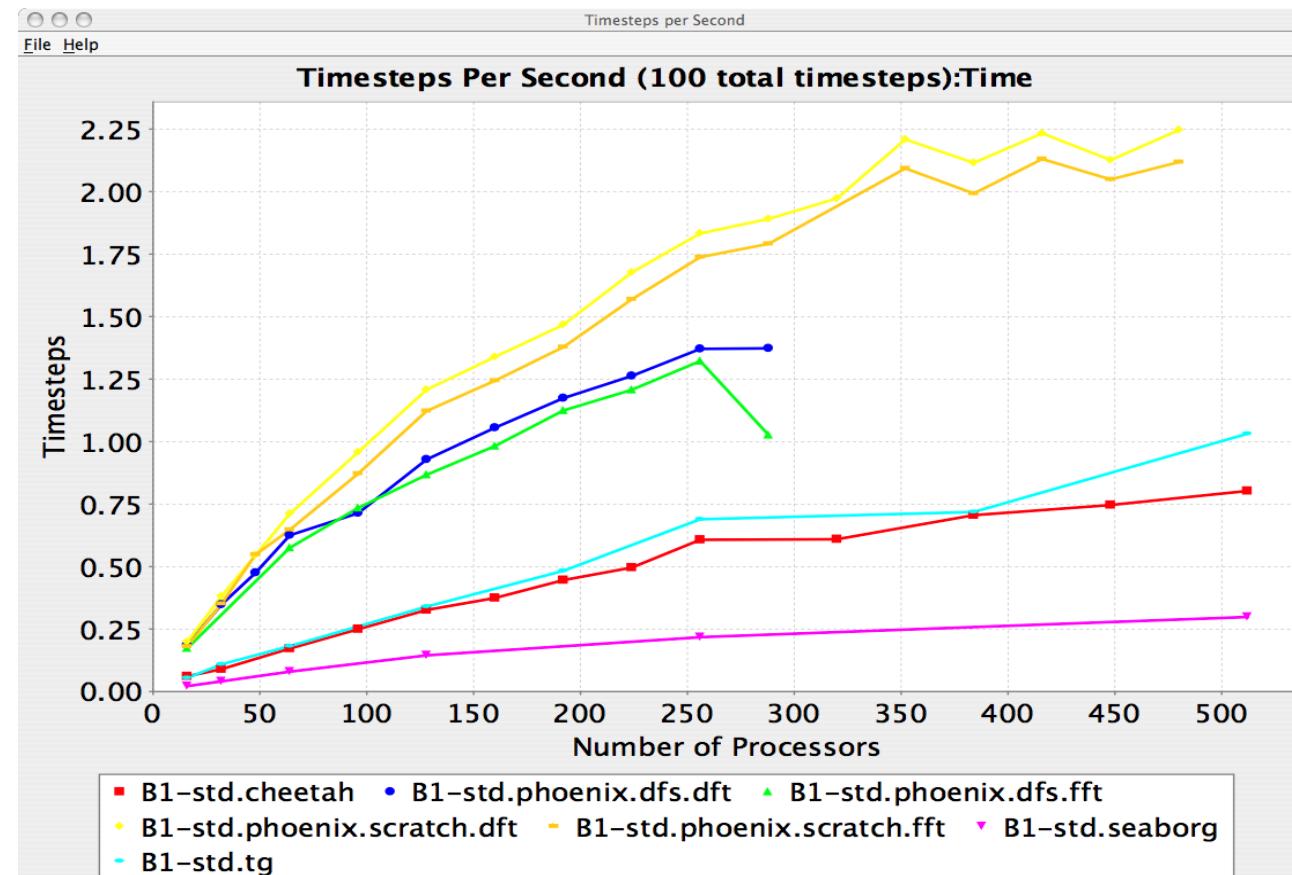
PerfExplorer - Relative Efficiency by Routine



PerfExplorer - Relative Speedup

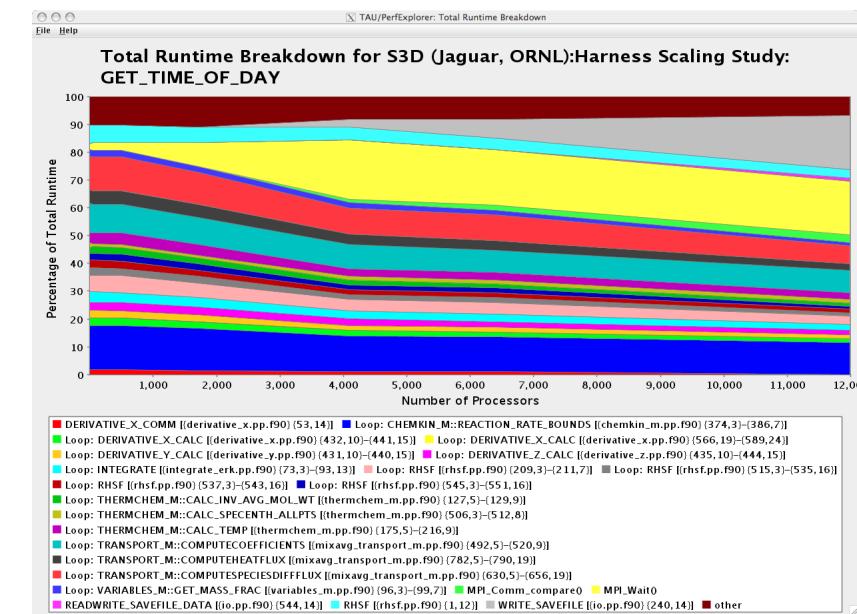
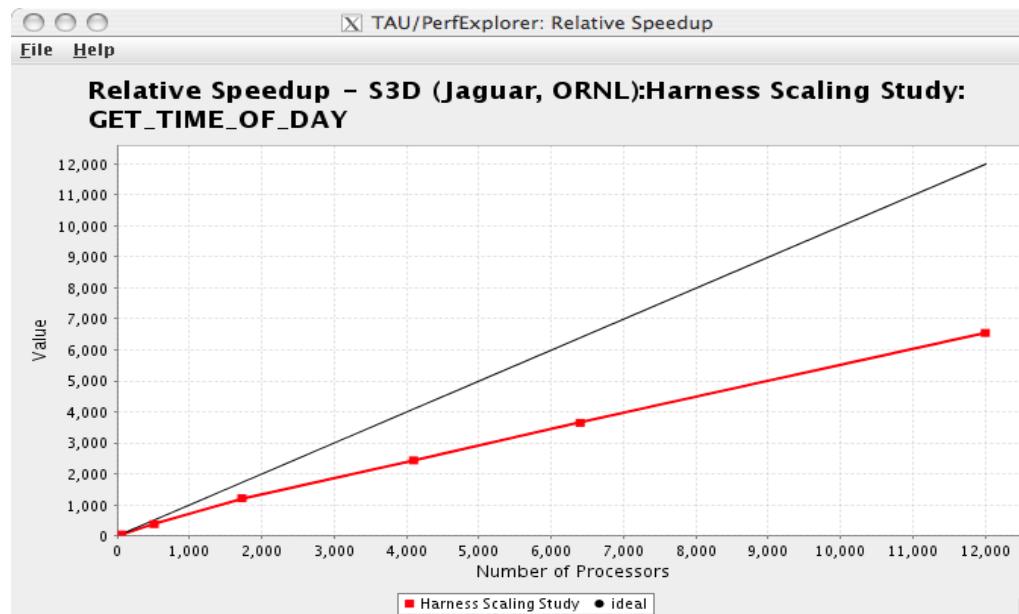


PerfExplorer - Timesteps Per Second

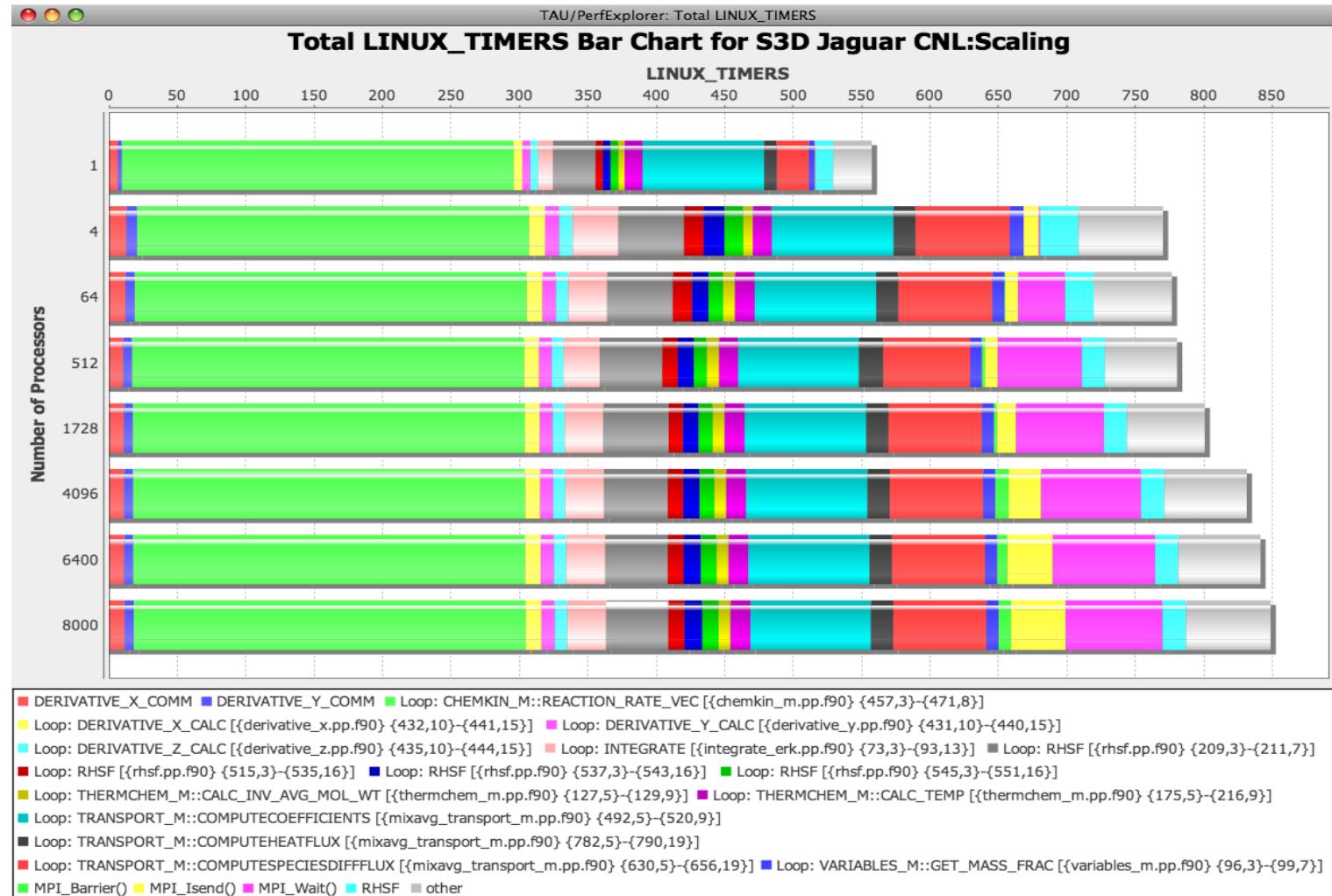


Evaluate Scalability

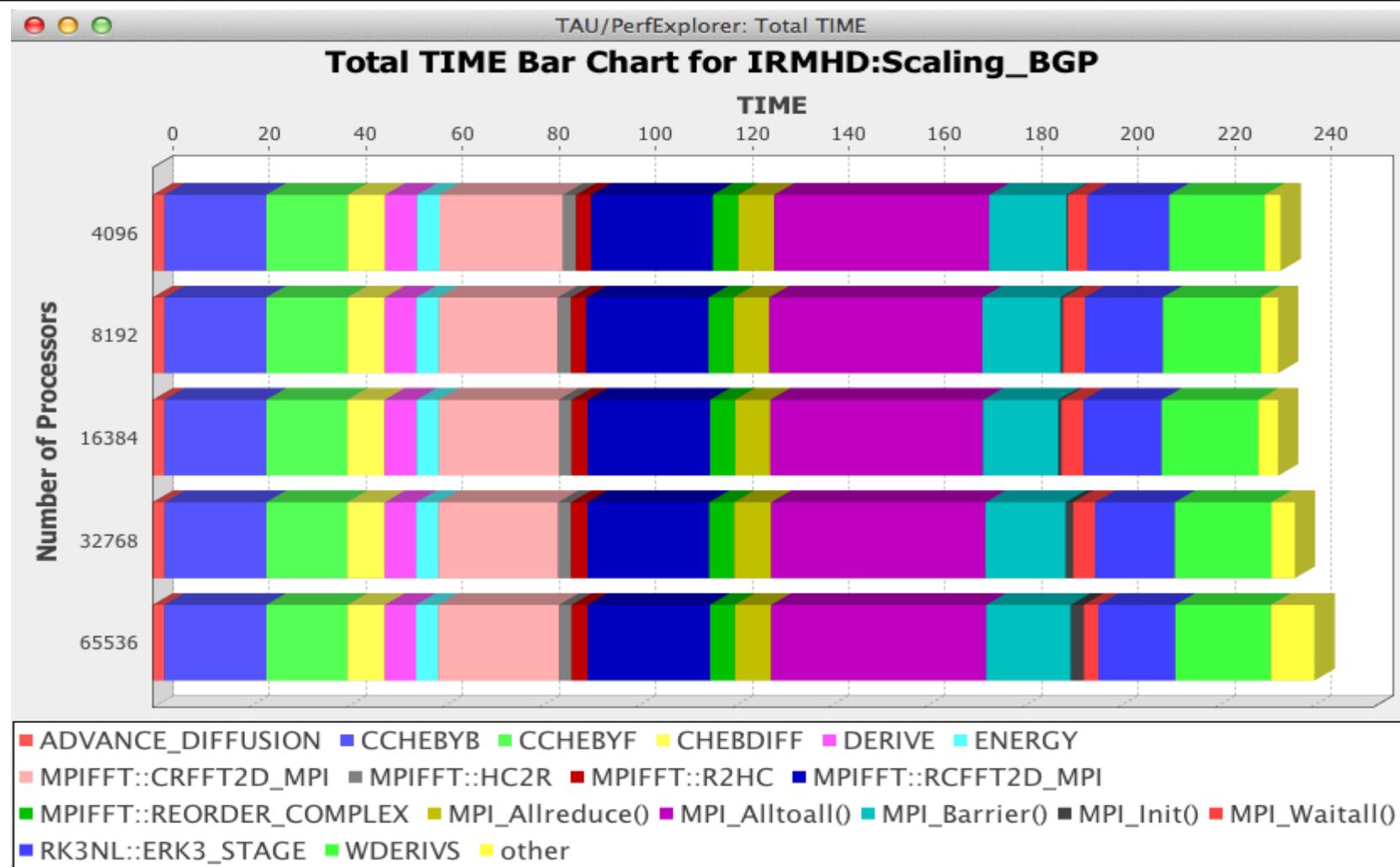
- Goal: How does my application scale? What bottlenecks occur at what core counts?
- Load profiles in taudb database and examine with PerfExplorer



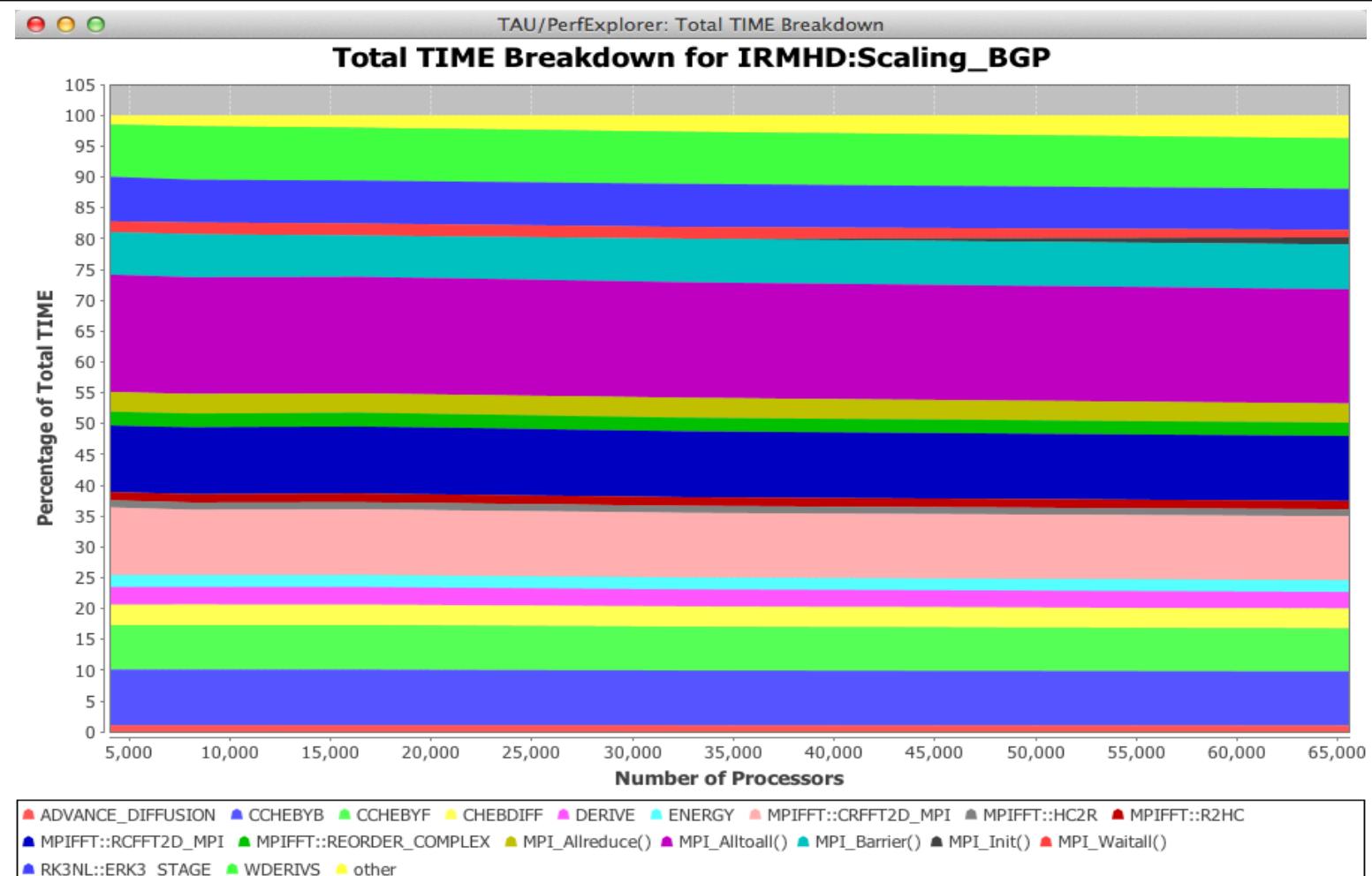
Evaluate Scalability



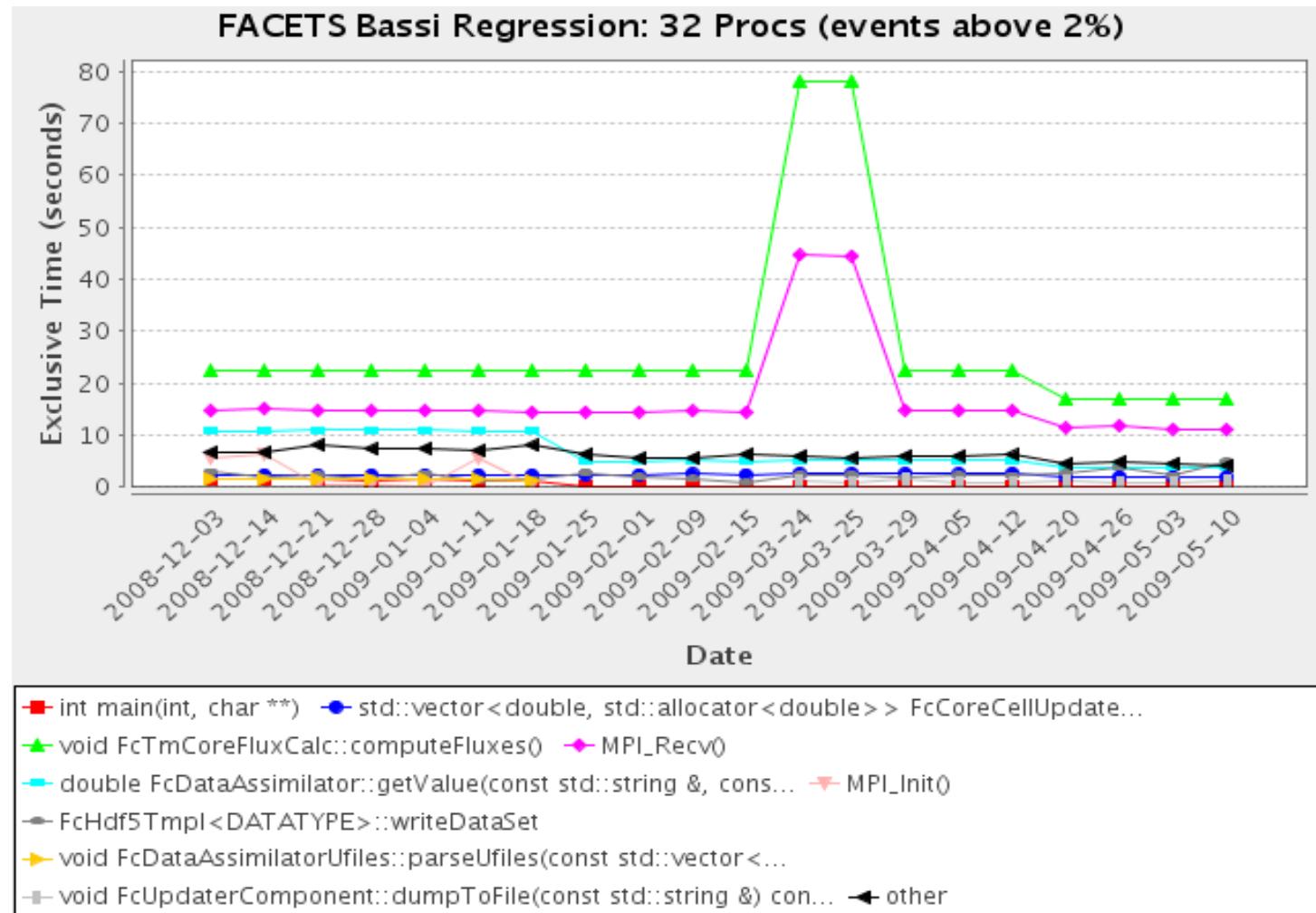
PerfExplorer



PerfExplorer

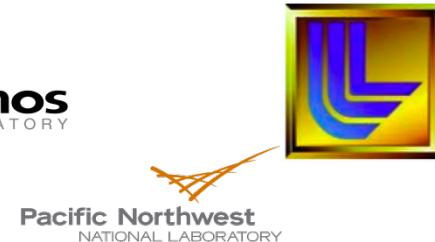


Performance Regression Testing



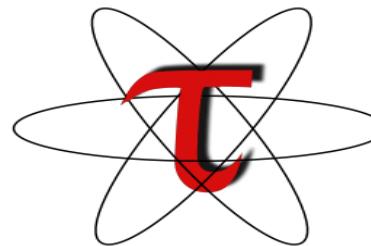
Support Acknowledgments

- U.S. Department of Energy (DOE)
 - Office of Science
 - PNNL, LBL, ORNL
 - ASC/NNSA, Tri-labs (LLNL, LANL, SNL)
- U.S. Department of Defense (DoD)
 - HPC Modernization Office (HPCMO)
- NSF Software Development for Cyberinfrastructure (SDCI)
- Juelich Supercomputing Center, NIC
- Argonne National Laboratory
- T.U. Dresden
- ParaTools, Inc.



ParaTools

Download TAU from U. Oregon



<http://tau.uoregon.edu>

<http://www.hpclinux.com> [LiveDVD, OVA]

Free download, open source, BSD license