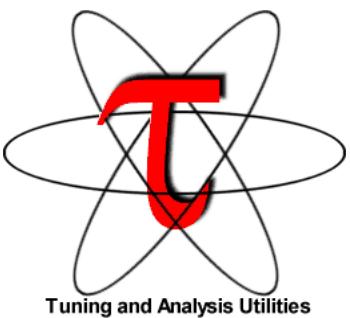


## TAU Performance System®

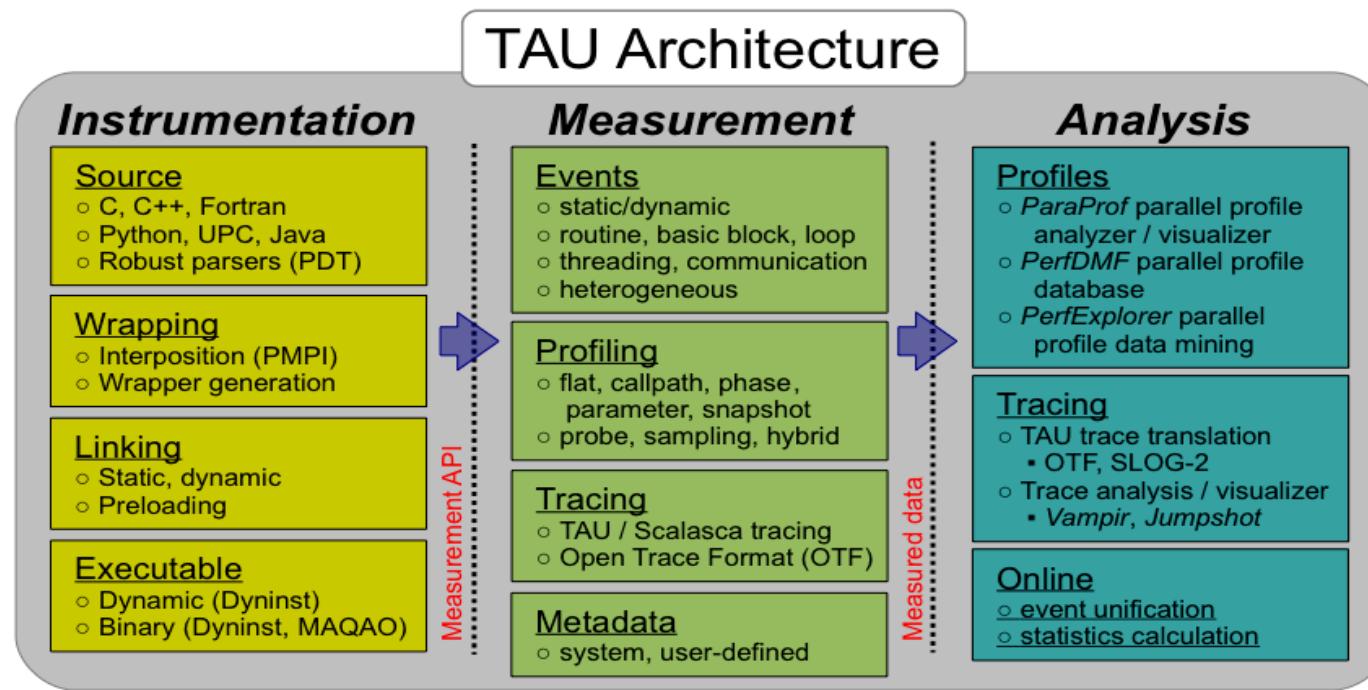
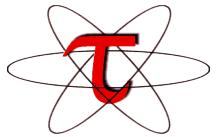


Sameer Shende  
[sameer@cs.uoregon.edu](mailto:sameer@cs.uoregon.edu)  
University of Oregon  
<http://tau.uoregon.edu>



# 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, ROCm, HIP
  - 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
  - Native generation of OTF2 traces (TAU\_TRACE=1, TAU\_TRACE\_FORMAT=otf2)
  - Callsite instrumentation with profiles and traces (TAU\_CALLSITE=1)
  - 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 TR6 (OMPT, tau\_exec –T ompt,tr6)
  - CUDA CUPTI, OpenCL (tau\_exec -T cupti -cupti)

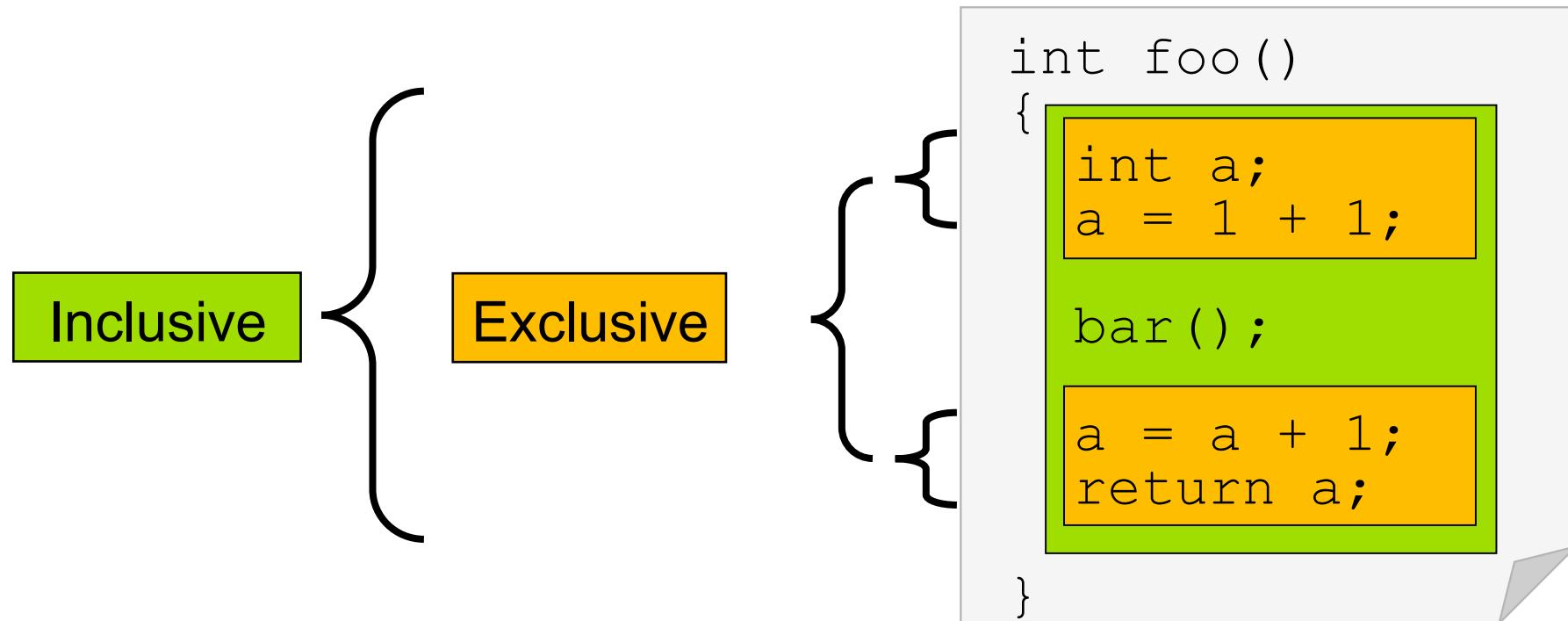
# Application Performance Engineering using TAU

---

- 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?

# Inclusive vs. Exclusive values

- Inclusive
  - Information of all sub-elements aggregated into single value
- Exclusive
  - Information cannot be subdivided further

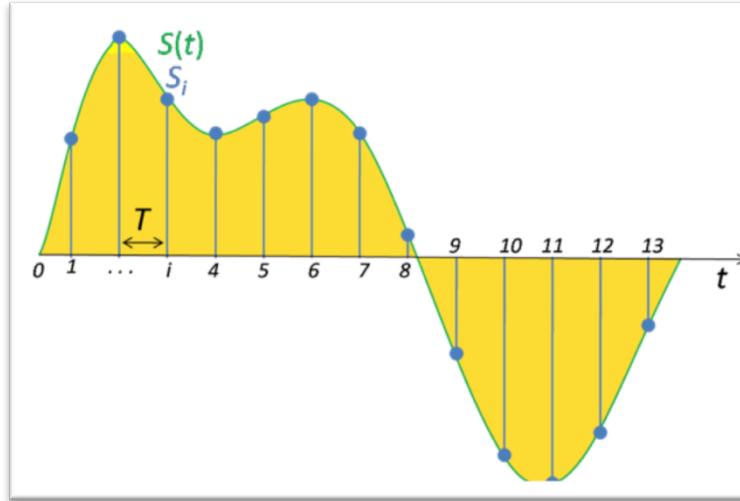


# Performance Data Measurement

## Direct via Probes

```
Call START('potential')
// code
Call STOP('potential')
```

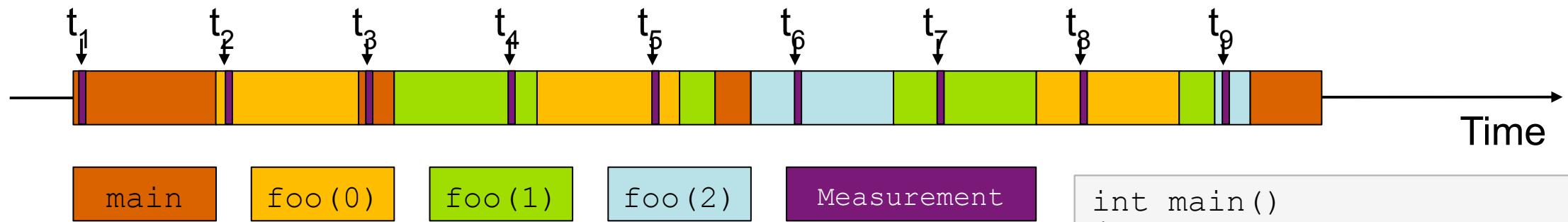
## Indirect via Sampling



- Exact measurement
- Fine-grain control
- Calls inserted into code

- No code modification
- Minimal effort
- Relies on debug symbols (**-g**)

# Event-Based Sampling (EBS)



- Running program is periodically interrupted to take measurement
  - Timer interrupt, OS signal, or HWC overflow
  - Service routine examines return-address stack
  - Addresses are mapped to routines using symbol table information
- Statistical inference of program behavior
  - Not very detailed information on highly volatile metrics
  - Requires long-running applications
- Works with unmodified executables (`tau_exec -ebs`)

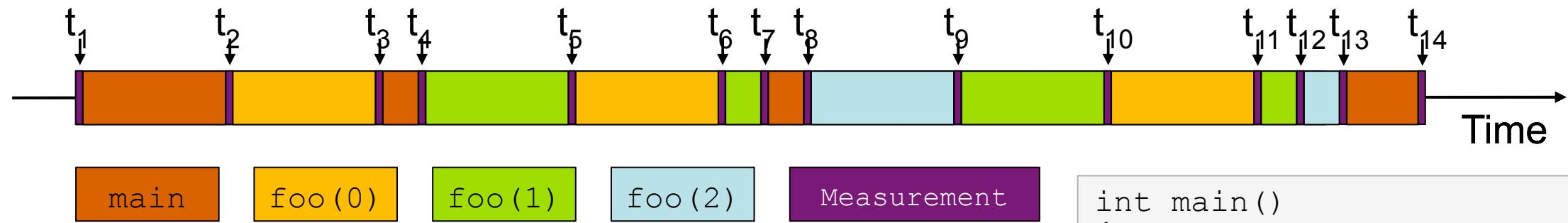
```
int main()
{
    int i;

    for (i=0; i < 3; i++)
        foo(i);

    return 0;
}

void foo(int i)
{
    if (i > 0)
        foo(i - 1);
}
```

# Instrumentation

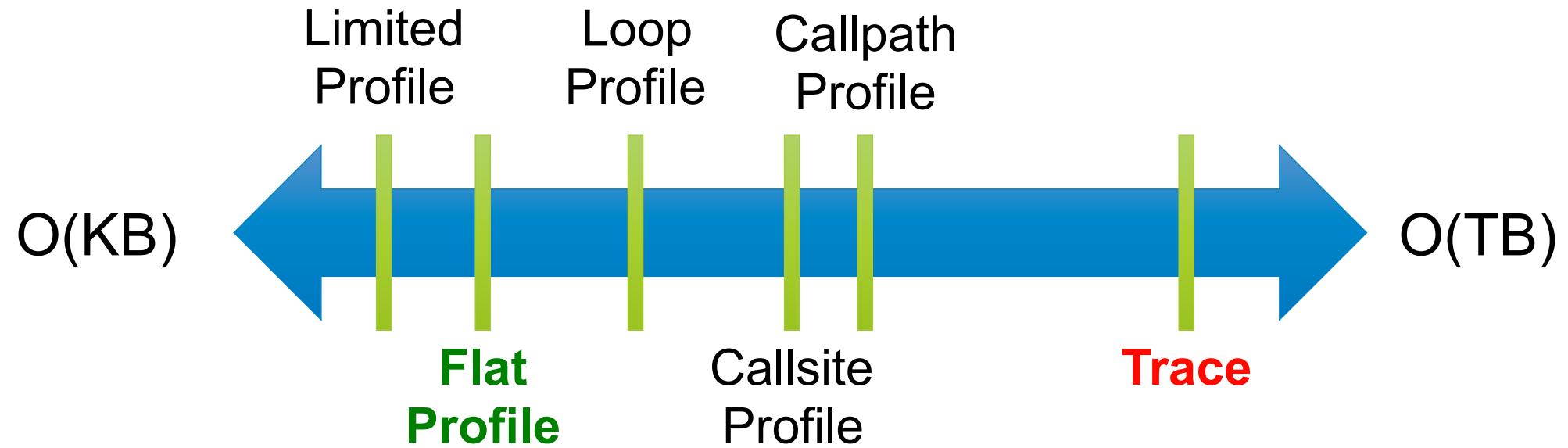


- Measurement code is inserted such that every event of interest is captured directly
  - Can be done in various ways
- Advantage:
  - Much more detailed information
- Disadvantage:
  - Processing of source-code / executable necessary
  - Large relative overheads for small functions

```
int main()
{
    int i;
    TAU_START("main");
    for (i=0; i < 3; i++)
        foo(i);
    TAU_STOP("main");
    return 0;
}

void foo(int i)
{
    TAU_START("foo");
    if (i > 0)
        foo(i - 1);
    TAU_STOP("foo");
}
```

# How much data do you want?



# Types of Performance Profiles

---

- *Flat* profiles

- Metric (e.g., time) spent in an event
  - Exclusive/inclusive, # of calls, child calls, ...

- *Callpath* profiles

- Time spent along a calling path (edges in callgraph)
  - “*main=> f1 => f2 => MPI\_Send*”
  - Set the `TAU_CALLPATH` and `TAU_CALLPATH_DEPTH` environment variables

- *Callsite* profiles

- Time spent along in an event at a given source location
  - Set the `TAU_CALLSITE` environment variable

- *Phase* profiles

- Flat profiles under a phase (nested phases allowed)
  - Default “main” phase
  - Supports static or dynamic (e.g. per-iteration) phases

# Using TAU's Runtime Preloading Tool: tau\_exec

---

- Preload a wrapper that intercepts the runtime system call and substitutes with another
  - MPI
  - OpenMP
  - POSIX I/O
  - Memory allocation/deallocation routines
  - Wrapper library for an external package
- No modification to the binary executable!
- Enable other TAU options (communication matrix, OTF2, event-based sampling)
- Add tau\_exec before the name of the binary
  - srun tau\_exec ./a.out
  - srun tau\_exec -T ompt,tr6,mpi,papi -ompt ./a.out

# 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:

```
mpirun -np 2 tau_exec -T icpc,ompt,mpi -ompt ./a.out  
aprun -n 2 tau_exec -io ./a.out
```

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

```
aprun -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)
```

```
mpirun -np 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:

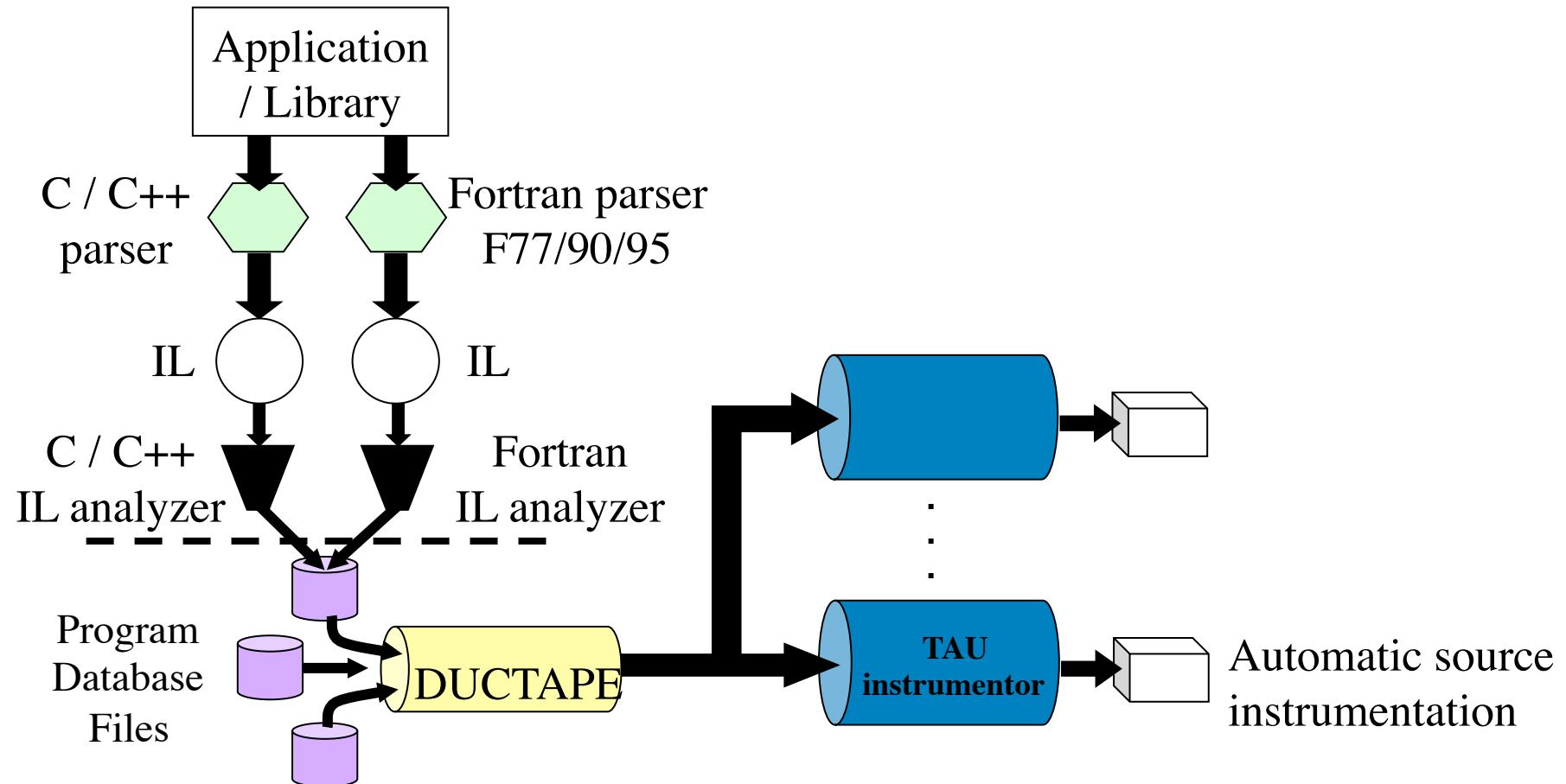
Detects memory leaks, checks for invalid alignment, and checks for array overflow. This is exactly like setting TAU\_TRACK\_MEMORY\_LEAKS=1 and TAU\_MEMDBG\_PROTECT\_ABOVE=1 and running with -memory

- tau\_exec can enable event based sampling while launching the executable using the **-ebs** flag!

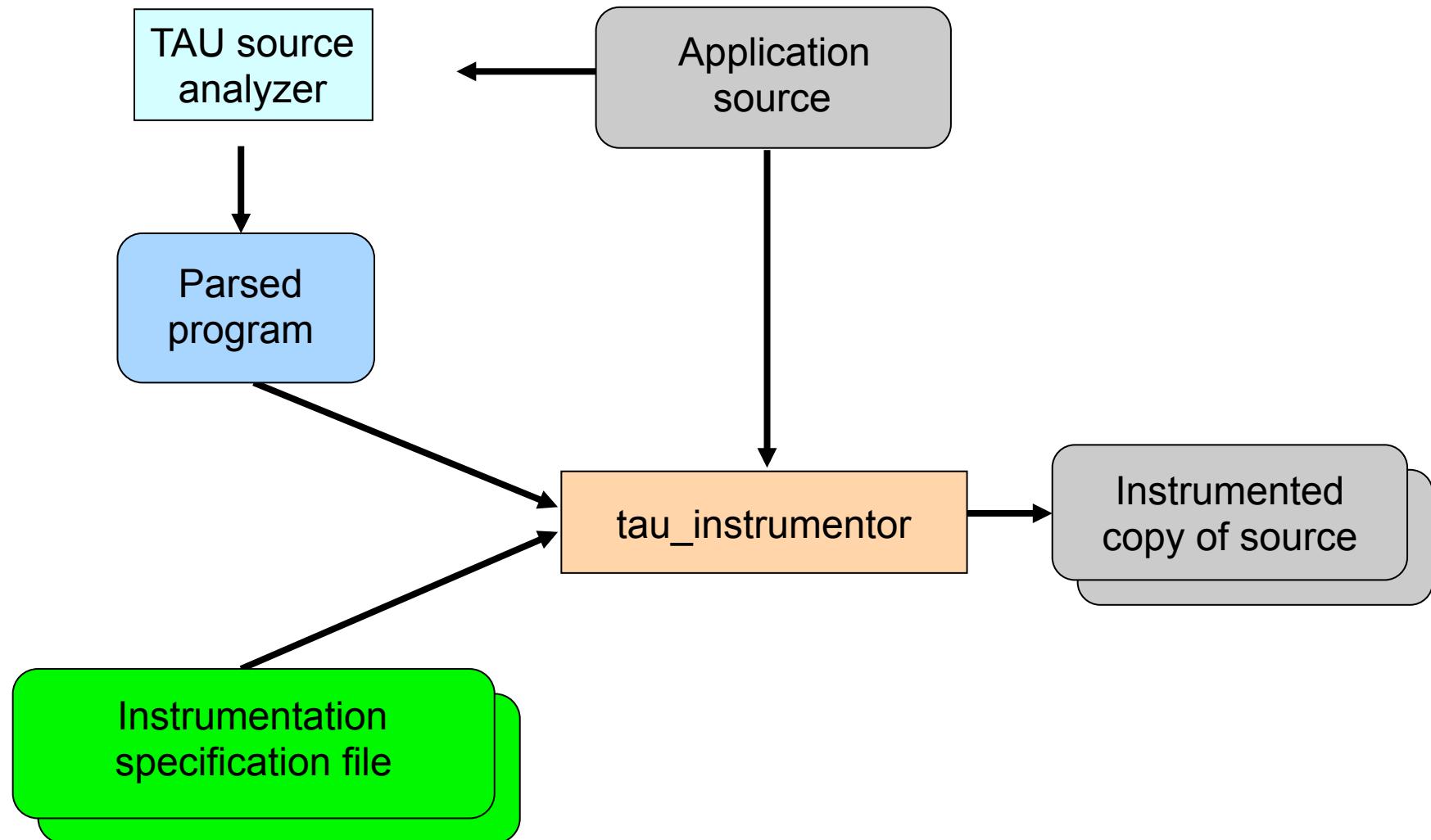
---

# TAU's Source Instrumentation

# TAU's Static Analysis System: Program Database Toolkit (PDT)



# Automatic Source Instrumentation using PDT



# Installing and Configuring TAU

---

- **Installing PDT:**

- wget http://tau.uoregon.edu/pdt.tgz
  - ./configure; make ; make install

- **Installing TAU:**

- wget http://tau.uoregon.edu/tau.tgz
  - ./configure **-ompt=download-tr6 -c++=mpiicpc -cc=mpicc -fortran=mpiifort -mpi -bfd=download -pdt=<dir>**  
-papi=<dir> ...
  - make install; export PATH=<taudir>/x86\_64/bin:\$PATH

- **Using TAU for source instrumentation:**

- export TAU\_MAKEFILE=<taudir>/x86\_64/lib/Makefile.tau-<TAGS>
  - make CC=tau\_cc.sh CXX=tau\_cxx.sh F90=tau\_f90.sh

# 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 PDT

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

```
% source /home/ri-sshende/tau.bashrc  
% export TAU_MAKEFILE=$TAU/Makefile.tau-gnu-mpi-pdt-openmp-opari  
% 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:

```
% mpif90 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)
```

# Makefiles for TAU Compiler and Runtime Options (Isambard)

```
% source /home/ri-sshende/tau.bashrc
% echo $TAU
/home/ri-sshende/pkgs/tau_latest/craycnl/lib
% ls $TAU/Makefile.*
tau_latest/craycnl/lib/Makefile.tau-gnu-mpi-pdt
tau_latest/craycnl/lib/Makefile.tau-gnu-mpi-pdt-openmp-opari
tau_latest/craycnl/lib/Makefile.tau-gnu-mpi-pthread-pdt
tau_latest/craycnl/lib/Makefile.tau-gnu-ompt-tr6-mpi-pdt-openmp
tau_latest/craycnl/lib/Makefile.tau-gnu-papi-ompt-tr6-mpi-pdt-openmp
tau_latest/craycnl/lib/Makefile.tau-allinea-ompt-tr6-mpi-pdt-openmp
```

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

## Makefile.tau-gnu-mpi-pdt-openmp-opari

- Supports MPI instrumentation & PDT for automatic source instrumentation

```
% export TAU_MAKEFILE=$TAU/Makefile.tau-gnu-mpi-pdt-openmp-opari
% tau_f90.sh matmult.f90 -o matmult
% aprun -n 16 ./matmult
% paraprof
```

# Makefiles for TAU Compiler and Runtime Options (Archer)

```
% source /home/y14/y14/shende/tau.bashrc
% echo $TAU
/home/y14/y14/shende/pkgs/tau_latest/craycn1/lib
% ls $TAU/Makefile.*
/home/y14/y14/shende/pkgs/tau_latest/craycn1/lib/Makefile.tau-gnu-mpi
/home/y14/y14/shende/pkgs/tau_latest/craycn1/lib/Makefile.tau-gnu-papi-mpi
/home/y14/y14/shende/pkgs/tau_latest/craycn1/lib/Makefile.tau-gnu-papi-mpi-openmp-opari
/home/y14/y14/shende/pkgs/tau_latest/craycn1/lib/Makefile.tau-gnu-papi-mpi-pdt-openmp-opari
/home/y14/y14/shende/pkgs/tau_latest/craycn1/lib/Makefile.tau-gnu-papi-mpi-pthread
/home/y14/y14/shende/pkgs/tau_latest/craycn1/lib/Makefile.tau-gnu-papi-ompt-tr6-mpi-pdt-openmp
```

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

## Makefile.tau-gnu-mpi-pdt-openmp-opari

- Supports MPI instrumentation & PDT for automatic source instrumentation

```
% export TAU_MAKEFILE=$TAU/Makefile.tau-gnu-papi-mpi-pdt-openmp-opari
% tau_f90.sh matmult.f90 -o matmult
% aprun -n 16 ./matmult
% paraprof
```

# Configuration tags for tau\_exec

```
% ./configure -pdt=<dir> -mpi -papi=<dir>; make install  
Creates in $TAU:  
Makefile.tau-papi-mpi-pdt (Configuration parameters in stub makefile)  
shared-papi-mpi-pdt/libTAU.so
```

```
% ./configure -pdt=<dir> -mpi; make install creates  
Makefile.tau-mpi-pdt  
shared-mpi-pdt/libTAU.so
```

To explicitly choose preloading of shared-<options>/libTAU.so change:

```
% aprun -n 256 ./a.out      to  
% aprun -n 256 tau_exec -T <comma_separated_options> ./a.out
```

```
% aprun -n 256 tau_exec -T papi,mpi,pdt ./a.out
```

Preloads \$TAU/shared-papi-mpi-pdt/libTAU.so

```
% aprun -n 256 tau_exec -T papi ./a.out
```

Preloads \$TAU/shared-papi-mpi-pdt/libTAU.so by matching.

```
% aprun -n 256 tau_exec -T papi,mpi,pdt -s ./a.out
```

Does not execute the program. Just displays the library that it will preload if executed without the **-s** option.

NOTE: -mpi configuration is selected by default. Use -T serial for Sequential programs.

# Simplifying TAU's usage (`tau_exec`)

---

- Uninstrumented execution
  - % aprun -n 16 ./a.out
- Track MPI performance
  - % aprun -n 16 **tau\_exec** ./a.out
- Track OpenMP, and MPI performance (MPI enabled by default)
  - % export TAU\_OMPT\_SUPPORT\_LEVEL=full; export TAU\_OMPT\_RESOLVE\_ADDRESS\_EAGERLY=1
  - % aprun -n 16 **tau\_exec -T mpi,pdt,ompt,papi -ompt** ./a.out
- Track memory operations
  - % export TAU\_TRACK\_MEMORY\_LEAKS=1
  - % mpirun -np 16 **tau\_exec -memory\_debug** ./a.out (bounds check)
- Use event based sampling (compile with `-g`)
  - % mpirun -np 16 **tau\_exec -ebs** ./a.out
  - Also `-ebs_source=<PAPI_COUNTER>` `-ebs_period=<overflow_count>`
- Load wrapper interposition library
  - % mpirun -np 16 **tau\_exec -loadlib=<path/libwrapper.so>** ./a.out
- **Track GPGPU operations (-rocm, -opencl, -cuhti, -cuhti -um, -openacc):**
  - % mpirun -np 16 **tau\_exec -cuhti** ./a.out

---

# TAU hands-on exercises

# TAU tutorial exercise 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 (*Isambard*)

---

- Setup preferred program environment compilers

```
% source /home/ri-sshende/tau.bashrc

# If you have previous performance data from Score-P, you may view it with TAU's paraprof
% paraprof profile.cubex &
```

## Local Installation (*Archer*)

---

- Setup preferred program environment compilers

```
% source /home/y14/y14/shende/tau.bashrc

# If you have previous performance data from Score-P, you may view it with TAU's paraprof
% paraprof profile.cubex &
```

## 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 (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
  - Proposed hands-on setup on Stampede2:
    - 2 compute nodes with 1 Intel Xeon Phi 7250 CPU (Knights Landing, KNL) each
    - 32 MPI processes with 4 OpenMP threads each
  - bt-mz\_C.16 should run in less than 30 seconds
- Use `tau_exec` with **dynamic** executables

# NPB-MZ-MPI / BT: config/make.def : Use default compiler!

```
# SITE- AND/OR PLATFORM-SPECIFIC DEFINITIONS.  
#  
#-----  
# Configured for generic MPI with GCC compiler  
#-----  
#COMPILER = -homp          # Cray/CCE compiler  
COMPILER = -fopenmp      # GNU/GCC compiler  
#COMPILER = -qopenmp       # Intel compiler  
  
...  
#-----  
# The Fortran compiler used for MPI programs  
#-----  
# ftn = ftn # OpenMPI with Intel compiler  
ftn = ftn  
# Alternative variant to perform instrumentation  
# ftn = tau_f90.sh  
# PREP is a generic preposition macro for instrumentation preparation  
#ftn = $(PREP) ftn -f77=ifort  
#ftn = scorep ...  
  
FLINKFLAGS = $(FFLAGS) -dynamic  
...
```

Default (no instrumentation)

Uncomment ftn and comment Score-P's compiler wrapper to generate uninstrumented binary!  
Add -dynamic flag to link line.

# 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=C NPROCS=16 *  
*****
```

- Type “make” for instructions

# Building an NPB-MZ-MPI Benchmark

```
% make bt-mz CLASS=C NPROCS=16
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 16 C
make[2]: Entering directory `../BT-MZ'
ftn -c -O3 -fopenmp          bt.f
[...]
Ftn -c -O3 -fopenmp          mpi_setup.f
cd ..;/common; ftn -c -O3 -fopenmp      print_results.f
cd ..;/common; ftn -c -O3 -fopenmp      timers.f
ftn -O3 -fopenmp -o -dynamic ..;/bin/bt-mz_C.16 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_C.16
make[1]: Leaving directory `BT-MZ'
```

- Specify the benchmark configuration
  - benchmark name: **bt-mz**, lu-mz, sp-mz
  - the number of MPI processes: **NPROCS=16**
  - the benchmark class (S, W, A, B, C, D, E): **CLASS=C**

Shortcut: % **make suite**

# NPB-MZ-MPI / BT reference execution on Isambard

```
% cd bin  
% cp ..../jobscript/isambard/run.pbs .  
% cat run.pbs  
% qsub run.pbs  
% less mzmpibt.o<job_id>  
  
NAS Parallel Benchmarks (NPB3.3-MZ-MPI) - BT-MZ MPI+OpenMP Benchmark  
Number of zones: 16 x 16  
Iterations: 200 dt: 0.000100  
Number of active processes: 32  
Total number of threads: 128 ( 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 = 22.34
```

- 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 reference execution on Archer

```
% cd bin  
% cp ..../jobscript/archer/run.pbs .  
% cat run.pbs  
% qsub run.pbs  
% less mzmpibt.o<job_id>  
  
NAS Parallel Benchmarks (NPB3.3-MZ-MPI) - BT-MZ MPI+OpenMP Benchmark  
Number of zones: 16 x 16  
Iterations: 200 dt: 0.000100  
Number of active processes: 32  
Total number of threads: 128 ( 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 = 22.34
```

- 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 with TAU's Event-Based Sampling on Isambard (tau\_exec -T openmp-gnu -ebs)

```
% cd bin  
% cp /home/ri-sshende/isambard/tau.pbs .  
% qsub tau.pbs  
% cat mzmpibt.o<job_id>  
NAS Parallel Benchmarks (NPB3.3-MZ-MPI) - BT-MZ MPI+OpenMP Benchmark  
Number of zones: 16 x 16  
Iterations: 200 dt: 0.000100  
Number of active processes: 32  
Total number of threads: 128 ( 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 = 23  
% pprof -a -m | more  
% paraprof &  
% paraprof --pack bt_ebs.ppk  
<Copy file over to desktop using scp>  
% paraprof bt_ebs &
```

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

tau\_exec launches un-instrumented binaries!

# NPB-MZ-MPI / BT with TAU's Event-Based Sampling on Archer (tau\_exec -T openmp-gnu -ebs)

```
% cd bin  
% cp /home/ri-sshende/isambard/tau.pbs .  
% qsub tau.pbs  
% cat mzmpibt.o<job_id>  
NAS Parallel Benchmarks (NPB3.3-MZ-MPI) - BT-MZ MPI+OpenMP Benchmark  
Number of zones: 16 x 16  
Iterations: 200 dt: 0.000100  
Number of active processes: 32  
Total number of threads: 128 ( 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 = 23  
% pprof -a -m | more  
% paraprof &  
% paraprof --pack bt_ebs.ppk  
<Copy file over to desktop using scp>  
% paraprof bt_ebs &
```

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

tau\_exec launches un-instrumented binaries!

# Source Instrumentation: Compiler and MPI modules on Isambard

- Select appropriate compiler / MPI combination (if not done automatically)

```
source /home/ri-sshende/tau.bashrc
```

- Copy tutorial sources to your scratch directory

```
% tar zxvf /projects/bristol/vi-hps/tutorial/NPB3.3-MZ-MPI.tar.gz
% cd NPB3.3-MZ-MPI; make clean;
echo $TAU_MAKEFILE
/home/ri-sshende/pkgs/tau_latest/craycnl/lib/Makefile.tau-gnu-mpi-pdt-openmp-opari
% make suite MPIF77=tau_f77.sh
% cd bin.tau
% cp ../jobscript/archer/run.pbs .
% qsub run.pbs
% pprof -a -m | more
% paraprof --pack bt.ppk
<SCP to your laptop/desktop>
% paraprof bt.ppk &
```

# Compiler and MPI modules on Archer

- Select appropriate compiler / MPI combination (if not done automatically)

```
source /home/y14/y14/shende/tau.bashrc
```

- Copy tutorial sources to your scratch directory

```
% tar zxvf /work/y14/shared/tutorial/NPB3.3-MZ-MPI.tar.gz
% cd NPB3.3-MZ-MPI; make clean
% echo $TAU_MAKEFILE
/home/y14/y14/shende/pkgs/tau_latest/craycnl/lib/Makefile.tau-gnu-papi-mpi-pdt-openmp-opari

% make suite MPIF77=tau_f77.sh
% cd bin.tau
% cp ..../jobscript/archer/run.pbs .
% qsub run.pbs
% pprof -a -m | more
% paraprof --pack bt.ppk
<SCP to your laptop/desktop>
% paraprof bt.ppk &
```

# Compile-Time Options for TAU's compiler scripts (e.g., tau\_f90.sh)

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 <i>–iowrapper</i> )
-optTrackGOMP	Enable tracking GNU OpenMP runtime layer (used without <i>–opari</i> )
-optMemDbg	Enable runtime bounds checking (see TAU_MEMDBG_* env vars)
-optKeepFiles	Does not remove intermediate .pdb and .inst.* files
-optPreProcess	Preprocess sources (OpenMP, Fortran) before instrumentation
-optTauSelectFile=<file>	Specify selective instrumentation file for <i>tau_instrumentor</i>
-optTauWrapFile=<file>	Specify path to <i>link_options.tau</i> generated by <i>tau_gen_wrapper</i>
-optHeaderInst	Enable Instrumentation of headers
-optTrackUPCR	Track UPC runtime layer routines (used with <i>tau_upc.sh</i> )
-optLinking=""	Options passed to the linker. Typically \$(TAU_MPI_FLIBS) \$(TAU_LIBS) \$(TAU_CXXLIBS)
-optCompile=""	Options passed to the compiler. Typically \$(TAU_MPI_INCLUDE) \$(TAU_INCLUDE) \$(TAU_DEFS)
-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

-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

See tau\_compiler.sh for a full list of TAU\_OPTIONS.

...

# Compiling Fortran Codes with TAU

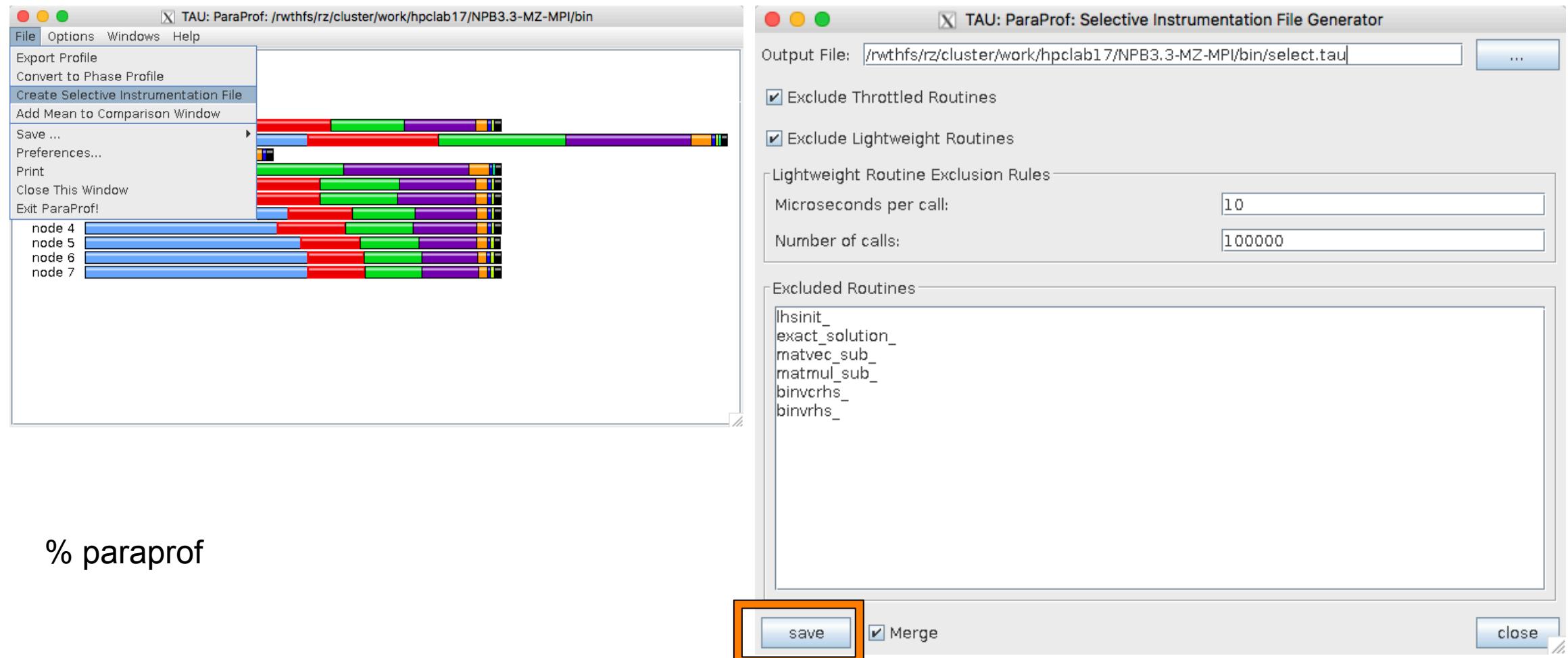
---

- If your Fortran code uses free format in .f files (fixed is default for .f), you may use:  
`% export TAU_OPTIONS=‘-optPdtF95Opts=“-R free” -optVerbose’`
- To use the compiler based instrumentation instead of PDT (source-based):  
`% export TAU_OPTIONS=‘-optComInst -optVerbose’`
- If your Fortran code uses C preprocessor directives (#include, #ifdef, #endif):  
`% export TAU_OPTIONS=‘-optPreProcess -optVerbose -optDetectMemoryLeaks’`
- To use an instrumentation specification file:  
`% export TAU_OPTIONS=‘-optTauSelectFile=select.tau -optVerbose -optPreProcess’`  
`% cat select.tau`  
BEGIN\_INSTRUMENT\_SECTION  
loops routine="#"  
# this statement instruments all outer loops in all routines. # is wildcard as well as comment in first column.  
END\_INSTRUMENT\_SECTION

# TAU's Selective Instrumentation (Filter) File Format

```
% export TAU_OPTIONS='-optTauSelectFile=/path/to/select.tau ...'  
% cat select.tau  
BEGIN_INCLUDE_LIST  
int main#  
int dgemm#  
END_INCLUDE_LIST  
BEGIN_FILE_INCLUDE_LIST  
Main.c  
Blas/*.f77  
END_FILE_INCLUDE_LIST  
# replace INCLUDE with EXCLUDE list for excluding routines/files  
  
BEGIN_INSTRUMENT_SECTION  
loops routine="foo"  
loops routine="int main#"  
END_INSTRUMENT_SECTION  
% export TAU_SELECT_FILE=select.tau      (to use at runtime)
```

# Create a Selective Instrumentation File, Re-instrument, Re-run



% paraprof

## Re-instrument with TAU

---

- Select appropriate compiler / MPI combination (if not done automatically)

```
source /home/ri-sshende/tau.bashrc (Isambard)
source /home/y14/y14/shende/tau.bashrc (Archer)
```

- Choose path to selective instrumentation file saved by paraprof (copy if necessary):

```
% export TAU_OPTIONS=' -optTauSelectFile=/path/to/select.tau -optVerbose'
% echo $TAU_MAKEFILE
.../pkgs/tau_latest/craycnl/lib/Makefile.tau-gnu-papi-mpi-pdt-openmp-opari
% make clean
% make suite MPIF77=tau_f77.sh
% cd bin.tau
% cp ../jobscript/archer/run.pbs .
% qsub run.pbs
% pprof -a -m | more
% paraprof --pack bt.ppk
<SCP to your laptop/desktop>
% paraprof bt.ppk &
```

# TAU's 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_CALLSITE	0	Setting to 1 enables callsite profiling that shows where an instrumented function was called. Also compatible with tracing.
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

---

Environment Variable	Default	Description
TAU_TRACE	0	Setting to 1 turns on tracing
TAU_TRACE_FORMAT	Default	Setting to “otf2” turns on TAU’s native OTF2 trace generation (configure with –otf=download)
TAU_EBS_UNWIND	0	Setting to 1 turns on unwinding the callstack during sampling (use with tau_exec –ebs or TAU_SAMPLING=1)
TAU_EBS_RESOLUTION	line	Setting to “function” or “file” changes the sampling resolution to function or file level respectively.
TAU_TRACK_LOAD	0	Setting to 1 tracks system load on the node
TAU_SELECT_FILE	Default	Setting to a file name, enables selective instrumentation based on exclude/include lists specified in the file.
TAU_OMPT_SUPPORT_LEVEL	basic	Setting to “full” improves resolution of OMPT TR6 regions on threads 1.. N-1. Also, “lowoverhead” option is available.
TAU_OMPT_RESOLVE_ADDRESS_EAGERLY	0	Setting to 1 is necessary for event based sampling to resolve addresses with OMPT TR6 (-ompt=download-tr6)

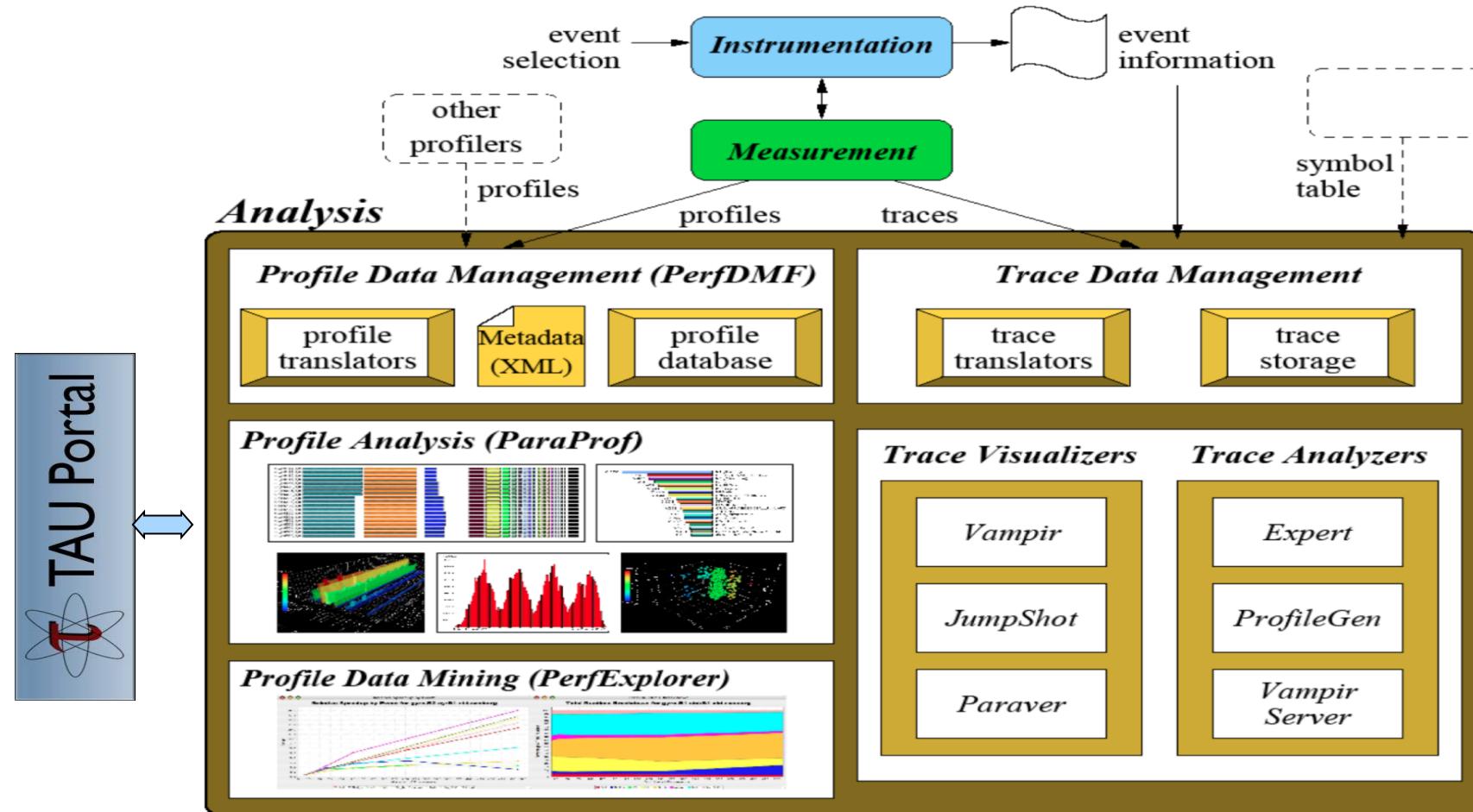
# Runtime Environment Variables

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

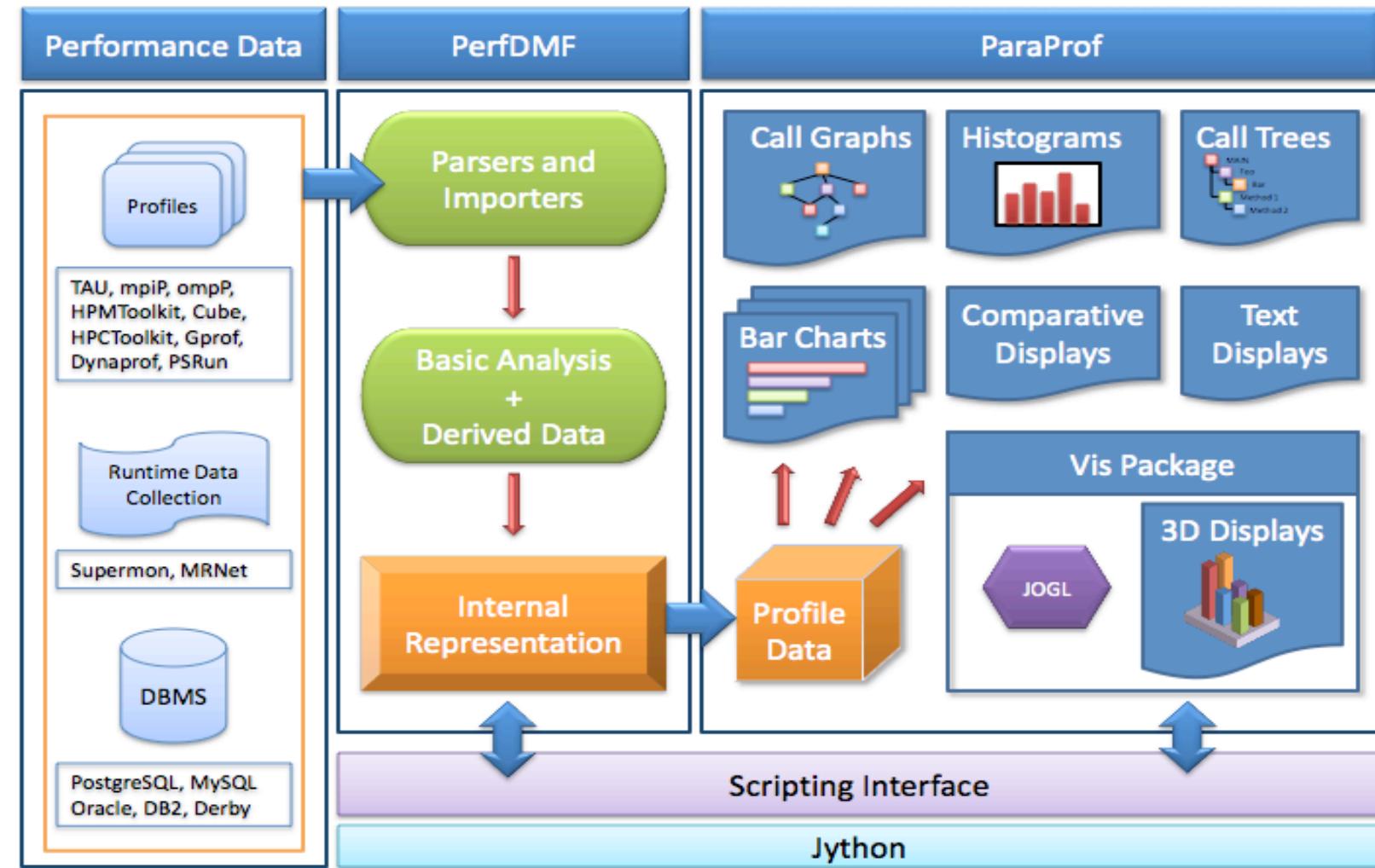
---

# TAU's Analysis Tools: ParaProf

# TAU Analysis



# ParaProf Profile Analysis Framework



# 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 of applications under 'Standard Applications'. Under 'Default App' > 'Default Exp', there is a file named 'bt\_ompt.ppk' which is highlighted with a yellow dot. A blue arrow points from the word 'Metric' in the text above to this file. To the right of the tree view is a large table titled 'TAU: ParaProf Manager' containing various experimental parameters.

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 Manager Widow: scout.cubex

TAU: ParaProf Manager

File Options Help

Applications

- Standard Applications
- Default App
- Default Exp
- scout.cubex
  - Time
  - Wait at Barrier
  - Barrier Completion
  - Late Sender
  - Late Sender => Messages in Wrong Order
  - Late Sender => Messages in Wrong Order => Messages from different sources
  - Late Sender => Messages in Wrong Order => Messages from same source
  - Late Receiver
  - Early Reduce
  - Early Scan
  - Late Broadcast
  - Wait at N x N
  - N x N Completion
  - Management
  - Management => Fork
  - P2P send synchronizations
  - P2P send synchronizations => Late Receivers
  - P2P recv synchronizations
  - P2P recv synchronizations => Late Senders
  - P2P recv synchronizations => Late Senders => Messages in Wrong Order
  - Collective synchronizations
  - P2P send communications
  - P2P send communications => Late Receivers
  - P2P recv communications
  - P2P recv communications => Late Senders
  - P2P recv communications => Late Senders => Messages in Wrong Order
  - Collective exchange communications
  - Collective communications as source
  - Collective communications as destination
  - P2P bytes sent
  - P2P bytes received
  - Collective bytes outgoing
  - Collective bytes incoming
  - RMA bytes received
  - RMA bytes put

Metrics in the profile

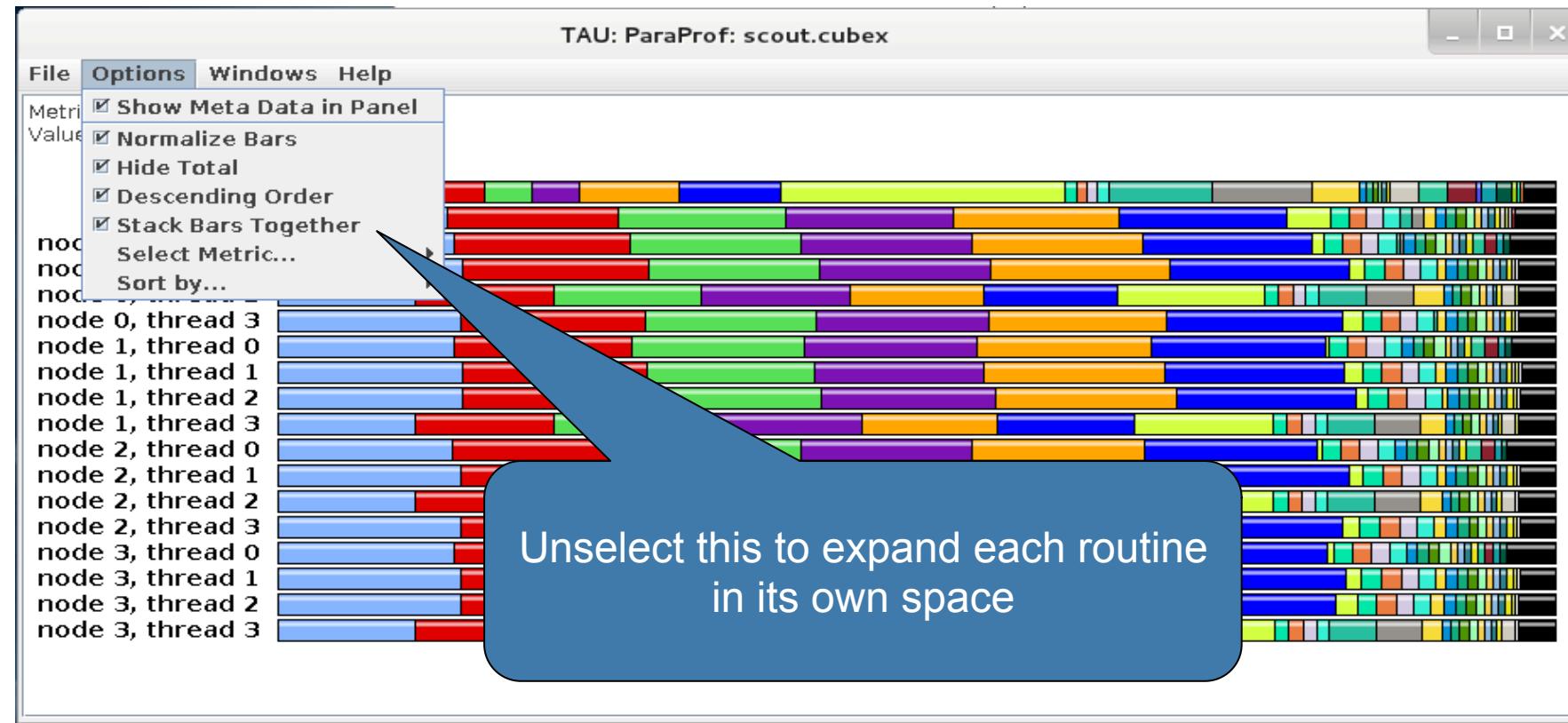
# Paraprof main window



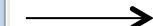
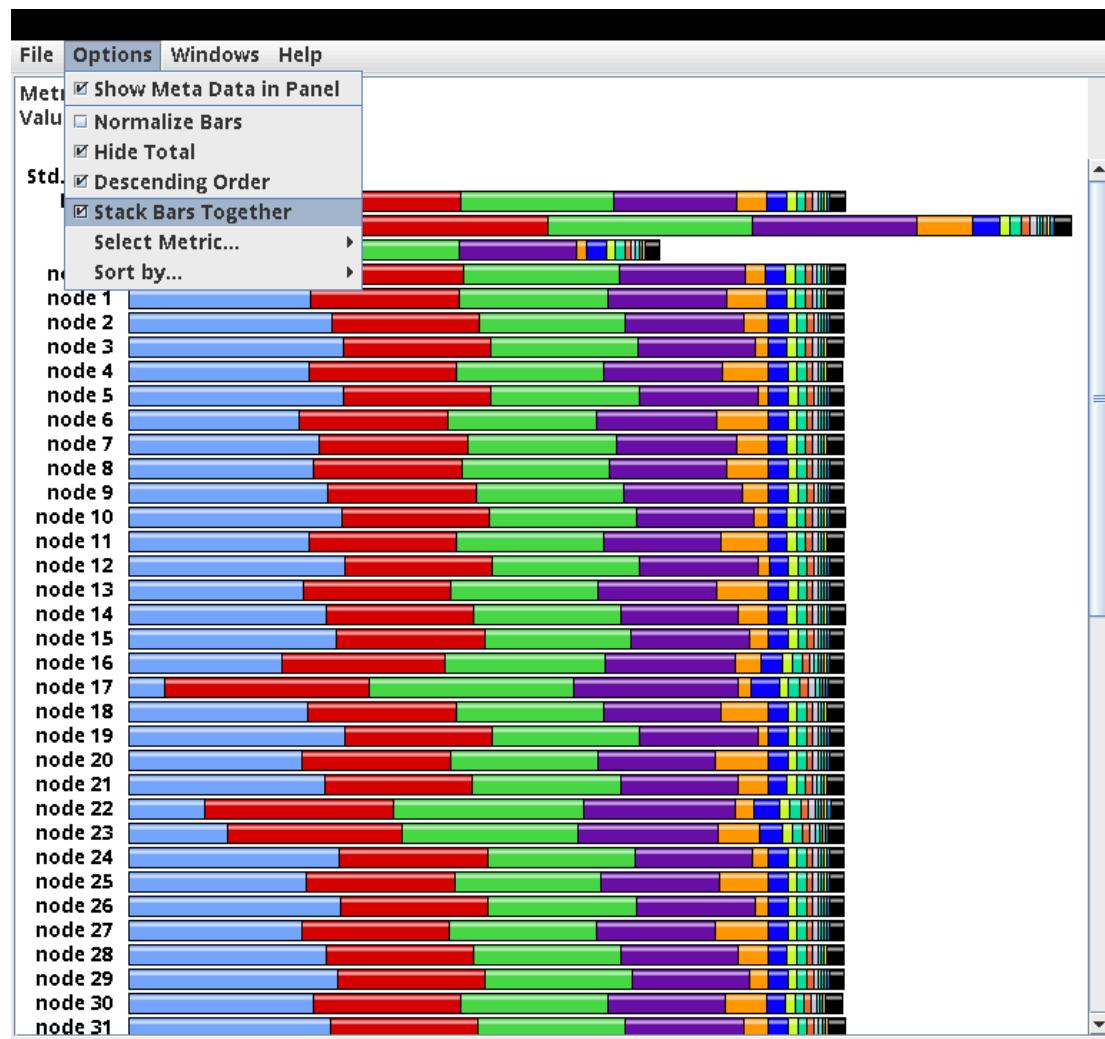
Colors represent code regions

Options -> uncheck Stack Bars Together

# Paraprof main window



# ParaProf Profile Browser



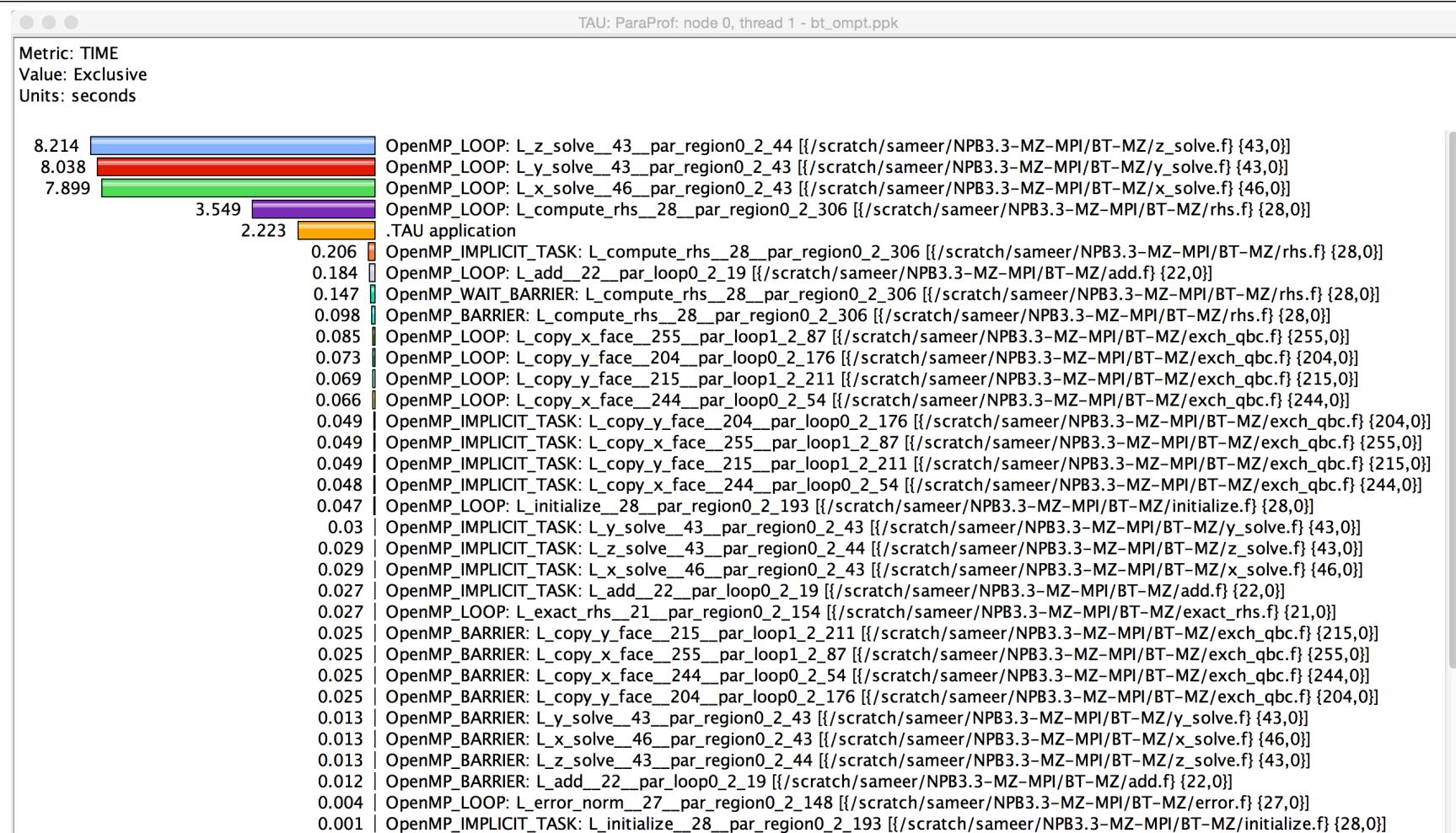
# Paraprof main window



Left/right  
click here

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

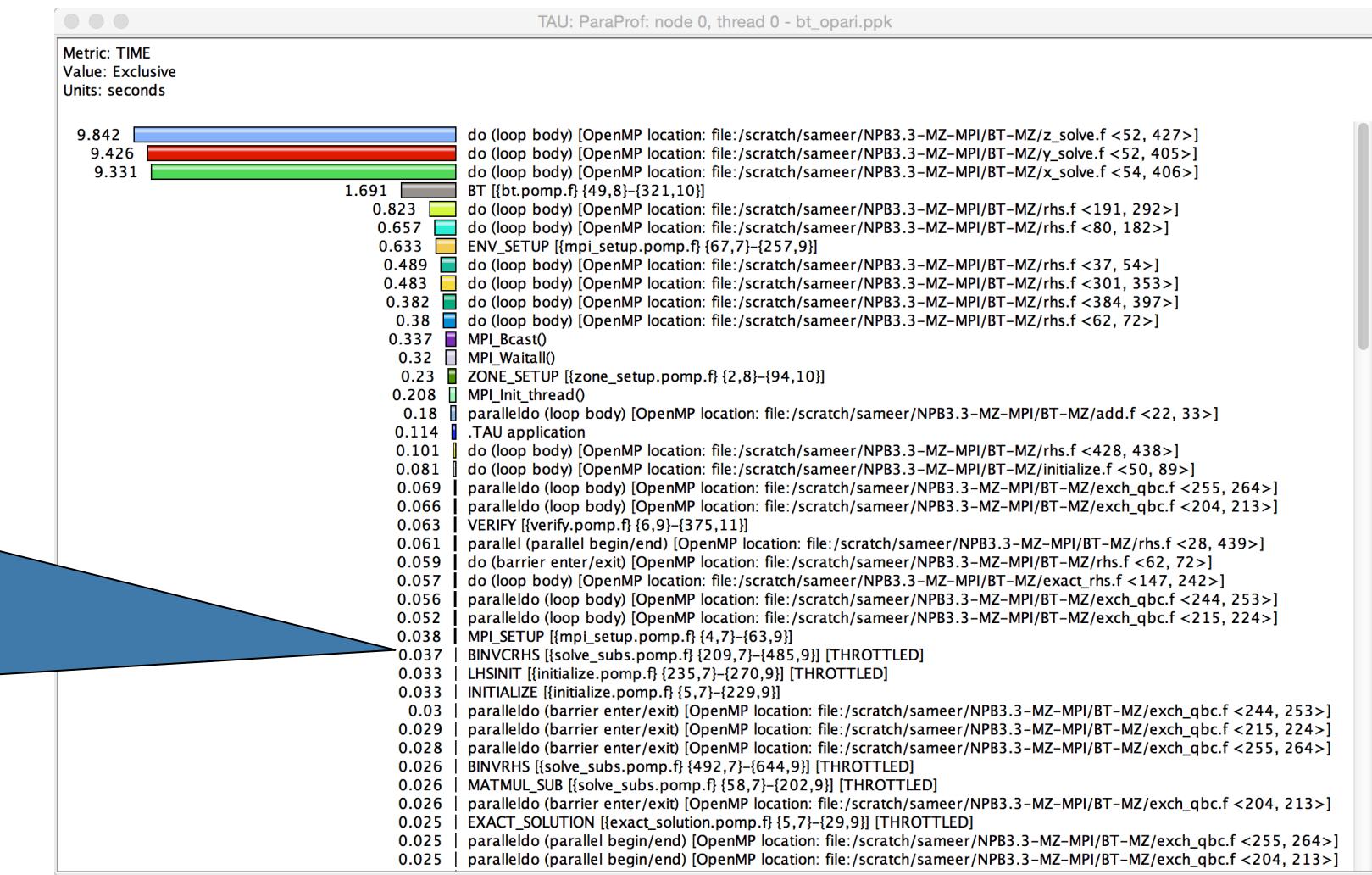
# Paraprof node window (function barchart window)



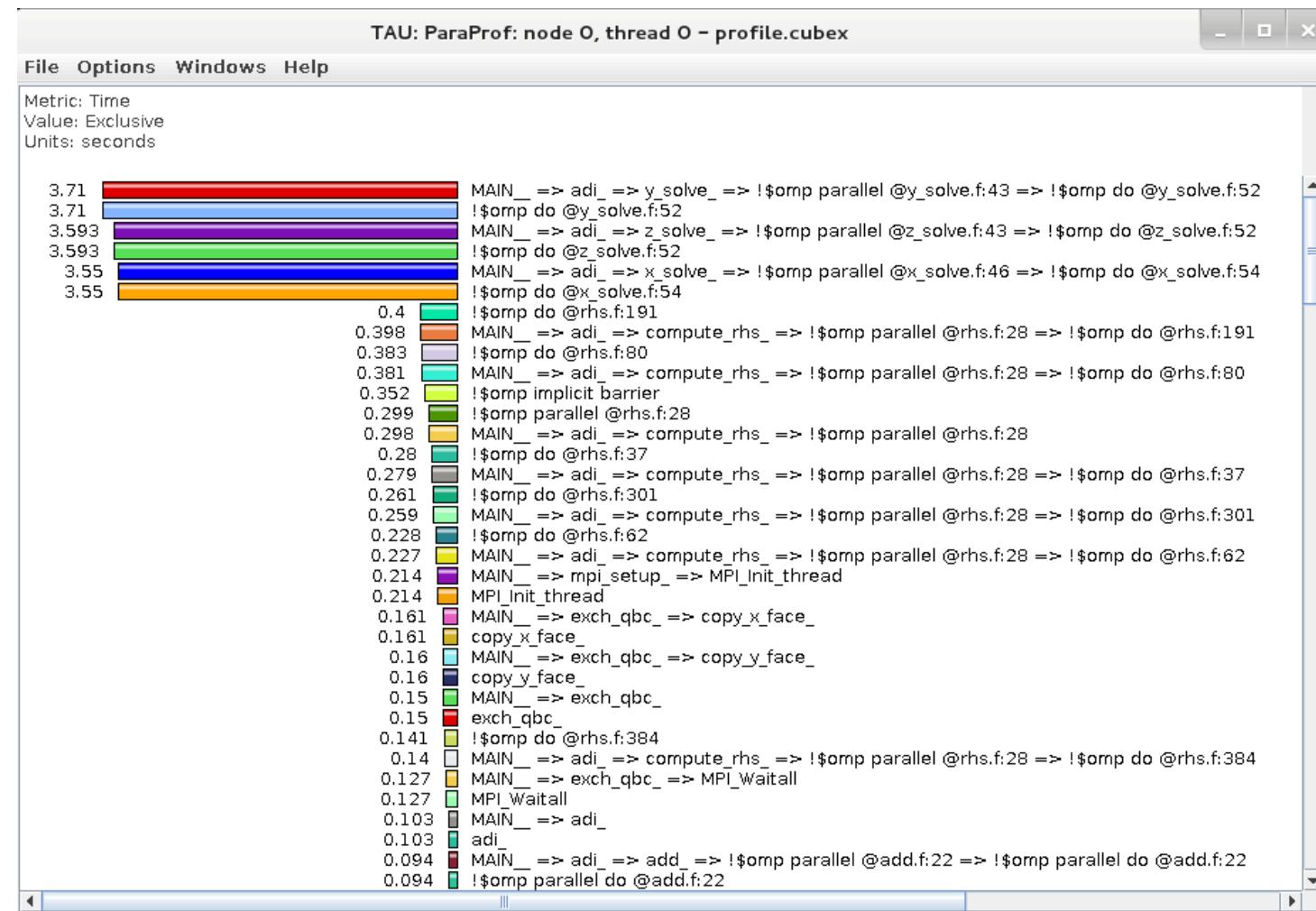
Exclusive time spent in each code region (OpenMP loop) is shown here for MPI rank 0 thread 1

# Instrumenting Source Code with PDT and Opari

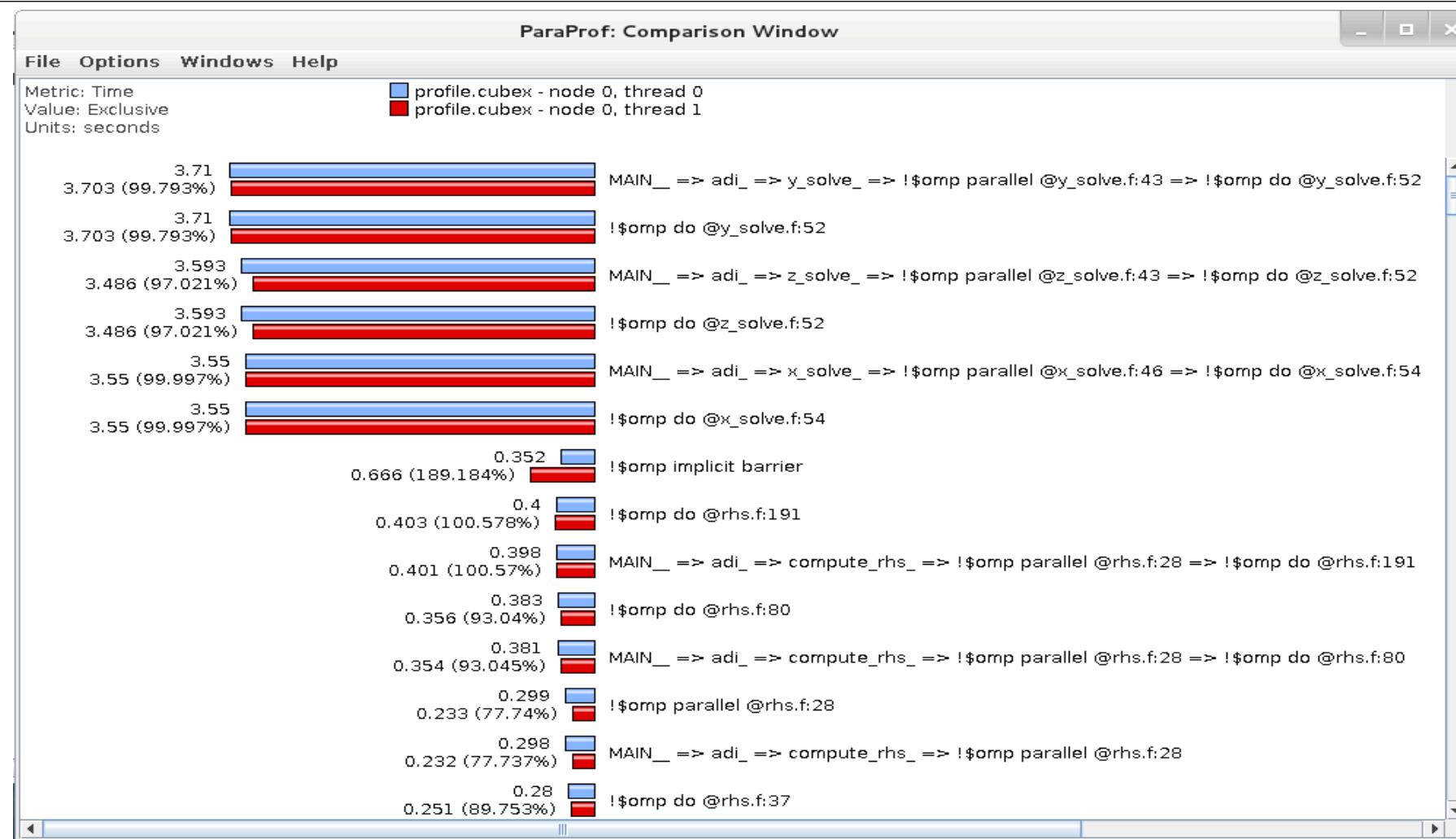
Frequently executing lightweight routines are automatically throttled at runtime. Reduces runtime dilation.



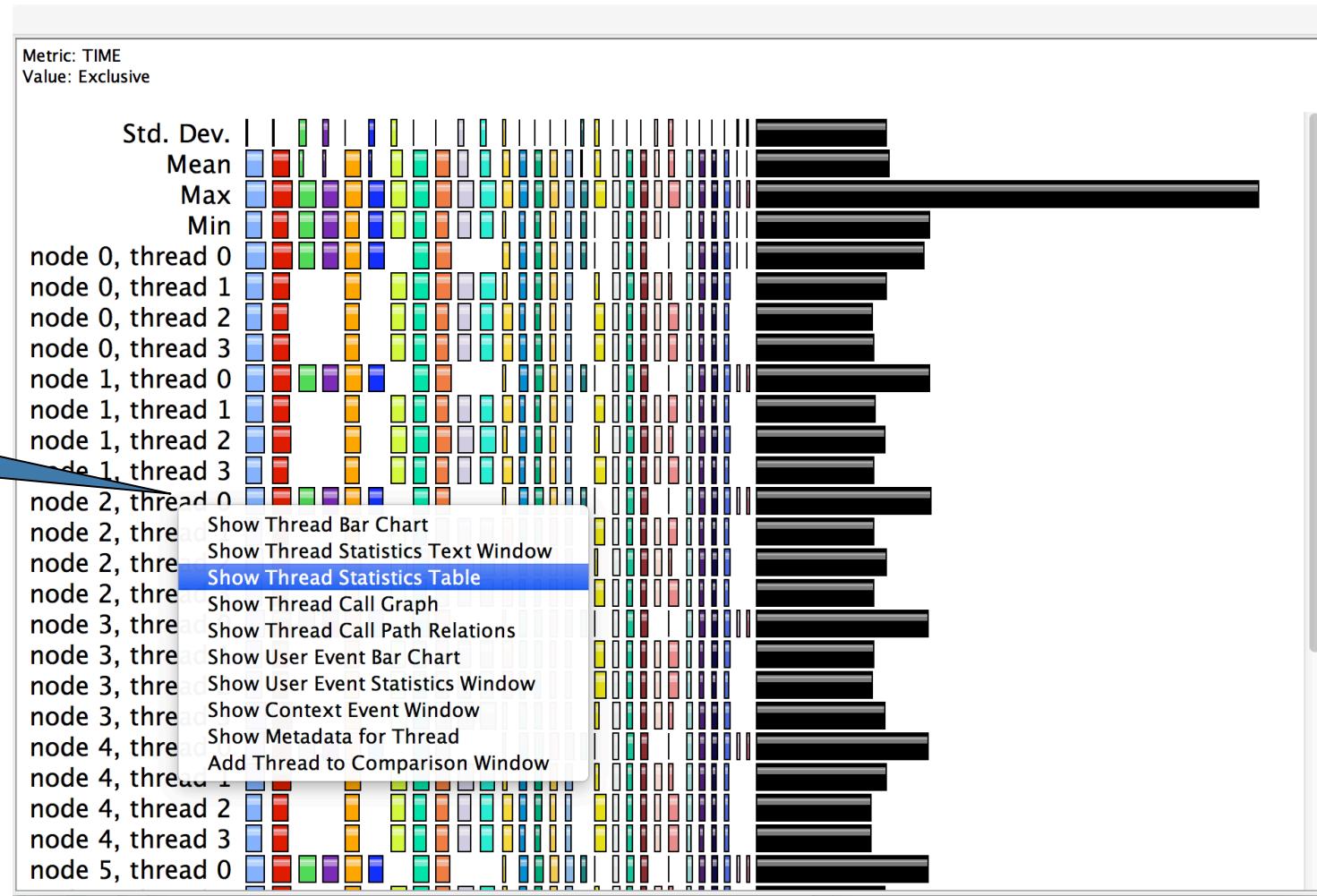
# ParaProf: Node view in a callpath profile



# ParaProf: Add thread to comparison window



# Paraprof Thread Statistics Table with TAU\_SAMPLING=1



# ParaProf: Thread Statistics Table

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

File Options Windows Help

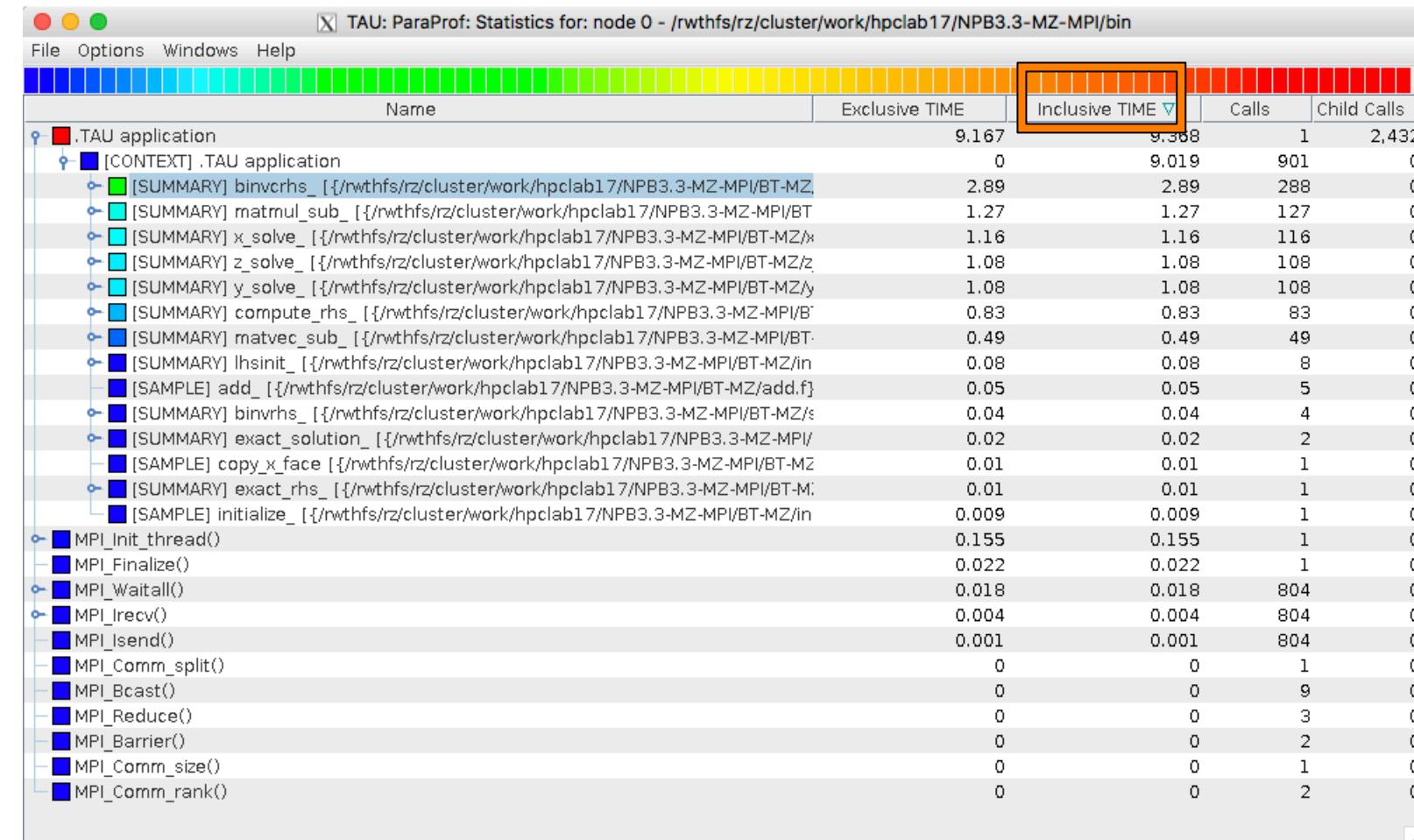
Time

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

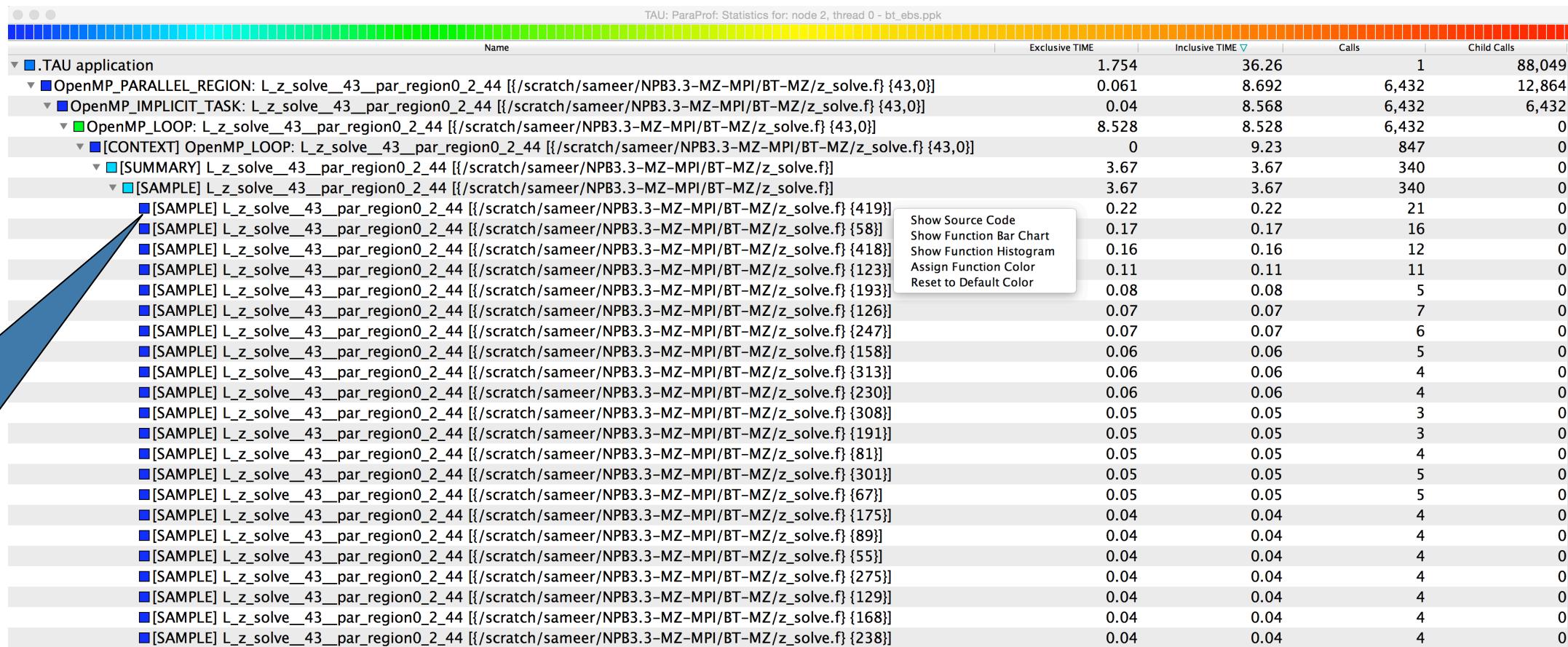
Click to sort by a given metric, drag and move to rearrange columns

# ParaProf

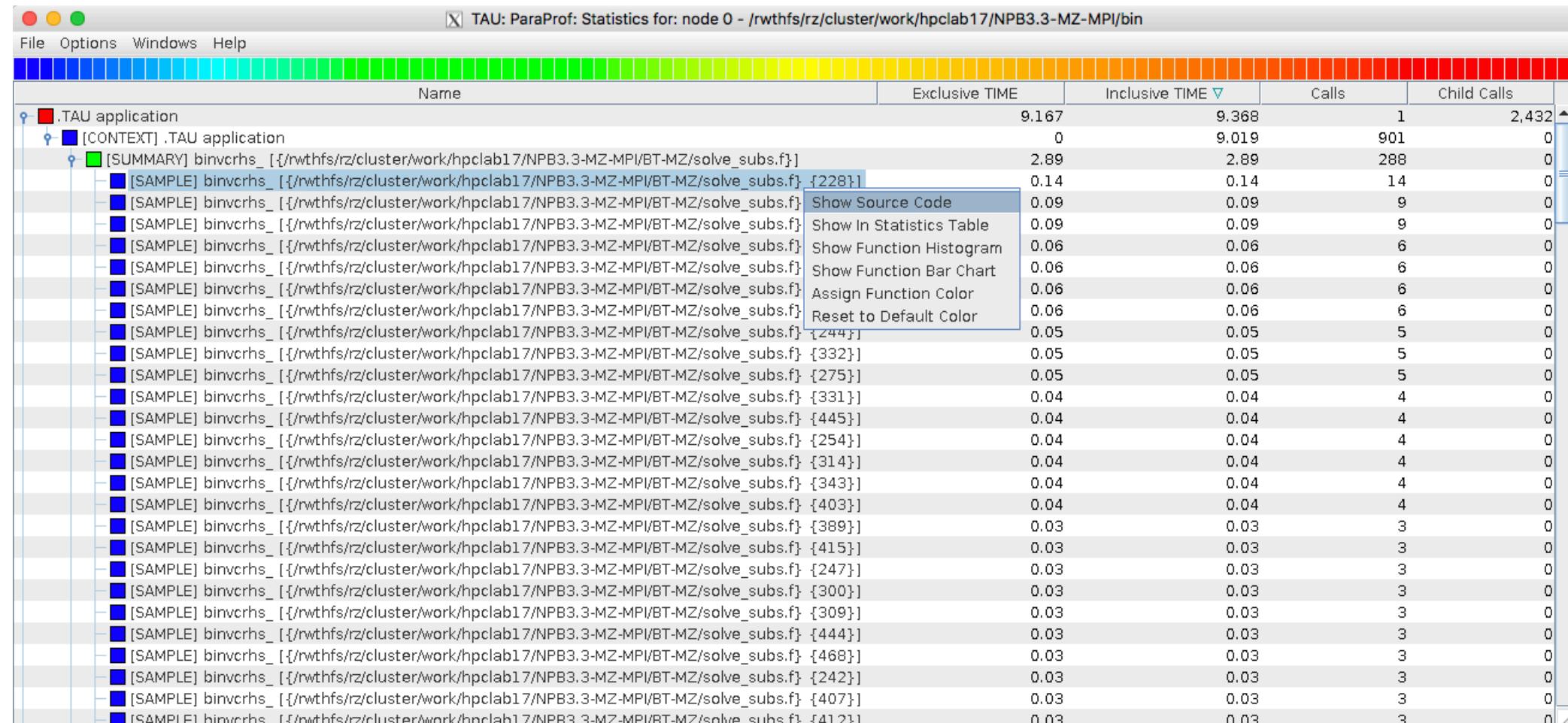
- Click on Columns:
  - to sort by incl time
- Open binvcrhs
- Click on Sample



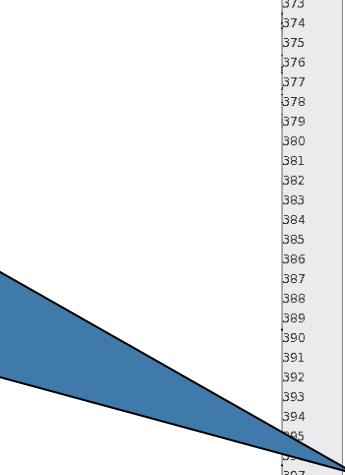
# Paraprof Thread Statistics Table



# ParaProf



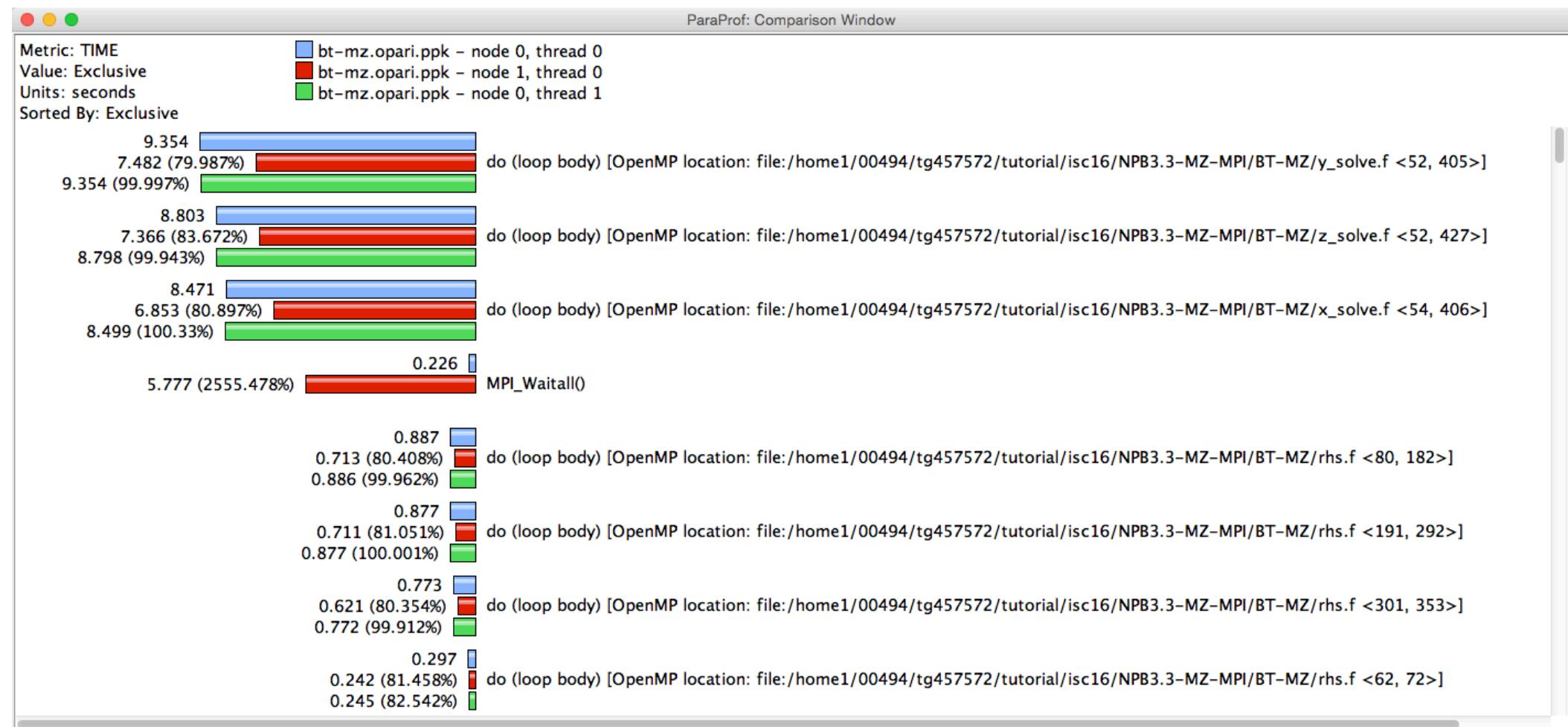
# Statement Level Profiling with TAU



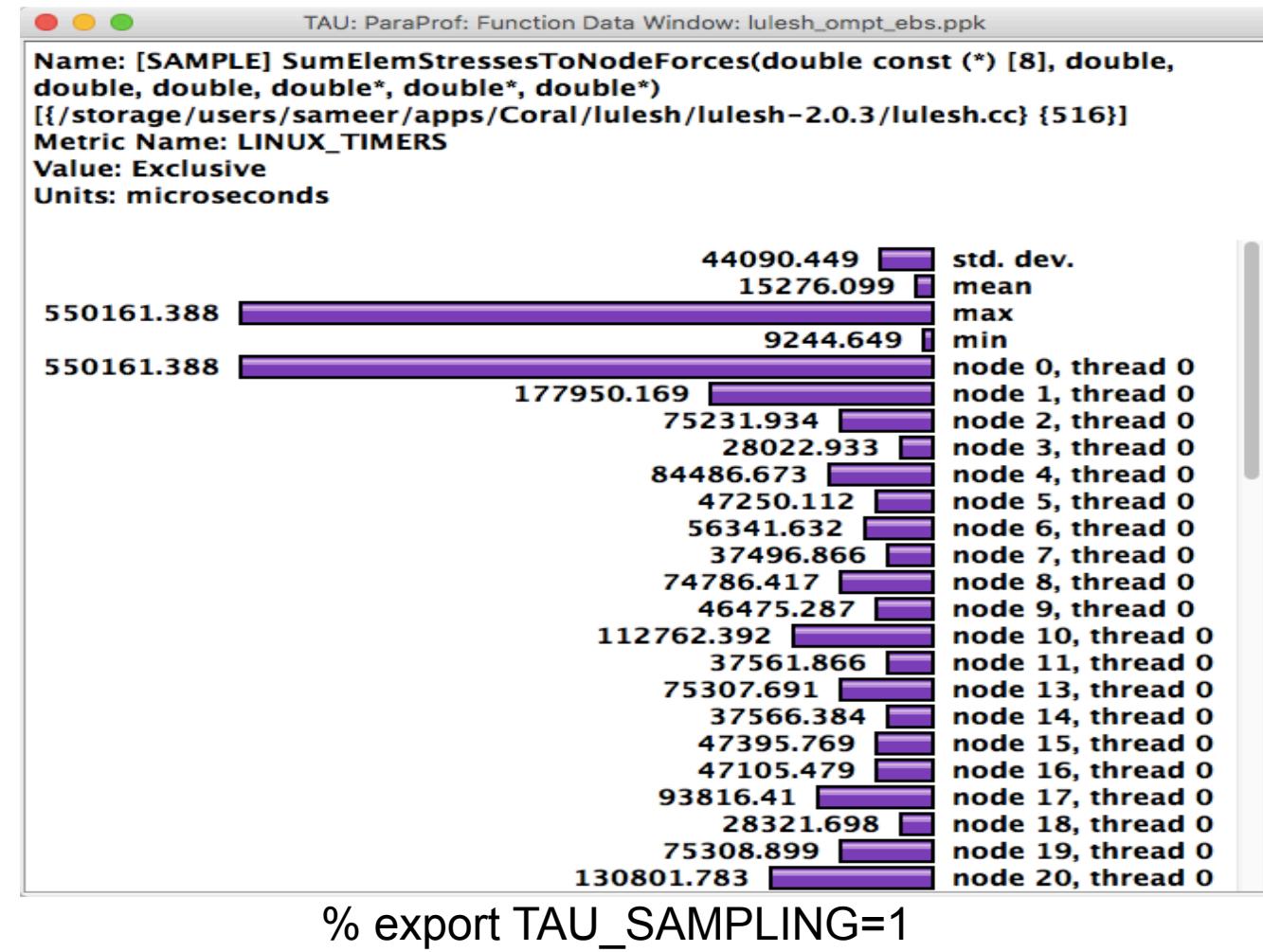
```
File Help
353     call matmul_sub(lhs(1,1,aa,i),
354             >           lhs(1,1,cc,i-1),
355             >           lhs(1,1,bb,i))
356
357
358 C-----c
359 c   multiply c(i,j,k) by b_inverse and copy back to c
360 c   multiply rhs(1,j,k) by b_inverse(1,j,k) and copy to rhs
361 C-----c
362     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
388 C-----c
389 c   back solve: if last cell, then generate U(isize)=rhs(isize)
390 c   else assume U(isize) is loaded in un pack backsolve_info
391 c   so just use it
392 c   after call u(istart) will be sent to next cell
393 C-----c
394
395     do i=isize-1,0,-1
396       do m=1,BLOCK_SIZE
397         do n=1,BLOCK_SIZE
398           rhs(m,i,j,k) = rhs(m,i,j,k)
399             >           - lhs(m,n,cc,i)*rhs(n,i+1,j,k)
400
401       enddo
402     enddo
403   enddo
```

Source location where samples are taken.  
Compute intensive region.

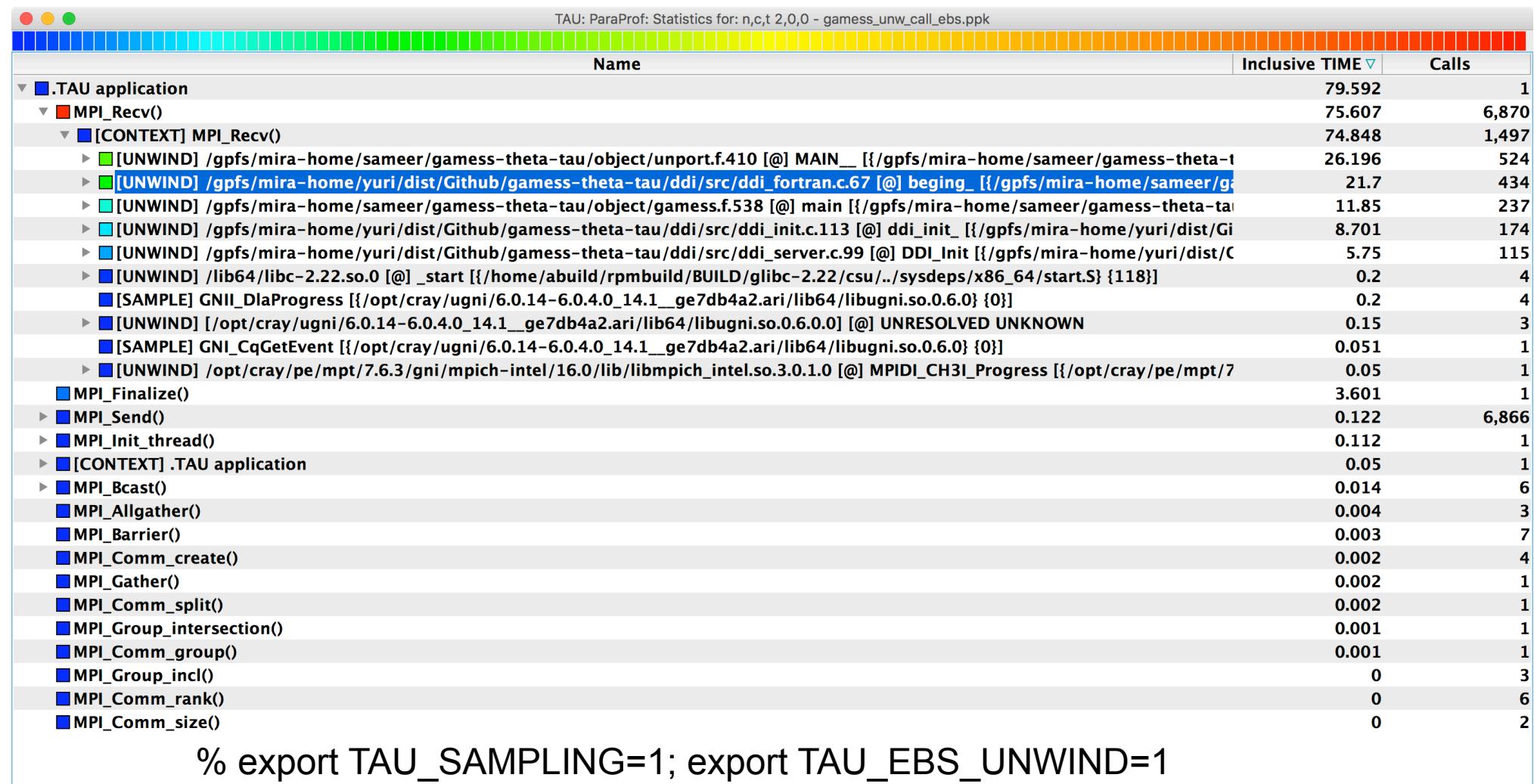
# ParaProf Comparison Window



# TAU – Event Based Sampling (EBS)



# Examples: Callstack Sampling in TAU



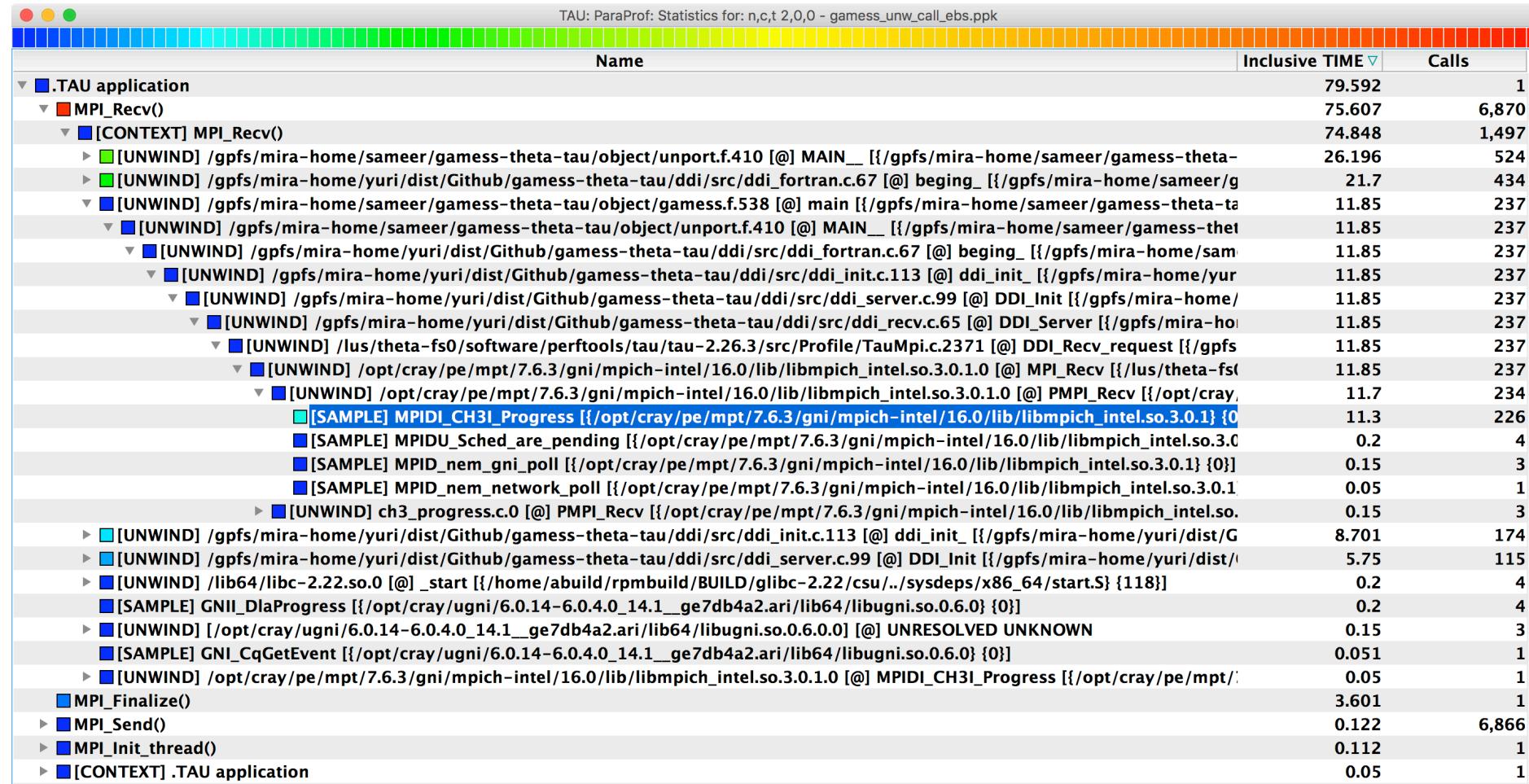
# UNWINDING CALLSTACKS

TAU: ParaProf: Statistics for: n,c,t 2,0,0 - gamess\_unw\_call\_ebs.ppk

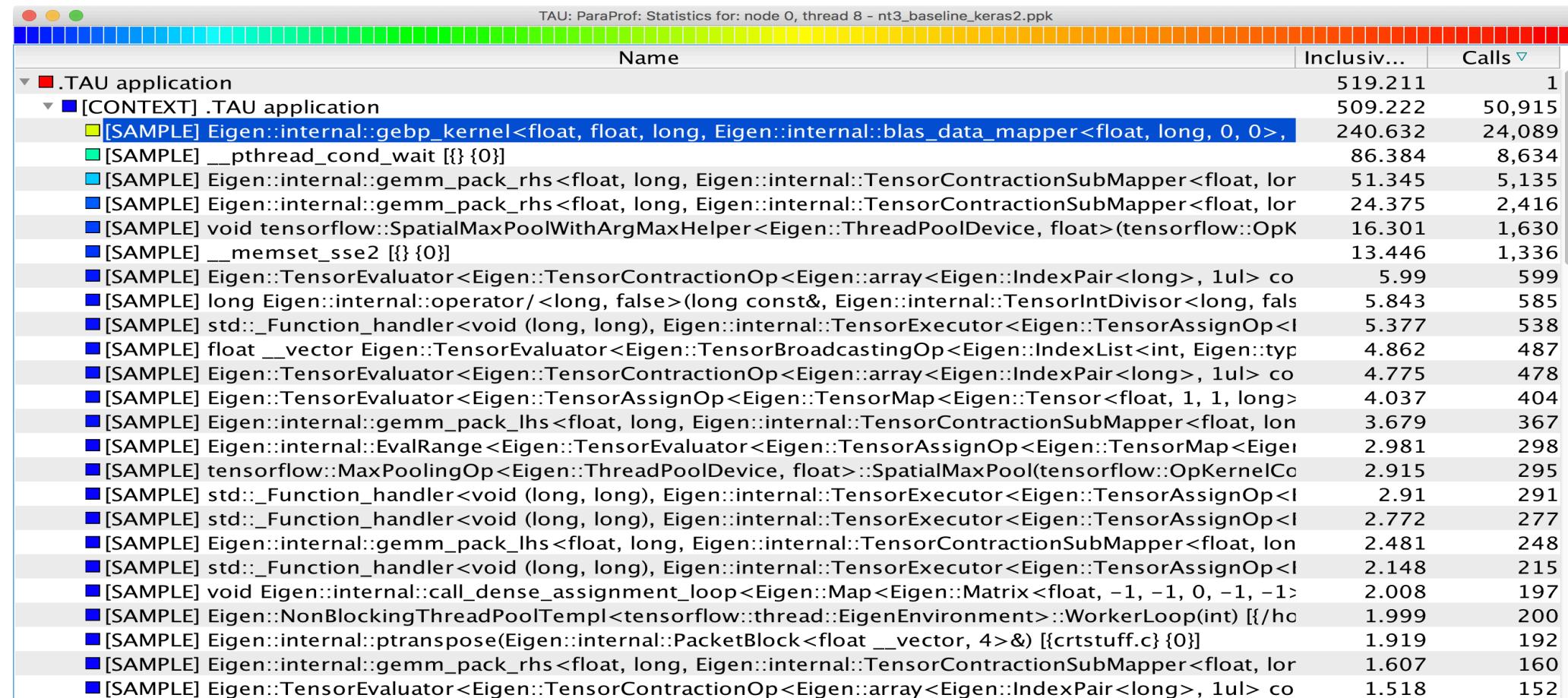
	Name	Inclusive TIME ▼	Calls
▼ .TAU application		79.592	1
▼ MPI_Recv()		75.607	6,870
▼ [CONTEXT] MPI_Recv()		74.848	1,497
► [UNWIND] /gpfs/mira-home/sameer/gamess-theta-tau/object/unport.f.410 [@] MAIN__ [{/gpfs/mira-home/sameer/gamess-theta-tau}]	26.196	524	
▼ [UNWIND] /gpfs/mira-home/yuri/dist/Github/gamess-theta-tau/ddi/src/ddi_fortran.c.67 [@] begin_ [{/gpfs/mira-home/sameer/gamess-theta-tau}]	21.7	434	
▼ [UNWIND] /gpfs/mira-home/yuri/dist/Github/gamess-theta-tau/ddi/src/ddi_init.c.113 [@] ddi_init_ [{/gpfs/mira-home/yuri/dist/Github/gamess-theta-tau/ddi/src/ddi_init.c.113}]	21.7	434	
▼ [UNWIND] /gpfs/mira-home/yuri/dist/Github/gamess-theta-tau/ddi/src/ddi_server.c.99 [@] DDI_Init [{/gpfs/mira-home/yuri/dist/Github/gamess-theta-tau/ddi/src/ddi_server.c.99}]	21.7	434	
▼ [UNWIND] /lus/theta-fs0/software/perf-tools/tau/tau-2.26.3/src/Profile/TauMpi.c.2371 [@] DDI_Recv_request [{/lus/theta-fs0/software/perf-tools/tau/tau-2.26.3/src/Profile/TauMpi.c.2371}]	21.7	434	
▼ [UNWIND] /opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.so.3.0.1.0 [@] MPI_Recv [{/opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.so.3.0.1.0}]	21.7	434	
▼ [UNWIND] /opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.so.3.0.1.0 [@] PMPI_Recv [{/opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.so.3.0.1.0}]	21.7	434	
▼ [UNWIND] /opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.so.3.0.1.0 [@] MPIIDI_CH3I_Progress [{/opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.so.3.0.1.0}]	21.45	429	
▼ [UNWIND] /opt/cray/ugni/6.0.14-6.0.4.0_14.1_ge7db4a2.ari/lib64/libugni.so.0.6.0.0 [@] MPID_nem_gni_poll [{/opt/cray/ugni/6.0.14-6.0.4.0_14.1_ge7db4a2.ari/lib64/libugni.so.0.6.0.0}]	15.95	319	
■ [SAMPLE] GNI_SmsgGetNextWTag [{/opt/cray/ugni/6.0.14-6.0.4.0_14.1_ge7db4a2.ari/lib64/libugni.so.0.6.0}]	10.349	207	
■ [SAMPLE] GNI_CqGetEvent [{/opt/cray/ugni/6.0.14-6.0.4.0_14.1_ge7db4a2.ari/lib64/libugni.so.0.6.0} {0}]	5.6	112	
► [UNWIND] gni_poll.c.0 [@] MPID_nem_gni_poll [{/opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.so.3.0.1.0}]	5.25	105	
► [UNWIND] /opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.so.3.0.1.0 [@] MPID_nem_gni_poll [{/opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.so.3.0.1.0}]	0.25	5	
► [UNWIND] UNRESOLVED [@] MPIIDI_CH3I_Progress [{/opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.so.3.0.1.0}]	0.25	5	
► [UNWIND] /gpfs/mira-home/sameer/gamess-theta-tau/object/gamess.f.538 [@] main [{/gpfs/mira-home/sameer/gamess-theta-tau}]	11.85	237	
► [UNWIND] /gpfs/mira-home/yuri/dist/Github/gamess-theta-tau/ddi/src/ddi_init.c.113 [@] ddi_init_ [{/gpfs/mira-home/yuri/dist/Github/gamess-theta-tau/ddi/src/ddi_init.c.113}]	8.701	174	
► [UNWIND] /gpfs/mira-home/yuri/dist/Github/gamess-theta-tau/ddi/src/ddi_server.c.99 [@] DDI_Init [{/gpfs/mira-home/yuri/dist/Github/gamess-theta-tau/ddi/src/ddi_server.c.99}]	5.75	115	
► [UNWIND] /lib64/libc-2.22.so.0 [@] _start [{/home/abuild/rpmbuild/BUILD/glibc-2.22/csu/../sysdeps/x86_64/start.S} {118}]	0.2	4	
■ [SAMPLE] GNII_DlaProgress [{/opt/cray/ugni/6.0.14-6.0.4.0_14.1_ge7db4a2.ari/lib64/libugni.so.0.6.0} {0}]	0.2	4	
► [UNWIND] /opt/cray/ugni/6.0.14-6.0.4.0_14.1_ge7db4a2.ari/lib64/libugni.so.0.6.0.0 [@] UNRESOLVED UNKNOWN	0.15	3	
■ [SAMPLE] GNI_CqGetEvent [{/opt/cray/ugni/6.0.14-6.0.4.0_14.1_ge7db4a2.ari/lib64/libugni.so.0.6.0} {0}]	0.051	1	
► [UNWIND] /opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.so.3.0.1.0 [@] MPIIDI_CH3I_Progress [{/opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.so.3.0.1.0}]	0.05	1	
■ [MPI_Finalize()]	3.601	1	
► [MPI_Send()]	0.122	6,866	
► [MPI_Init_thread()]	0.112	1	
► [CONTEXT] .TAU application	0.05	1	

```
% export TAU_SAMPLING=1; export TAU_EBS_UNWIND=1
```

# UNWINDING CALLSTACKS

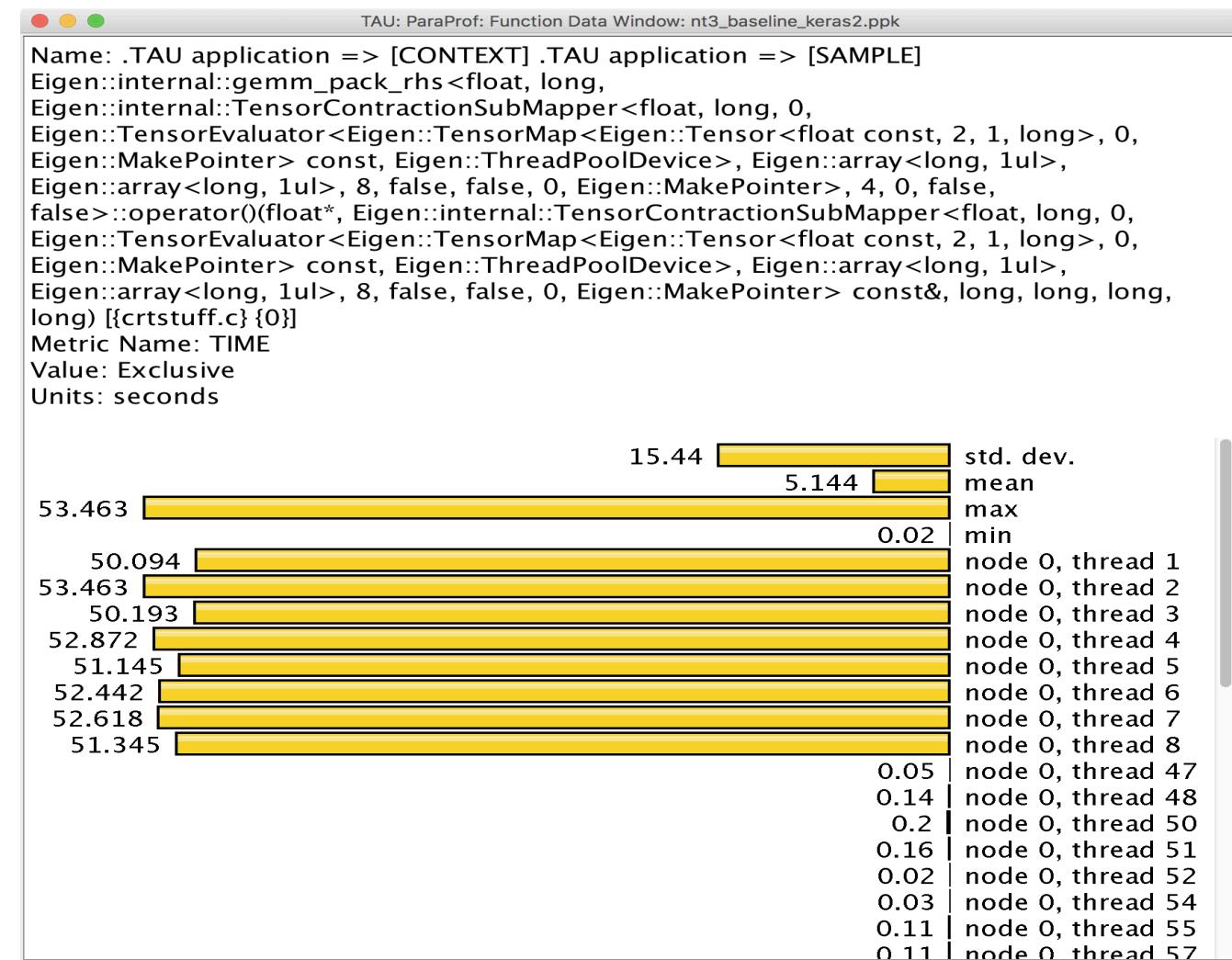


# **Deep Learning: Tensorflow**

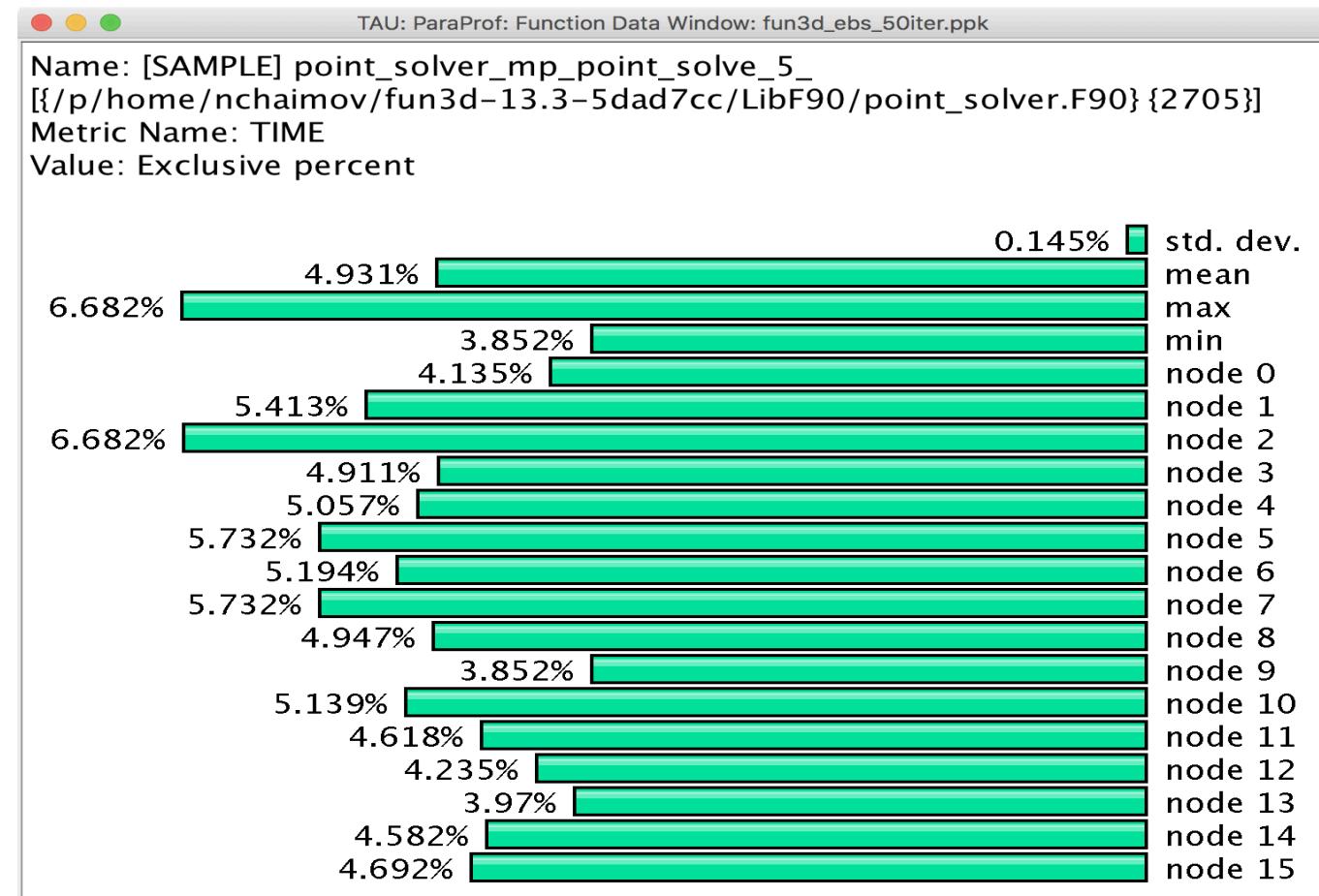


```
% tau_python -ebs nt3_baseline_keras2.py (CANDLE)
```

# Sampling Tensorflow



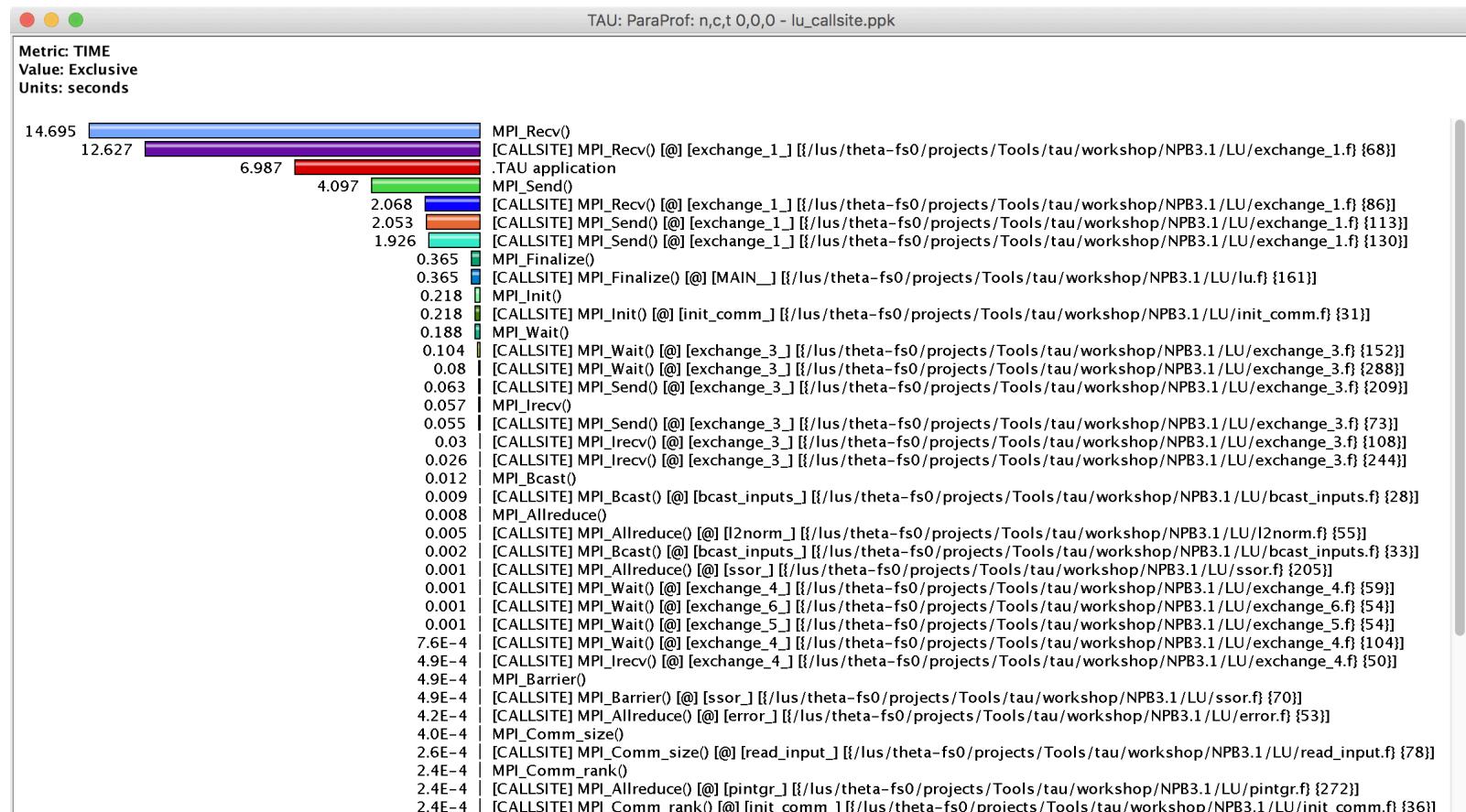
# Event Based Sampling (EBS)



Uninstrumented!

% aprun -n 16 tau\_exec **-ebs** a.out

# Callsite Profiling and Tracing



% export TAU\_CALLSITE=1

# CALLPATH THREAD RELATIONS WINDOW

TAU: ParaProf: Call Path Data n,c,t, 2,0,0 - gamess\_unw\_call\_ebs.ppk

Metric Name: TIME  
Sorted By: Inclusive  
Units: seconds

	Exclusive	Inclusive	Calls/Tot.Calls	Name[id]
-->	0.121	79.592	1	.TAU application
	0.002	0.002	1/1	MPI_Gather()
	0.004	0.004	3/3	MPI_Allgather()
	0.122	0.122	6866/6866	MPI_Send()
	0.002	0.002	1/1	MPI_Comm_split()
	8.9E-5	8.9E-5	2/2	MPI_Comm_size()
	4.6E-4	4.6E-4	3/3	MPI_Group_incl()
	75.607	75.607	6870/6870	MPI_Recv()
	0.002	0.002	4/4	MPI_Comm_create()
	9.5E-5	9.5E-5	6/6	MPI_Comm_rank()
	5.4E-4	5.4E-4	1/1	MPI_Comm_group()
	0.003	0.003	7/7	MPI_BARRIER()
	0.112	0.112	1/1	MPI_Init_thread()
	6.3E-4	6.3E-4	1/1	MPI_Group_intersection()
	0	0.05	1/1	[CONTEXT] .TAU application
	3.601	3.601	1/1	MPI_Finalize()
	0.014	0.014	6/6	MPI_Bcast()
-->	75.607	75.607	6870/6870	.TAU application
-->	75.607	75.607	6870	MPI_Recv()
-->	0	74.848	1497/1497	[CONTEXT] MPI_Recv()
-->	0	74.848	1497	MPI_Recv()
-->	0	8.701	174/1371	[CONTEXT] MPI_Recv()
-->	0	26.196	524/763	[UNWIND] /gpfs/mira-home/yuri/dist/Github/gamess-theta-tau/ddi/src/ddi_init.c.113 [0] ddi_in
-->	0.2	0.2	4/138	[UNWIND] /gpfs/mira-home/sameer/gamess-theta-tau/object/unport.f.410 [0] MAIN_ [{/gpfs/mir
-->	0	5.75	115/1484	[SAMPLE] GNIIL_DlaProgress [{/opt/cray/ugni/6.0.14-6.0.4.0_14.1_ge7db4a2.ari/lib64/libugni.
-->	0	0.2	4/5	[UNWIND] /gpfs/mira-home/yuri/dist/Github/gamess-theta-tau/ddi/src/ddi_server.c.99 [0] DDI_
-->	0	11.85	237/239	[UNWIND] /lib64/libc-2.22.so.0 [0] _start [{/home/abuild/rpmbuild/BUILD/glibc-2.22/csu/.s
-->	0.051	0.051	1/273	[UNWIND] /gpfs/mira-home/sameer/gamess-theta-tau/object/gamess.f.538 [0] main [{/gpfs/mira-
-->	0	0.05	1/1197	[SAMPLE] GNI_CqGetEvent [{/opt/cray/ugni/6.0.14-6.0.4.0_14.1_ge7db4a2.ari/lib64/libugni.so
-->	0	0.15	3/7	[UNWIND] /opt/cray/pe/mpit/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.so.3.0.1.0 [0] MPID:
-->	0	21.7	434/1197	[UNWIND] /opt/cray/ugni/6.0.14-6.0.4.0_14.1_ge7db4a2.ari/lib64/libugni.so.0.6.0.0) [0] UNI
-->				[UNWIND] /gpfs/mira-home/yuri/dist/Github/gamess-theta-tau/ddi/src/ddi_fortran.c.67 [0] beg:

# CALLPATH THREAD RELATIONS WINDOW

TAU: ParaProf: Call Path Data n,c,t, 2,0,0 - gamess\_unw\_call\_ebs.ppk

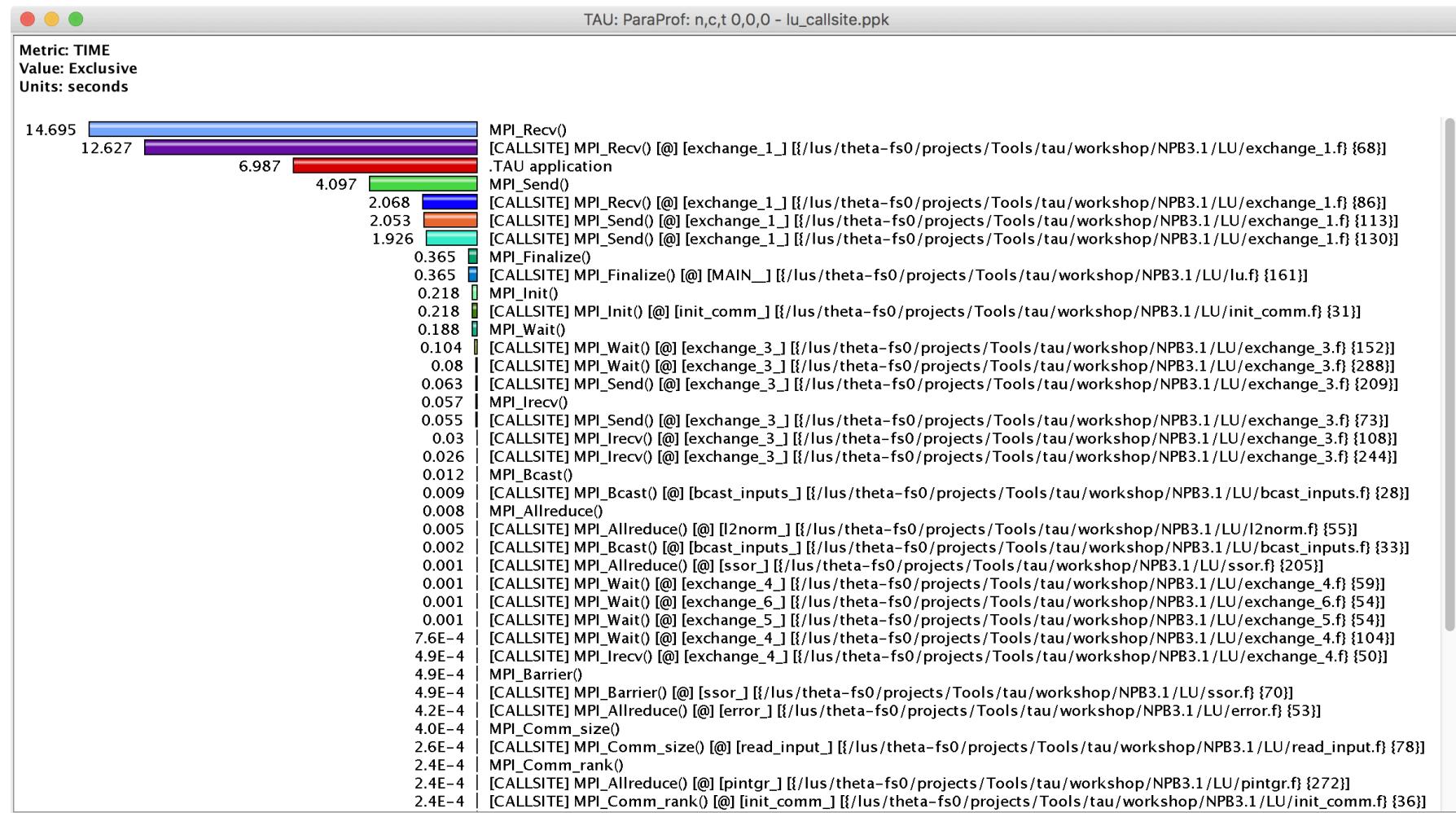
Metric Name: TIME  
Sorted By: Exclusive  
Units: seconds

	Exclusive	Inclusive	Calls/Tot.Calls	Name[id]
-->	75.607	75.607	6870/6870	.TAU application
-->	75.607	75.607	6870	MPI_Recv()
-->	0	74.848	1497/1497	[CONTEXT] MPI_Recv()
-->	0.15	0.15	3/444	[UNWIND] /opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.so.3.0.1.0 [0] PMPI_Recv
-->	22.046	22.046	441/444	[UNWIND] /opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.so.3.0.1.0 [0] MPIDI_CH3I
-->	22.196	22.196	444	[SAMPLE] MPID_nem_gni_poll [{/opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.so.3
-->	5.6	5.6	112/273	[UNWIND] /opt/cray/ugni/6.0.14-6.0.4.0_14.1_ge7db4a2.ari/lib64/libugni.so.0.6.0.0 [0] MPID_nem_6
-->	0.051	0.051	1/273	[CONTEXT] MPI_Recv()
-->	7.651	7.651	153/273	[UNWIND] /opt/cray/ugni/6.0.14-6.0.4.0_14.1_ge7db4a2.ari/lib64/libugni.so.0.6.0.0 [0] MPID_nem_6
-->	0.35	0.35	7/273	[UNWIND] [/opt/cray/ugni/6.0.14-6.0.4.0_14.1_ge7db4a2.ari/lib64/libugni.so.0.6.0.0] [0] UNRESOLV
-->	13.652	13.652	273	[SAMPLE] GNI_CqGetEvent [{/opt/cray/ugni/6.0.14-6.0.4.0_14.1_ge7db4a2.ari/lib64/libugni.so.0.6.0.0]
-->	11.3	11.3	226/226	[UNWIND] /opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.so.3.0.1.0 [0] PMPI_Recv
-->	11.3	11.3	226	[SAMPLE] MPIDI_CH3I_Progress [{/opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.so
-->	10.349	10.349	207/207	[UNWIND] /opt/cray/ugni/6.0.14-6.0.4.0_14.1_ge7db4a2.ari/lib64/libugni.so.0.6.0.0 [0] MPID_nem_6
-->	10.349	10.349	207	[SAMPLE] GNI_SmsgGetNextWTag [{/opt/cray/ugni/6.0.14-6.0.4.0_14.1_ge7db4a2.ari/lib64/libugni.so.0.6.0.0]
-->	0.2	0.2	4/138	[CONTEXT] MPI_Recv()
-->	6.701	6.701	134/138	[UNWIND] /opt/cray/ugni/6.0.14-6.0.4.0_14.1_ge7db4a2.ari/lib64/libugni.so.0.6.0.0 [0] GNI_CqGet
-->	6.901	6.901	138	[SAMPLE] GNII_DlaProgress [{/opt/cray/ugni/6.0.14-6.0.4.0_14.1_ge7db4a2.ari/lib64/libugni.so.0.6.0.0]
-->	5.25	5.25	105/109	[UNWIND] gni_poll.c.0 [0] MPID_nem_gni_poll [{/opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/l1]
-->	0.2	0.2	4/109	[UNWIND] gni_poll.c.0 [0] MPIDI_CH3I_Progress [{/opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/
-->	5.45	5.45	109	[SAMPLE] MPID_nem_gni_check_localCQ [{gni_poll.c} {0}]
-->	3.601	3.601	1/1	.TAU application
-->	3.601	3.601	1	MPI_Finalize()

# ParaProf: Callpath Thread Relations Window

TAU: ParaProf: Call Path Data n,c,t, O,O,O - scout.cubex				
	File	Options	Windows	Help
<b>Metric Name:</b> Time <b>Sorted By:</b> Exclusive <b>Units:</b> seconds				
-->	0.04	0.04	32/32	!\$omp parallel @initialize.f:28 !\$omp do @initialize.f:50
-->	0.03	2.536	3232/3232	compute_rhs_
9.8E-4	0.03	2.536	3232	!\$omp parallel @rhs.f:28 !\$omp master @rhs.f:424
0.225	9.8E-4	3232/3232	!\$omp do @rhs.f:62	
0.002	0.228	3232/3232	!\$omp master @rhs.f:74	
0.002	0.002	3232/3232	!\$omp master @rhs.f:293	
0.199	0.002	3232/3232	!\$omp do @rhs.f:384	
0.002	0.199	3232/3232	!\$omp master @rhs.f:183	
0.343	0.002	3232/3232	!\$omp do @rhs.f:37	
0.016	0.343	3232/3232	!\$omp do @rhs.f:372	
0.014	0.016	3232/3232	!\$omp do @rhs.f:413	
0.609	0.014	3232/3232	!\$omp do @rhs.f:191	
0.36	0.609	3232/3232	!\$omp do @rhs.f:301	
0.583	0.36	3232/3232	!\$omp do @rhs.f:80	
0.019	0.583	3232/3232	!\$omp do @rhs.f:400	
0.006	0.019	3232/51680	!\$omp implicit barrier	
0.069	0.006	3232/51680	!\$omp do @rhs.f:428	
0.015	0.069	3232/3232	!\$omp do @rhs.f:359	
-->	0.021	0.029	6432/6432	!\$omp parallel @exch_qbc.f:215
0.021	0.021	6432	!\$omp parallel do @exch_qbc.f:215	
0.007	0.007	6432/51680	!\$omp implicit barrier	
-->	0.02	0.021	6432/6432	!\$omp parallel @exch_qbc.f:255
0.02	0.02	6432	!\$omp parallel do @exch_qbc.f:255	
0.013	0.02	6432/51680	!\$omp implicit barrier	

# Callsite Profiling and Tracing (TAU\_CALLSITE=1)



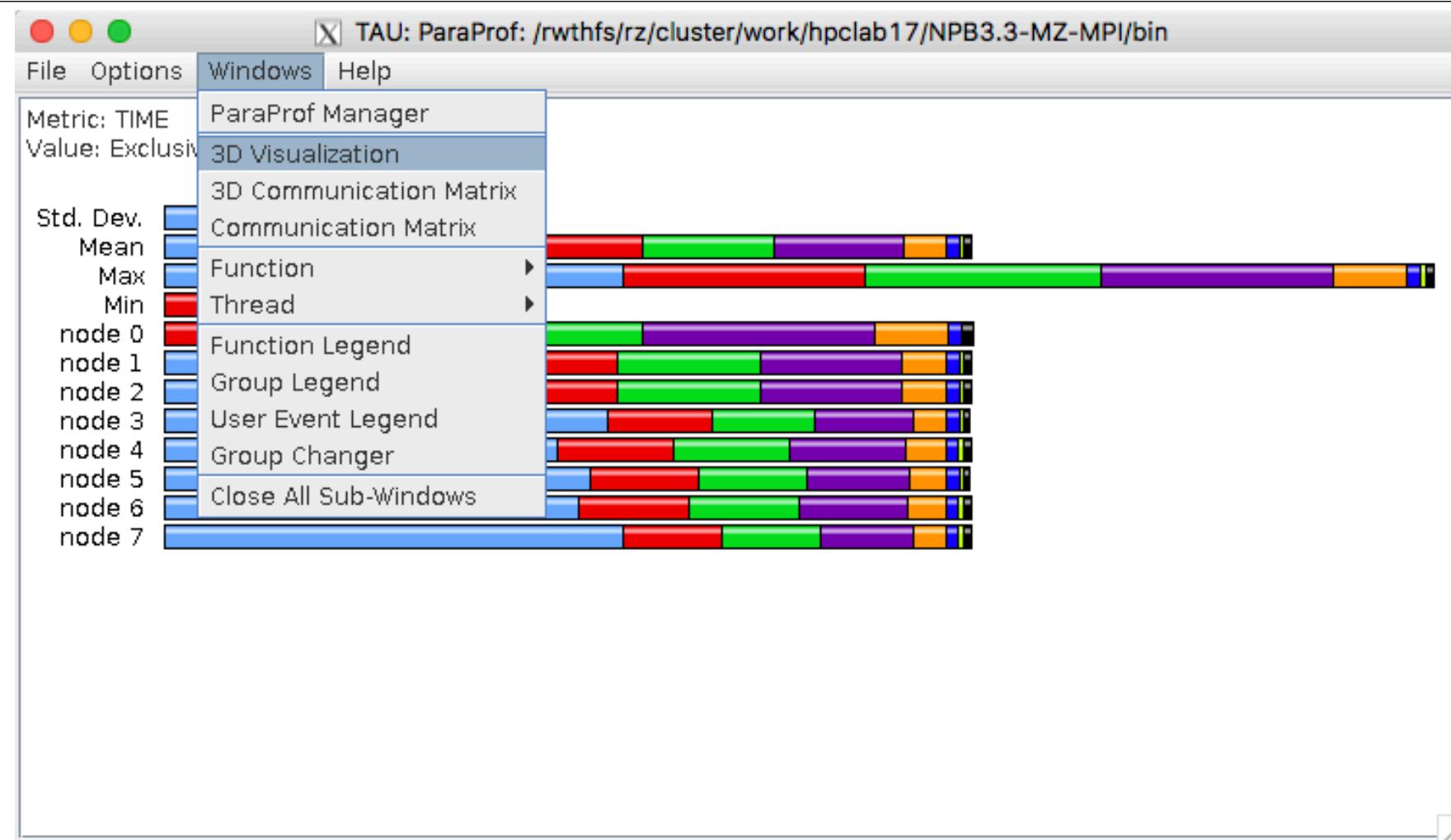
# TAU – Context Events

Name	Total	MeanValue	NumSamples	MinValue	MaxValue	Std. Dev.
.TAU application						
► read()						
► fopen64()						
► fclose()						
▼ OurMain()						
malloc size	25,235	1,097.174	23	11	12,032	2,851.143
free size	22,707	1,746.692	13	11	12,032	3,660.642
▼ OurMain [{wrapper.py}{3}]						
► read()						
malloc size	3,877	323.083	12	32	981	252.72
free size						122
► fopen64()						
► fclose()						
▼ <module> [{obe.py}{8}]						
▼ writeRestartData [{samarcInterface.py}{145}]						
▼ samarcWriteRestartData						
▼ write()						
WRITE Bandwidth (MB/s) <file="samarc/restore.00002/nodes.00004/proc.00001">	74.565	117	0	2,156.889	246.386	
WRITE Bandwidth (MB/s) <file="samarc/restore.00001/nodes.00004/proc.00001">	77.594	117	0	1,941.2	228.366	
WRITE Bandwidth (MB/s)	76.08	234	0	2,156.889	237.551	
Bytes Written <file="samarc/restore.00002/nodes.00004/proc.00001">	2,097,552	17,927.795	117	1	1,048,576	133,362.946
Bytes Written <file="samarc/restore.00001/nodes.00004/proc.00001">	2,097,552	17,927.795	117	1	1,048,576	133,362.946
Bytes Written	4,195,104	17,927.795	234	1	1,048,576	133,362.946
► open64()						

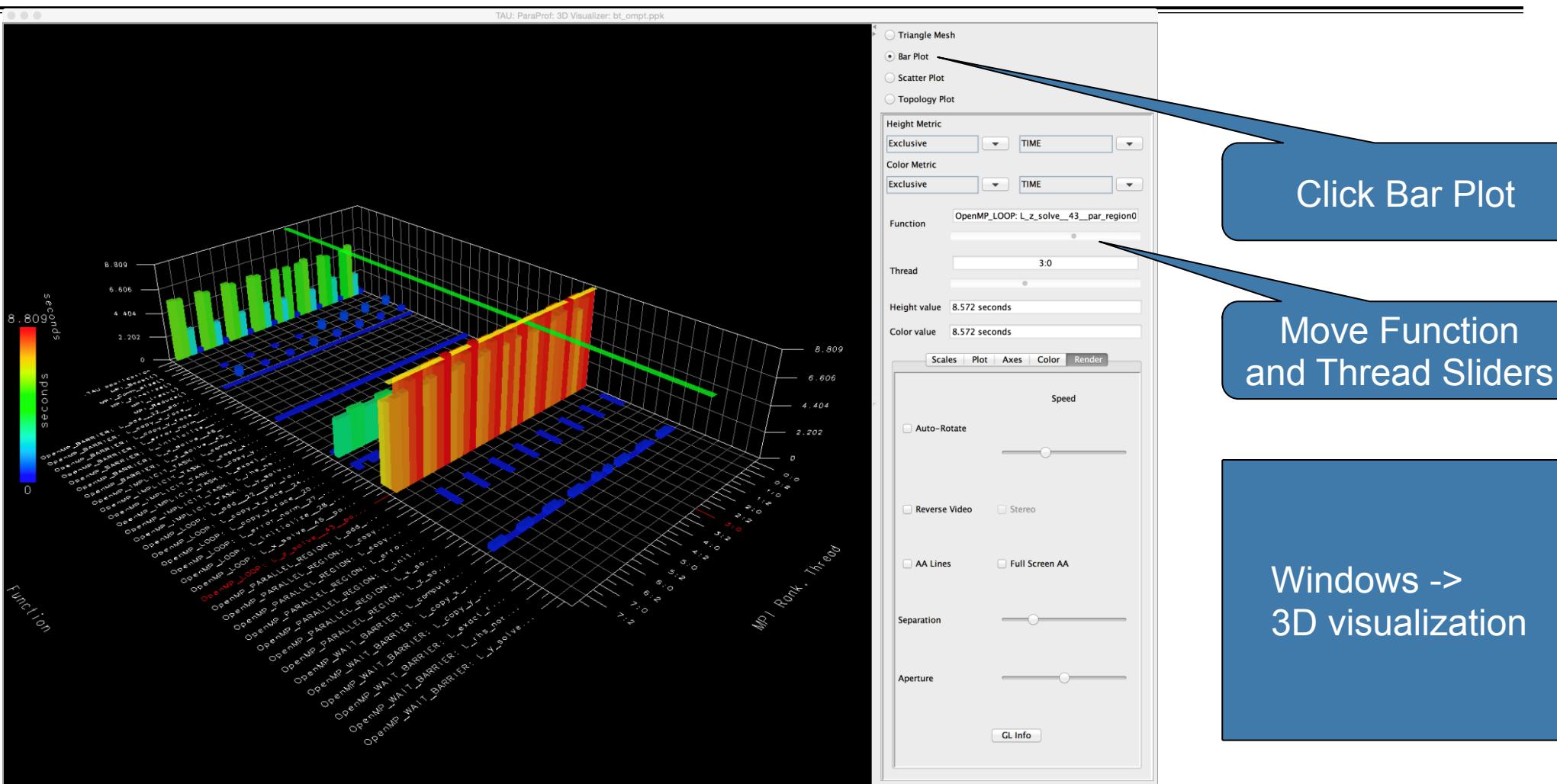
Write bandwidth per file

Bytes written to each file

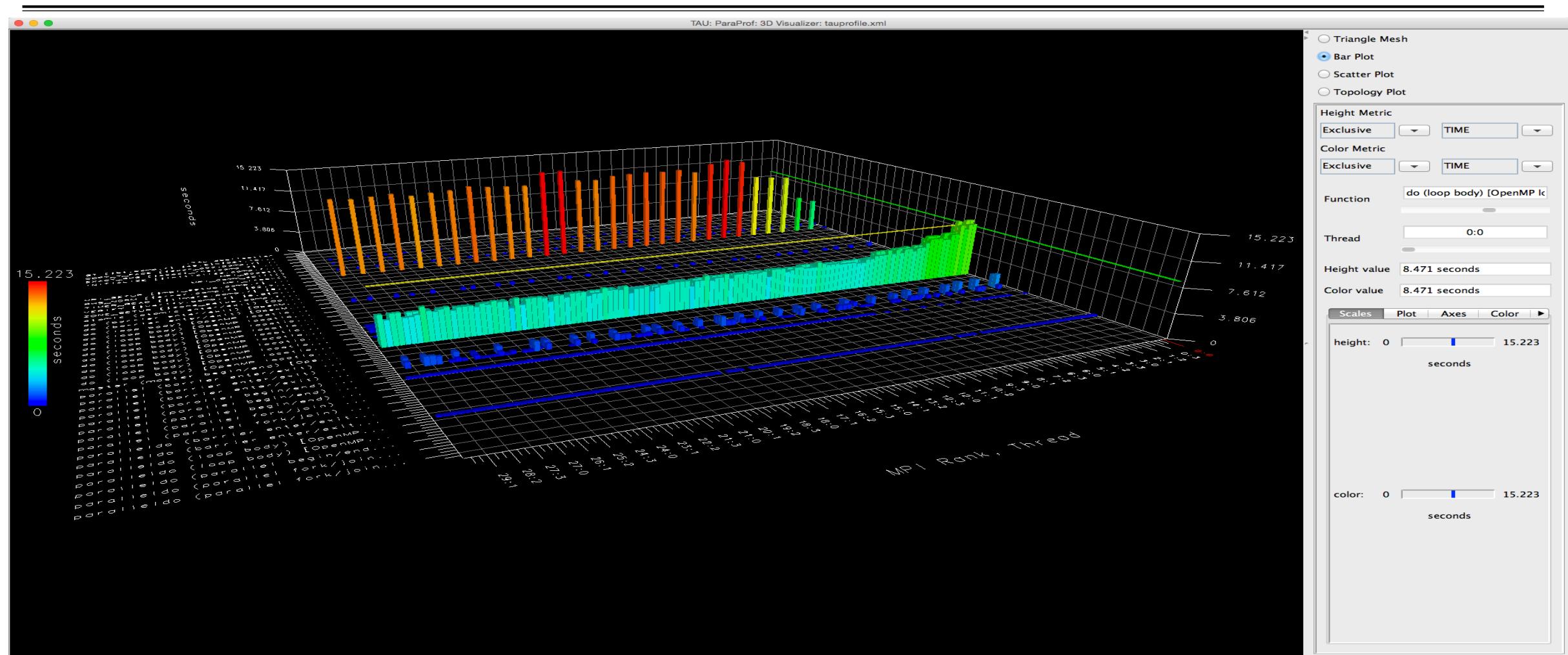
# ParaProf with Optimized Instrumentation



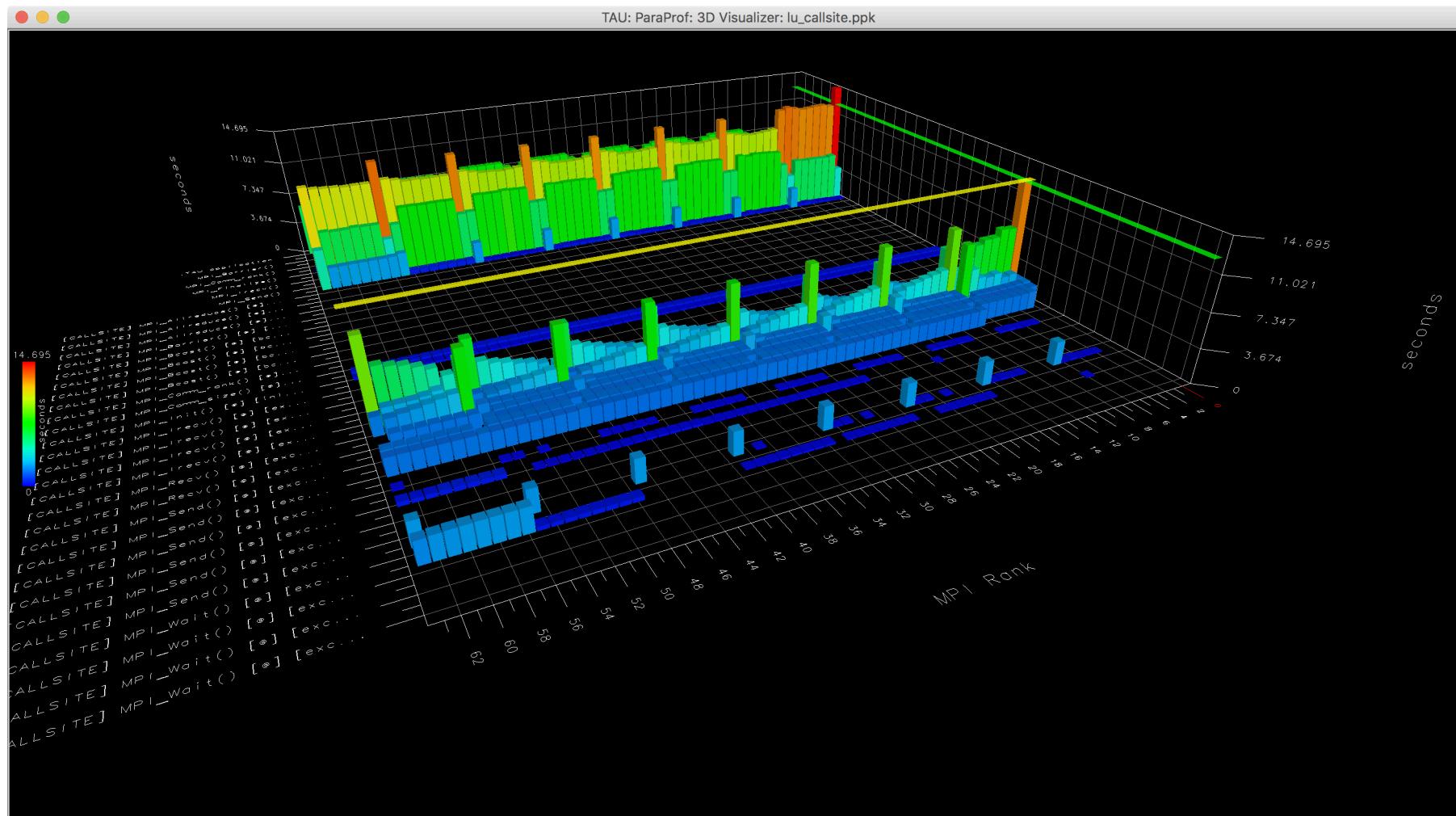
# Paraprof 3D visualization window



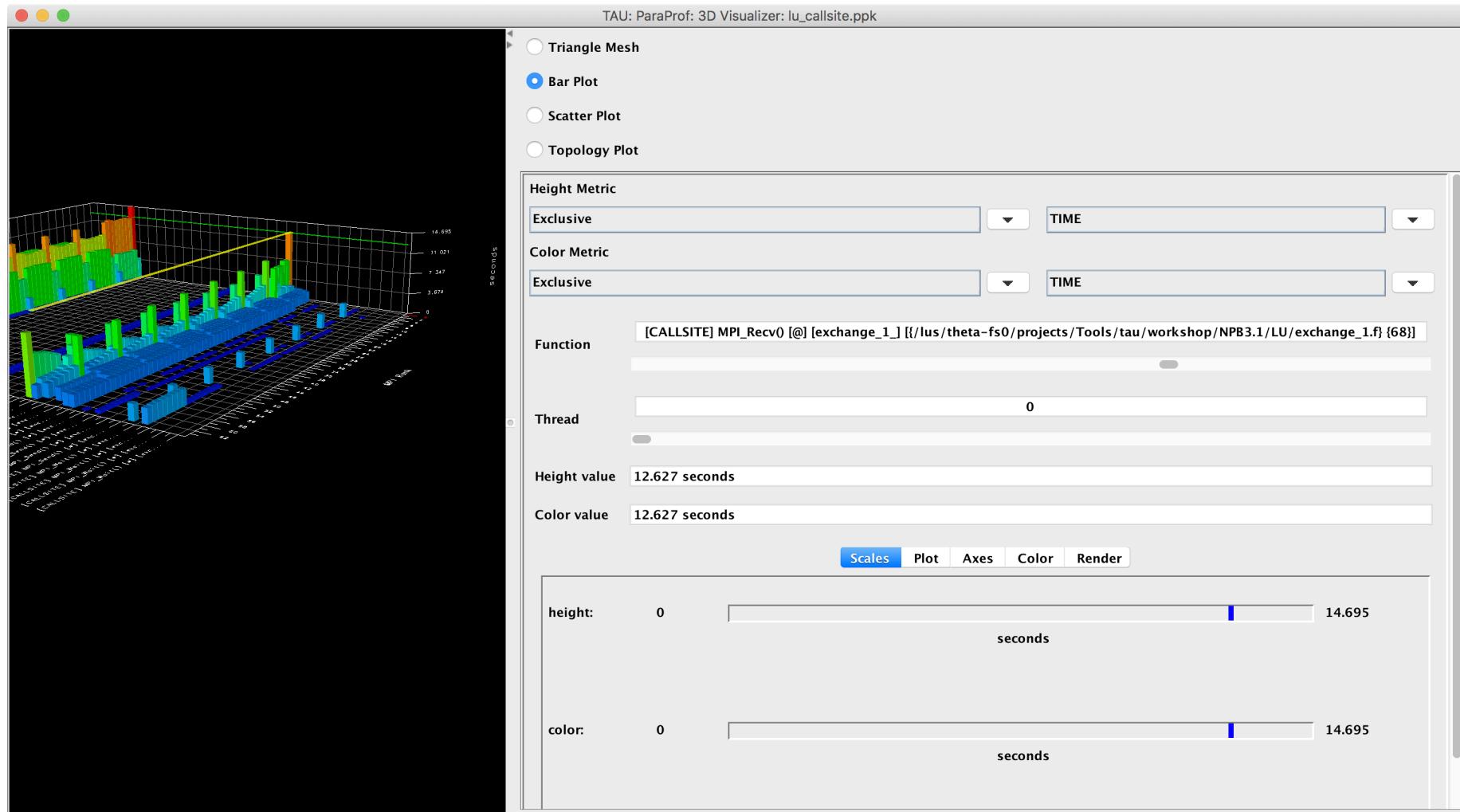
# ParaProf: 3D Visualization Window Showing Entire Profile



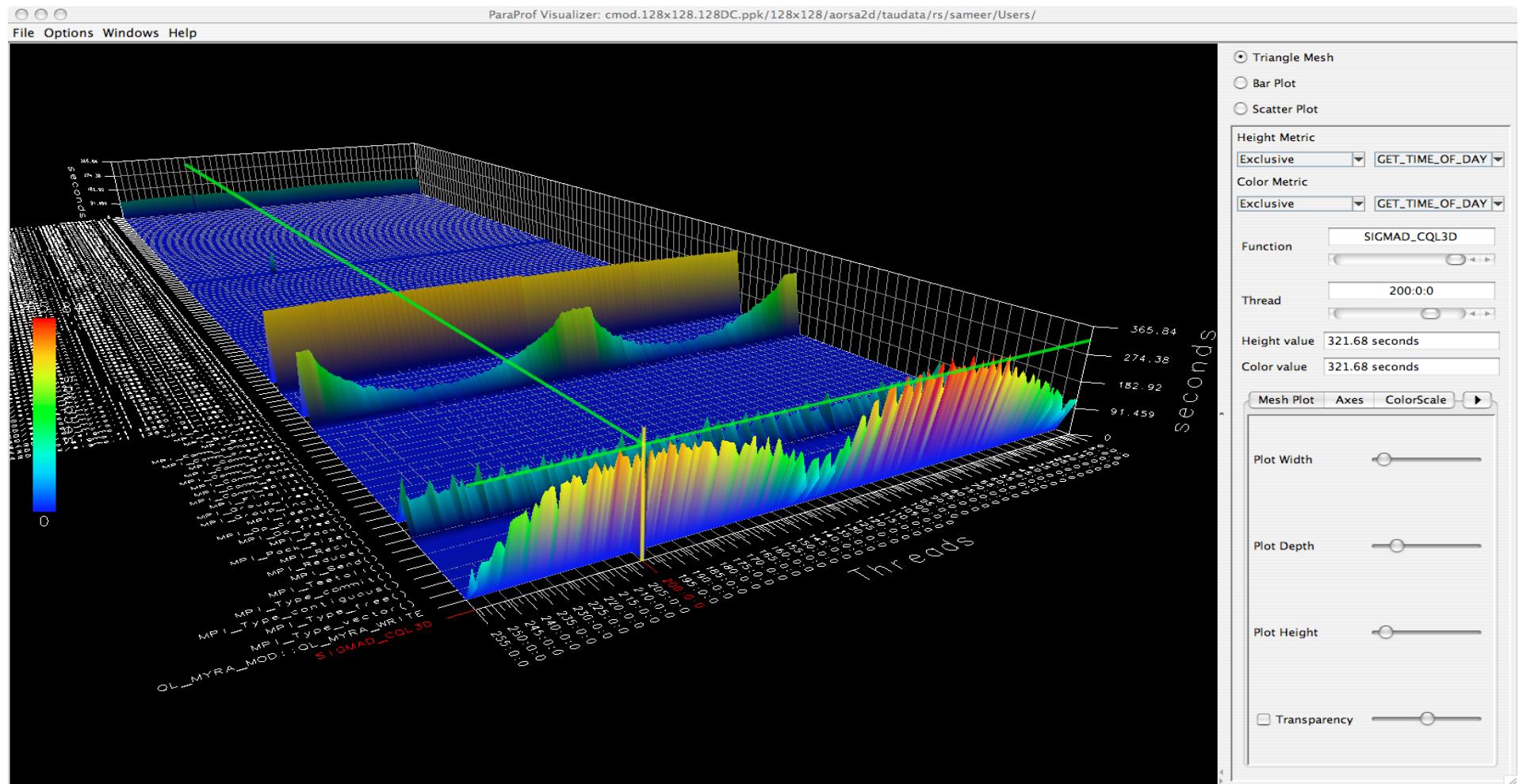
# Callsite Profiling and Tracing



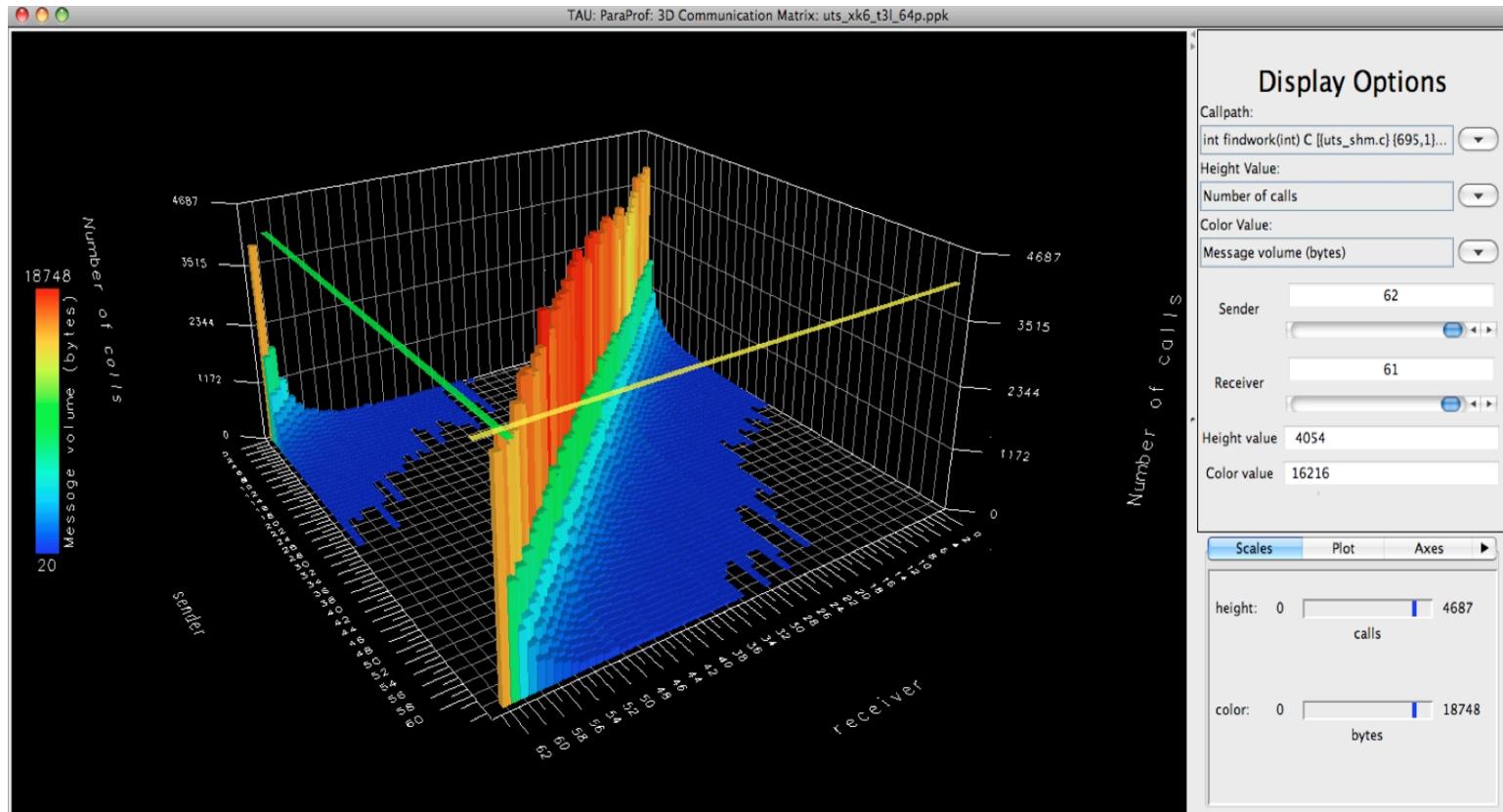
# Callsite Profiling and Tracing



# Parallel Profile Visualization: ParaProf

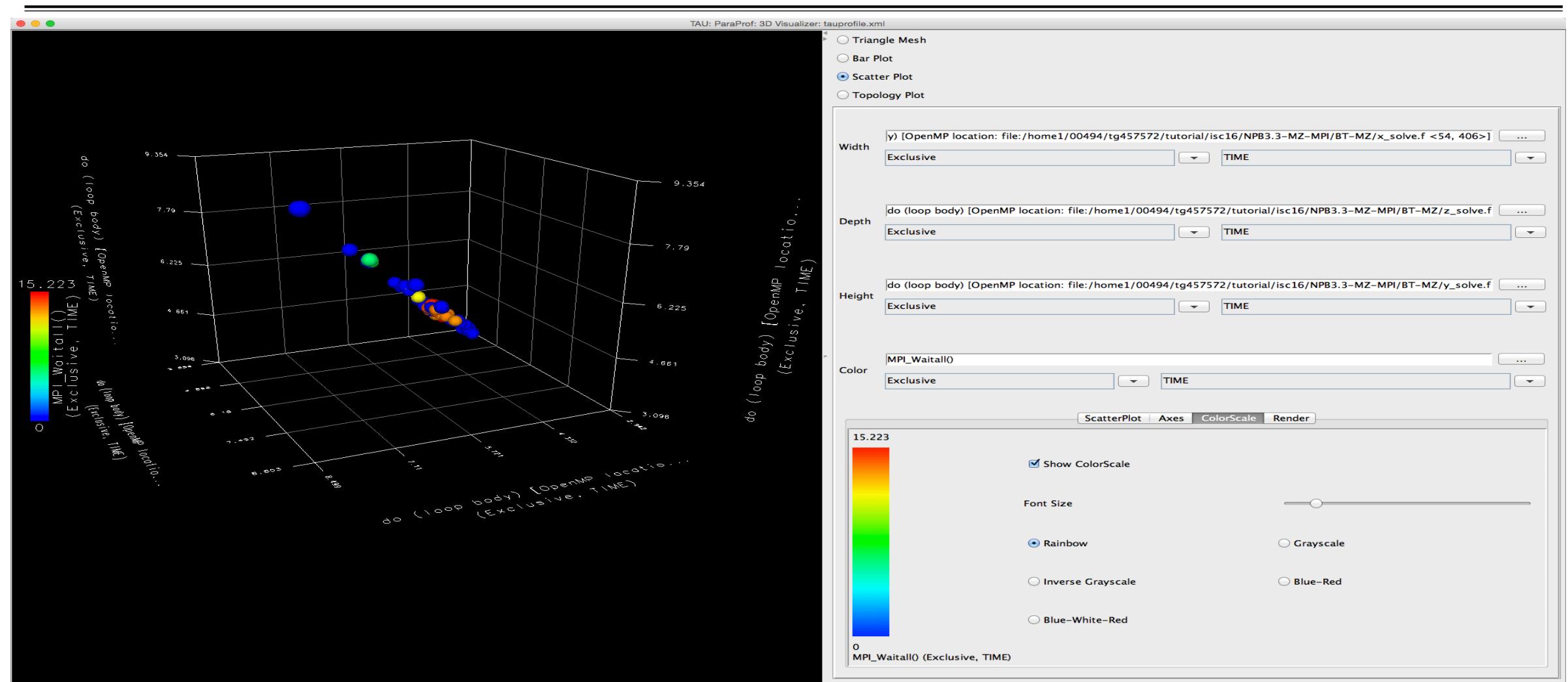


# ParaProf 3D Communication Matrix

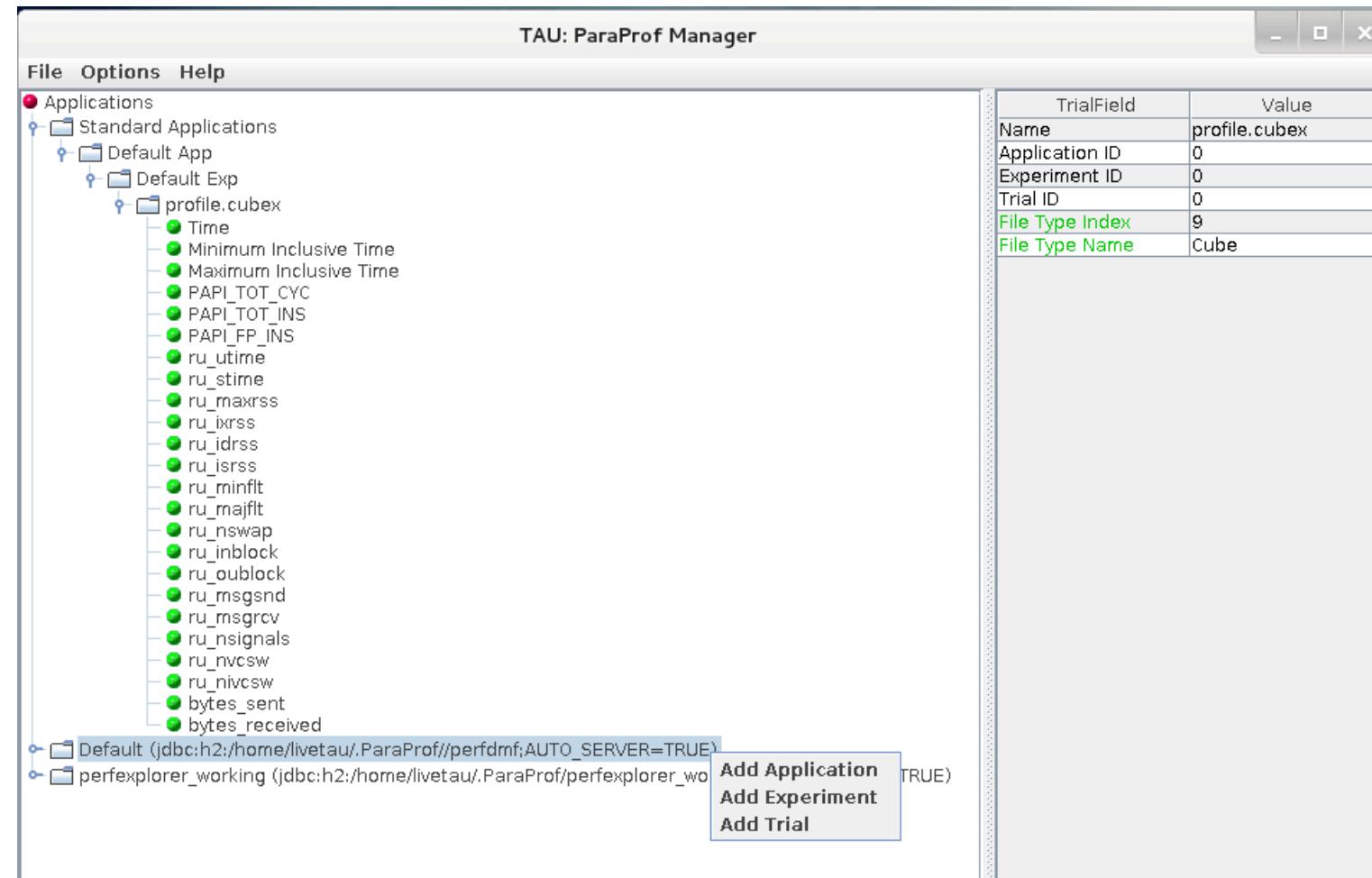


% export TAU\_COMM\_MATRIX=1

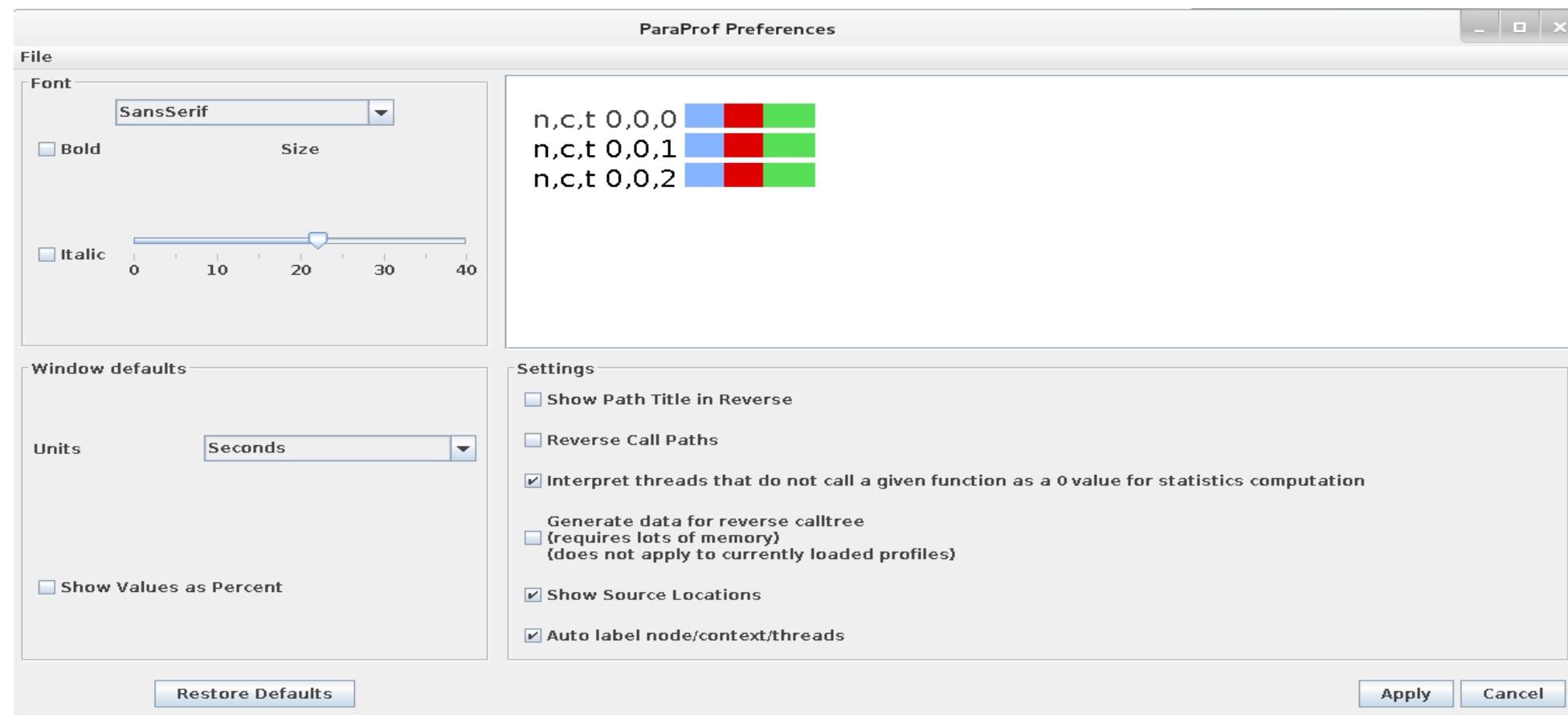
# ParaProf: 3D Scatter Plot



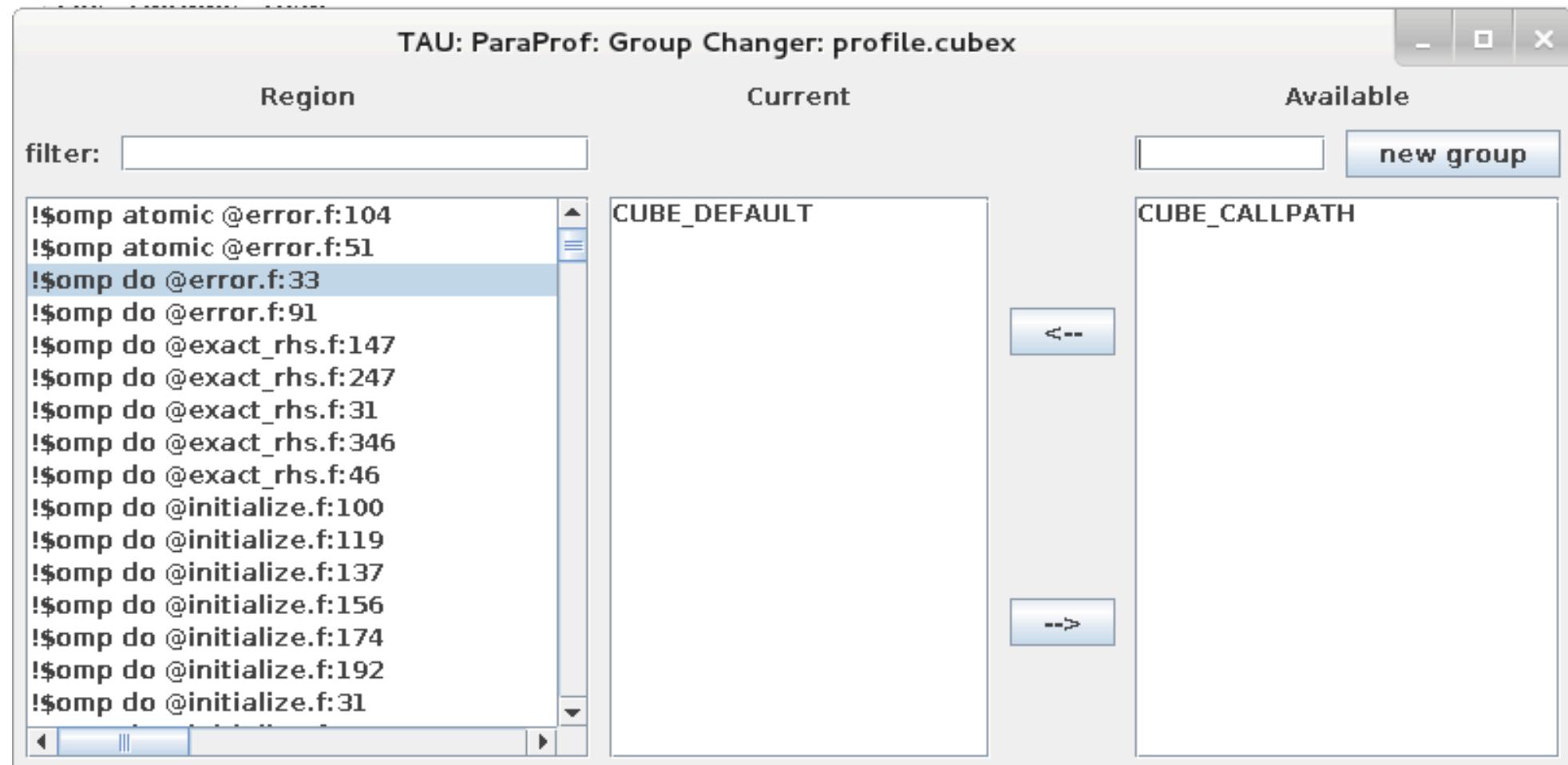
# ParaProf: Score-P Profile Files, Database



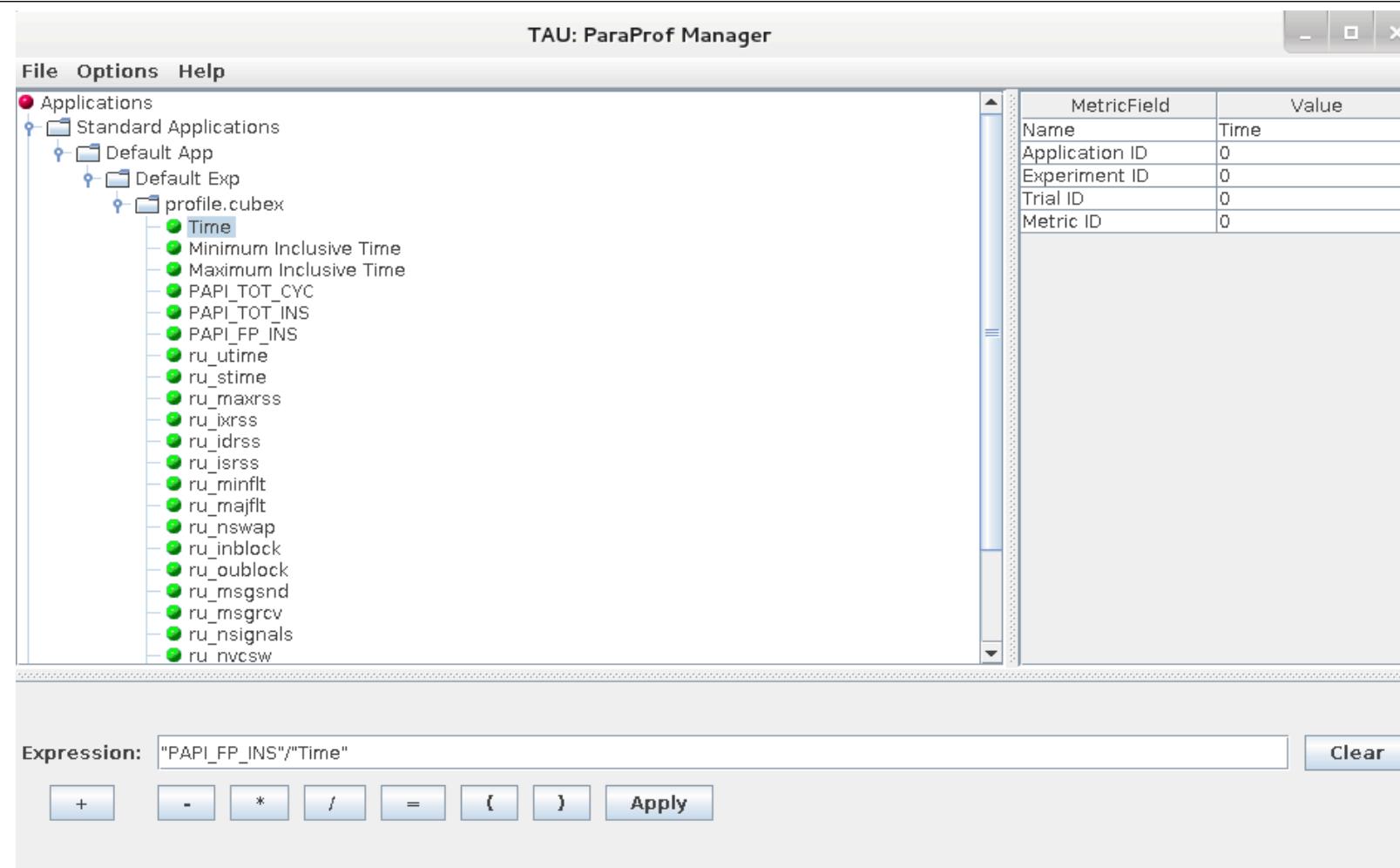
# ParaProf: File Preferences Window



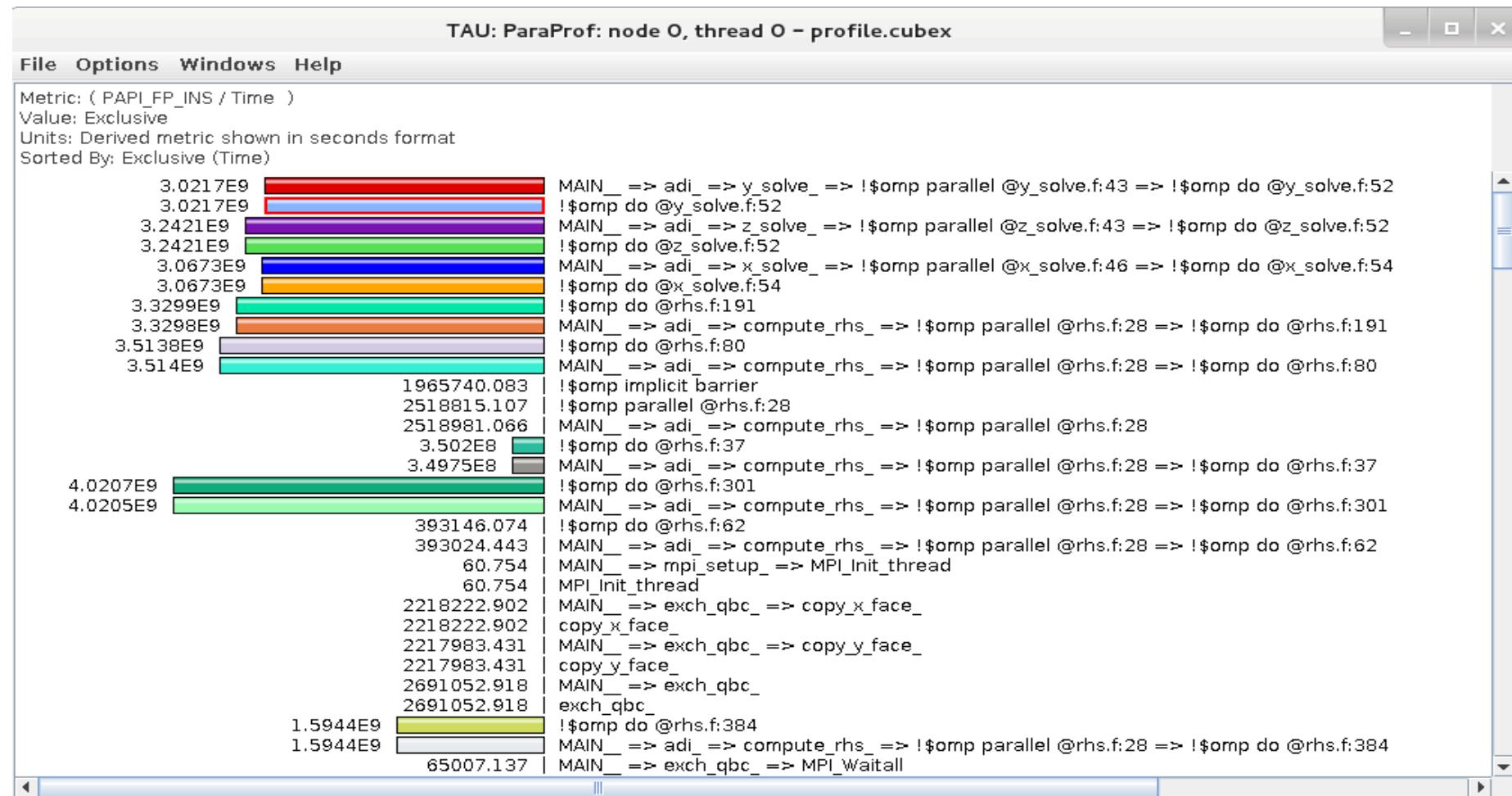
# ParaProf: Group Changer Window



# ParaProf: Derived Metric Panel in Manager Window



# Sorting Derived FLOPS metric by Exclusive Time



# Performance Research Lab, University of Oregon, Eugene, USA



# Support Acknowledgments

- US Department of Energy (DOE)
  - Office of Science contracts
  - SciDAC, LBL contracts
  - LLNL-LANL-SNL ASC/NNSA contract
  - Battelle, PNNL contract
  - ANL, ORNL contract
- Department of Defense (DoD)
  - PETTT, HPCMP
- National Science Foundation (NSF)
  - Glassbox, SI-2
  - NASA
  - CEA, France
- Partners:
  - University of Oregon
  - ParaTools, Inc., ParaTools, SAS
  - The Ohio State University
  - University of Tennessee, Knoxville
  - T.U. Dresden, GWT
  - Juelich Supercomputing Center

ParaTools



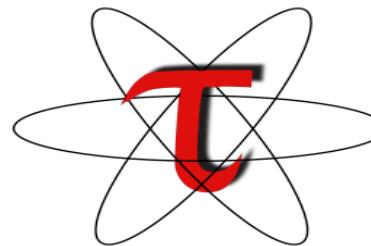


## Acknowledgement

This research was supported by the Exascale Computing Project (17-SC-20-SC), a collaborative effort of two U.S. Department of Energy organizations (Office of Science and the National Nuclear Security Administration) responsible for the planning and preparation of a capable exascale ecosystem, including software, applications, hardware, advanced system engineering, and early testbed platforms, in support of the nation's exascale computing imperative.

## Download TAU from U. Oregon

---



<http://tau.uoregon.edu>

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

<https://e4s.io> [Containers for Extreme-Scale Scientific Software Stack]

**Free download, open source, BSD license**