# Parallel Performance Optimization Using TAU Training and Workshop

**Sameer Shende**

*ParaTools, Inc. and University of Oregon.*

*sameer@paratools.com*

*http://www.paratools.com/TAU*

**Download slides and workshop tarball from:**

# http://tau.uoregon.edu/tau_llnl19.pdf
**Quartz.llnl.gov (or Lassen): /usr/global/tools/tau/training/workshop.tgz**

**Serrano.sandia.gov: /projects/tau/workshop.tgz**

**LANL: /turquoise/usr/projects/packages/tau/workshop.tgz**

**ParaTools**

UNIVERSITY OF OREGON

# TAU: Quickstart Guide

**Setup:**

- `% module load tau (Lasssen, LLNL)`

- `% source /usr/global/tools/tau/training/tau.bashrc (Quartz,LLNL)`

**Profiling:**

`MPI: % mpirun –np 16 tau_exec -ebs ./a.out`

- `MPI+OpenMP: % export TAU_OMPT_SUPPORT_LEVEL=full;`
  `% mpirun –np 16 tau_exec –T ompt,v5 –ompt ./a.out`

- `Pthread:   % mpirun –np 16 tau_exec –T mpi,pthread –ebs ./a.out`

- `Python+MPI+EBS % mpirun –np 16 tau_python –ebs ./a.py`

`Analysis: % pprof –a –m | more;  % paraprof (GUI)`

**Tracing:**

- `Vampir: MPI: % export TAU_TRACE=1; export TAU_TRACE_FORMAT=otf2`
  `% mpirun –np 16 tau_exec ./a.out; vampir traces.otf2 &`

- `Chrome: % export TAU_TRACE=1; mpirun –np 64 tau_exec ./a.out`

  `% tau_treemerge.pl;`

`% tau_trace2json tau.trc tau.edf –chrome –ignoreatomic –o app.json`

  `Chrome browser: chrome://tracing    (Load -> app.json)`

## ParaTools

UNIVERSITY OF OREGON

# Outline

**Day 1:**

**Morning:**

- Introduction to TAU
- Instrumentation: PDT, MPI, OpenMP OMPT, tau_exec
- Hands-on workshop examples using paraprof
- PAPI
- Hands-on using loop level instrumentation, PAPI

**Afternoon**
- Demonstration of analysis tools: Paraprof, TAUdb and PerfExplorer
- Vampir and Jumpshot
- Hands-on

**Breaks**
- 10:30am – 10:45am, noon – 1:30pm, 3pm – 3:15pm

**Day 2:**
**Individual Sessions:.**
- Applying performance evaluation tools to user codes

# TAU's Runtime Environment Variables

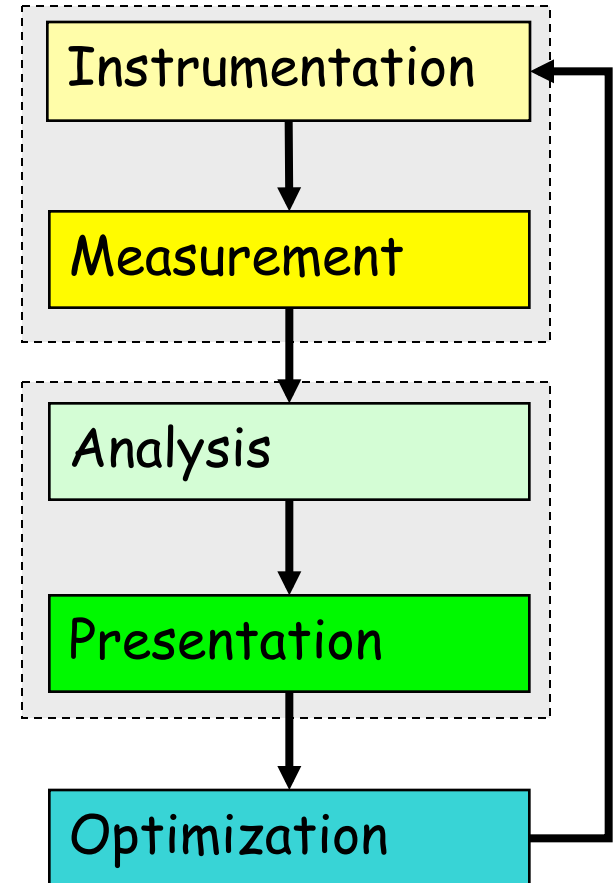| Environment Variable | Default | Description |
|---|---|---|
| TAU_TRACE | 0 | Setting to 1 turns on tracing |
| TAU_CALLPATH | 0 | Setting to 1 turns on callpath profiling |
| TAU_TRACK_MEMORY_FOOTPRINT | 0 | Setting to 1 turns on tracking memory usage by sampling periodically the resident set size and high water mark of memory usage |
| TAU_TRACK_LOAD | 0 | Tracks system load on a node (e.g., also seen in tools like w, top, uptime) |
| 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 (same as tau_exec –ebs). |
| TAU_EBS_RESOLUTION | Line | Setting to file \| function \| line will resolve addresses at the given resolution. |
| TAU_EBS_UNWIND | 0 | Setting to 1 enables callstack unwinding during 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_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>) |

ParaTools

UNIVERSITY OF OREGON

# Tutorial Goals

**This tutorial is an introduction to portable performance evaluation tools. You should leave here with a better understanding of…**

- Concepts and steps involved in performance evaluation
- Understanding key concepts in understanding code performance
- How to collect and analyze data from hardware performance counters (PAPI)
- How to instrument your programs with TAU
- Measurement options provided by TAU
- Environment variables used for choosing metrics, generating performance data
- How to use ParaProf, TAU's profile browser
- General familiarity with TAU use for Fortran, C++, C, and mixed language
- How to generate trace data in different formats

# Performance Optimization Cycle

- **Expose factors**
- **Collect performance data**
- **Calculate metrics**
- **Analyze results**
- **Visualize results**
- **Identify problems**
- **Tune performance**

ParaTools

UNIVERSITY OF OREGON

# TAU Performance System®

**http://tau.uoregon.edu**

- **Tuning and Analysis Utilities (20+ year project)**
- **Comprehensive performance profiling and tracing**
  - Integrated, scalable, flexible, portable
  - Targets all parallel programming/execution paradigms

- **Integrated performance toolkit**
  - Instrumentation, measurement, analysis, visualization
  - Widely-ported performance profiling / tracing system
  - Performance data management and data mining
  - Open source (BSD-style license)

- **Integrates with application frameworks**

ParaTools

UNIVERSITY OF OREGON

# Direct Performance Observation

**Execution *actions* exposed as *events***

- In general, actions reflect some execution state
  - presence at a code location or change in data
  - occurrence in parallelism context (thread of execution)
- Events encode actions for observation

**Observation is *direct***

- Direct instrumentation of program code (probes)
- Instrumentation invokes performance measurement
- Event measurement = performance data + context

**Performance experiment**

- Actual events + performance measurements

# Indirect Performance Observation

**Program code instrumentation is not used**

**Performance is observed indirectly**

- Execution is interrupted
  - can be triggered by different events
- Execution state is queried (sampled)
  - different performance data measured
- *Event-based sampling* (EBS)

**Performance attribution is inferred**

- Determined by execution context (state)
- Observation resolution determined by interrupt period
- Performance data associated with context for period

**ParaTools**

UNIVERSITY OF OREGON

# Understanding Application Performance using TAU

- **How much time** is spent in each application routine and outer *loops*? Within loops, what is the contribution of each *statement*?

- **How many instructions** are executed in these code regions? Floating point, Level 1 and 2 *data cache misses*, hits, branches taken?

- **What is the memory usage** of the code? When and where is memory allocated/de-allocated? Are there any memory leaks?

- **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?

# Profiling and Tracing

## Profiling



## Tracing



- Profiling and tracing

  **Profiling** shows you **how much** (total) time was spent in each routine

  **Tracing** shows you **when** the events take place on a timeline

ParaTools

UNIVERSITY OF OREGON

# Inclusive vs. Exclusive values

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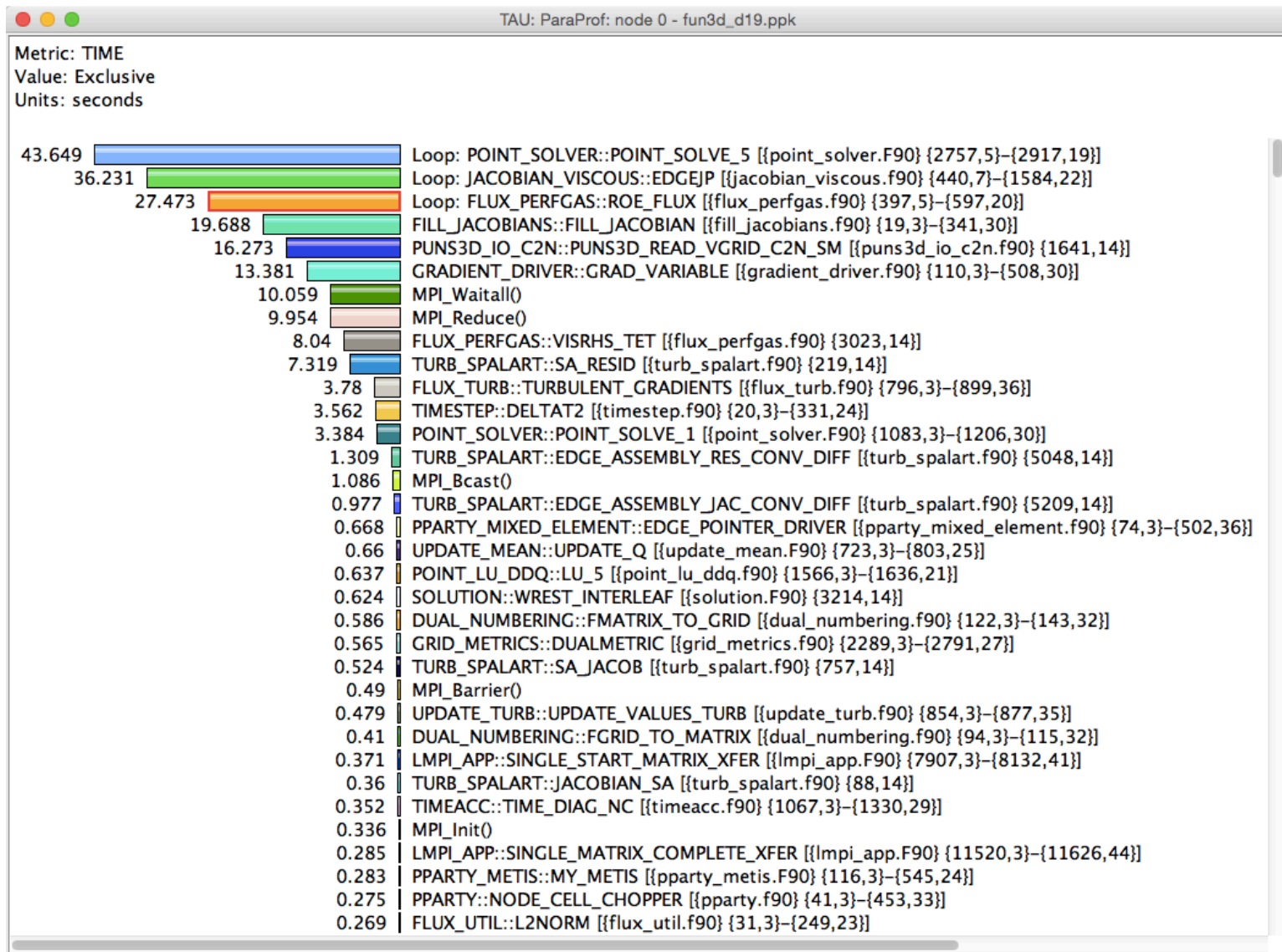    - Information of all sub-elements aggregated into single value
- **Exclusive**
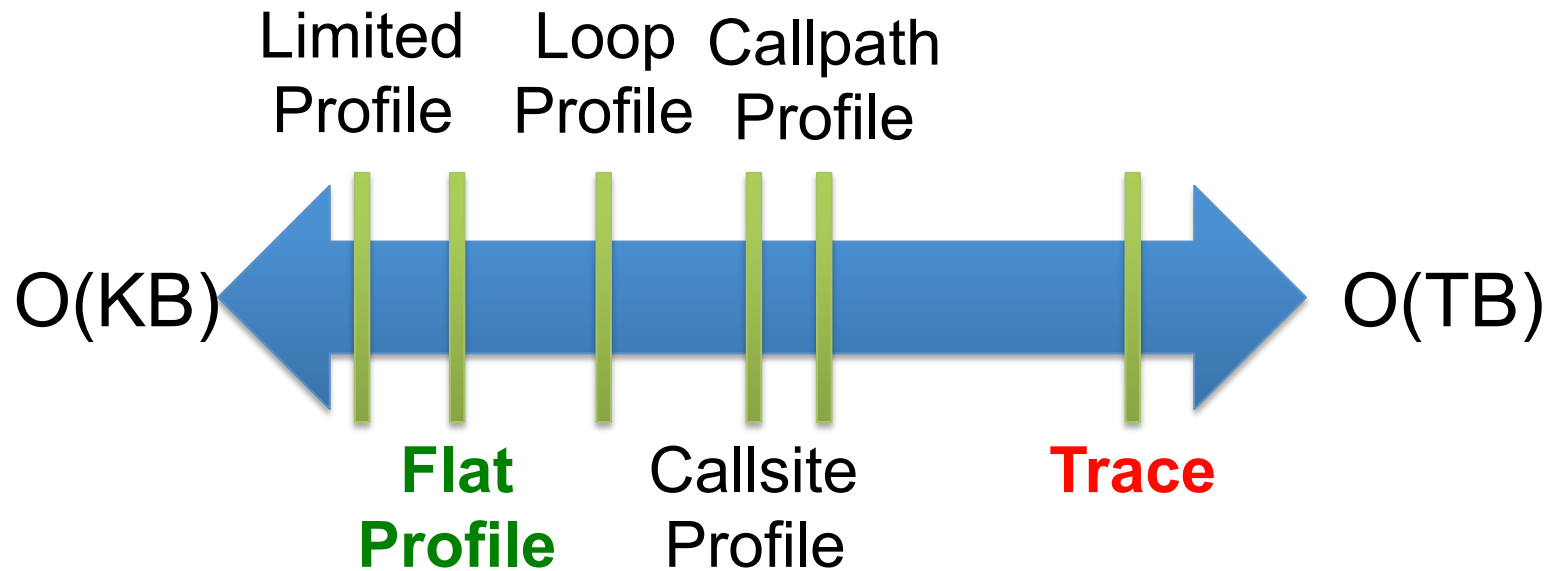    - Information cannot be subdivided further



```
int foo()
{
    int a;
    a = 1 + 1;

    bar();

    a = a + 1;
    return a;
}
```

Inclusive

Exclusive

**ParaTools**

UNIVERSITY OF OREGON

# Inclusive Measurements



TAU: ParaProf: node 0 - fun3d_d19.ppk

Metric: TIME
Value: Inclusive
Units: seconds

| Value | Function |
|---|---|
| 221.305 | .TAU application |
| 221.304 | NODET [{main.f90} {4,1}–{35,17}] |
| 197.989 | FLOW::ITERATE [{flow.F90} {1692,14}] |
| 195.577 | FLOW::STEP_SOLVER [{flow.F90} {1845,14}] |
| 195.569 | RELAX_STEADY::RELAX [{relax_steady.f90} {30,3}–{307,22}] |
| 61.927 | UPDATE_MEAN::UPDATE_LINEAR_SYSTEM_MEAN [{update_mean.F90} {195,3}–{275,42}] |
| 61.28 | UPDATE_MEAN::UPDATE_JACOBIAN_DRIVER_MEAN [{update_mean.F90} {460,3}–{505,44}] |
| 61.275 | UPDATE_MEAN::UPDATE_JACOBIAN [{update_mean.F90} {513,3}–{588,32}] |
| 61.258 | FILL_JACOBIANS::FILL_JACOBIAN [{fill_jacobians.f90} {19,3}–{341,30}] |
| 59.068 | GCR_SOLVE::GCR_SOLVER_QSET [{gcr_solve.f90} {47,3}–{415,32}] |
| 57.635 | GCR_SOLVE_UTIL::GCR_PRECONDITIONER_QSET [{gcr_solve_util.f90} {40,3}–{131,40}] |
| 57.152 | POINT_SOLVER::POINT_SOLVE [{point_solver.F90} {31,3}–{214,28}] |
| 56.882 | UPDATE_MEAN::UPDATE_RHS_MEAN [{update_mean.F90} {102,3}–{185,32}] |
| 54.402 | RELAX_MEAN::RELAX [{relax_mean.f90} {22,3}–{84,22}] |
| 53.103 | LINEARSOLVE_NODIVCHECK::NODIVCHECK_RELAX_Q [{linearsolve_nodivcheck.F90} {56,14}] |
| 52.867 | UPDATE_MEAN::RESIDUAL_S [{update_mean.F90} {42,3}–{94,27}] |
| 52.866 | FUN3D_RES_FLOW::RES_FLOW [{fun3d_res_flow.f90} {27,3}–{279,25}] |
| 52.756 | FLUX::RESIDUAL_COMPRESSIBLE [{flux.f90} {25,3}–{592,38}] |
| 52.747 | POINT_SOLVER::POINT_SOLVE_5 [{point_solver.F90} {2700,3}–{2921,30}] |
| 52.744 | Loop: POINT_SOLVER::POINT_SOLVE_5 [{point_solver.F90} {2757,5}–{2917,19}] |
| 36.232 | JACOBIAN_VISCOUS::VISCOUS_JACOBIAN [{jacobian_viscous.f90} {20,14}] |
| 36.231 | JACOBIAN_VISCOUS::EDGEJP [{jacobian_viscous.f90} {324,14}] |
| 36.231 | Loop: JACOBIAN_VISCOUS::EDGEJP [{jacobian_viscous.f90} {440,7}–{1584,22}] |
| 27.474 | FLUX_PERFGAS::INVISCID_FLUX_DRIVER [{flux_perfgas.f90} {37,14}] |
| 27.474 | FLUX_PERFGAS::ROE_FLUX [{flux_perfgas.f90} {236,14}] |
| 27.473 | Loop: FLUX_PERFGAS::ROE_FLUX [{flux_perfgas.f90} {397,5}–{597,20}] |
| 22.707 | FLOW::INITIALIZE_DATA [{flow.F90} {465,14}] |
| 22.694 | FLOW::INITIALIZE_DATA2 [{flow.F90} {663,14}] |
| 20.916 | PPARTY_PREPROCESSOR::PPARTY_PREPROCESS [{pparty_preprocessor.f90} {28,14}] |
| 16.726 | PPARTY_PREPROCESSOR::PPARTY_READ_GRID [{pparty_preprocessor.f90} {735,14}] |
| 16.726 | PUNS3D_IO_C2N::PUNS3D_READ_VGRID_C2N [{puns3d_io_c2n.f90} {1543,14}] |
| 16.657 | PUNS3D_IO_C2N::PUNS3D_READ_VGRID_C2N_SM [{puns3d_io_c2n.f90} {1641,14}] |
| 14.159 | GRADIENT_DRIVER::GRAD_VARIABLE [{gradient_driver.f90} {110,3}–{508,30}] |
| 13.852 | UPDATE_TURB::UPDATE_RHS_TURB [{update_turb.f90} {742,3}–{845,32}] |

ParaTools

UNIVERSITY OF OREGON

# Exclusive Time



TAU: ParaProf: node 0 - fun3d_d19.ppk

Metric: TIME
Value: Exclusive
Units: seconds

| Time | Function |
|------|----------|
| 43.649 | Loop: POINT_SOLVER::POINT_SOLVE_5 [{point_solver.F90} {2757,5}–{2917,19}] |
| 36.231 | Loop: JACOBIAN_VISCOUS::EDGEJP [{jacobian_viscous.f90} {440,7}–{1584,22}] |
| 27.473 | Loop: FLUX_PERFGAS::ROE_FLUX [{flux_perfgas.f90} {397,5}–{597,20}] |
| 19.688 | FILL_JACOBIANS::FILL_JACOBIAN [{fill_jacobians.f90} {19,3}–{341,30}] |
| 16.273 | PUNS3D_IO_C2N::PUNS3D_READ_VGRID_C2N_SM [{puns3d_io_c2n.f90} {1641,14}] |
| 13.381 | GRADIENT_DRIVER::GRAD_VARIABLE [{gradient_driver.f90} {110,3}–{508,30}] |
| 10.059 | MPI_Waitall() |
| 9.954 | MPI_Reduce() |
| 8.04 | FLUX_PERFGAS::VISRHS_TET [{flux_perfgas.f90} {3023,14}] |
| 7.319 | TURB_SPALART::SA_RESID [{turb_spalart.f90} {219,14}] |
| 3.78 | FLUX_TURB::TURBULENT_GRADIENTS [{flux_turb.f90} {796,3}–{899,36}] |
| 3.562 | TIMESTEP::DELTAT2 [{timestep.f90} {20,3}–{331,24}] |
| 3.384 | POINT_SOLVER::POINT_SOLVE_1 [{point_solver.F90} {1083,3}–{1206,30}] |
| 1.309 | TURB_SPALART::EDGE_ASSEMBLY_RES_CONV_DIFF [{turb_spalart.f90} {5048,14}] |
| 1.086 | MPI_Bcast() |
| 0.977 | TURB_SPALART::EDGE_ASSEMBLY_JAC_CONV_DIFF [{turb_spalart.f90} {5209,14}] |
| 0.668 | PPARTY_MIXED_ELEMENT::EDGE_POINTER_DRIVER [{pparty_mixed_element.f90} {74,3}–{502,36}] |
| 0.66 | UPDATE_MEAN::UPDATE_Q [{update_mean.F90} {723,3}–{803,25}] |
| 0.637 | POINT_LU_DDQ::LU_5 [{point_lu_ddq.f90} {1566,3}–{1636,21}] |
| 0.624 | SOLUTION::WREST_INTERLEAF [{solution.F90} {3214,14}] |
| 0.586 | DUAL_NUMBERING::FMATRIX_TO_GRID [{dual_numbering.f90} {122,3}–{143,32}] |
| 0.565 | GRID_METRICS::DUALMETRIC [{grid_metrics.f90} {2289,3}–{2791,27}] |
| 0.524 | TURB_SPALART::SA_JACOB [{turb_spalart.f90} {757,14}] |
| 0.49 | MPI_Barrier() |
| 0.479 | UPDATE_TURB::UPDATE_VALUES_TURB [{update_turb.f90} {854,3}–{877,35}] |
| 0.41 | DUAL_NUMBERING::FGRID_TO_MATRIX [{dual_numbering.f90} {94,3}–{115,32}] |
| 0.371 | LMPI_APP::SINGLE_START_MATRIX_XFER [{lmpi_app.F90} {7907,3}–{8132,41}] |
| 0.36 | TURB_SPALART::JACOBIAN_SA [{turb_spalart.f90} {88,14}] |
| 0.352 | TIMEACC::TIME_DIAG_NC [{timeacc.f90} {1067,3}–{1330,29}] |
| 0.336 | MPI_Init() |
| 0.285 | LMPI_APP::SINGLE_MATRIX_COMPLETE_XFER [{lmpi_app.F90} {11520,3}–{11626,44}] |
| 0.283 | PPARTY_METIS::MY_METIS [{pparty_metis.F90} {116,3}–{545,24}] |
| 0.275 | PPARTY::NODE_CELL_CHOPPER [{pparty.f90} {41,3}–{453,33}] |
| 0.269 | FLUX_UTIL::L2NORM [{flux_util.f90} {31,3}–{249,23}] |

ParaTools

UNIVERSITY OF OREGON

# How much data do you want?

Limited Profile

Loop Profile

Callpath Profile

O(KB) ⟵⟶ O(TB)

**Flat Profile**

Callsite Profile

**Trace**

# Performance Data Measurement

## Direct via Probes

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

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

## Indirect via Sampling



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

ParaTools

UNIVERSITY OF OREGON

# Sampling



Time

| main | foo(0) | foo(1) | foo(2) | Measurement |

**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**

```
int main()
{
  int i;

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

  return 0;
}

void foo(int i)
{

  if (i > 0)
    foo(i – 1);

}
```

ParaTools

UNIVERSITY OF OREGON

# Instrumentation



| main | foo(0) | foo(1) | foo(2) | Measurement |
|------|--------|--------|--------|-------------|

**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;
  Start("main");
  for (i=0; i < 3; i++)
    foo(i);
  Stop("main");
  return 0;
}

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

ParaTools

UNIVERSITY OF OREGON

# TAU Architecture and Workflow

## TAU Architecture

### Instrumentation

**Source**
- C, C++, Fortran
- Python, UPC, Java
- Robust parsers (PDT)

**Wrapping**
- Interposition (PMPI)
- Wrapper generation

**Linking**
- Static, dynamic
- Preloading

**Executable**
- Dynamic (Dyninst)
- Binary (Dyninst, MAQAO)

*Measurement API*

### Measurement

**Events**
- static/dynamic
- routine, basic block, loop
- threading, communication
- heterogeneous

**Profiling**
- flat, callpath, phase, parameter, snapshot
- probe, sampling, hybrid

**Tracing**
- TAU / Scalasca tracing
- Open Trace Format (OTF)

**Metadata**
- system, user-defined

*Measured data*

### Analysis

**Profiles**
- *ParaProf* parallel profile analyzer / visualizer
- *PerfDMF* parallel profile database
- *PerfExplorer* parallel profile data mining

**Tracing**
- TAU trace translation
  - OTF, SLOG-2
- Trace analysis / visualizer
  - *Vampir, Jumpshot*

**Online**
- event unification
- statistics calculation

# TAU's Support for Runtime Systems

*MPI*

- PMPI profiling interface
- MPI_T tools interface using performance and control variables

*Pthread*

- Captures time spent in routines per thread of execution

*OpenMP*

- OMPT tools interface to track salient OpenMP runtime events
- Opari source rewriter
- Preloading wrapper OpenMP runtime library when OMPT is not supported

*OpenACC*

- OpenACC instrumentation API
- Track data transfers between host and device (per-variable)
- Track time spent in kernels

# TAU's Support for Runtime Systems (contd.)

*OpenCL*

- OpenCL profiling interface
- Track timings of kernels

*CUDA*

- Cuda Profiling Tools Interface (CUPTI)
- Track data transfers between host and GPU
- Track access to uniform shared memory between host and GPU

*ROCm*

- Rocprofiler and Roctracer instrumentation interfaces
- Track data transfers and kernel execution between host and GPU

*Kokkos*

- Kokkos profiling API
- Push/pop interface for region, kernel execution interface

*Python*

- Python interpreter instrumentation API
- Tracks Python routine transitions as well as Python to C transitions

**ParaTools**

UNIVERSITY OF OREGON

# Examples of Multi-Level Instrumentation

***MPI + OpenMP***

- MPI_T + PMPI + OMPT may be used to track MPI and OpenMP

***MPI + CUDA***

- PMPI + CUPTI interfaces

***OpenCL + ROCm***

- Rocprofiler + OpenCL instrumentation interfaces

***Kokkos + OpenMP***

- Kokkos profiling API + OMPT to transparently track events

***Kokkos + pthread + MPI***

- Kokkos + pthread wrapper interposition library + PMPI layer

***Python + CUDA***

- Python + CUPTI + pthread profiling interfaces (e.g., Tensorflow, PyTorch)

***MPI + OpenCL***

- PMPI + OpenCL profiling interfaces

**ParaTools**

UNIVERSITY OF OREGON

# Simplifying the use of TAU!

## Uninstrumented code:

- % make

- % mpirun –np 64 ./a.out

## With TAU using event based sampling (EBS):

- % mpirun –np 64 tau_exec <span style="color:red">–ebs</span> ./lu.B.64

- % paraprof          (GUI)

- % pprof –a | more

**NOTE:**

- Requires dynamic executables (-dynamic link flag on Cray XC systems).

- Source code should be compiled with –g for access to symbol table.

- Replace mpirun with srun on Quartz, LLNL or your appropriate launch command.

ParaTools

UNIVERSITY OF OREGON

# TAU Execution Command (tau_exec)

**Uninstrumented execution**

- % mpirun -np 256  ./a.out

**Track GPU operations**

- % mpirun –np 256  tau_exec –rocm ./a.out
- % mpirun –np 256  tau_exec –cupti ./a.out
- % mpirun –np 256  tau_exec –cupti  -um ./a.out  (for Unified Memory)
- % mpirun –np 256 tau_exec –opencl ./a.out
- % mpirun –np 256 tau_exec –openacc ./a.out

**Track MPI performance**

- % mpirun -np 256  tau_exec ./a.out

**Track I/O, and MPI performance (MPI enabled by default)**

- % mpirun -np 256  tau_exec -io  ./a.out

**Track OpenMP and MPI execution (using OMPT for Intel v19)**

- % export TAU_OMPT_SUPPORT_LEVEL=full;
  % export TAU_OMPT_RESOLVE_ADDRESS_EAGERLY=1
- % mpirun –np 256  tau_exec –T ompt,v5,mpi  -ompt  ./a.out

**Track memory operations**

- % export TAU_TRACK_MEMORY_LEAKS=1
- % mpirun –np 256 tau_exec –memory_debug ./a.out (bounds check)

**Use event based sampling (compile with –g)**

- % mpirun –np 256 tau_exec –ebs ./a.out
- Also  -ebs_source=<PAPI_COUNTER> -ebs_period=<overflow_count>
       -ebs_resolution=<file | function | line>

## ParaTools

# 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

ParaTools

UNIVERSITY OF OREGON

# ParaProf Profile Browser



Click "node X" next to see details

# Python Instrumentation

| Name | Exclusive TAUGPU_TIME | Inclusive TAUGPU_TIME ▽ | Calls | Child Calls |
|---|---|---|---|---|
| ▼ ■ .TAU application | 0.575 | 182.783 | 1 | 6 |
| ▼ ■ <module> [{micro_benchmarking_pytorch.py}{1}] | 0.002 | 182.151 | 1 | 13 |
| ▼ ■ main [{micro_benchmarking_pytorch.py}{81}] | 0.002 | 168.702 | 1 | 1 |
| ▼ ■ run_benchmarking [{micro_benchmarking_pytorch.py}{40}] | 0.006 | 168.7 | 1 | 40 |
| ▼ ■ forwardbackward [{micro_benchmarking_pytorch.py}{33}] | 0.002 | 155.924 | 22 | 110 |
| ▼ ■ backward [{tensor.py}{79}] | 0.001 | 106.141 | 22 | 22 |
| ▼ ■ backward [{__init__.py}{38}] | 0.001 | 106.14 | 22 | 88 |
| ▼ ■ run_backward | 106.135 | 106.135 | 22 | 3 |
| ■ pthread_create | 0 | 0 | 3 | 0 |
| ▸ ■ _make_grads [{__init__.py}{20}] | 0.001 | 0.004 | 22 | 110 |
| ■ isinstance | 0 | 0 | 22 | 0 |
| ■ len | 0 | 0 | 22 | 0 |
| ▼ ■ __call__ [{module.py}{485}] | 0 | 49.77 | 22 | 110 |
| ▼ ■ forward [{container.py}{95}] | 0 | 49.768 | 22 | 66 |
| ▼ ■ __call__ [{module.py}{485}] | 0.001 | 49.767 | 44 | 220 |
| ▼ ■ forward [{resnet.py}{151}] | 0.003 | 49.765 | 22 | 484 |
| ▼ ■ __call__ [{module.py}{485}] | 0.006 | 49.759 | 220 | 1,100 |
| ▼ ■ forward [{container.py}{95}] | 0.002 | 45.622 | 88 | 440 |
| ▼ ■ __call__ [{module.py}{485}] | 0.007 | 45.616 | 352 | 1,760 |
| ▼ ■ forward [{resnet.py}{78}] | 0.071 | 45.598 | 352 | 6,600 |
| ▼ ■ __call__ [{module.py}{485}] | 0.07 | 45.495 | 3,256 | 16,280 |
| ▼ ■ forward [{conv.py}{319}] | 0.017 | 29.675 | 1,056 | 3,168 |
| ■ conv2d | 29.648 | 29.648 | 1,056 | 0 |
| ■ __getattr__ [{module.py}{523}] | 0.01 | 0.01 | 2,112 | 0 |
| ▸ ■ forward [{container.py}{95}] | 0.002 | 9.401 | 88 | 264 |
| ▸ ■ forward [{batchnorm.py}{59}] | 0.262 | 6.097 | 1,056 | 9,504 |

% tau_python ./foo.py

ParaTools

UNIVERSITY OF OREGON

# Identifying Wait States Using EBS



TAU: ParaProf: Statistics for: node 0, thread 0 – nt3_baseline_keras2.ppk

| Name | Inclusive ... | Calls ▽ |
|---|---|---|
| ■ _do_call [{session.py}{1348}] | 512.135 | 82 |
| ■ _run_fn [{session.py}{1317}] | 512.134 | 82 |
| ▼ ■ TF_Run | 512.093 | 82 |
| ▼ ■ [CONTEXT] TF_Run | 512.173 | 51,211 |
| ■ [SAMPLE] __pthread_cond_wait [{} {0}] | 511.273 | 51,123 |
| ■ [SAMPLE] tensorflow::TensorBuffer* tensorflow::(anonymous namespace)::FromProtoFiel( | 0.42 | 42 |
| ■ [SAMPLE] __memcpy_ssse3_back [{} {0}] | 0.28 | 28 |
| ■ [SAMPLE] _int_free [{malloc.c} {0}] | 0.03 | 3 |
| ■ [SAMPLE] __GI___libc_malloc [{} {0}] | 0.02 | 2 |
| ■ [SAMPLE] std::basic_string<char, std::char_traits<char>, std::allocator<char> >::assign(s | 0.02 | 1 |
| ■ [SAMPLE] google::protobuf::internal::MapField<tensorflow::NodeDef::NodeDef_AttrEntry, | 0.02 | 1 |
| ■ [SAMPLE] __exchange_and_add [{/home/msarahan/miniconda2/conda-bld/compilers_lin | 0.01 | 1 |
| ■ [SAMPLE] void google::protobuf::internal::RepeatedPtrFieldBase::MergeFromInnerLoop<g | 0.01 | 1 |
| ■ [SAMPLE] google::protobuf::internal::ArenaStringPtr::Destroy(std::basic_string<char, std:: | 0.01 | 1 |
| ■ [SAMPLE] std::_Hashtable<std::basic_string<char, std::char_traits<char>, std::allocator<( | 0.01 | 1 |
| ■ [SAMPLE] std::_Hashtable<tensorflow::Node*, std::pair<tensorflow::Node* const, tensorf | 0.01 | 1 |
| ■ [SAMPLE] std::basic_string<char, std::char_traits<char>, std::allocator<char> >::_Rep::_I | 0.01 | 1 |
| ■ [SAMPLE] std::_Hash_bytes(void const*, unsigned long, unsigned long) [{/home/msaraha | 0.01 | 1 |
| ■ [SAMPLE] std::basic_string<char, std::char_traits<char>, std::allocator<char> >::_M_cop | 0.01 | 1 |
| ■ [SAMPLE] std::basic_string<char, std::char_traits<char>, std::allocator<char> >::size() c( | 0.01 | 1 |
| ■ [SAMPLE] PyObject_Malloc [{/home/nwani/m2u/conda-bld/python_1500576437846/wor | 0.01 | 1 |
| ■ [SAMPLE] jemalloc_free [{/home/nchaimov/candle/anaconda3/lib/python3.6/site-packag | 0.01 | 1 |

% tau_python -ebs ./foo.py

**ParaTools**

UNIVERSITY OF OREGON

# Kokkos and OpenMP Instrumentation



% tau_exec −ompt ./a.out

ParaTools

UNIVERSITY OF OREGON

# Kokkos Instrumentation with OpenMP



% tau_exec −ompt ./a.out

# Vampir [TU Dresden] Timeline



% export TAU_TRACE_FORMAT=otf2
% tau_exec –ompt ./a.out
% vampir traces.otf2 &

# TAU − Flat Profile



TAU: ParaProf: node 102 - cesm.fideal.f09.n240.pdt.callPath.4_7_16.ppk

**Metric: TIME**
**Value: Exclusive**
**Units: seconds**

| Value | Function |
|---|---|
| 1.424 | CD_CORE [{cd_core.F90} {8,2}-{1591,28}] |
| 1.403 | TP_CORE::XTPV [{tp_core.F90} {294,2}-{509,20}] |
| 1.247 | PHYS_GRID::CREATE_CHUNKS [{phys_grid.F90} {3680,4}-{4265,31}] |
| 1.203 | TP_CORE::YMIST [{tp_core.F90} {1105,2}-{1241,21}] |
| 1.053 | TP_CORE::TP2D [{tp_core.F90} {166,2}-{285,20}] |
| 0.882 | TE_MAP [{te_map.F90} {7,4}-{1057,27}] |
| 0.874 | MAPZ_MODULE::PPM2M [{mapz_module.F90} {675,2}-{953,21}] |
| 0.83 | M_MERGESORTS::DSORT_::MERGESORT_ [{m_MergeSorts.F90} {548,1}-{569,25}] |
| 0.812 | TP_CORE::FYPPM [{tp_core.F90} {1249,2}-{1427,21}] |
| 0.748 | SW_CORE::D_SW [{sw_core.F90} {542,2}-{1416,20}] |
| 0.7 | TP_CORE::XMIST [{tp_core.F90} {517,2}-{568,21}] |
| 0.619 | PHYS_GRID::ASSIGN_CHUNKS [{phys_grid.F90} {4648,4}-{4826,31}] |
| 0.564 | MAPZ_MODULE::MAP1_PPM [{mapz_module.F90} {200,3}-{350,24}] |
| 0.556 | PHYS_GRID::FIND_TWIN [{phys_grid.F90} {4430,4}-{4644,27}] |
| 0.524 | GEOPK [{geopk.F90} {36,7}-{149,26}] |
| 0.52 | GHOSTMODULE::GHOSTREGULAR3D [{ghostmodule.F90} {634,7}-{808,35}] |
| 0.515 | SW_CORE::C_SW [{sw_core.F90} {63,2}-{532,20}] |
| 0.454 | MOD_COMM::GA_GET4D_R8 [{mod_comm.F90} {1827,7}-{1881,32}] |
| 0.449 | TP_CORE::FXPPM [{tp_core.F90} {576,2}-{697,21}] |
| 0.429 | MOD_COMM::GA_PUT4D_R8 [{mod_comm.F90} {1711,7}-{1776,32}] |
| 0.415 | VPASSM [{fft99.F90} {937,7}-{1300,27}] |
| 0.413 | M_MERGESORTS::ISORT_::MERGE_ [{m_MergeSorts.F90} {337,1}-{370,21}] |
| 0.397 | TP_CORE::LMPPM [{tp_core.F90} {779,2}-{889,21}] |
| 0.385 | MOD_COMM::MP_SENDIRR [{mod_comm.F90} {2495,7}-{2838,31}] |
| 0.382 | M_MERGESORTS::ISORT_::MERGESORT_ [{m_MergeSorts.F90} {314,1}-{335,25}] |
| 0.381 | MAPZ_MODULE::STEEPZ [{mapz_module.F90} {1209,2}-{1285,22}] |
| 0.37 | MCT_MOD::MCT_AVECT_VECMULT [{mct_mod.F90} {988,1}-{1120,32}] |
| 0.366 | SHR_CONST_MOD::SHR_CONST_ISSPVAL [{shr_const_mod.F90} {67,3}-{78,32}] |
| 0.34 | MOD_COMM::MP_RECVIRR [{mod_comm.F90} {2846,7}-{3040,31}] |
| 0.323 | PHYS_GRID::PHYS_GRID_INIT [{phys_grid.F90} {311,3}-{930,31}] |

ParaTools

http://tau.uoregon.edu/tau_llnl19.pdf

UNIVERSITY OF OREGON

# ParaProf Thread Statistics Table



Right click over "node X" and choose
Show Thread Statistics Table

# TAU − Callsite Profiling



% export TAU_CALLSITE=1

# Callsite Profiling and Tracing



% export TAU_CALLSITE=1

# Callsite Profiling and Tracing



TAU: ParaProf: 3D Visualizer: lu_callsite.ppk

# Callsite Profiling and Tracing

# TAU − Callstack Sampling



% export TAU_SAMPLING=1; export TAU_EBS_UNWIND=1

# TAU – Event Based Sampling (EBS)



% export TAU_SAMPLING=1

# TAU − Callpath Profiling



% export TAU_CALLPATH=1; export TAU_CALLPATH_DEPTH=100

# TAU Atomic Events

| Name ▽ | Total | NumSamples | MaxValue | MinValue | MeanValue | Std. Dev. |
|---|---|---|---|---|---|---|
| Bytes Written <file=stdout> | 911 | 62 | 21 | 1 | 14.694 | 7.441 |
| Bytes Written <file=pipe> | 22 | 22 | 1 | 1 | 1 | 0 |
| Bytes Written <file=Process_Output/VelRsdl.dat> | 7,826 | 100 | 302 | 76 | 78.26 | 22.487 |
| Bytes Written <file=Process_Output/MomRsdl.dat> | 7,826 | 100 | 302 | 76 | 78.26 | 22.487 |
| Bytes Written <file=Process_Output/MassRsdl.dat> | 11,325 | 100 | 435 | 110 | 113.25 | 32.337 |
| Bytes Written <file=Grid_Output/bodyBndry.dat> | 9,724 | 5 | 8,192 | 4 | 1,944.8 | 3,174.201 |
| Bytes Written <file=/home/sameer/apps/sukra/RotCFD_Regression/case_catalog/UNS2D/N/ | 45 | 1 | 45 | 45 | 45 | 0 |
| Bytes Written <file=./Restarts/Restart_History//NACA0012_LargeGrid_00010.Rst> | 44,619,720 | 5,484 | 8,192 | 4 | 8,136.346 | 640.325 |
| Bytes Written <file=./Restarts/Restart_History//NACA0012_LargeGrid_00005.Rst> | 44,619,720 | 5,484 | 8,192 | 4 | 8,136.346 | 640.325 |
| Bytes Written <file=./Restarts//NACA0012_LargeGrid.Rst> | 44,619,720 | 5,484 | 8,192 | 4 | 8,136.346 | 640.325 |
| Bytes Written <file=./Process_Output/TurbRsdl.dat> | 4,271 | 72 | 224 | 57 | 59.319 | 19.544 |
| Bytes Written <file=./Process_Output/Solver.out> | 2,039 | 13 | 797 | 43 | 156.846 | 191.359 |
| Bytes Written <file=./Field_Solutions/Solution_History/NACA0012_LargeGrid_00010.Sln> | 4,356,976 | 534 | 8,192 | 4 | 8,159.131 | 501.319 |
| Bytes Written <file=./Field_Solutions/Solution_History/NACA0012_LargeGrid_00005.Sln> | 4,356,976 | 534 | 8,192 | 4 | 8,159.131 | 501.319 |
| Bytes Written <file=./Field_Solutions/NACA0012_LargeGrid.Sln> | 4,356,976 | 534 | 8,192 | 4 | 8,159.131 | 501.319 |
| Bytes Written <file=./Body_Pressure/NACA0012_LargeGrid_00010_body.Prs> | 65,986 | 9 | 8,190 | 1,300 | 7,331.778 | 2,133.204 |
| Bytes Written <file=./Body_Pressure/NACA0012_LargeGrid_00005_body.Prs> | 65,986 | 9 | 8,190 | 1,300 | 7,331.778 | 2,133.204 |
| Bytes Written <file=./Body_Pressure/FrcMnt.out> | 1,497 | 3 | 1,185 | 108 | 499 | 486.656 |
| Bytes Written | 147,107,546 | 18,550 | 8,192 | 1 | 7,930.326 | 1,420.552 |

**ParaTools**

# TAU – Context Events

# 3D Visualization in ParaProf

# TAU − ParaProf 3D Visualization



% paraprof app.ppk
Windows -> 3D Visualization -> Bar Plot (right pane)

ParaTools

UNIVERSITY OF OREGON

# TAU − 3D Communication Window



% export TAU_COMM_MATRIX=1; mpirun ⋯ tau_exec ./a.out
% paraprof  app.ppk;     Windows -> 3D Communication Matrix

**ParaTools**

UNIVERSITY OF OREGON

# TAU − 3D Scatter Plot



% paraprof  app.ppk;    Windows -> 3D Visualization -> Scatter Plot

47

# Examples

# Setting up Accounts, Examples at Sandia

```
On Serrano.sandia.gov
% tar zxf /projects/tau/workshop.tgz
% cd workshop
% source /projects/tau/tau.bashrc   (or tau.cshrc) for OpenMPI
% echo $TAU_MAKEFILE_BASE
% cat README
And try the examples.


OR:
See http://tau.uoregon.edu/workshop_summer19.tgz  (cat README)
```

ParaTools

UNIVERSITY OF OREGON

# Setting up Accounts, Examples at LANL

```
On Grizzly gr-fe.lanl.gov:
% tar zxf /usr/projects/packages/tau/workshop.tgz
% cat README
% module load friendly-testing
% module load intel
% module load openmpi
% module load tau/2.28.2
% echo $TAU_MAKEFILE_BASE
/usr/projects/hpcsoft/toss3/grizzly/tau/2.28.2/x86_64/lib/
Makefile.tau-intel-18.0.5-openmpi-2.1.2
% echo $TAU_LTAG
intel-18.0.5-openmpi-2.1.2
% export TAU_MAKEFILE=$TAU_MAKEFILE_BASE-icpc-mpi-pdt
OR:
See http://tau.uoregon.edu/workshop_summer19.tgz  (cat README)
```

# Setting up Accounts, Examples at LLNL

```
On quartz.llnl.gov:
% tar zxf /usr/global/tools/tau/training/workshop.tgz
% cat workshop/README workshop/handson.txt
% source /usr/global/tools/tau/training/tau.bashrc
% ls $TAU/Makefile*


On Lassen.llnl.gov:
% module load tau


OR:
See http://tau.uoregon.edu/workshop_llnl19.tgz  (cat README)
```

**ParaTools**

UNIVERSITY OF OREGON

# Simplifying the use of TAU!

## Uninstrumented code:

- % make

- % mpirun –np 256 ./a.out

## With TAU:

- % mpirun –np 256 tau_exec ./a.out

- % paraprof

**For more information at the statement level:**

- % mpirun –np 256 tau_exec –ebs ./a.out

- % paraprof

**To rewrite the binary to instrument individual functions (using MAQAO):**

- % tau_rewrite a.out a.inst; mpirun –np 256 ./a.inst   (beta)

- % paraprof

**ParaTools**

UNIVERSITY OF OREGON

# TAU for Heterogeneous Measurement

**Multiple performance perspectives**

**Integrate Host-GPU support in TAU measurement framework**

- Enable use of each measurement approach
- Include use of PAPI and CUPTI
- Provide profiling and tracing support

**Tutorial**

- Use TAU library wrapping of libraries
- Use `tau_exec` to work with binaries
  ```
  % ./a.out    (uninstrumented)
  % tau_exec –T <configuration tags> –cupti ./a.out
  % paraprof
  ```

**ParaTools**

UNIVERSITY OF OREGON

# Kokkos profiling in TAU



Kokkos kernel

Use tau_exec to launch an uninstrumented binary that uses Kokkos to expose the Kokkos runtime profiling interface.

# OpenACC with PGI compilers

| Name | Exclusive... | Inclusive... ▽ | Calls | Child... |
|---|---|---|---|---|
| ▼ 🟨 .TAU application | 4.982 | 9.443 | 1 | 5,168 |
|   ▼ 🟦 openacc_enqueue_upload bench_staggeredleapfrog2 [{/storage | 0.694 | 3.35 | 3,700 | 29,867 |
|     🟩 cuMemcpyHtoDAsync_v2 | 2.47 | 2.47 | 3,700 | 0 |
|     🟦 cuEventRecord | 0.06 | 0.06 | 7,400 | 0 |
|     🟦 cuDeviceGetCount | 0.032 | 0.032 | 7,401 | 0 |
|     🟦 cuEventElapsedTime | 0.031 | 0.031 | 3,700 | 0 |
|     🟦 cuCtxSynchronize | 0.031 | 0.031 | 3,700 | 0 |
|     🟦 cuEventSynchronize | 0.028 | 0.028 | 3,700 | 0 |
|     🟦 cuDeviceGetAttribute | 0.002 | 0.002 | 249 | 0 |
|     🟦 cuDeviceGetName | 0 | 0 | 3 | 0 |
|     🟦 cuDeviceTotalMem_v2 | 0 | 0 | 3 | 0 |
|     🟦 cuCtxGetDevice | 0 | 0 | 1 | 0 |
|     🟦 cuEventCreate | 0 | 0 | 2 | 0 |
|     🟦 cuDeviceGet | 0 | 0 | 3 | 0 |
|     🟦 cuDriverGetVersion | 0 | 0 | 1 | 0 |
|     🟦 cuCtxGetCurrent | 0 | 0 | 2 | 0 |
|     🟦 cuCtxGetApiVersion | 0 | 0 | 1 | 0 |
|     🟦 cuInit | 0 | 0 | 1 | 0 |
|   ▼ 🟦 openacc_enqueue_download bench_staggeredleapfrog2 [{/stora | 0.116 | 0.556 | 600 | 5,407 |
|     🟦 cuMemcpyDtoHAsync_v2 | 0.405 | 0.405 | 600 | 0 |
|     🟦 cuEventRecord | 0.013 | 0.013 | 1,800 | 0 |
|     🟦 cuDeviceGetCount | 0.005 | 0.005 | 1,200 | 0 |

TAU: ParaProf: Statistics for: node 0, thread 0 - benchADM_acc2.ppk

% configure –c++=pgCC –cc=pgcc –fortran=pgi …
% tau_exec –T pgi  -openacc -cupti  ./a.out

ParaTools

http://tau.uoregon.edu/tau_llnl19.pdf

UNIVERSITY OF OREGON

# Tracking OpenACC Data Transfers



TAU: ParaProf: Context Events for: node 0, thread 0 - benchADM_acc2.ppk

| Name △ | Total | NumSamples | MaxValue | MinValue | MeanValue | Std. Dev. |
|---|---|---|---|---|---|---|
| ▼ .TAU application | | | | | | |
| ▼ openacc_enqueue_upload bench_staggeredleapfrog2 [{/st | | | | | | |
| ▼ cuMemcpyHtoDAsync_v2 | | | | | | |
| [GROUP=MAX_MARKER] Bytes copied from Host to De | 512,000 | 1 | 512,000 | 512,000 | 512,000 | 0 |
| Bytes copied from Host to Device | 973,016,000 | 3,700 | 512,000 | 120 | 262,977.297 | 255,846.506 |
| ▼ openacc_enqueue_download bench_staggeredleapfrog2 [{ | | | | | | |
| ▼ cuMemcpyDtoHAsync_v2 | | | | | | |
| Bytes copied from Device to Host | 307,200,000 | 600 | 512,000 | 512,000 | 512,000 | 0 |
| Bytes copied from Device to Host | 307,200,000 | 600 | 512,000 | 512,000 | 512,000 | 0 |
| Bytes copied from Host to Device | 973,016,000 | 3,700 | 512,000 | 120 | 262,977.297 | 255,846.506 |
| [GROUP=MAX_MARKER] Bytes copied from Host to Device | 512,000 | 1 | 512,000 | 512,000 | 512,000 | 0 |

% configure –c++=pgCC –cc=pgcc –fortran=pgi …
% tau_exec –T pgi  -openacc -cupti  ./a.out
Context events show extent of variation

**ParaTools**

UNIVERSITY OF OREGON

# Measuring Memory Footprint



**User Event Window: /Users/sameer/rs/taudata/coral/MCBenchmark/data/memory**

Name: Peak Memory Usage Resident Set Size (VmHWM) (KB)
Value Type: Max Value

| Value | | Label |
|---|---|---|
| 531376 | | node 0, thread 0 |
| 171897.25 | | Std. Dev. |
| 78212 | | Mean |
| 31476 | | node 1, thread 0 |
| 31436 | | node 3, thread 0 |
| 31404 | | node 2, thread 0 |

% export TAU_TRACK_MEMORY_FOOTPRINT=1
Paraprof:
Right click on a node -> Show Context Event Window -> see memory events

# ParaProf 3D Profile Browser

# TAU Performance System®

**http://tau.uoregon.edu**

- **Tuning and Analysis Utilities (20+ year project)**
- **Comprehensive performance profiling and tracing**
  - Integrated, scalable, flexible, portable
  - Targets all parallel programming/execution paradigms

- **Integrated performance toolkit**
  - Instrumentation, measurement, analysis, visualization
  - Widely-ported performance profiling / tracing system
  - Performance data management and data mining
  - Open source (BSD-style license)

- **Integrates with application frameworks**

ParaTools

UNIVERSITY OF OREGON

# What does TAU support?

C/C++  CUDA  UPC  OpenCL  Python

Fortran  GPI

OpenACC  Java  MPI

pthreads

Intel MIC  OpenMP

Intel  GNU

LLVM  PGI  Cray  Sun

MPC

AIX

Insert yours here

Linux  Windows

BlueGene  Fujitsu  ARM64

NVIDIA  Power 8  OS X

ParaTools

UNIVERSITY OF OREGON

# Profiling and Tracing

## Profiling



## Tracing



- Profiling and tracing

  **Profiling** shows you **how much** (total) time was spent in each routine

  **Tracing** shows you **when** the events take place on a timeline

# Inclusive vs. Exclusive Measurements

- Performance with respect to code regions
- Exclusive measurements for region only
- Inclusive measurements includes child regions

# Inclusive vs. Exclusive Measurements



Inclusive time

Exclusive time

http://tau.uoregon.edu/tau_llnl19.pdf

# TAU Architecture and Workflow

# TAU Architecture and Workflow

**Instrumentation**: **Add probes to perform measurements**

- Source code instrumentation using pre-processors and compiler scripts
- Wrapping external libraries (I/O, MPI, Memory, CUDA, OpenCL, pthread)
- Rewriting the binary executable

**Measurement:** **Profiling or tracing using various metrics**

- Direct instrumentation (Interval events measure exclusive or inclusive duration)
- Indirect instrumentation (Sampling measures statement level contribution)
- Throttling and runtime control of low-level events that execute frequently
- Per-thread storage of performance data
- Interface with external packages (e.g. PAPI hw performance counter library)

**Analysis**: **Visualization of profiles and traces**

- 3D visualization of profile data in paraprof or perfexplorer tools
- Trace conversion & display in external visualizers (Vampir, Jumpshot, ParaVer)

# Instrumentation

**Direct and indirect performance observation**

- Instrumentation invokes performance measurement
- Direct measurement with *probes*
- Indirect measurement with periodic sampling or hardware performance counter overflow interrupts
- Events measure performance data, metadata, context, etc.

**User-defined events**

- *Interval* (start/stop) events to measure exclusive & inclusive duration
- *Atomic events* take measurements at a single point
  - Measures total, samples, min/max/mean/std. deviation statistics
- *Context events* are atomic events with executing context
  - Measures above statistics for a given calling path

# Instrumentation

**Add hooks in the code to perform measurements**

**Source instrumentation using a preprocessor**

- Add timer start/stop calls in a copy of the source code.
- Use Program Database Toolkit (PDT) for parsing source code.
- Requires recompiling the code using TAU shell scripts (tau_cc.sh, tau_f90.sh)
- Selective instrumentation (filter file) can reduce runtime overhead and narrow instrumentation focus.

**Compiler-based instrumentation**

- Use system compiler to add a special flag to insert hooks at routine entry/exit.
- Requires recompiling using TAU compiler scripts (tau_cc.sh, tau_f90.sh…)

**Runtime preloading of TAU's Dynamic Shared Object (DSO)**

- No need to recompile code! Use **aprun tau_exec ./app** with options.
- Requires dynamic executable (link using −**dynamic** on Cray).

**Para𝒯ools**

UNIVERSITY OF OREGON

# Examples

ParaTools

UNIVERSITY OF OREGON

# Using TAU

**TAU supports several measurement and thread options**

Phase profiling, profiling with hardware counters, MPI library, CUDA…

Each measurement configuration of TAU corresponds to a unique stub makefile 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:

 **(or module load tau… )**

**% export TAU_MAKEFILE=$TAU/Makefile.tau-icpc-papi-mpi-pdt**

**% export TAU_OPTIONS= '-optVerbose …' (see tau_compiler.sh )**

**% export PATH=$TAUDIR/x86_64/bin:$PATH**

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)**

**Para Tools**

UNIVERSITY OF OREGON

# Choosing TAU_MAKEFILE

```
At LLNL (Quartz):
% source /usr/global/tools/tau/training/tau.bashrc
% echo $TAU
/usr/global/tools/tau/training//tau_latest/x86_64/lib
% ls $TAU/Makefile*
/usr/global/tools/tau/training//tau_latest/x86_64/lib/Makefile.tau-icpc-mpi-pdt
/usr/global/tools/tau/training//tau_latest/x86_64/lib/Makefile.tau-icpc-papi-mpi-pdt-openmp-opari
/usr/global/tools/tau/training//tau_latest/x86_64/lib/Makefile.tau-icpc-papi-mpi-pthread-pdt
/usr/global/tools/tau/training//tau_latest/x86_64/lib/Makefile.tau-icpc-papi-mpi-pthread-python-pdt
/usr/global/tools/tau/training//tau_latest/x86_64/lib/Makefile.tau-icpc-papi-ompt-v5-mpi-pdt-openmp
/usr/global/tools/tau/training//tau_latest/x86_64/lib/Makefile.tau-icpc-papi-ompt-v5-pdt-openmp
/usr/global/tools/tau/training//tau_latest/x86_64/lib/Makefile.tau-icpc-pdt
```

## For an MPI+F90 application with Intel MPI, you may choose
```
$TAU/Makefile.tau-icpc-mpi-pdt
```

- Supports MPI instrumentation & PDT for automatic source instrumentation

```
% export TAU_MAKEFILE=$TAU/Makefile.tau-icpc-mpi-pdt

% tau_f90.sh matrix.f90 -o matrix
OR with build systems:
% make CC=tau_cc.sh CXX=tau_cxx.sh F90=tau_f90.sh
% cmake -DCMAKE_Fortran_COMPILER=tau_f90.sh
        -DCMAKE_C_COMPILER=tau_cc.sh -DCMAKE_CXX_COMPILER=tau_cxx.sh
% <ALLOCATE a NODE>
% mpirun -np 256 ./matrix    (lrun on lassen/srun on Quartz)
% paraprof
```

**ParaTools**

UNIVERSITY OF OREGON

# 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:**
```
% mpirun -np 256 ./a.out       to
% mpirun -np 256  tau_exec –T <comma_separated_options> ./a.out
```

```
% mpirun -np 256  tau_exec –T papi,mpi,pdt ./a.out
```
**Preloads $TAU/shared-papi-mpi-pdt/libTAU.so**
```
% mpirun -np 256  tau_exec –T papi ./a.out
```
**Preloads $TAU/shared-papi-mpi-pdt/libTAU.so by matching.**
```
% mpirun –np 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.**

# 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
- **Dyninst [U. Wisconsin, U. Maryland]**:

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

- **MAQAO [Intel Exascale Labs, UVSQ]**:

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

- **Pebil [SDSC]**:

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

- Execute the application to get measurement data:

```
% mpirun -np 4  ./a.inst
```

**ParaTools**

UNIVERSITY OF OREGON

# Binary Rewriting Instrumentation

```
% mpif90 –g matmult.f90 –o matmult
% tau_rewrite matmult matmult.i


Or use a selective instrumentation file (include/exclude lists)
% tau_rewrite –f select.tau –T icpc,papi \
        ./matmult -o matmult.i
% mpirun -np 256   ./matmult.i
% paraprof
```

ParaTools

UNIVERSITY OF OREGON

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



Application / Library

C / C++ parser

Fortran parser F77/90/95

IL

IL

C / C++ IL analyzer

Fortran IL analyzer

Program Database Files

DUCTAPE

TAU instrumentor

Automatic source instrumentation

ParaTools

http://tau.uoregon.edu/tau_llnl19.pdf

UNIVERSITY OF OREGON

# Automatic Source Instrumentation using PDT

# Selective Instrumentation File

```
% export TAU_OPTIONS='-optTauSelectFile=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

BEGIN_INSTRUMENT_SECTION
loops routine="foo"
loops routine="int main#"
END_INSTRUMENT_SECTION
% export TAU_SELECT_FILE=select.tau    (to use at runtime)
```

ParaTools

UNIVERSITY OF OREGON

# Automatic Instrumentation

- **Use TAU's compiler wrappers**
  - Simply replace `CXX` with `tau_cxx.sh, etc.`
  - Automatically instruments source code, links with TAU libraries.
- **Use `tau_cc.sh` for C, `tau_f90.sh` for Fortran, `tau_upc.sh` for UPC, etc.**

```
Before
CXX = mpicxx
F90 = mpif90
CXXFLAGS =
LIBS = -lm
OBJS = f1.o f2.o f3.o … fn.o

app: $(OBJS)
     $(CXX) $(LDFLAGS) $(OBJS) -o $@
     $(LIBS)
.cpp.o:
     $(CXX) $(CXXFLAGS) -c $<
```

```
After
CXX = tau_cxx.sh
F90 = tau_f90.sh
CXXFLAGS =
LIBS = -lm
OBJS = f1.o f2.o f3.o … fn.o

app: $(OBJS)
     $(CXX) $(LDFLAGS) $(OBJS) -o $@
     $(LIBS)
.cpp.o:
     $(CXX) $(CXXFLAGS) -c $<
```

**ParaTools**

UNIVERSITY OF OREGON

# Routine Level Profile

**How much time is spent in each application routine?**

ParaTools

UNIVERSITY OF OREGON

# Generating a flat profile with MPI

```
% export TAU_MAKEFILE=$TAU_MAKEFILE_BASE-icpc-papi-mpi-pdt
% make F90=tau_f90.sh
Or
% tau_f90.sh matmult.f90
; mpirun –np 16   ./a.out
% paraprof
To view. To view the data locally on the workstation,
% paraprof --pack app.ppk
  Move the app.ppk file to your desktop.
% paraprof app.ppk


Click on the "node 0" label to see profile for that node. Right
click to see other options. Windows -> 3D Visualization for 3D
window.
```

# ParaProf 3D Profile Browser

# ParaProf

http://tau.uoregon.edu/tau_llnl19.pdf

# ParaProf 3D Topology Display

# ParaProf Comparison Window



Comparing Rank 0 with 5.
Right click on "node 5" -> Add node to comparison window

# Event Based Sampling in TAU



% export TAU_MAKEFILE=$TAU_MAKEFILE_BASE-icpc-papi-mpi-pdt
% make CC=tau_cc.sh CXX=tau_cxx.sh
% export TAU_SAMPLING=1
% mpirun –np 256 ./a.out
% paraprof

# Mixed MPI and OpenMP Instrumentation



Options -> Uncheck "Stack Bars Together"

# Opari OpenMP Instrumentation, Sampling



% export TAU_MAKEFILE=$TAU_MAKEFILE_BASE-icpc-papi-mpi-pdt-opari-openmp
% make CC=tau_cc.sh CXX=tau_cxx.sh
% export TAU_SAMPLING=1; export OMP_NUM_THREADS=16
% mpirun –np 256 ./a.out
% paraprof

ParaTools

# TAU's support for OMPT TR6

| Name | Exclusive... | Inclusive ...▽ | Calls | Child Calls |
|---|---|---|---|---|
| ■ void ljDestroy(BasePotential **) C [{ljForce.c} {99,1}–{108,1}] | 0 | 0 | 1 | 0 |
| ▶ ■ OpenMP_Parallel_Region advanceVelocity [{/turquoise/users/sameer/workshop/CoMD/src–openmp/timestep.c} {71, 0}] | | | | |
| ▼ ■ OpenMP_Parallel_Region redistributeAtoms [{/turquoise/users/sameer/workshop/CoMD/src–openmp/timestep.c} {152, 0}] | | | | |
| ▼ ■ [CONTEXT] OpenMP_Parallel_Region redistributeAtoms [{/turquoise/users/sameer/workshop/CoMD/src–openmp/timestep.c} {152, 0}] | 0 | 12.958 | 201 | 0 |
| ■ [SAMPLE] __lll_unlock_wake [{interp.c} {0}] | 11.321 | 11.321 | 185 | 0 |
| ■ [SAMPLE] __pthread_cond_signal [{interp.c} {0}] | 1.636 | 1.636 | 16 | 0 |
| ▼ ■ OpenMP_Sync_Region_Barrier L_ljForce_172__par_loop1_2_2 [{/turquoise/users/sameer/workshop/CoMD/src–openmp/ljForce.c} {172, 0}] | | | | |
| ▼ ■ [CONTEXT] OpenMP_Sync_Region_Barrier L_ljForce_172__par_loop1_2_2 [{/turquoise/users/sameer/workshop/CoMD/src–openmp/ljForce.c} {172, 0}] | 0 | 1.577 | 100 | 0 |
| ■ [SAMPLE] __GI_sched_yield [{interp.c} {0}] | 1.02 | 1.02 | 74 | 0 |
| ■ [SAMPLE] __kmp_hardware_timestamp [{eqtf2.c} {0}] | 0.31 | 0.31 | 15 | 0 |
| ■ [SAMPLE] __kmp_barrier [{eqtf2.c} {0}] | 0.247 | 0.247 | 11 | 0 |
| ▶ ■ OpenMP_Sync_Region_Barrier L_setTemperature_218__par_loop1_2_4 [{/turquoise/users/sameer/workshop/CoMD/src–openmp/initAtoms.c} {218, 0}] | | | | |
| ▶ ■ OpenMP_Sync_Region_Barrier advancePosition [{/turquoise/users/sameer/workshop/CoMD/src–openmp/timestep.c} {85, 0}] | | | | |
| ▶ ■ OpenMP_Sync_Region_Barrier advanceVelocity [{/turquoise/users/sameer/workshop/CoMD/src–openmp/timestep.c} {71, 0}] | | | | |
| ▶ ■ OpenMP_Sync_Region_Barrier kineticEnergy [{/turquoise/users/sameer/workshop/CoMD/src–openmp/timestep.c} {107, 0}] | | | | |

TAU: ParaProf: Statistics for: node 1, thread 0 – tauprofile.xml

Configure TAU with –ompt=download (without –opari)
% export TAU_MAKEFILE=$TAU_MAKEFILE_BASE-icpc-papi-ompt-tr6-mpi-pdt-openmp
% export TAU_OMPT_SUPPORT_LEVEL=full
% export TAU_OMPT_RESOLVE_ADDRESS_EAGERLY=1
% make CC=tau_cc.sh CXX=tau_cxx.sh
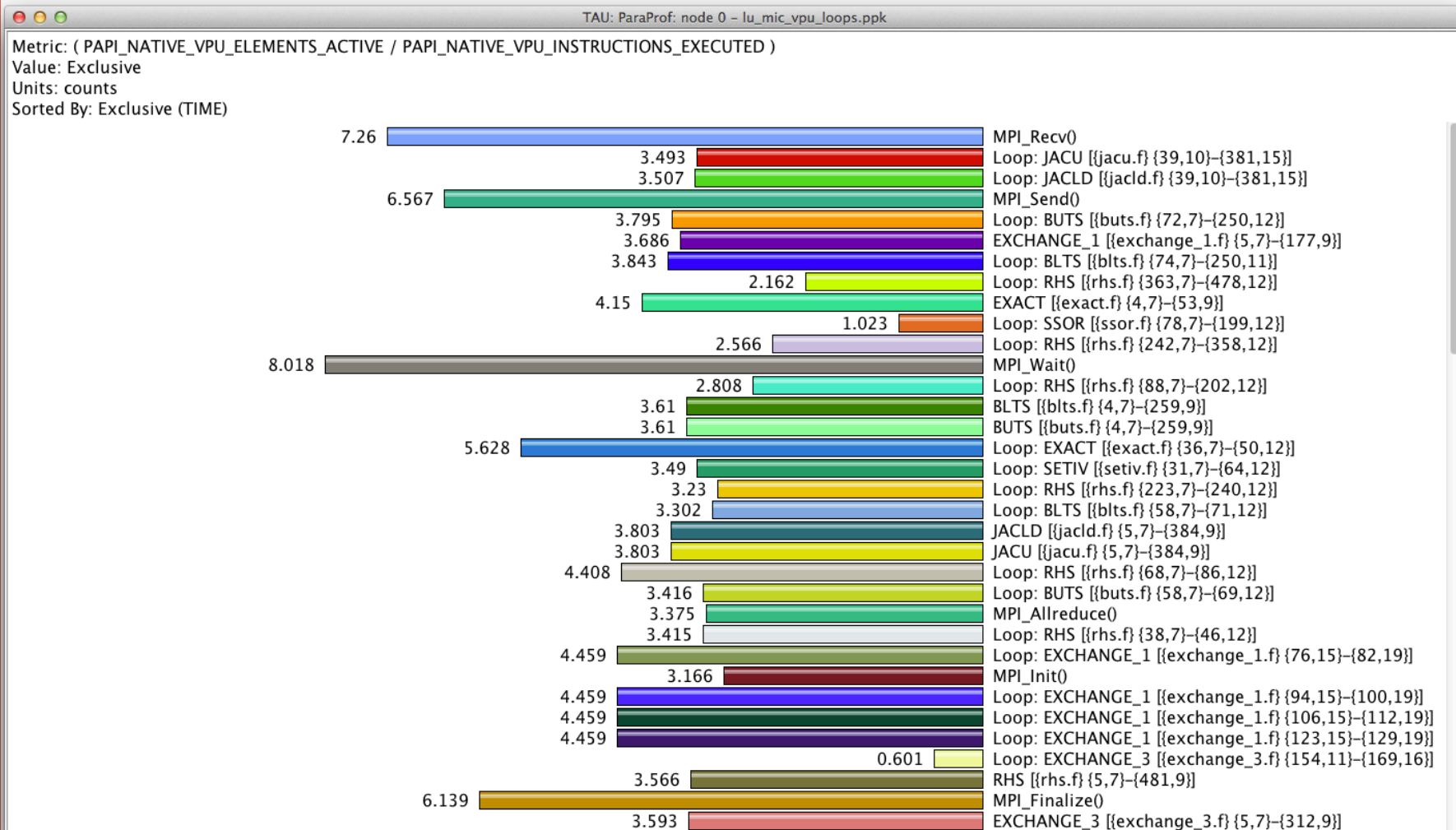% export OMP_NUM_THREADS=16
% mpirun –np 256 tau_exec –T ompt,tr6,papi,pdt  –ompt –ebs ./a.out
% paraprof
NOTE: Instrumentation is at the source, MPI, and OpenMP levels with sampling

ParaTools

UNIVERSITY OF OREGON

# Evaluating Extent of Vectorization on MIC



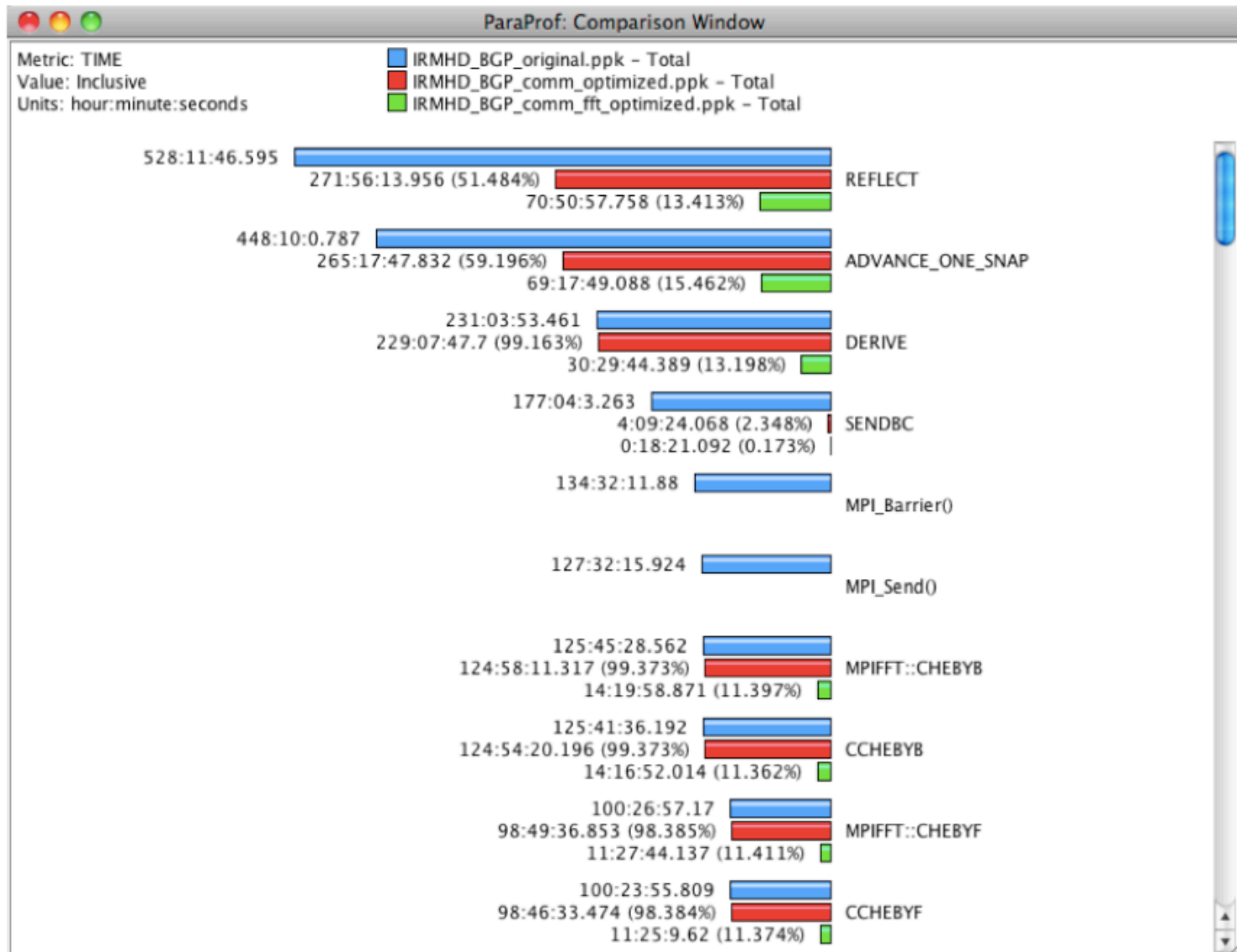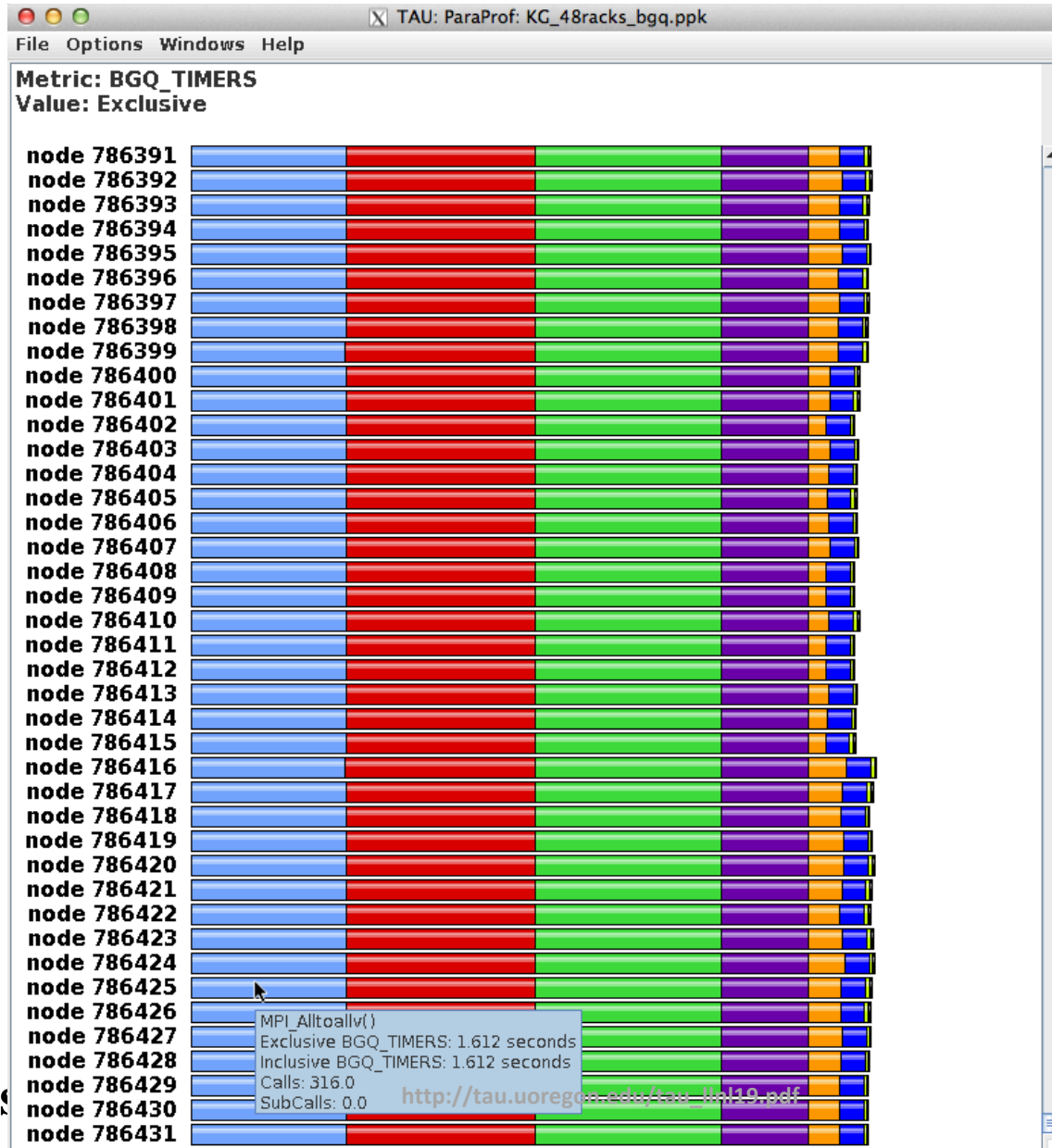% export TAU_MAKEFILE=$TAUROOT/mic_linux/lib/Makefile.tau-papi-mpi-pdt
% export TAU_METRICS=TIME,
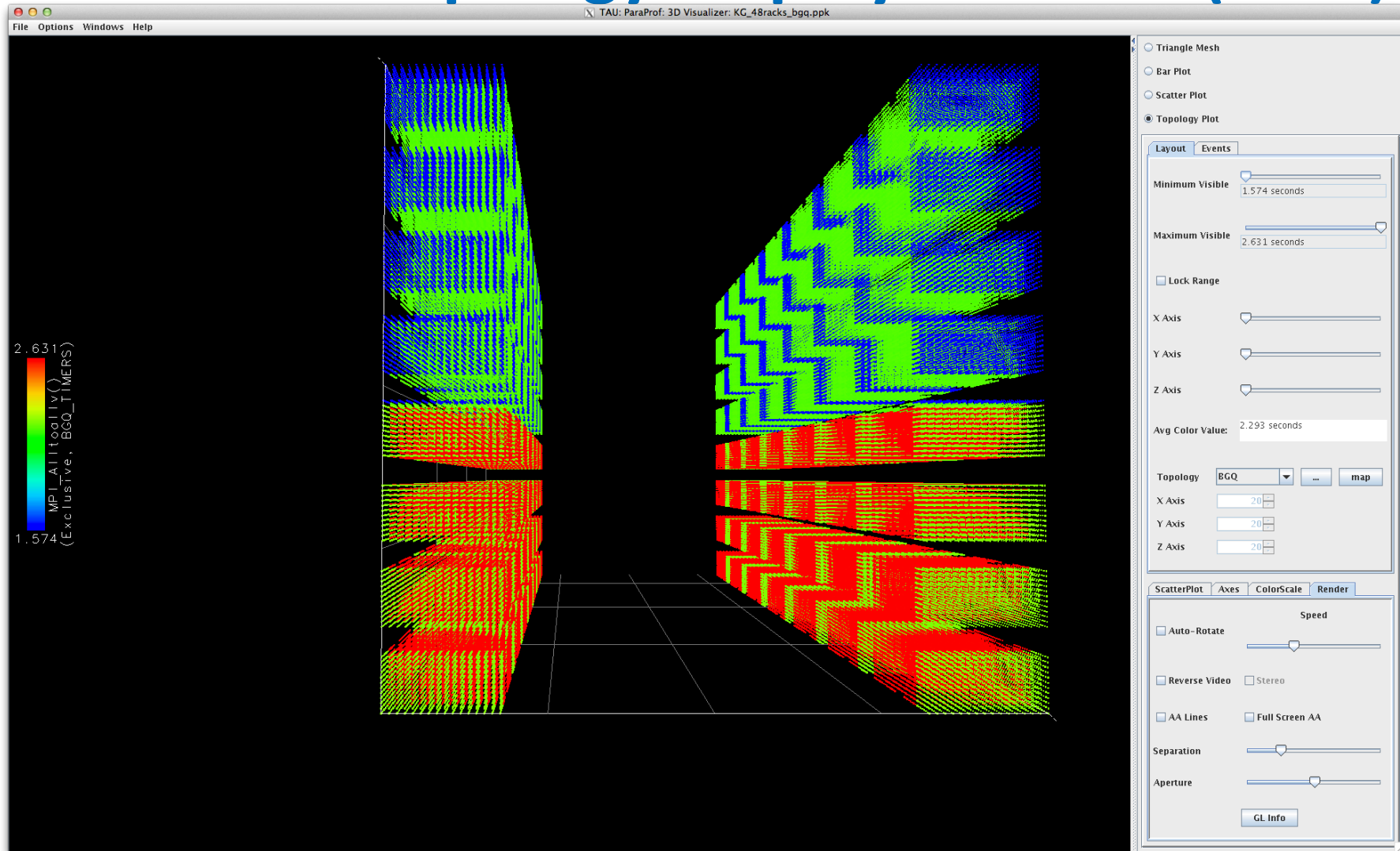        PAPI_NATIVE_VPU_ELEMENTS_ACTIVE,PAPI_NATIVE_VPU_INSTRUCTIONS_EXECUTED
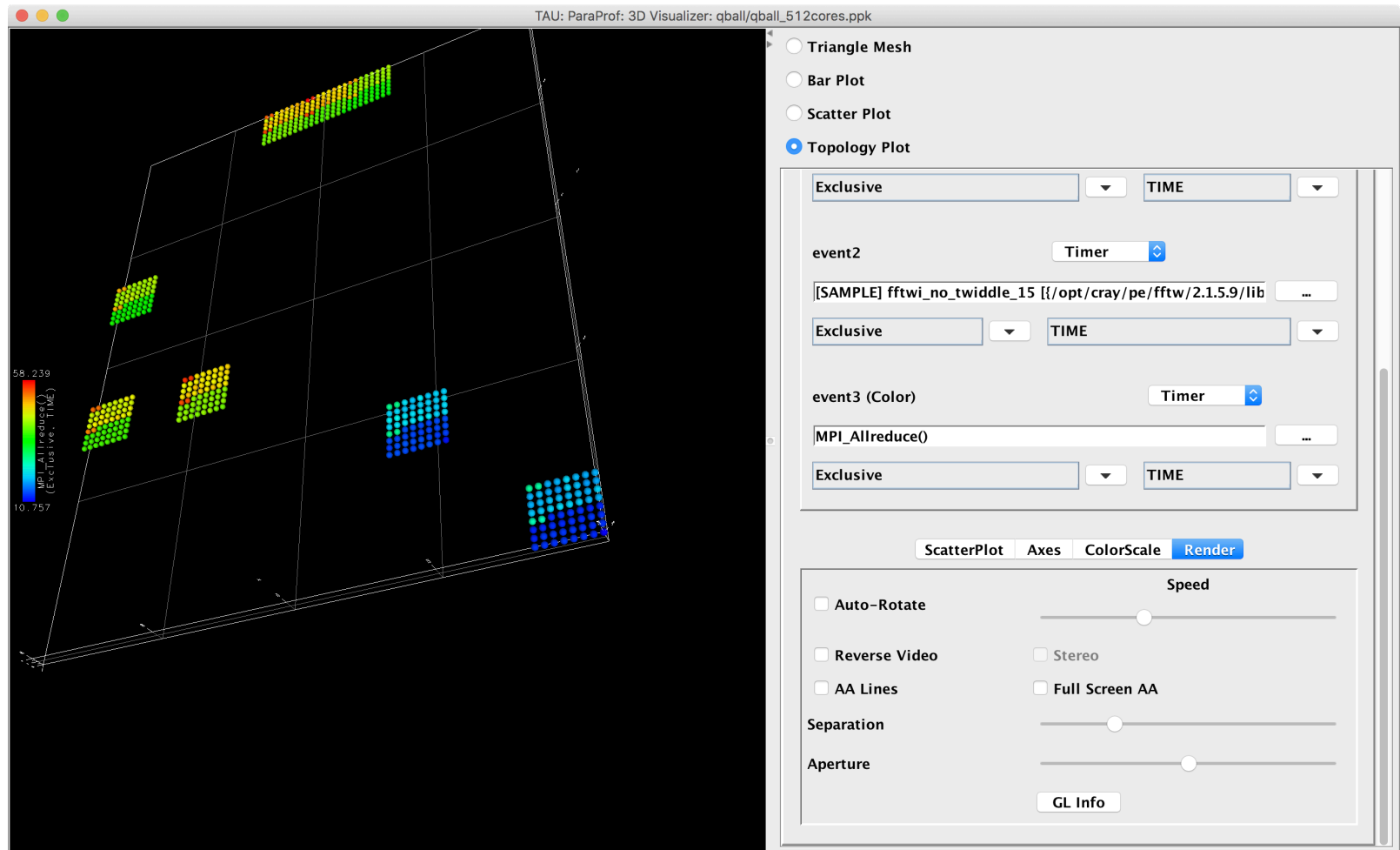
# ParaProf Comparison Window

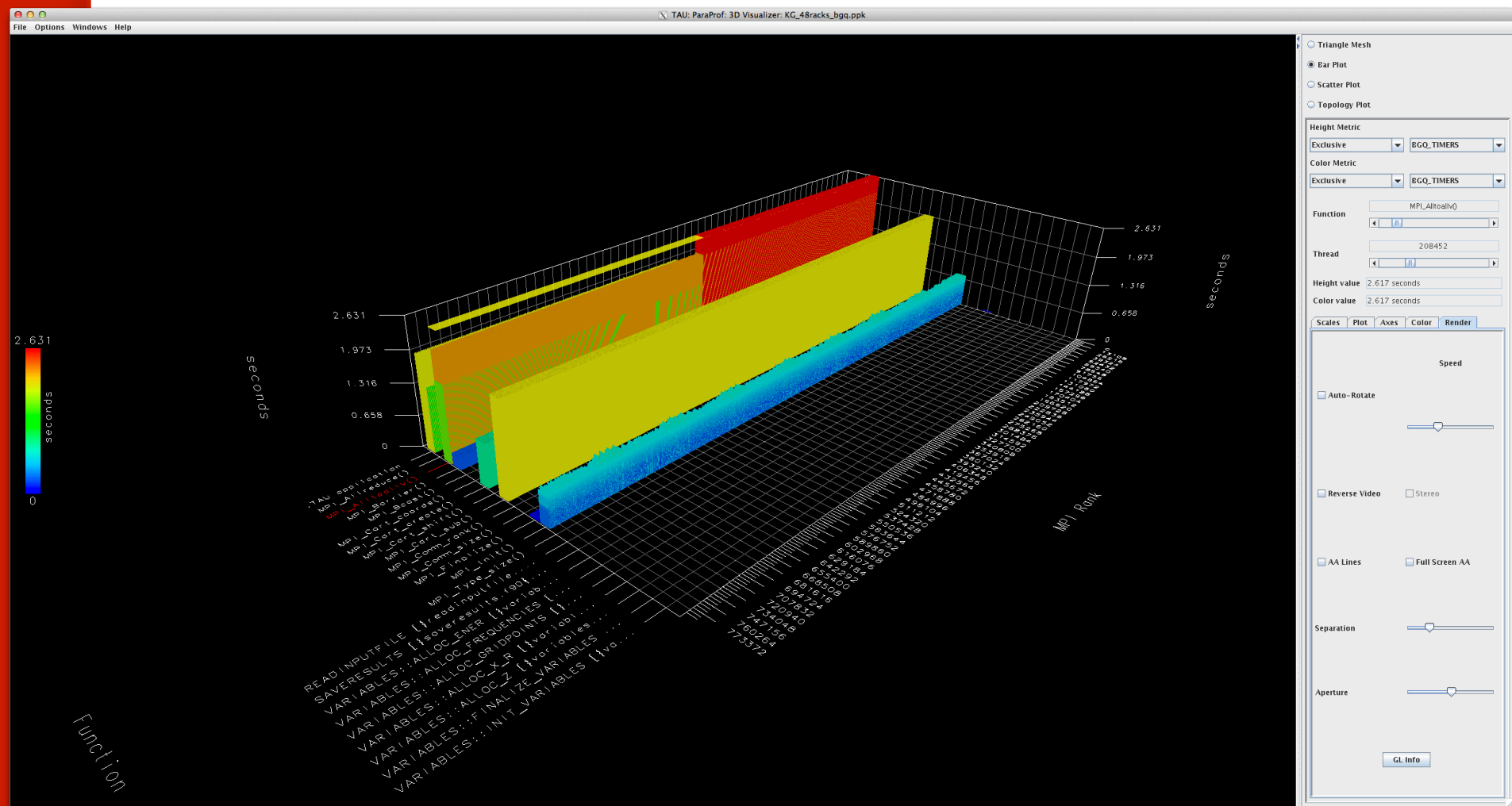http://tau.uoregon.edu/tau_llnl19.pdf

# TAU's ParaProf Profile Browser

# ParaProf's Topology Display Window (BGQ)

# ParaProf Topology Display
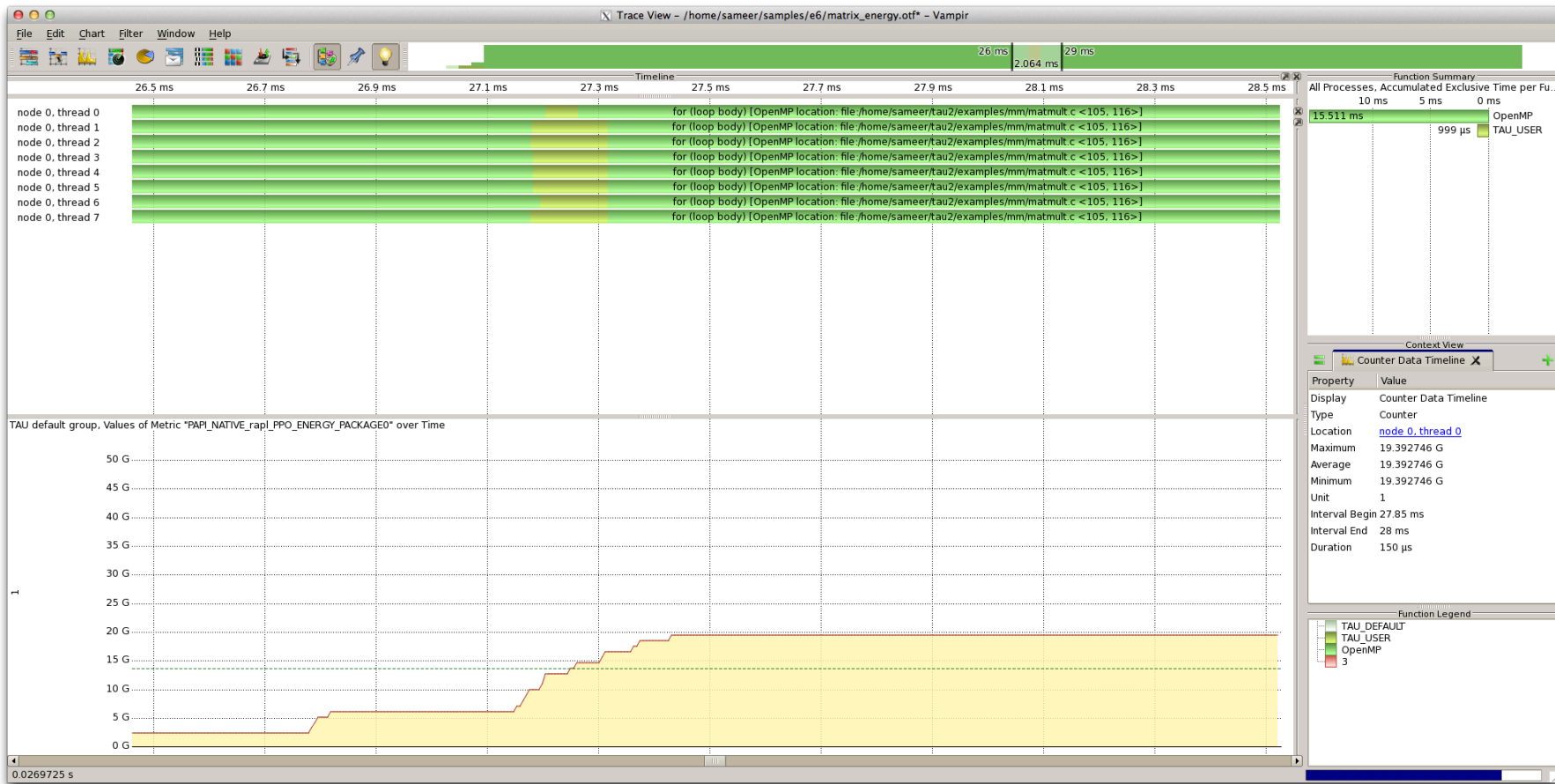
# ParaProf's Scalable 3D Visualization (BGQ)



786,432 ranks

# ParaProf Histogram Display
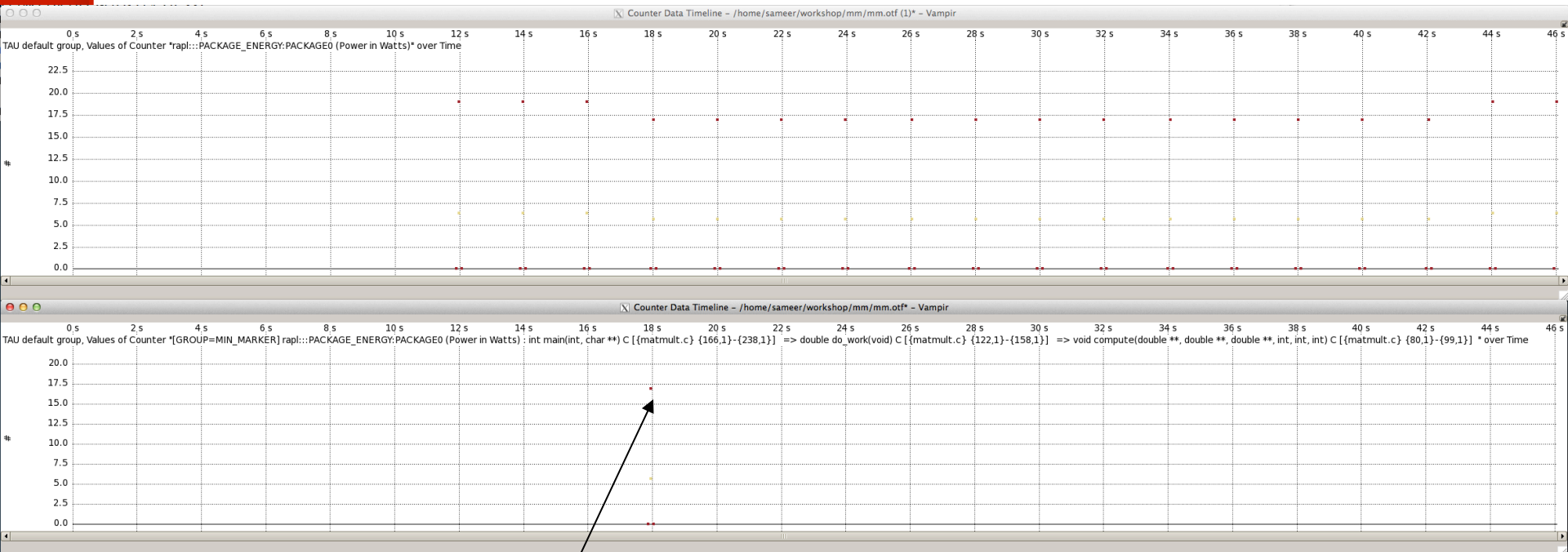
# Tracing Energy Usage with TAU and Vampir

http://tau.uoregon.edu/tau_llnl19.pdf

# TAU's Runtime Merging of Profile Data



| TrialField | Value |
|---|---|
| Name | KG_48racks_bgq.ppk |
| Application ID | 0 |
| Experiment ID | 0 |
| Trial ID | 0 |
| BGQ Block Thread ID | 3342328 |
| BGQ Coords | (188,108,272) |
| BGQ DDR Size (MB) | 16384 |
| BGQ Job ID | 221109 |
| BGQ Node ID | 77280 |
| BGQ Node Name | R2b-M0-N0f-J00 <8,12,16,16,2> |
| BGQ Period | (0,0,0) |
| BGQ Physical HW Thread ID | 0 |
| BGQ Physical Processor ID | 15 |
| BGQ Process Count | 16 |
| BGQ Processor Core ID | 15 |
| BGQ Processor Count | 4 |
| BGQ Processor ID | 60 |
| BGQ Processor Thread ID | 0 |
| BGQ Rank | 786431 |
| BGQ Size | (8,12,16,16,2,64) |
| BGQ tCoord | 15 |
| CPU MHz | 1600.000000MHz |
| CPU Type | A2 (Blue Gene/Q) |
| CWD | /gpfs/mira-fs0/projects/MiraBootCamp2013/KG/tau/49152 |
| Command Line | /gpfs/mira-fs0/projects/MiraBootCamp2013/KG/tau/49152/./Kg |
| Executable | /gpfs/mira-fs0/projects/MiraBootCamp2013/KG/tau/49152/./Kg |
| File Type Index | 0 |
| File Type Name | ParaProf Packed Profile |
| Hostname | Q2H-I3-J03.mira.i2b |
| Local Time | 2013-05-24T19:20:06+00:00 |
| MPI Processor Name | Task 786431 of 786432 (7,11,15,15,1,15)  R2B-M0-N15-J00 |
| Memory Size | 16718464 kB |
| Node Name | Q2H-I3-J03.mira.i2b |
| OS Machine | BGQ |
| OS Name | CNK |
| OS Release | 2.6.32-279.14.1.bgq.el6_V1R2M0_26.ppc64 |
| OS Version | 1 |
| Starting Timestamp | 1369423205614897 |
| TAU Architecture | bgq |
| TAU Config |  -BGQTIMERS -arch=bgq -pdt=/home/projects/tau/pdt_latest -pdt_c++=xlC -mpi -papi=/soft/perftools/tau/papi_latest -bfd=/home/projects/tau/tau2/bgq/binutils-2.20 -iowrapper |
| TAU Makefile | /soft/perftools/tau/tau-2.22.2p1/bgq/lib/Makefile.tau-bgqtimers-papi-mpi-pdt |
| TAU MetaData Merge Time | 0.000545 seconds |
| TAU Profile Merge Time | 47.34 seconds |
| TAU Unification Time | 0.01323 seconds |
| TAU Version | 2.22.2 |
| TAU_CALLPATH | off |
| TAU_CALLPATH_DEPTH | 2 |
| TAU_CALLSITE_LIMIT | 1 |

% export TAU_PROFILE_FORMAT=merged
It took ~48 seconds to merge and write profiles from 786,432 ranks

**ParaTools**

# Marker Events: Tracing



When an atomic event exceeds the max or min value by a threshold (say 20%), a marker context event is triggered to record the callstack.

# Marker events show sudden spikes

| Name △ | MaxValue | MinValue | NumSamples | MeanValue | Std. Dev. | ... |
|---|---|---|---|---|---|---|
| ▼ 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:::PP0_ENERGY:PACKAGE0 (Power in Watts) | 100.138 | 96.266 | 24 | 98.206 | 1.13 | |
| [GROUP=MAX_MARKER] rapl:::PP0_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:::PP0_ENERGY:PACKAGE0 (Power in Watts) | 75.096 | 75.096 | 1 | 75.096 | 0 | |
| [GROUP=MIN_MARKER] rapl:::PP0_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:::PP0_ENERGY:PACKAGE0 (Power in Watts) | 103.981 | 102.21 | 6 | 102.769 | 0.584 | |
| [GROUP=MAX_MARKER] rapl:::PP0_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:::PP0_ENERGY:PACKAGE0 (Power in Watts) | 103.981 | 75.096 | 36 | 98.208 | 4.466 | |
| rapl:::PP0_ENERGY:PACKAGE1 (Power in Watts) | 102.615 | 79.646 | 36 | 94.872 | 3.662 | |

TAU: ParaProf: Context Events for: node 2 – rapl_marker_16p.ppk

% export TAU_EVENT_THRESHOLD 0.5

ParaTools

UNIVERSITY OF OREGON

# Generating a loop level profile

```
% export TAU_MAKEFILE=$TAU_MAKEFILE_BASE-icpc-papi-mpi-pdt
% export TAU_OPTIONS='-optTauSelectFile=select.tau –optVerbose'
% cat select.tau
  BEGIN_INSTRUMENT_SECTION
  loops routine="#"
  END_INSTRUMENT_SECTION


% module load tau
% make F90=tau_f90.sh
(Or edit Makefile and change F90=tau_f90.sh)


% paraprof --pack app.ppk
  Move the app.ppk file to your desktop.


% paraprof app.ppk
Also, you may use export TAU_SELECT_FILE=select.tau to
filter events at runtime!
```
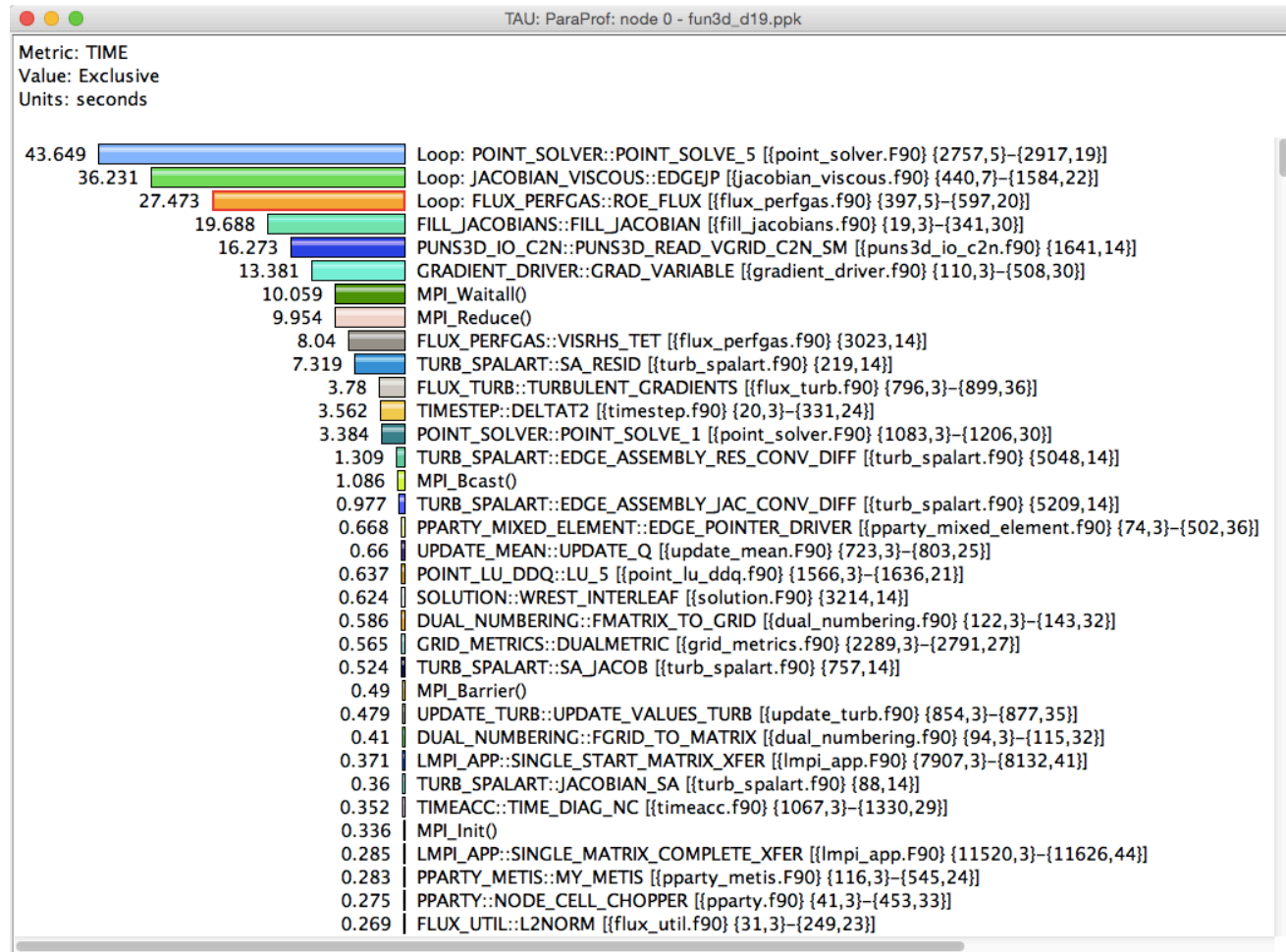
**ParaTools**

# Loop Level Instrumentation

**Goal: What loops account for the most time? How much?**

**Flat profile with wallclock time with loop instrumentation:**

# Tools: PAPI

ParaTools

# PAPI

- Consistent interface to hardware performance counters
- Preset and native events
- Higher level tools use PAPI
  - TAU, Score-P, HPCToolkit, VampirTrace, Open|SpeedShop
- API, Library, Tools
  - papi_avail - shows the list of preset counters
  - papi_native_avail – shows the list of native processor specific counters
  - papi_event_chooser – allows you to find a list of compatible events
  - papi_decode – shows how a preset event is comprised of native events
  - Components – perf, RAPL for energy, network counters
- *From University of Tenessee, Knoxville*
  - http://icl.cs.utk.edu/papi

# PAPI's Preset Cache Events on KNL

```
●●●        tg457572@c455-073.stampede2:~/tmp/mm — ssh stampede — 84×47
c455-073[knl](24)$ papi_avail --ca
Available PAPI preset and user defined events plus hardware information.
-------------------------------------------------------------------------------
PAPI Version            : 5.5.1.0
Vendor string and code  : GenuineIntel (1)
Model string and code   : Intel(R) Xeon Phi(TM) CPU 7250 @ 1.40GHz (87)
CPU Revision            : 1.000000
CPUID Info              : Family: 6  Model: 87  Stepping: 1
CPU Max Megahertz       : 1600
CPU Min Megahertz       : 1000
Hdw Threads per core    : 4
Cores per Socket        : 68
Sockets                 : 1
NUMA Nodes              : 1
CPUs per Node           : 272
Total CPUs              : 272
Running in a VM         : no
Number Hardware Counters : 5
Max Multiplex Counters  : 384
-------------------------------------------------------------------------------


===============================================================================
  PAPI Preset Events
===============================================================================
    Name         Code    Deriv Description (Note)
PAPI_L1_DCM  0x80000000  No    Level 1 data cache misses
PAPI_L1_ICM  0x80000001  No    Level 1 instruction cache misses
PAPI_L1_TCM  0x80000006  Yes   Level 1 cache misses
PAPI_L2_TCM  0x80000007  No    Level 2 cache misses
PAPI_TLB_DM  0x80000014  No    Data translation lookaside buffer misses
PAPI_L1_LDM  0x80000017  No    Level 1 load misses
PAPI_L2_LDM  0x80000019  No    Level 2 load misses
PAPI_STL_ICY 0x80000025  No    Cycles with no instruction issue
PAPI_BR_UCN  0x8000002a  Yes   Unconditional branch instructions
PAPI_BR_CN   0x8000002b  No    Conditional branch instructions
PAPI_BR_TKN  0x8000002c  No    Conditional branch instructions taken
PAPI_BR_NTK  0x8000002d  Yes   Conditional branch instructions not taken
PAPI_BR_MSP  0x8000002e  No    Conditional branch instructions mispredicted
PAPI_TOT_INS 0x80000032  No    Instructions completed
PAPI_LD_INS  0x80000035  No    Load instructions
PAPI_SR_INS  0x80000036  No    Store instructions
PAPI_BR_INS  0x80000037  No    Branch instructions
PAPI_RES_STL 0x80000039  No    Cycles stalled on any resource
PAPI_TOT_CYC 0x8000003b  No    Total cycles
PAPI_LST_INS 0x8000003c  Yes   Load/store instructions completed
PAPI_L1_DCA  0x80000040  Yes   Level 1 data cache accesses
PAPI_L1_ICH  0x80000049  No    Level 1 instruction cache hits
```

% papi_avail –ca

ParaTools

UNIVERSITY OF OREGON

# PAPI's Native Events for powercap



% papi_native_avail

# Profiling with multiple counters

```
% export TAU_MAKEFILE=$TAU_MAKEFILE_BASE-icpc-papi-mpi-pdt
% export TAU_OPTIONS='-optTauSelectFile=select.tau –optVerbose'
% cat select.tau
  BEGIN_INSTRUMENT_SECTION
  loops routine="#"
  END_INSTRUMENT_SECTION
% make F90=tau_f90.sh

% export TAU_METRICS=TIME,PAPI_TOT_CYC,PAPI_L1_DCM
% mpirun -np 4  ./matmult
% paraprof --pack app.ppk
  Move the app.ppk file to your desktop.
% paraprof app.ppk
  Choose Options -> Show Derived Panel -> Click PAPI_TOT_CYC,
 Click "/",  Click TIME, Apply, Choose new metric by double
clicking.
```

# Computing FLOPS per loop

**Goal: What is the execution rate of my loops in MFLOPS?**

**Flat profile with PAPI_FP_INS and time with loop instrumentation:**

Metric: PAPI_FP_INS / GET_TIME_OF_DAY
Value: Exclusive
Units: Derived metric shown in microseconds format

| Value | Loop |
|---|---|
| 770.699 | Loop: MULTIPLY_MATRICES [{matmult.f90} {31,9}–{36,14}] |
| 223.39 | Loop: INITIALIZE [{matmult.f90} {10,9}–{14,14}] |
| 223.24 | Loop: INITIALIZE [{matmult.f90} {17,9}–{21,14}] |
| 171.855 | Loop: MAIN [{matmult.f90} {71,9}–{74,14}] |
| 170.862 | Loop: MAIN [{matmult.f90} {112,9}–{115,14}] |
| 122.96 | Loop: MAIN [{matmult.f90} {117,9}–{128,14}] |
| 37.549 | MULTIPLY_MATRICES |
| 21.367 | INITIALIZE |
| 13.795 | Loop: MAIN [{matmult.f90} {86,9}–{106,14}] |
| 11 | MPI_Comm_size() |
| 8.935 | Loop: MAIN [{matmult.f90} {77,9}–{84,14}] |
| 1.131 | MPI_Send() |
| 0.794 | MPI_Comm_rank() |
| 0.647 | MPI_Bcast() |
| 0.355 | MPI_Recv() |
| 0.171 | MPI_Barrier() |
| 0.115 | MPI_Finalize() |
| 0.023 | MAIN |

ParaTools

UNIVERSITY OF OREGON

# Generate a Callpath Profile

```
% export TAU_MAKEFILE=$TAU_MAKEFILE_BASE-icpc-papi-mpi-pdt

% make F90=tau_f90.sh

(Or edit Makefile and change F90=tau_f90.sh)



% export TAU_CALLPATH=1

% export TAU_CALLPATH_DEPTH=100
(truncates all calling paths to a specified depth)

% mpirun -np 4    ./a.out

% paraprof --pack app.ppk

  Move the app.ppk file to your desktop.

% paraprof app.ppk

(Windows -> Thread -> Call Graph)
```

ParaTools

UNIVERSITY OF OREGON

# Callpath Profiling: FUN3D



| Name | Exclusive... ▽ | Inclusive... | Calls | Child... |
|---|---|---|---|---|
| ▼ ■ .TAU application | 0.001 | 221.305 | 1 | 1 |
| ▼ ■ NODET [{main.f90} {4,1}–{35,17}] | 0 | 221.304 | 1 | 105 |
| ▶ ■ FLOW::INITIALIZE_PROJECT [{flow.F90} {366,14}] | 0 | 0.517 | 1 | 9 |
| ▼ ■ FLOW::ITERATE [{flow.F90} {1692,14}] | 0 | 197.989 | 100 | 500 |
| ▶ ■ FLOW::STEP_POST [{flow.F90} {2098,14}] | 0.001 | 2.394 | 100 | 1,202 |
| ▼ ■ FLOW::STEP_SOLVER [{flow.F90} {1845,14}] | 0.001 | 195.577 | 100 | 702 |
| ▼ ■ RELAX_STEADY::RELAX [{relax_steady.f90} {30,3}–{307,22}] | 0.049 | 195.569 | 100 | 800 |
| ▶ ■ UPDATE_TURB::UPDATE_VALUES_TURB [{update_turb.f90} {854,3}–{877,35}] | 0.479 | 0.737 | 100 | 300 |
| ▶ ■ RELAX_TURB::RELAX [{relax_turb.f90} {22,3}–{68,22}] | 0.024 | 4.77 | 100 | 300 |
| ▼ ■ RELAX_MEAN::RELAX [{relax_mean.f90} {22,3}–{84,22}] | 0.002 | 54.402 | 100 | 300 |
| ▶ ■ WU_DEFS::TIMES [{wu_defs.f90} {59,3}–{174,22}] | 0.003 | 0.065 | 200 | 200 |
| ▼ ■ GCR_SOLVE::GCR_SOLVER_QSET [{gcr_solve.f90} {47,3}–{415,32}] | 0.002 | 54.334 | 100 | 801 |
| ▶ ■ GCR_UTIL::RES_RMS_QSET [{gcr_util.f90} {375,3}–{395,29}] | 0.001 | 0.15 | 100 | 100 |
| ▶ ■ GCR_UTIL::MATRIX_TO_GRID_RES [{gcr_util.f90} {313,3}–{336,35}] | 0.001 | 0.536 | 100 | 100 |
| ▶ ■ GCR_UTIL::MATRIX_TO_GRID_DQ [{gcr_util.f90} {282,3}–{305,34}] | 0.001 | 0.195 | 100 | 100 |
| ▶ ■ GCR_UTIL::GRID_TO_MATRIX_RES [{gcr_util.f90} {344,3}–{367,35}] | 0 | 0.341 | 100 | 100 |
| ▼ ■ GCR_SOLVE_UTIL::GCR_PRECONDITIONER_QSET [{gcr_solve_util.f90} {40,3}–{131,40}] | 0 | 53.104 | 100 | 100 |
| ▼ ■ LINEARSOLVE_NODIVCHECK::NODIVCHECK_RELAX_Q [{linearsolve_nodivcheck.F90} {56,14}] | 0.008 | 53.103 | 100 | 4,900 |
| ▶ ■ WU_DEFS::TIMES [{wu_defs.f90} {59,3}–{174,22}] | 0.02 | 0.34 | 3,200 | 3,200 |
| ▼ ■ POINT_SOLVER::POINT_SOLVE [{point_solver.f90} {31,3}–{214,28}] | 0.004 | 52.751 | 1,500 | 1,500 |
| ▼ ■ POINT_SOLVER::POINT_SOLVE_5 [{point_solver.F90} {2700,3}–{2921,30}] | 0.003 | 52.747 | 1,500 | 1,500 |
| ▼ ■ Loop: POINT_SOLVER::POINT_SOLVE_5 [{point_solver.F90} {2757,5}–{2917,19}] | 43.649 | 52.744 | 1,500 | 36,000 |
| ▶ ■ LMPI_APP::SINGLE_START_MATRIX_XFER [{lmpi_app.F90} {7907,3}–{8132,41}] | 0.271 | 0.512 | 18,000 | 85,500 |
| ▼ ■ LMPI_APP::SINGLE_MATRIX_COMPLETE_XFER [{lmpi_app.F90} {11520,3}–{11626,44}] | 0.228 | 8.583 | 18,000 | 30,000 |
| ▼ ■ LMPI::LMPI_WAITALL [{lmpi.F90} {20175,3}–{20200,29}] | 0.139 | 8.355 | 30,000 | 30,000 |
| ■ MPI_Waitall() | 8.217 | 8.217 | 30,000 | 0 |
| ▶ ■ LMPI::INTEGR_SCALAR_REDUCE [{lmpi.F90} {4584,3}–{4611,37}] | 0 | 0.002 | 100 | 100 |
| ▶ ■ LINEAR_SPECTRAL::SET_FIELD_POINTS [{linear_spectral.f90} {173,3}–{184,33}] | 0 | 0.002 | 100 | 200 |

```
% export TAU_CALLPATH=1
% export TAU_CALLPATH_DEPTH=100
```

ParaTools

http://tau.uoregon.edu/tau_llnl19.pdf

UNIVERSITY OF OREGON

108

# ParaProf Function Window

# ParaProf Callpath Thread Relations Window



Shows the contribution of parents and children for each routine (marked by an arrow)
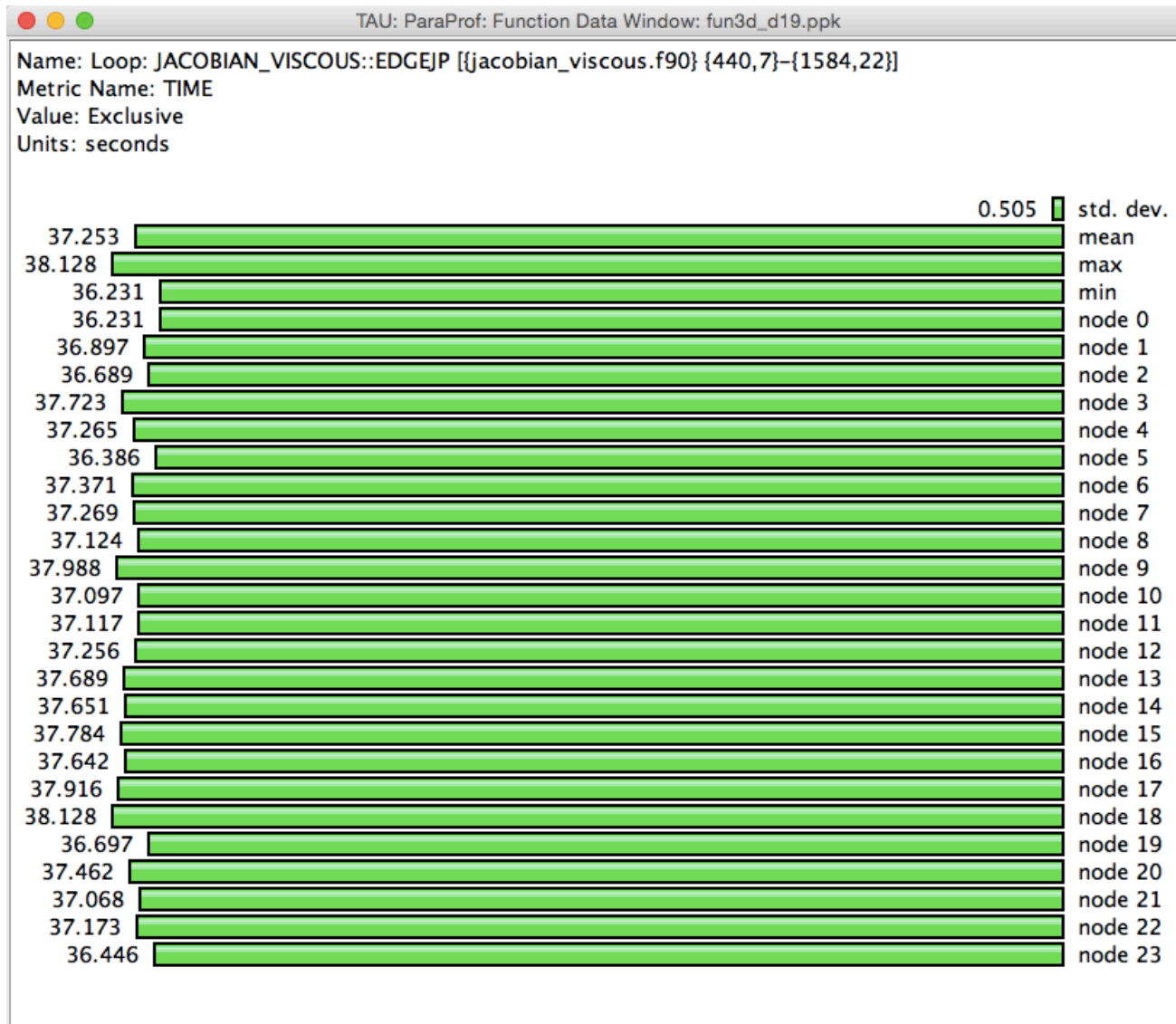
ParaTools

UNIVERSITY OF OREGON

# ParaProf Callpath Thread Relations Window



TAU: ParaProf: Call Path Data n,c,t, 13,0,0 - fun3d_d19.ppk

Metric Name: TIME
Sorted By: Exclusive
Units: seconds

| Exclusive | Inclusive | Calls/Tot.Calls | Name[id] |
|---|---|---|---|
| 45.642 | 52.774 | 1500/1500 | POINT_SOLVER::POINT_SOLVE_5 [{point_solver.F90} {2700,3}-{2921,30}] |
| 45.642 | 52.774 | 1500 | Loop: POINT_SOLVER::POINT_SOLVE_5 [{point_solver.F90} {2757,5}-{2917,19}] |
| 0.299 | 6.259 | 18000/30100 | LMPI_APP::SINGLE_MATRIX_COMPLETE_XFER [{lmpi_app.F90} {11520,3}-{11626,44}] |
| 0.6 | 0.873 | 18000/30100 | LMPI_APP::SINGLE_START_MATRIX_XFER [{lmpi_app.F90} {7907,3}-{8132,41}] |
| 37.689 | 37.689 | 78/78 | JACOBIAN_VISCOUS::EDGEJP [{jacobian_viscous.f90} {324,14}] |
| 37.689 | 37.689 | 78 | Loop: JACOBIAN_VISCOUS::EDGEJP [{jacobian_viscous.f90} {440,7}-{1584,22}] |
| 28.431 | 28.431 | 100/100 | FLUX_PERFGAS::ROE_FLUX [{flux_perfgas.f90} {236,14}] |
| 28.431 | 28.431 | 100 | Loop: FLUX_PERFGAS::ROE_FLUX [{flux_perfgas.f90} {397,5}-{597,20}] |
| 0.003 | 0.003 | 1197/27956 | LMPI::CHARACTER_BCAST [{lmpi.F90} {3100,3}-{3136,32}] |
| 0.542 | 0.542 | 489/27956 | LMPI::INTEGR_VECTOR_BCAST [{lmpi.F90} {3196,3}-{3232,36}] |
| 0.033 | 0.033 | 3/27956 | LMPI::INTEG8_SCALAR_BCAST [{lmpi.F90} {3331,3}-{3367,36}] |
| 0.005 | 0.005 | 10/27956 | LMPI::DOUBLE_MATRIX_BCAST [{lmpi.F90} {4140,3}-{4176,36}] |
| 16.724 | 16.724 | 31/27956 | LMPI::INTEGR_MATRIX_BCAST [{lmpi.F90} {3240,3}-{3276,36}] |
| 0.032 | 0.032 | 1040/27956 | LMPI::DOUBLE_SCALAR_BCAST [{lmpi.F90} {4051,3}-{4087,36}] |
| 1.48 | 1.48 | 5371/27956 | LMPI::DOUBLE_VECTOR_BCAST [{lmpi.F90} {4096,3}-{4132,36}] |
| 1.5E-5 | 1.5E-5 | 2/27956 | LMPI::LOGICL_VECTOR_BCAST [{lmpi.F90} {3736,3}-{3772,36}] |
| 0.002 | 0.002 | 2/27956 | LMPI::DOUBLE_TENSOR_BCAST [{lmpi.F90} {4185,3}-{4214,36}] |
| 0.013 | 0.013 | 48/27956 | LMPI::INTEG8_VECTOR_BCAST [{lmpi.F90} {3376,3}-{3412,36}] |
| 6.1E-4 | 6.1E-4 | 45/27956 | PPARTY_METIS::MY_METIS [{pparty_metis.F90} {116,3}-{545,24}] |
| 5.481 | 5.481 | 18941/27956 | LMPI::INTEGR_SCALAR_BCAST [{lmpi.F90} {3151,3}-{3187,36}] |
| 0.243 | 0.243 | 777/27956 | LMPI::LOGICL_SCALAR_BCAST [{lmpi.F90} {3691,3}-{3727,36}] |
| 24.557 | 24.557 | 27956 | MPI_Bcast() |
| 20.045 | 61.19 | 78/78 | UPDATE_MEAN::UPDATE_JACOBIAN [{update_mean.F90} {513,3}-{588,32}] |
| 20.045 | 61.19 | 78 | FILL_JACOBIANS::FILL_JACOBIAN [{fill_jacobians.f90} {19,3}-{341,30}] |
| 1.4E-4 | 1.4E-4 | 78/78 | SOURCE::SOURCE_JACOBIAN [{source.f90} {93,3}-{168,32}] |
| 0.006 | 2.491 | 3822/16665 | LMPI::LMPI_CONDITIONAL_STOP [{lmpi.F90} {611,3}-{672,38}] |
| 0.003 | 0.003 | 3822/8622 | BC_NAMES::BC_HAS_PRESSURE_CLOSURE [{bc_names.f90} {1618,3}-{1693,38}] |
| 0.008 | 0.008 | 7644/17444 | BC_NAMES::ELEMENT_BASED_BC [{bc_names.f90} {1390,3}-{1439,31}] |
| 3.2E-4 | 37.689 | 78/78 | JACOBIAN_VISCOUS::VISCOUS_JACOBIAN [{jacobian_viscous.f90} {20,14}] |
| 0.443 | 0.445 | 78/123 | TIMEACC::TIME_DIAG_NC [{timeacc.f90} {1067,3}-{1330,29}] |

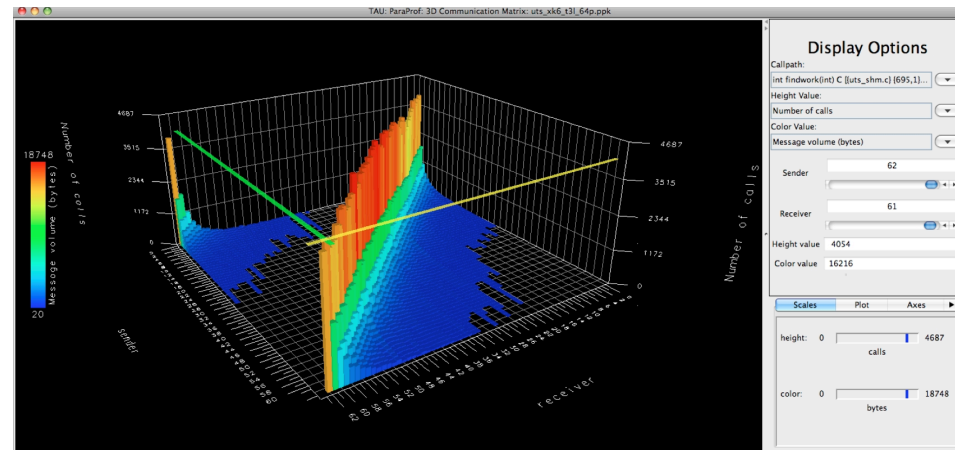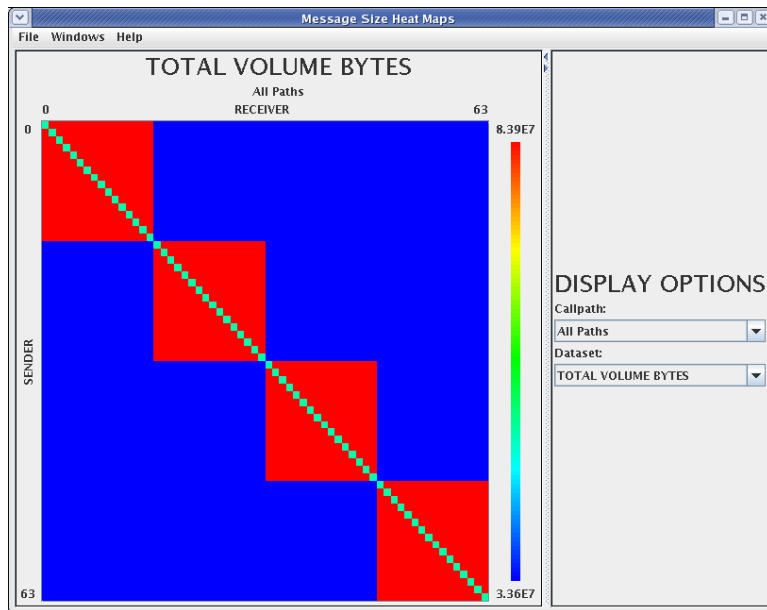# ParaProf Function Window

# Generating Communication Matrix

```
% export TAU_MAKEFILE=$TAU_MAKEFILE_BASE-icpc-papi-mpi-pdt
% make F90=tau_f90.sh
(Or edit Makefile and change F90=tau_f90.sh)

% export TAU_COMM_MATRIX=1
% mpirun -np 4    ./a.out

% paraprof
(Windows -> Communication Matrix)
(Windows -> 3D Communication Matrix)
```
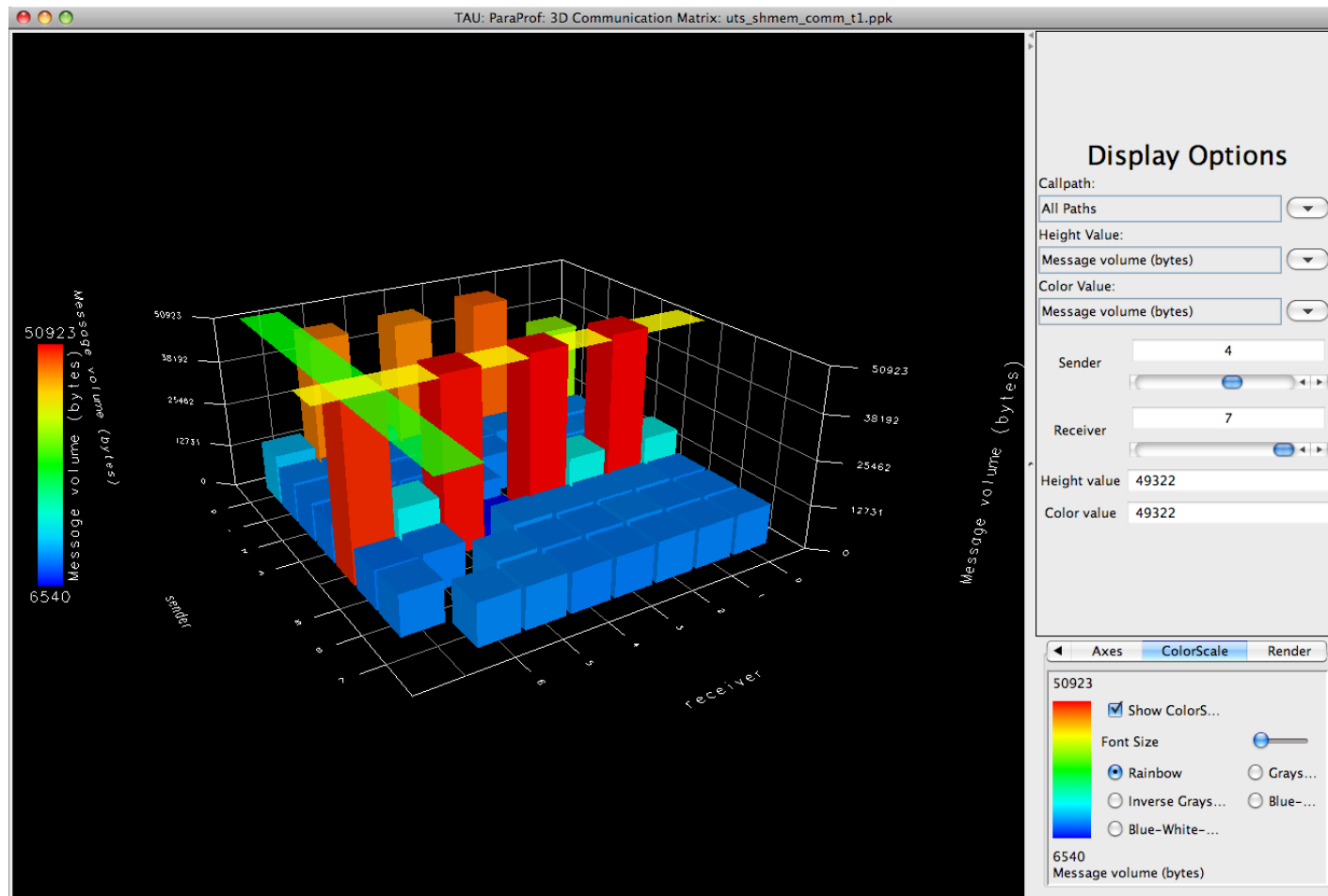
ParaTools

UNIVERSITY OF OREGON

# Communication Matrix Display

**Goal: What is the volume of inter-process communication? Along which calling path?**

# SHMEM Communication Matrix

# Compiler-based Instrumentation

- Compiler automatically **emits instrumentation calls** in the object code instead of parsing the source code using PDT

- To enable: export TAU_OPTIONS="-optCompInst"

- Configure TAU with "-bfd=download" for best results

ParaTools

O | UNIVERSITY OF OREGON

# Use Compiler-Based Instrumentation

```
% export TAU_MAKEFILE=$TAU_MAKEFILE_BASE-icpc-papi-mpi-pdt
% export TAU_OPTIONS='-optCompInst -optQuiet'

% make CC=tau_cc.sh CXX=tau_cxx.sh F90=tau_f90.sh

NOTE: You may also use the short-hand scripts taucc, tauf90,
taucxx instead of specifying TAU_OPTIONS and using the
traditional tau_<cc,cxx,f90>.sh scripts. These scripts use
compiler-based instrumentation by default.

% make CC=taucc CXX=taucxx F90=tauf90

% mpirun -np 4 ./a.out
% paraprof --pack app.ppk
  Move the app.ppk file to your desktop.
% paraprof app.ppk
```

**ParaTools**

UNIVERSITY OF OREGON

# Compiler-based Instrumentation

```
% export TAU_MAKEFILE=$TAU_MAKEFILE_BASE-icpc-papi-mpi-pdt
% export TAU_OPTIONS='-optCompInst -optQuiet'

% make CC=tau_cc.sh CXX=tau_cxx.sh F90=tau_f90.sh

NOTE: You may also use the short-hand scripts taucc, tauf90,
taucxx instead of specifying TAU_OPTIONS and using the
traditional tau_<cc,cxx,f90>.sh scripts. These scripts use
compiler-based instrumentation by default.

% make CC=taucc CXX=taucxx F90=tauf90

% mpirun -np 4 ./a.out
% paraprof --pack app.ppk
  Move the app.ppk file to your desktop.
% paraprof app.ppk
```

ParaTools

# Jumpshot Trace Visualizer in TAU

ParaTools

UNIVERSITY OF OREGON

# Tracing Communication in Jumpshot



% export TAU_MAKEFILE=$TAU_MAKEFILE_BASE-icpc-papi-mpi-pdt
% cmake –DCMAKE_CXX_COMPILER=tau_cxx.sh; make –j 8
% export TAU_TRACE=1
% mpirun –np 16 ./a.out ; tau_treemerge.pl; tau2slog2 tau.trc tau.edf –o a.slog2
% jumpshot a.slog2 &

# Performance Analysis

# Tools: Vampir

ParaTools

UNIVERSITY OF OREGON

# Vampir

**Alternative and supplement to automatic analysis**

**Show dynamic run-time behavior graphically at any level of detail**

**Provide statistics and performance metrics**

**Timeline charts**

- **Show application activities and communication along a time axis**

**Summary charts**

- **Provide quantitative results for the currently selected time interval**

- **Commercial trace visualization tool**

*From TU Dresden, Germany*

**http://www.vampir.eu**

**ParaTools**

UNIVERSITY OF OREGON

# Vampir − Trace Visualization

# Vampir − Trace Visualization

# Stencil2D Parallel Profile / Trace in Vampir

# Vampir – TAU's Kokkos Profiling Interface

# Tools: Jumpshot

ParaTools

UNIVERSITY OF OREGON

# Jumpshot

- **Open source alternative to Vampir**

- **Developed by Argonne National Laboratory**

- **Packaged with TAU**

**Timeline charts**

- **Show application activities and communication along a time axis**

- **Shows boxes within boxes to show nesting of events**

**ParaTools**

UNIVERSITY OF OREGON

# Jumpshot

# Generating Event Traces

```
% export TAU_MAKEFILE=$TAU/Makefile.tau-icpc-mpi-pdt
% make F90=tau_f90.sh
(Or edit Makefile and change F90=tau_f90.sh)

For Jumpshot:
% export TAU_TRACE=1
% mpirun -np 64    ./a.out

% tau_treemerge.pl
% tau_treemerge; tau2slog2 tau.trc tau.edf –o app.slog2;
% jumpshot app.slog2 &

For Vampir:
% export TAU_TRACE_FORMAT=otf2
# TAU's native OTF2 trace generation capability!
% mpirun -np 64   ./a.out
% vampir traces.otf2 &

For ParaVer:
% tau_convert –paraver tau.trc tau.edf app.prv; paraver app.prv
```

# Chrome Browser

```
% export TAU_MAKEFILE=$TAU/Makefile.tau-icpc-mpi-pdt
% make F90=tau_f90.sh
(Or edit Makefile and change F90=tau_f90.sh)


For Chrome:
% export TAU_TRACE=1
% mpirun -np 64    ./a.out


% tau_treemerge.pl
% tau_trace2json tau.trc tau.edf –chrome –ignoreatomic –o app.json
Copy app.json to your laptop and launch Chrome browser and in address:
chrome://tracing
Load -> app.json
```

ParaTools

UNIVERSITY OF OREGON

# Chrome Browser

# Tags in tau_exec and other tools

```
% cd $TAU; ls Makefile.*
Makefile.tau-icpc-papi-mpi-pdt
% mpirun -np 4 ./matrix
% tau_exec –T icpc,mpi,pdt ./a.out

Chooses Makefile.tau-icpc-mpi,pdt and associated libraries.

% tau_exec –T serial,pdt ./a.out
Chooses Makefile.tau-pdt or the shortest Makefile name without –mpi.


-T <list_of_tags> is used in several TAU tools:
```

- tau_run
  - tau_python
- tau_rewrite
- tau_exec
- tau_gen_wrapper

# Three Instrumentation Techniques for Wrapping External Libraries

**Pre-processor based substitution by re-defining a call (e.g., read)**

- Tool defined header file with same name *<unistd.h>* takes precedence

- Header redefines a routine as a different routine using macros

- Substitution: *read()* substituted by preprocessor as *tau_read()* at callsite

**Preloading a library at runtime**

- Library preloaded (*LD_PRELOAD* env var in Linux) in the address space of executing application intercepts calls from a given library

- Tool's wrapper library defines *read(),* gets address of global *read()* symbol (dlsym), internally calls timing calls around call to global read

**Linker based substitution**

- Wrapper library defines *__wrap_read* which calls *__real_read* and linker is passed *-Wl,-wrap,read* to substitute all references to read from application's object code with the *__wrap_read* defined by the tool

# Preprocessor based substitution

**Pre-processor based substitution by re-defining a call**

- Compiler replaces read() with tau_read() in the body of the source code

**Advantages:**

- Simple to instrument
  - Preprocessor based replacement
  - A header file redefines the calls
  - No special linker or runtime flags required

**Disadvantages**

- Only works for C & C++ for replacing calls in the body of the code.
- Incomplete instrumentation: fails to capture calls in uninstrumented libraries (e.g., libhdf5.a)

ParaTools

UNIVERSITY OF OREGON

# Linker based substitution

**Linker based substitution**

- Wrapper library defines __wrap_read which calls __real_read and linker is passed -Wl,-wrap, read

**Advantages**

- Tool can intercept all references to a given call
- Works with static as well as dynamic executables
- No need to recompile the application source code, just re-link the application objects and libraries with the tool wrapper library

**Disadvantages**

- Wrapping an entire library can lengthen the linker command line with multiple –Wl,-wrap,<func> arguments. It is better to store these arguments in a file and pass the file to the linker
- Approach does not work with un-instrumented binaries

ParaTools

http://tau.uoregon.edu/tau_llnl19.pdf

O | UNIVERSITY OF OREGON

# tau_gen_wrapper

**Automates creation of wrapper libraries using TAU**

**Input:**

- header file (foo.h)
- library to be wrapped (/path/to/libfoo.a)
- technique for wrapping
    - Preprocessor based redefinition (-d)
    - Runtime preloading (-r)
    - Linker based substitution (-w: default)
- Optional selective instrumentation file (-f select)
    - Exclude list of routines, or
    - Include list of routines

**Output:**

- wrapper library
- optional *link_options.tau* file (-w), pass –optTauWrapFile=<file> in TAU_OPTIONS environment variable

ParaTools

UNIVERSITY OF OREGON

# Design of wrapper generator (tau_gen_wrapper)

*tau_gen_wrapper* **shell script:**

- parses source of header file using static analysis tool Program Database Toolkit (PDT)

- Invokes *tau_wrap*, a tool that generates

    - instrumented wrapper code,

    - an optional *link_options.tau* file (for linker-based substitution, -w)

    - Makefile for compiling the wrapper interposition library

- Builds the wrapper library using make

**Use TAU_OPTIONS environment variable to pass location of link_options.tau file using**
```
% export TAU_OPTIONS='-optTauWrapFile=<path/to/
link_options.tau> -optVerbose'
```

**Use *tau_exec –loadlib=<wrapperlib.so>* to pass location of wrapper library for preloading based substitution**

**ParaTools**

UNIVERSITY OF OREGON

# tau_wrap

ParaTools

http://tau.uoregon.edu/tau_llnl19.pdf

UNIVERSITY OF OREGON

# Using POSIX I/O wrapper library

**Setting environment variable TAU_OPTIONS=-optTrackIO links in TAU's wrapper interposition library using linker-based substitution**

**Instrumented application generates bandwidth, volume data**

**Workflow:**

- % export TAU_OPTIONS= '-optTrackIO –optVerbose'
- % export TAU_MAKEFILE=$TAU_MAKEFILE_BASE-icpc-papi-mpi-pdt
- % make CC=tau_cc.sh CXX=tau_cxx.sh F90=tau_f90.sh
- % mpirun –np 8 ./a.out
- % paraprof

**Get additional data regarding individual arguments by setting environment variable TAU_TRACK_IO_PARAMS=1 prior to running**

**ParaTools**

UNIVERSITY OF OREGON

# Preloading a wrapper library

**Preloading a library at runtime**

- Tool defines read(), gets address of global read() symbol (dlsym), internally calls timing calls around call to global read
- *tau_exec* tool uses this mechanism to intercept library calls

**Advantages**

- No need to re-compile or re-link the application source code
- Drop-in replacement library implemented using LD_PRELOAD environment variable under Linux, Cray CNL, IBM BG/P CNK, Solaris…

**Disadvantages**

- Only works with dynamic executables. Default compilation mode under Cray XE6 and IBM BG/P is to use static executables
- Not all operating systems support preloading of dynamic shared objects (DSOs)

ParaTools

UNIVERSITY OF OREGON

# Profiling Python using tau_python



TAU: ParaProf: n,c,t 0,0,0 – flat.ppk

Metric: Time
Value: Exclusive
Units: seconds

| 5.473 | dot |
|---|---|
| 2.548 | standard_normal |
| 0.152 | MachAr::__init__ [{/usr/local/PET/pkgs/python-2.5.1/lib/python2.5/site-packages/numpy/lib/machar.py} {50}] |
| 0.112 | _parse [{/usr/local/PET/pkgs/python-2.5.1/lib/python2.5/sre_parse.py} {385}] |
| 0.105 | _compile [{/usr/local/PET/pkgs/python-2.5.1/lib/python2.5/sre_compile.py} {38}] |
| 0.056 | any [{/usr/local/PET/pkgs/python-2.5.1/lib/python2.5/site-packages/numpy/core/fromnumeric.py} {540}] |
| 0.048 | append |
| 0.046 | len |
| 0.046 | Tokenizer::__next [{/usr/local/PET/pkgs/python-2.5.1/lib/python2.5/sre_parse.py} {188}] |
| 0.045 | any |
| 0.035 | <module> [{/usr/local/PET/pkgs/python-2.5.1/lib/python2.5/site-packages/numpy/linalg/linalg.py} {7}] |
| 0.033 | SubPattern::getwidth [{/usr/local/PET/pkgs/python-2.5.1/lib/python2.5/sre_parse.py} {146}] |
| 0.033 | _optimize_charset [{/usr/local/PET/pkgs/python-2.5.1/lib/python2.5/sre_compile.py} {213}] |
| 0.028 | Tokenizer::get [{/usr/local/PET/pkgs/python-2.5.1/lib/python2.5/sre_parse.py} {207}] |
| 0.023 | _parse_sub [{/usr/local/PET/pkgs/python-2.5.1/lib/python2.5/sre_parse.py} {307}] |
| 0.022 | add_newdoc [{/usr/local/PET/pkgs/python-2.5.1/lib/python2.5/site-packages/numpy/lib/function_base.py} {1154}] |
| 0.019 | abs |
| 0.019 | _compile_charset [{/usr/local/PET/pkgs/python-2.5.1/lib/python2.5/sre_compile.py} {184}] |
| 0.016 | <module> [{/usr/local/PET/pkgs/python-2.5.1/lib/python2.5/site-packages/numpy/core/__init__.py} {2}] |
| 0.015 | <module> [{/usr/local/PET/pkgs/python-2.5.1/lib/python2.5/site-packages/numpy/core/numeric.py} {1}] |
| 0.014 | SubPattern::__len__ [{/usr/local/PET/pkgs/python-2.5.1/lib/python2.5/sre_parse.py} {132}] |
| 0.011 | _mk_bitmap [{/usr/local/PET/pkgs/python-2.5.1/lib/python2.5/sre_compile.py} {264}] |
| 0.011 | Tokenizer::match [{/usr/local/PET/pkgs/python-2.5.1/lib/python2.5/sre_parse.py} {201}] |
| 0.01 | <module> [{/usr/local/PET/pkgs/python-2.5.1/lib/python2.5/site-packages/numpy/core/defmatrix.py} {1}] |
| 0.009 | SubPattern::__getitem__ [{/usr/local/PET/pkgs/python-2.5.1/lib/python2.5/sre_parse.py} {136}] |
| 0.009 | _simple [{/usr/local/PET/pkgs/python-2.5.1/lib/python2.5/sre_compile.py} {360}] |
| 0.008 | SubPattern::append [{/usr/local/PET/pkgs/python-2.5.1/lib/python2.5/sre_parse.py} {144}] |
| 0.007 | min |
| 0.007 | <module> [{/usr/local/PET/pkgs/python-2.5.1/lib/python2.5/ctypes/__init__.py} {4}] |
| 0.007 | extend_all [{/usr/local/PET/pkgs/python-2.5.1/lib/python2.5/site-packages/numpy/core/numeric.py} {83}] |
| 0.006 | ord |
| 0.005 | <module> [{/usr/local/PET/pkgs/python-2.5.1/lib/python2.5/site-packages/numpy/__init__.py} {17}] |
| 0.005 | <module> [{/usr/local/PET/pkgs/python-2.5.1/lib/python2.5/site-packages/numpy/core/ma.py} {9}] |
| 0.005 | <module> [{/usr/local/PET/pkgs/python-2.5.1/lib/python2.5/site-packages/numpy/core/numerictypes.py} {76}] |
| 0.005 | SubPattern::__getslice__ [{/usr/local/PET/pkgs/python-2.5.1/lib/python2.5/sre_parse.py} {140}] |
| 0.004 | <module> [{/usr/local/PET/pkgs/python-2.5.1/lib/python2.5/site-packages/numpy/lib/index_tricks.py} {3}] |
| 0.004 | <module> [{/usr/local/PET/pkgs/python-2.5.1/lib/python2.5/site-packages/numpy/lib/__init__.py} {1}] |
| 0.004 | <module> [{/usr/local/PET/pkgs/python-2.5.1/lib/python2.5/site-packages/numpy/random/__init__.py} {2}] |
| 0.004 | <module> [{/usr/local/PET/pkgs/python-2.5.1/lib/python2.5/opcode.py} {5}] |
| 0.004 | _extendLine [{/usr/local/PET/pkgs/python-2.5.1/lib/python2.5/site-packages/numpy/core/arrayprint.py} {270}] |
| 0.004 | _formatArray [{/usr/local/PET/pkgs/python-2.5.1/lib/python2.5/site-packages/numpy/core/arrayprint.py} {278}] |

# Profiling Python codes

```
On Quartz:
% source /usr/global/tools/tau/training/tau.bashrc
% tar xf /usr/global/tools/tau/training/workshop.tgz
% cd workshop/py-c++-f90
% make clean; make
% mpirun -np 4 python ./samarcrun.py
With TAU:
% mpirun -np 4 tau_python ./samarcrun.py
% paraprof
```

**ParaTools**

UNIVERSITY OF OREGON

# TAU Commander

**Universal tool or integrated toolkit**

**Unbiased, accurate measurements**

- File I/O: serial and parallel
- Communication: inter- and intra-node
- Memory: allocation and access
- CPU: vectorization, cache utilization, etc.

**Minimal overhead**

- Provide multiple measurement methods
- Focus on one performance aspect at a time

**Easy to use**

- Intuitive, systematic, and well documented
- Easy to understand and configure

ParaTools

UNIVERSITY OF OREGON

# TAU Commander's Approach

**Say where you're going, not how to get there**

**Experiments give context to the user's actions**

- Defines desired metrics and measurement approach
- Defines operating environment
- Establishes a baseline for error checking



VS.

ParaTools

UNIVERSITY OF OREGON

# T-A-M Model for Performance Engineering

**Target**
- Installed software
- Available compilers
- Host architecture/ OS

**Application**
- MPI, OpenMP, CUDA, OpenACC, etc.

**Measurement**
- Profile, trace, or both
- Sample, source inst...

Measurement

Application

Target

**Experiment =
(Target, Application,
Measurement)**

# Which platform is best for my application?

- **Many targets:**
  - **Different MPI implementations**
  - **Different CPU architectures**
  - **GPU vs MIC**
  - **Cray vs SGI**
- **One measurement**
- **One application**

Measurement

Application

Target 0 ... Target N

ParaTools

UNIVERSITY OF OREGON

# What are the performance characteristics of my application?

- **One target**
- **Many measurements:**
  - **File I/O**
  - **Communication**
  - **Memory allocation**
  - **Performance counters**
  - **Vectorization**
- **One application**

Application

Measurement 0

Measurement N

Target 0

ParaTools

UNIVERSITY OF OREGON

# How well does my target perform various tasks?

- **One target**
- **One measurement**
- **Many applications:**
  - **Compute bound**
    - **Dense LA**
  - **Memory bound**
    - **Sparse LA**
    - **Graph**
  - **Scaling**
    - **Thread-level**
    - **Process-level**

Measurement

Application 0 . . . Application N

Target 0

**ParaTools**

UNIVERSITY OF OREGON

# Getting Started with TAU Commander

1. **tau** init −mpi −compilers Intel
2. **tau** mpif90 *.f90 -g -o foo
3. **tau** srun −n 64 ./foo
4. **tau** show

Just put `**tau**` in front of everything and see what happens.

This works on any supported system, even if TAU is not installed or has not been configured appropriately.

TAU and all its dependencies will be downloaded and installed if required.

**ParaTools**

UNIVERSITY OF OREGON

# TAU Commander Online Help

# Step 1: Initialize TAU Project

```
$ tau initialize --mpi --compilers Intel

$ tau init --mpi  --compilers Intel
```

- Creates a new project configuration using defaults
- Project files exist in a directory named ".tau"
- Like git, all directories below the directory containing the ".tau" directory can access the project
  - E.g. `tau dashboard` works in miniapp1/baseline

**ParaTools**

UNIVERSITY OF OREGON

# Project Initialization



Compiler detection

Project initialization

Download and install PDT

TAU installation progress

http://tau.uoregon.edu/tau_llnl19.pdf

ParaTools

UNIVERSITY OF OREGON

# Project Dashboard (`tau dashboard`)

# Step 2: Use `tau` to compile



```
1  CC = tau cc
2  LD = $(CC)
3  DEBUGFLAGS = -g -p -O0 -DDEBUG
4  OPTFLAGS = -O3 -DNDEBUG -xCORE-AVX2
5  CFLAGS += -Wall -Wextra -std=c99 #$(OPTFLAGS)
6  LDLIBS += -lrt -lm
7  LDFLAGS +=
```

Prepend `**tau**` command to compiler command

Compile as normal

```
[jlinford@cori09 ~/workspace/openshmem17/applications/ISx/SHMEM $ make optimized
tau cc -Wall -Wextra -std=c99  -O3 -DNDEBUG -xCORE-AVX2 -D SCALING_OPTION=1 -c pcg_basic.c -o obj/pcg_basic.o_s
[TAU] Cray SHMEM C compiler '/opt/cray/pe/craype/2.5.7/bin/cc' wraps
[TAU]     '/opt/intel/compilers_and_libraries_2017.2.174/linux/bin/intel64/icc'
[TAU] TAU_MAKEFILE=/global/project/projectdirs/m88/jlinford/taucmdr-test/system/tau/./tau-2.26.2/craycnl/lib/Makefile.ta
u-intel-3f5a233a-shmem-pdt
[TAU] TAU_OPTIONS=-optNoCompInst -optLinkOnly -optQuiet
[TAU] tau_cc.sh -g -Wall -Wextra -std=c99 -O3 -DNDEBUG -xCORE-AVX2 -D SCALING_OPTION=1 -c pcg_basic.c -o
[TAU]     obj/pcg_basic.o_s
```

- TAU Commander constructs a new compilation command line.
  - May replace compiler commands with TAU's compiler wrapper scripts.
  - May set environment variables, parse configuration files, etc.
  - If no changes are required then nothing is changed.

**ParaTools**

UNIVERSITY OF OREGON

# Step 3: Use `tau` to run



Prepend `**tau**` command to command line

Application executes, possibly with tau_exec

New data is added to the performance database

# Step 4: Use `tau` to view data (`tau show`)



| Name | Exclusive TIME | Inclusive TIME ▽ | Calls | Child Calls |
|---|---|---|---|---|
| ▼ .TAU application | 0.306 | 1.347 | 1 | 326 |
| void shmem_init(void) C | 0.498 | 0.498 | 1 | 0 |
| void shmem_finalize(void) C | 0.462 | 0.462 | 1 | 0 |
| ▼ [CONTEXT] .TAU application | 0 | 0.32 | 6.406 | 0 |
| ▼ [SUMMARY] main [{/global/project/projectdirs/m88/jlinford/openshmem17/applications/ISx/SHMEM/isx.c}] | 0.262 | 0.262 | 5.172 | 0 |
| [SAMPLE] main [{/global/project/projectdirs/m88/jlinford/openshmem17/applications/ISx/SHMEM/isx.c} {497}] | 0.149 | 0.149 | 2.922 | 0 |
| [SAMPLE] main [{/global/project/projectdirs/m88/jlinford/openshmem17/applications/ISx/SHMEM/isx.c} {374}] | 0.041 | 0.041 | 0.812 | 0 |
| [SAMPLE] main [{/global/project/projectdirs/m88/jlinford/openshmem17/applications/ISx/SHMEM/isx.c} {376}] | 0.036 | 0.036 | 0.719 | 0 |
| [SAMPLE] main [{/global/project/projectdirs/m88/jlinford/openshmem17/applications/ISx/SHMEM/isx.c} {378}] | 0.011 | 0.011 | 0.219 | 0 |
| [SAMPLE] main [{/global/project/projectdirs/m88/jlinford/openshmem17/applications/ISx/SHMEM/isx.c} {260}] | 0.008 | 0.008 | 0.172 | 0 |
| [SAMPLE] main [{/global/project/projectdirs/m88/jlinford/openshmem17/applications/ISx/SHMEM/isx.c} {381}] | 0.005 | 0.005 | 0.094 | 0 |
| [SAMPLE] main [{/global/project/projectdirs/m88/jlinford/openshmem17/applications/ISx/SHMEM/isx.c} {476}] | 0.004 | 0.004 | 0.078 | 0 |
| [SAMPLE] main [{/global/project/projectdirs/m88/jlinford/openshmem17/applications/ISx/SHMEM/isx.c} {379}] | 0.004 | 0.004 | 0.078 | 0 |
| [SAMPLE] main [{/global/project/projectdirs/m88/jlinford/openshmem17/applications/ISx/SHMEM/isx.c} {380}] | 0.004 | 0.004 | 0.078 | 0 |
| ▶ [SUMMARY] pcg32_boundedrand_r [{/global/project/projectdirs/m88/jlinford/openshmem17/applications/ISx/SHMEM/pcg_basic.c}] | 0.057 | 0.057 | 1.203 | 0 |
| [SAMPLE] __close_nocancel [{/home/abuild/rpmbuild/BUILD/glibc-2.19/nptl/../sysdeps/unix/syscall-template.S} {81}] | 0.001 | 0.001 | 0.016 | 0 |
| [SAMPLE] __wrap_shmem_n_pes [{/global/project/projectdirs/m88/jlinford/openshmem17/applications/ISx/SHMEM/bin/isx.strong} {0}] | 0.001 | 0.001 | 0.016 | 0 |
| void shmem_int_put(int *, const int *, size_t, int) C | 0.037 | 0.037 | 126 | 0 |
| long long shmem_longlong_fadd(long long *, long long, int) C | 0.018 | 0.018 | 128 | 0 |
| void *shmem_malloc(size_t) C | 0.015 | 0.015 | 16 | 0 |
| void shmem_barrier_all(void) C | 0.009 | 0.009 | 27 | 0 |
| void shmem_fcollect64(void *, const void *, size_t, int, int, int, long *) C | 0.001 | 0.001 | 7 | 0 |
| void shmem_collect32(void *, const void *, size_t, int, int, int, long *) C | 0 | 0 | 1 | 0 |
| void shmem_longlong_sum_to_all(long long *, const long long *, size_t, int, int, int, long long *, long *) C | 0 | 0 | 1 | 0 |
| int shmem_my_pe(void) C | 0 | 0 | 9 | 0 |
| void shmem_free(void *) C | 0 | 0 | 8 | 0 |
| int shmem_n_pes(void) C | 0 | 0 | 1 | 0 |
| ▼ [CALLSITE] void shmem_init(void) C | 0.996 | 0.996 | 2 | 0 |
| ▼ [CONTEXT] [CALLSITE] void shmem_init(void) C | 0 | 0.481 | 1.688 | 0 |
| [SAMPLE] __ioctl [{/home/abuild/rpmbuild/BUILD/glibc-2.19/misc/../sysdeps/unix/syscall-template.S} {81}] | 0.473 | 0.473 | 1.344 | 0 |
| [SAMPLE] _pmi_smp_barrier_join [{/usr/src/packages/BUILD/cray-pmi-5.0.10/src/pmi_core/smp_barrier.c} {70}] | 0.006 | 0.006 | 0.281 | 0 |
| [SAMPLE] Tau_lite_stop_timer [{/global/project/projectdirs/m88/jlinford/taucmdr-test/system/tau/tau-2.26.2/src/Profile/TauCAPI.cpp} { | 0.002 | 0.002 | 0.047 | 0 |
| [SAMPLE] _dmappi_sheap_alloc [{/home/abuild/rpmbuild/BUILD/cray-dmapp-7.1.1/src/dmapp_sheap.c} {318}] | 0.001 | 0.001 | 0.016 | 0 |

ParaTools

# Create a New Experiment

Select a new measurement to create a new experiment

```
[jlinford@nid00073 ~/workspace/openshmem17/applications/ISx/SHMEM $ tau select profile
[TAU] Created a new experiment 'cori09-ISx-profile'
[TAU] Installing TAU Performance System at
[TAU]      '/global/project/projectdirs/m88/jlinford/taucmdr-test/system/tau/./tau-2.26.2'
[TAU] Configuring TAU...
[TAU] Completed in 155.459 seconds
[TAU] Compiling and installing TAU...
[TAU] Completed in 48.596 seconds
[TAU] Checking installed files...
[TAU] Completed in 10.551 seconds
[TAU] Setting file permissions...
[TAU] Completed in 2.556 seconds
[TAU] Verifying TAU Performance System installation...
[TAU] Selected experiment 'cori09-ISx-profile'.
[TAU] Application rebuild required:
[TAU]   - source_inst changed from 'never' to 'automatic'
jlinford@nid00073 ~/workspace/openshmem17/applications/ISx/SHMEM $
```

ParaTools — ssh cori.nersc.gov — 120×17

TAU Performance System® automatically reconfigured and recompiled.

User advised that an application rebuild is required to use source-based instrumentation.

ParaTools

# TAU Commander

```
% which tau
% cd workshop/matmult
% tau init --mpi --compilers Intel
% make clean;
% make F90='tau mpif90'
% salloc –N 1 [args]
% tau mpirun –np 16 ./matmult
% tau show
And try the examples. Try:
% tau --help
% tau meas edit --help
```

ParaTools

UNIVERSITY OF OREGON

# TAU Commander: MPI + OpenMP

```
% which tau
% cd workshop/mm
% tau init --mpi --openmp ompt --compilers Intel
% make clean;
% mpicc -DTAU_MPI -DTAU_OPENMP *.c -g -qopenmp -o mm
% salloc -N 1 [args]
% tau mpirun -np 16 ./mm
% tau show
And try the examples. Try:
% tau --help
% tau meas edit --help
```

ParaTools

UNIVERSITY OF OREGON

# TAU and PDT for Source Instrumentation

```
% which tau
% cd workshop/matmult
% tau init --mpi --compilers Intel
% tau dash
% tau select profile
% make clean
% make F90='tau mpif90'
% salloc –N 1 [args]
% tau mpirun –np 16 ./matmult
% tau show
```
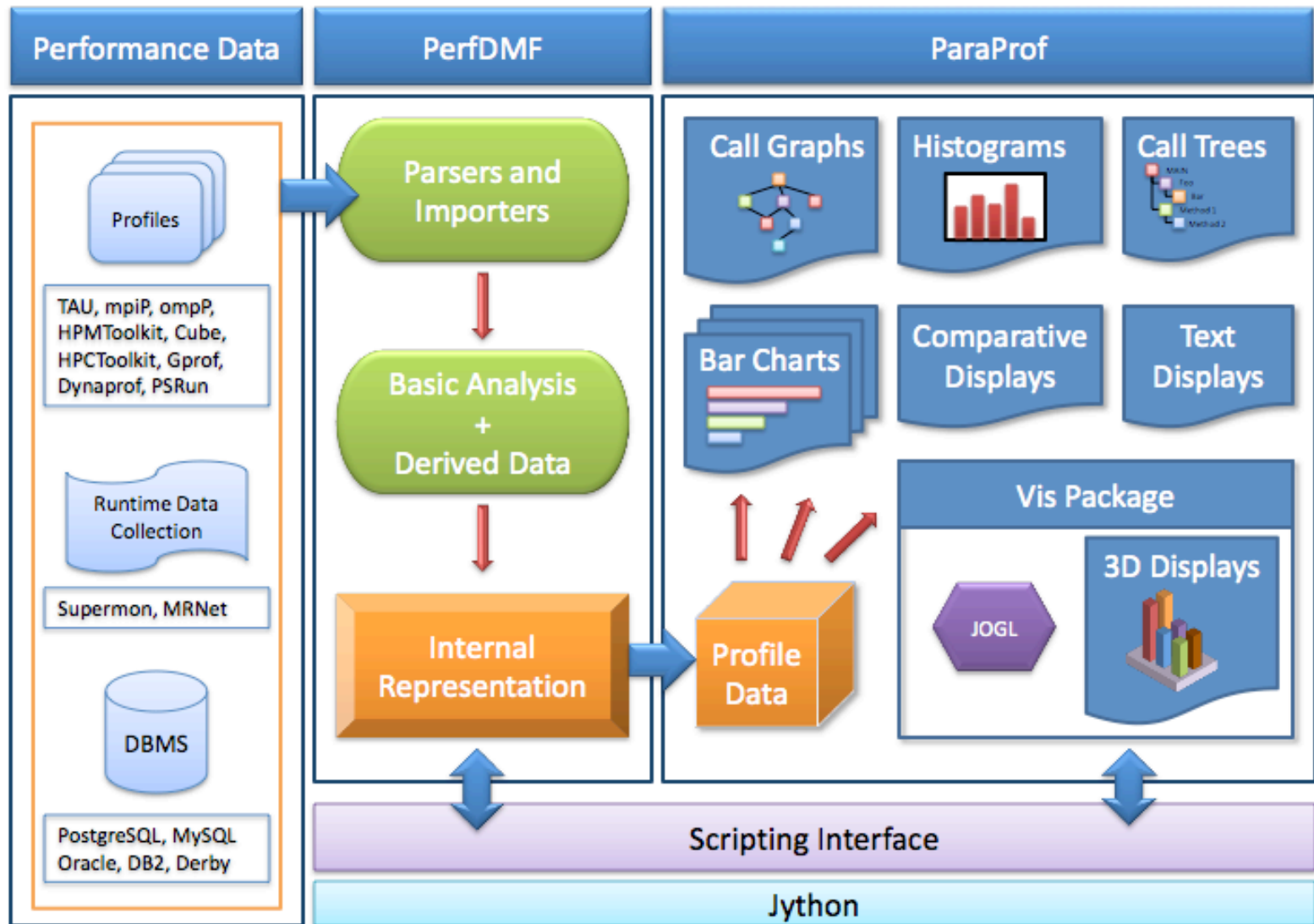
ParaTools

UNIVERSITY OF OREGON

# Selective Instrumentation File

```
% tau dash
% tau application edit <app_name> --select-file 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 (BEGIN_EXCLUDE_LIST/END···)

BEGIN_INSTRUMENT_SECTION
loops routine="foo"
loops routine="int main#"
END_INSTRUMENT_SECTION
% export TAU_SELECT_FILE=select.tau    (to use at runtime)
```

**Para**Tools

O UNIVERSITY OF OREGON

# Use Compiler-Based Instrumentation

```
% tau init --mpi --compilers Intel
% tau dash
% tau meas edit profile --compiler-inst always
% tau select profile
% make CC='tau mpicc'
% salloc –N 1 [args]
% tau mpirun –np 16 ./a.out
% tau show
% tau show --help
% tau show --profile-tools pprof | more
```

ParaTools

UNIVERSITY OF OREGON

# Generating Event Traces for Vampir

```
% cd workshop/matmult
% tau init --mpi --compilers Intel
% tau select trace
% make clean
% make F90='tau mpif90'
% salloc -N 1 [args]
% tau mpirun -np 16 ./matmult
% module load vampir
% tau show
```

ParaTools

UNIVERSITY OF OREGON

# Generating Event Traces for Jumpshot

```
% cd workshop/matmult
% tau init --mpi --compilers Intel
% tau meas edit profile --trace slog2
    (if it is profiling is being used in another experiment, you may
    have to delete it:
    tau experiment delete <exp_name>
    and retry

% make F90='tau mpif90'
% salloc -N 1 [args]
% tau mpirun -np 16 ./matmult
% tau show

NOTE: For OTF2 traces replace slog2 with otf2.
```

# ParaProf Profile Analysis Framework

ParaTools

UNIVERSITY OF OREGON

# Profile Snapshots in ParaProf

- Profile snapshots are profiles recorded at runtime
- Shows performance profile dynamics (all types allowed)

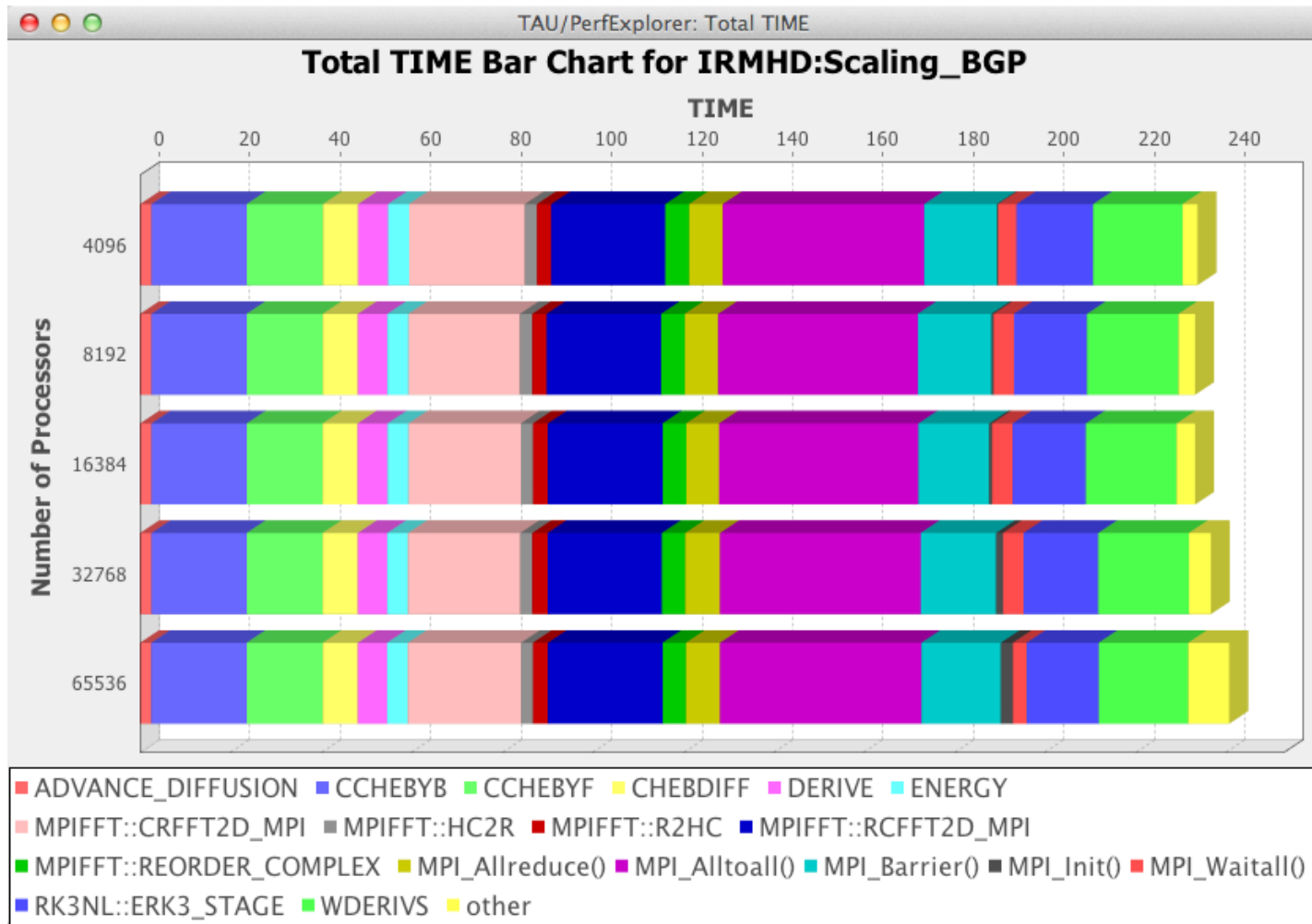# Profile Snapshot Views

## Percentage breakdown

## Only show main loop

# Snapshot Replay in ParaProf



All windows dynamically update

# PerfExplorer – Runtime Breakdown

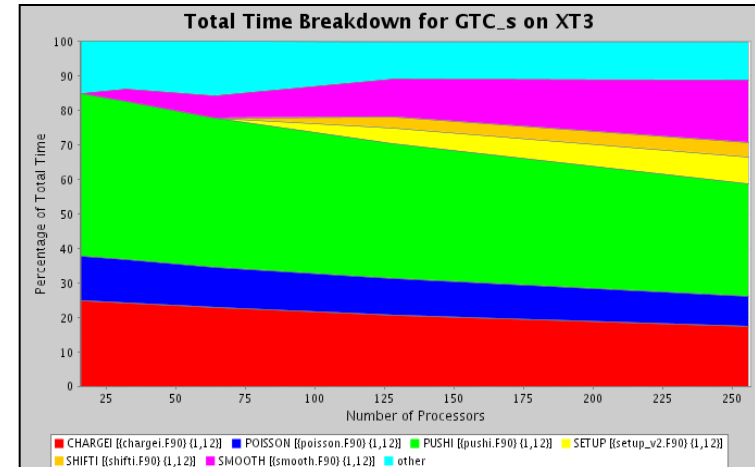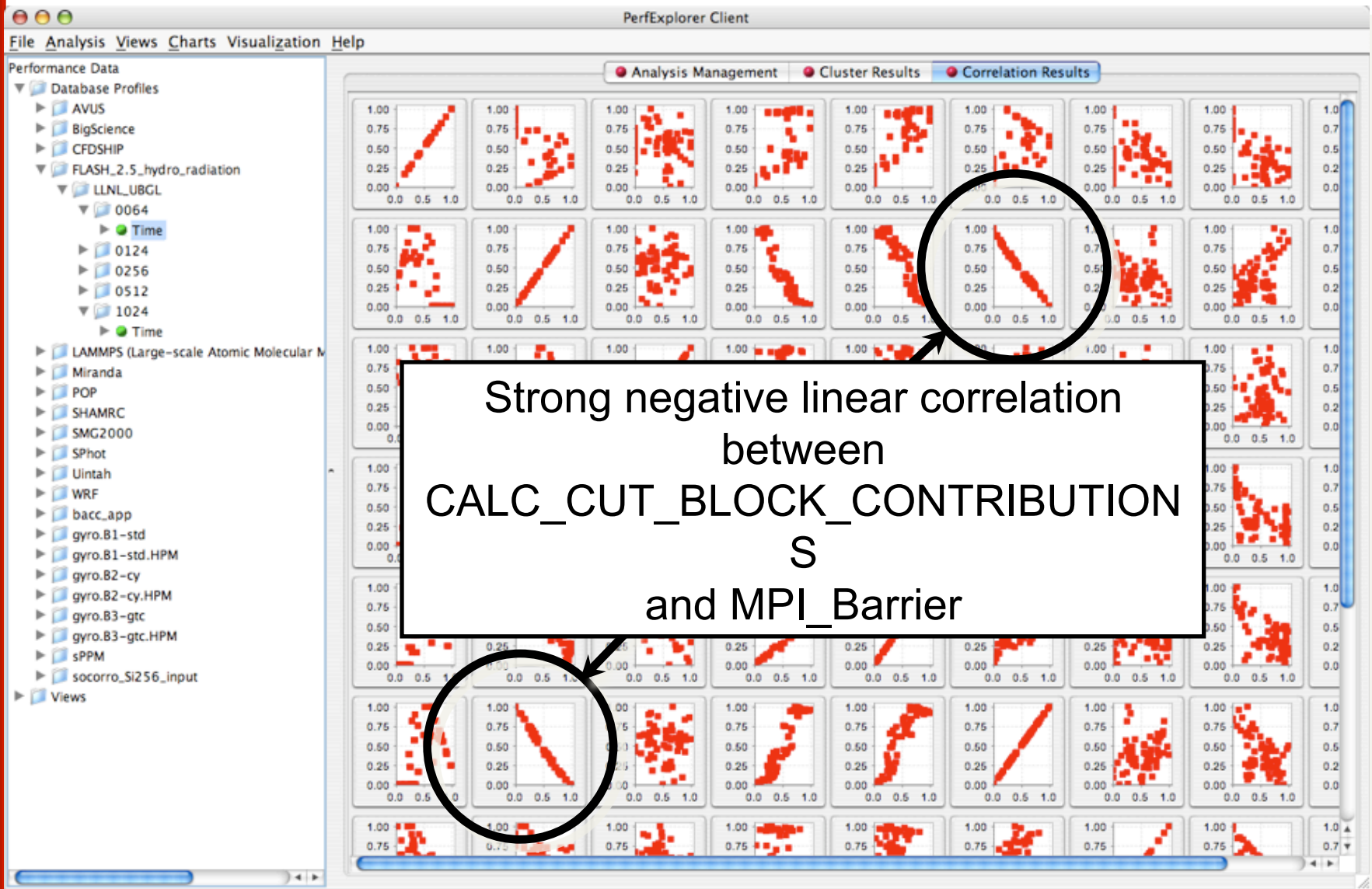# Evaluate Scalability

# Runtime Breakdown

# PerfExplorer – Relative Comparisons

**Total execution time**

**Timesteps per second**

**Relative efficiency**

**Relative efficiency per event**

**Relative speedup**

**Relative speedup per event**

**Group fraction of total**

**Runtime breakdown**

**Correlate events with total runtime**

**Relative efficiency per phase**

**Relative speedup per phase**
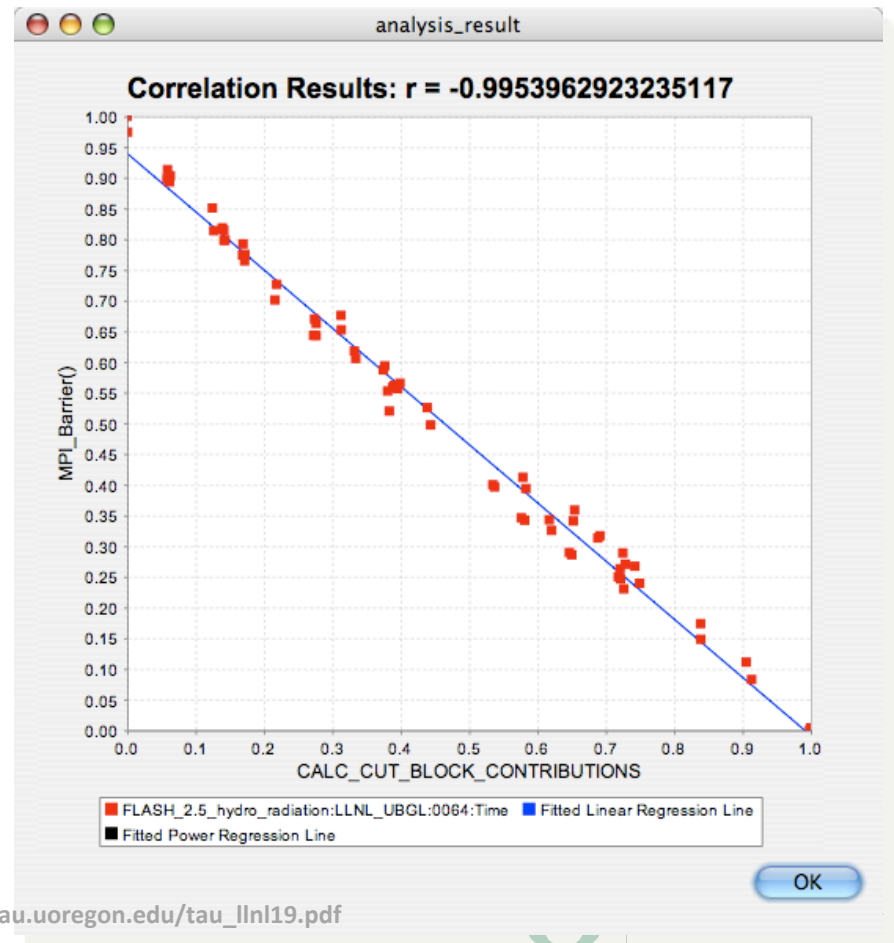
**Distribution visualizations**

# PerfExplorer – Correlation Analysis



Strong negative linear correlation
between
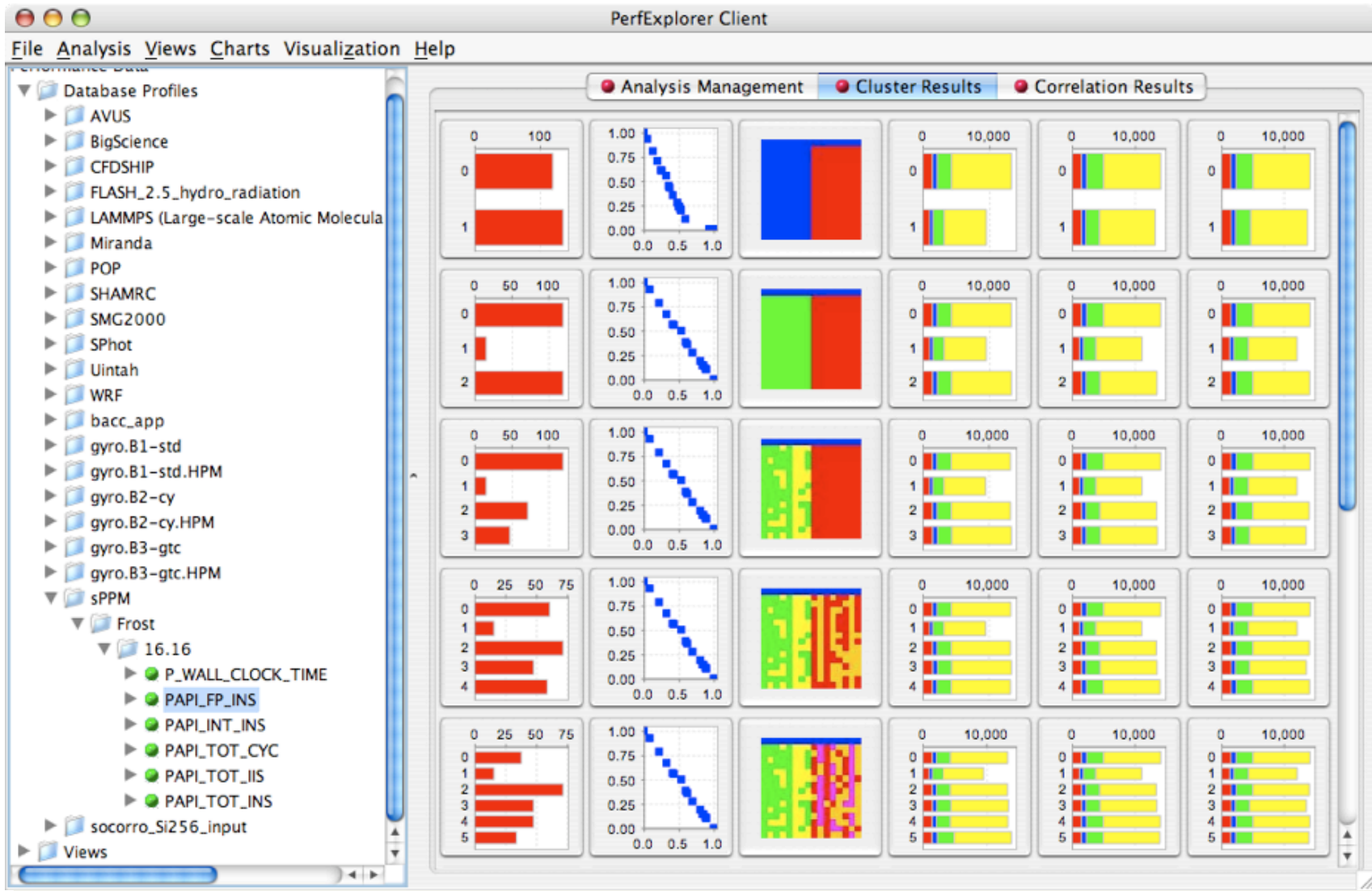CALC_CUT_BLOCK_CONTRIBUTIONS
and MPI_Barrier

# PerfExplorer – Correlation Analysis

**-0.995 indicates strong, negative relationship. As CALC_CUT_BLOCK_CONTRIBUTIONS() increases in execution time, MPI_Barrier() decreases**
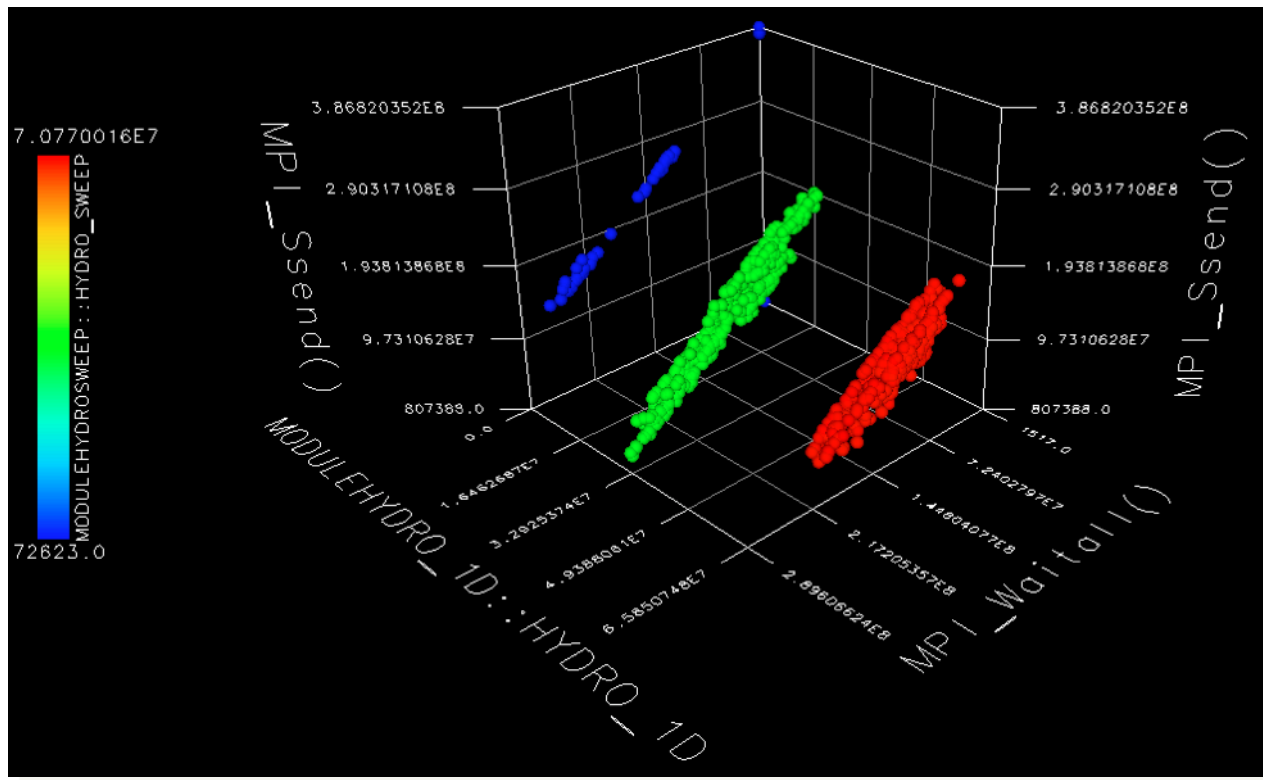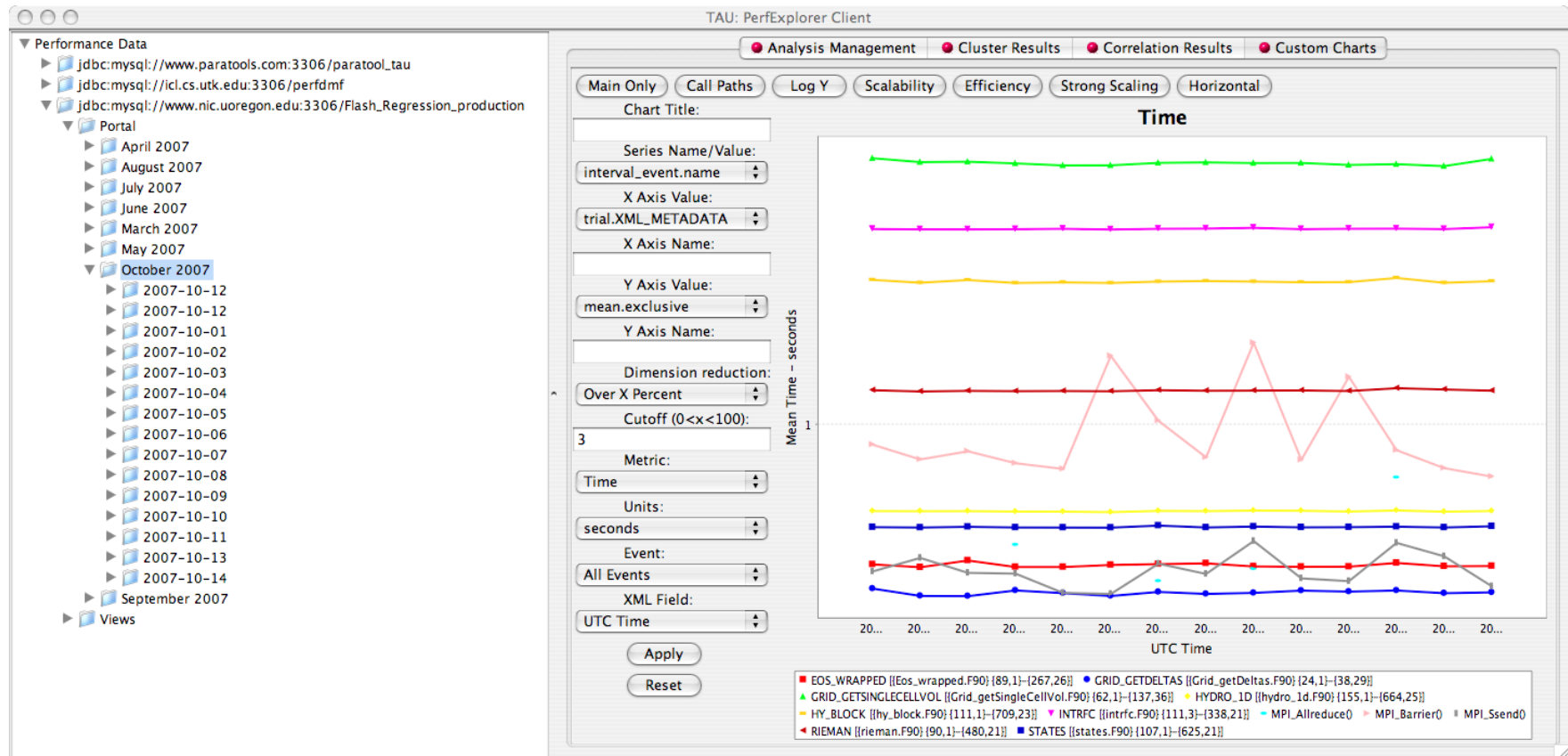
# PerfExplorer – Cluster Analysis

# PerfExplorer – Cluster Analysis

**Four significant events automatically selected**

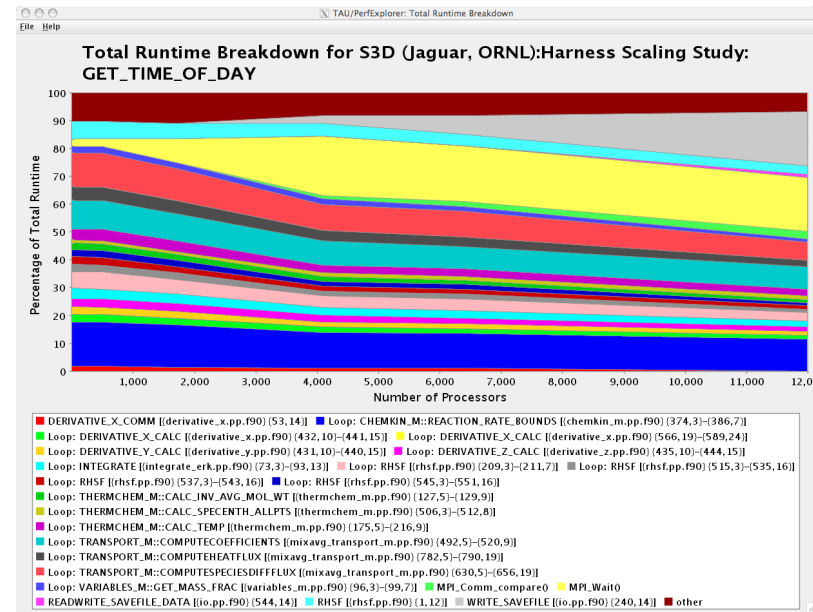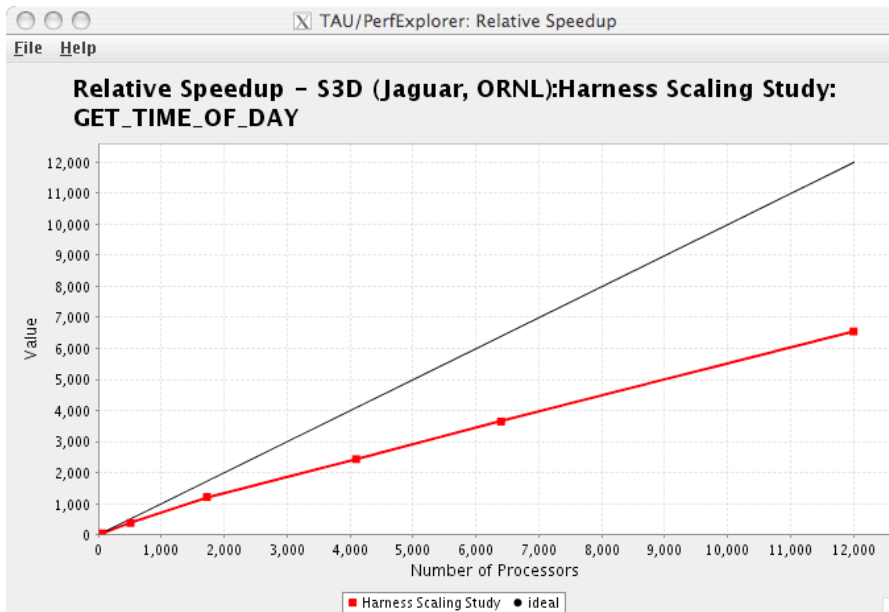**Clusters and correlations are visible**
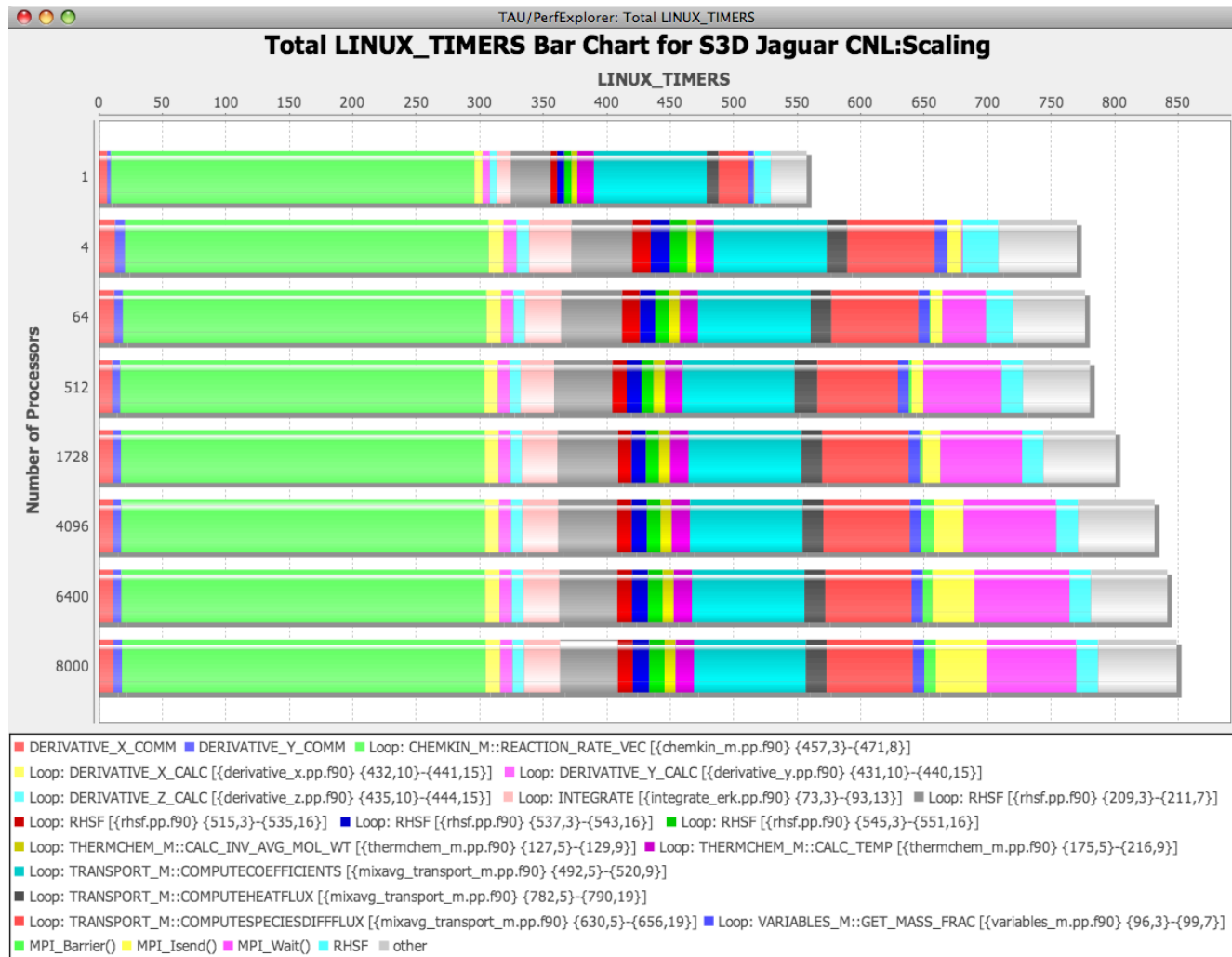
# PerfExplorer – Performance Regression

# Evaluate Scalability

**Goal: How does my application scale? What bottlenecks at what CPU counts?**

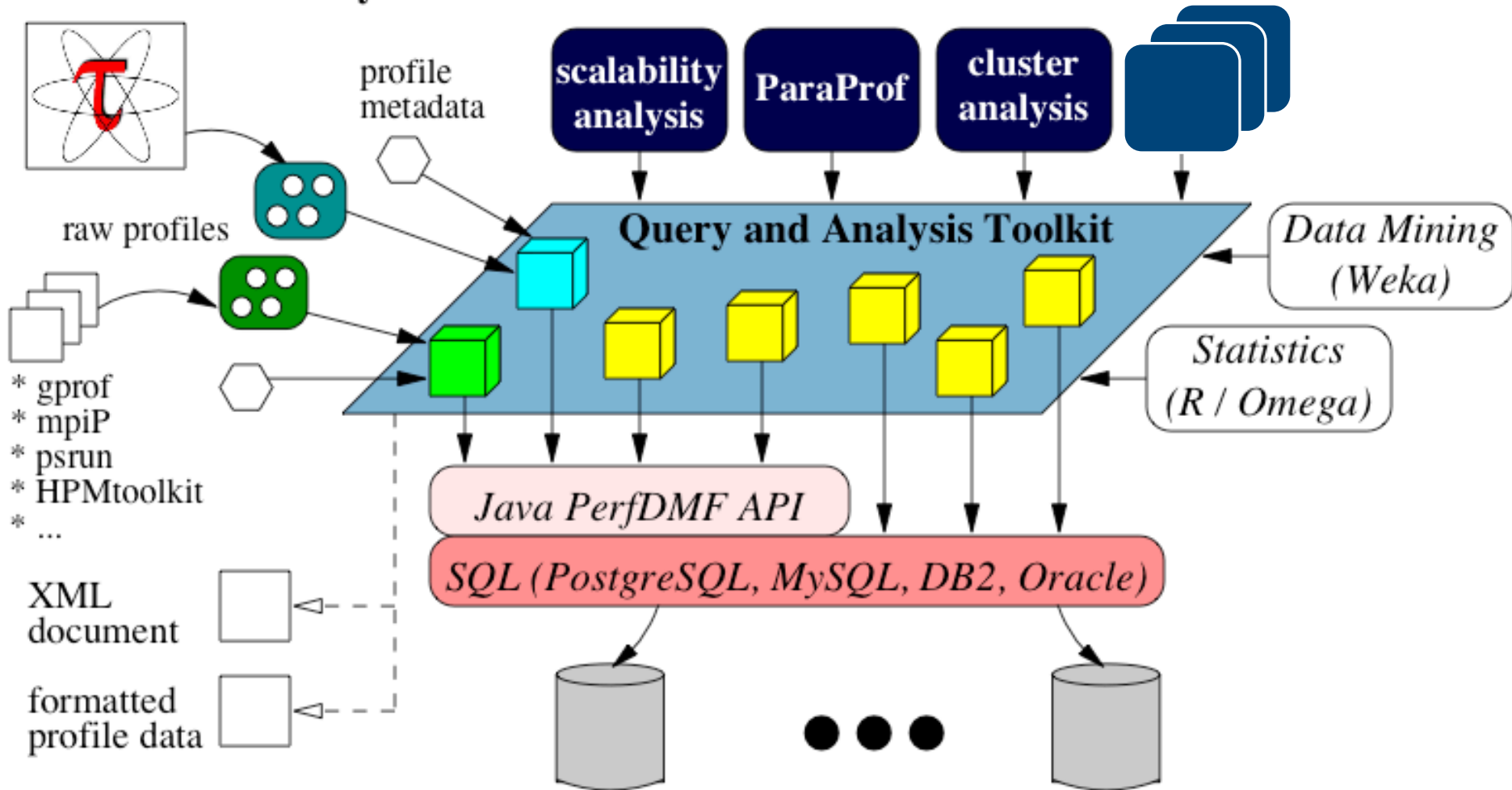**Load profiles in PerfDMF database and examine with PerfExplorer**

# Usage Scenarios: Evaluate Scalability

# TAUdb: Framework for Managing Performance Data

# Evaluate Scalability using PerfExplorer Charts

```
% export TAU_MAKEFILE=$TAU_MAKEFILE_BASE-icpc-papi-mpi-pdt
% make F90=tau_f90.sh
(Or edit Makefile and change F90=tau_f90.sh)
% qsub  run1p.job
% paraprof --pack 1p.ppk
% qsub  run2p.job …
% paraprof --pack 2p.ppk … and so on.
On your client:
% taudb_configure --create-default
% perfexplorer_configure
(Enter, y to load schema, defaults)
% paraprof
(load each trial: DB -> Add Trial -> Type (Paraprof Packed
Profile) -> OK, OR use taudb_loadtrial on the commandline)
% taudb_loadtrial –a App –x MyExp –n 4p 4p.ppk
% perfexplorer
(Charts -> Speedup)
OR:
wget http://tau.uoregon.edu/data.tgz; cat README in data
```

OREGON

# Multi-language Application Debugging

```
% export TAU_MAKEFILE=$TAU_MAKEFILE_BASE-icpc-papi-mpi-pdt
% export TAU_OPTIONS='-optMemDbg -optVerbose'
% make F90=tau_f90.sh CC=tau_cc.sh CXX=tau_cxx.sh

% export TAU_MEMDBG_PROTECT_ABOVE=1
% export TAU_MEMDBG_PROTECT_BELOW=1
% export TAU_MEMDBG_PROTECT_FREE=1
% mpirun -np 4  ./matmult
% paraprof
```

**ParaTools**

UNIVERSITY OF OREGON

# Multi-language Application Debugging

# Location of segmentation violation

ParaTools

UNIVERSITY OF OREGON

# Memory Leak Detection

```
% export TAU_MAKEFILE=$TAU_MAKEFILE_BASE-icpc-papi-mpi-pdt
% export TAU_OPTIONS='-optMemDbg -optVerbose'
% make F90=tau_f90.sh CC=tau_cc.sh CXX=tau_cxx.sh


% export TAU_TRACK_MEMORY_LEAKS=1
% mpirun -np 4   ./matmult
% paraprof
```

ParaTools

UNIVERSITY OF OREGON

# Multi-language Memory Leak Detection



| Name △ | Total | NumSamples | MaxValue | MinValue | MeanValue | Std. Dev. |
|---|---|---|---|---|---|---|
| Heap Allocate | 5,000,033 | 2 | 5,000,001 | 32 | 2,500,016.5 | 2,499,984.5 |
| Heap Allocate <file=simple.c, line=15> | 180 | 3 | 80 | 48 | 60 | 14.236 |
| Heap Allocate <file=simple.c, line=23> | 180 | 1 | 180 | 180 | 180 | 0 |
| Heap Free <file=simple.c, line=18> | 80 | 1 | 80 | 80 | 80 | 0 |
| Heap Free <file=simple.c, line=25> | 180 | 1 | 180 | 180 | 180 | 0 |
| Heap Memory Used (KB) | 4,884.829 | 8 | 4,883.196 | 0.047 | 610.604 | 1,614.888 |
| ▼ int foo(int) C [{simple.c} {36,1}–{44,1}] | | | | | | |
| ▼ int bar(int) C [{simple.c} {7,1}–{28,1}] | | | | | | |
| Heap Allocate <file=simple.c, line=23> | 180 | 1 | 180 | 180 | 180 | 0 |
| Heap Free <file=simple.c, line=25> | 180 | 1 | 180 | 180 | 180 | 0 |
| ▼ int g(int) C [{simple.c} {30,1}–{34,1}] | | | | | | |
| ▼ int bar(int) C [{simple.c} {7,1}–{28,1}] | | | | | | |
| Heap Allocate <file=simple.c, line=15> | 180 | 3 | 80 | 48 | 60 | 14.236 |
| Heap Free <file=simple.c, line=18> | 80 | 1 | 80 | 80 | 80 | 0 |
| MEMORY LEAK! Heap Allocate <file=simple.c, line=15> | 100 | 2 | 52 | 48 | 50 | 2 |
| ▼ int main(int, char **) C [{simple.c} {45,1}–{55,1}] | | | | | | |
| ▼ MPI_Finalize() | | | | | | |
| Heap Allocate | 5,000,033 | 2 | 5,000,001 | 32 | 2,500,016.5 | 2,499,984.5 |
| MEMORY LEAK! Heap Allocate | 5,000,033 | 2 | 5,000,001 | 32 | 2,500,016.5 | 2,499,984.5 |

TAU: ParaProf: Context Events for: node 0 – memleak.ppk

ParaTools

UNIVERSITY OF OREGON

# Extreme-scale Scientific Software Stack (E4S) https://e4s.io

- E4S is a community effort to provide open source software packages for developing, deploying, and running scientific applications on HPC platforms.

- E4S provides both source builds and containers of a broad collection of HPC software packages.

- E4S exists to accelerate the development, deployment and use of HPC software, lowering the barriers for HPC users.

- E4S provides containers and turn-key, from-source builds of 80+ popular HPC software packages:
    - MPI: MPICH and OpenMPI
    - Development tools: TAU, HPCToolkit, and PAPI
    - Math libraries: PETSc and Trilinos
    - Data and Viz tools: Adios, HDF5, and Paraview

**ParaTools**

UNIVERSITY OF OREGON

# Extreme-scale Scientific Software Stack (E4S) https://e4s.io

**Spack [http://spack.io] is the primary means for software delivery**

**SDKs: collection of related ECP ST products where coordination across package teams will improve usability and practicies, and foster community growth among teams that develop similar and complimentary capabilities. An SDK involves several products.**

**Containers of pre-built binaries of ECP ST products.**

**Container runtimes supported**

- Docker: Dockerhub: exascaleproject/sdk:AHM19
- Charliecloud
- Shifter
- Singularity
- Inception at NCAR

**VirtualBox Open Virtualization Appliance (OVA) image that contains these runtimes**

**MPI replacement strategies to use native network interconnect**

# Spack

- **E4S uses the Spack package manager for software delivery**
- **Spack provides the ability to specify versions of software packages that are and are not interoperable.**
- **Spack is a build layer for not only E4S software, but also a large collection of software tools and libraries outside of ECP ST.**
- **Spack supports achieving and maintaining interoperability between ST software packages.**

ParaTools

O | UNIVERSITY OF OREGON

# The Spack community is growing rapidly

**Spack simplifies HPC software for:**

- Users
- Developers
- Cluster installations
- The largest HPC facilities

**Spack is central to ECP's software strategy**

- Enable software reuse for developers and users
- Allow the facilities to consume the entire ECP stack

**The roadmap is packed with new features:**

- Building the ECP software distribution
- Better workflows for building containers
- Stacks for facilities
- Chains for rapid dev workflow
- Optimized binaries
- Better dependency resolution

**Visit spack.io**

github.com/spack/spack

**@spackpm**

ParaTools

UNIVERSITY OF OREGON

# E4S VirtualBox OVA image

**Contains all four container runtimes and the E4S Singularity image!**

- **Docker**
- **Singularity**
- **Shifter**
- **Charliecloud**

# Docker container of E4S

```
% docker pull exascaleproject/sdk:AHM19

Using USB stick or images from https://e4s.io:

% gunzip –c ecp.tgz | docker load
% docker images
Mount home directory:
% docker –i –v $HOME:$HOME –t exascaleproject/sdk:AHM19 /bin/bash
% which spack
% cp –r /usr/local/packages/ecp/demo . ; cd demo; cat README
```

ParaTools

UNIVERSITY OF OREGON

# Running MPI applications on other systems

- **Applications built with MPI in the E4S container can replace the MPI in the container with the system MPI!**

- **This allows fast inter-node communication using the native interconnect.**

- **Application and data are external to the E4S container.**

- **Programming models, compilers, runtime libraries, and tools are inside the container.**

- **We can replace MPI using the MPICH ABI compatibility layer.**

- **Goal: Build an MPI binary once and run it un-modified on all HPC Linux x86_64 clusters!**

# Using E4S Singularity Image on Quartz

**MVAPICH2 needs /lib. Mount it as /hostlib64 and add it to LD_LIBRARY_PATH**

```
% salloc –N 2 –p pdebug
% /usr/workspace/ecpsdk/apps/Zoltan/run_sing.sh
% paraprof &
% cat /usr/workspace/ecpsdk/apps/Zoltan/run_sing.sh
#!/bin/bash
srun –n 8 –N 2 --ntasks-per-node=4 singularity exec  -B /usr/workspace/ecpsdk:/usr/workspace/
ecpsdk -B /lib64:/hostlib64 -B $SLURM_SUBMIT_DIR:$SLURM_SUBMIT_DIR -B /usr/tce:/usr/tce /usr/
workspace/ecpsdk/ecp.simg /bin/bash -c ' . /etc/bashrc ; spack load   trilinos hypre parmetis
hdf5 metis openblas superlu zlib netcdf matio boost@1.66.0 scalapack suite-sparse tau; spack
unload openmpi mpich; export LD_LIBRARY_PATH=/usr/tce/packages/mvapich2/mvapich2-2.3-
intel-19.0.4/lib:$LD_LIBRARY_PATH:/hostlib64; /usr/workspace/ecpsdk/apps/Zoltan/Zoltan'
echo "Running with TAU:"
echo "----------------"
srun –n 8 –N 2 --ntasks-per-node=4 singularity exec  -B /usr/workspace/ecpsdk:/usr/workspace/
ecpsdk -B /lib64:/hostlib64 -B $SLURM_SUBMIT_DIR:$SLURM_SUBMIT_DIR -B /usr/tce:/usr/tce /usr/
workspace/ecpsdk/ecp.simg /bin/bash -c ' . /etc/bashrc ; spack load   trilinos hypre parmetis
hdf5 metis openblas superlu zlib netcdf matio boost@1.66.0 scalapack suite-sparse tau; spack
unload openmpi mpich; export LD_LIBRARY_PATH=/usr/tce/packages/mvapich2/mvapich2-2.3-
intel-19.0.4/lib:$LD_LIBRARY_PATH:/hostlib64; tau_exec -ebs /usr/workspace/ecpsdk/apps/Zoltan/
Zoltan'
```

ParaTools

# Replacing MPI using Shifter on Cori.nersc.gov

```
% shifterimg images
exascaleproject/sdk:AHM19 …
% To replace MPI with system MPI:
# salloc -N 2 -q interactive -t 00:30:00 \
--image=exascaleproject/sdk:AHM19 -C haswell -L SCRATCH

# cat ~sameer/run_shifter.sh
srun –n 32 shifter  -- /bin/bash  -c 'unset CRAYPE_VERSION;  . /
etc/bashrc ; spack load   trilinos hypre parmetis hdf5 metis
openblas superlu zlib netcdf matio boost@1.66.0 scalapack suite-
sparse tau; spack unload openmpi mpich; ./Zoltan'
# ~sameer/run_shifter.sh
```

ParaTools

# PRL, University of Oregon, Eugene



www.uoregon.edu

ParaTools

# Support Acknowledgments

**US Department of Energy (DOE)**

- ANL
- Office of Science contracts, ECP
- SciDAC, LBL contracts
- LLNL-LANL-SNL ASC/NNSA contract
- Battelle, PNNL and ORNL contract

**CEA, France**

**Department of Defense (DoD)**

- PETTT, HPCMP

**National Science Foundation (NSF)**

- SI2-SSI, Glassbox

**Intel Corporation**

**NASA**

**Partners:**

- University of Oregon
- The Ohio State University
- ParaTools, Inc.
- University of Tennessee, Knoxville
- T.U. Dresden, GWT
- Jülich Supercomputing Center

# Acknowledgment

ParaTools

UNIVERSITY OF OREGON

# Download TAU from U. Oregon



**http://tau.uoregon.edu**

**http://taucommander.com**

**http://www.hpclinux.com** **[OVA for VirtualBox]**
**https://e4s.io** **[ Extreme-Scale Scientific Software Stack, Containers for HPC]**

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

# Reference

ParaTools

UNIVERSITY OF OREGON

# Installing and Configuring TAU

- **Installing PDT:**
  - wget tau.uoregon.edu/pdt_lite.tgz
  - ./configure –prefix=<dir>; make ; make install

- **Installing TAU:**
  - wget tau.uoregon.edu/tau.tgz;   tar zxf tau.tgz; cd tau-2.<ver>
  - wget http://tau.uoregon.edu/ext.tgz ; tar xf ext.tgz
  - ./configure -bfd=download -pdt=<dir> -papi=<dir> -mpi –pthread –c++=mpicxx –cc=mpicc –fortran=mpif90 –dwarf=download –unwind=download –otf=download –iowrapper –papi=<dir>
  - make install

- **Using TAU:**
  - export TAU_MAKEFILE=<taudir>/x86_64/lib/Makefile.tau-<TAGS>
  - make CC=tau_cc.sh   CXX=tau_cxx.sh   F90=tau_f90.sh

**ParaTools**

UNIVERSITY OF OREGON

# Compile-Time Options

**Optional parameters for the TAU_OPTIONS environment variable:**
**% tau_compiler.sh**

| | |
|---|---|
| -optVerbose | Turn on verbose debugging messages |
| -optCompInst | Use compiler based instrumentation |
| -optNoCompInst | Do not revert to compiler instrumentation if source instrumentation fails. |
| -optTrackIO | Wrap POSIX I/O call and calculates vol/bw of I/O operations (Requires TAU to be configured with *–iowrapper*) |
| -optTrackGOMP | Enable tracking GNU OpenMP runtime layer (used without –opari) |
| -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 *tau_instrumentor* |
| -optTauWrapFile="<file>" | Specify path to *link_options.tau* generated by *tau_gen_wrapper* |
| -optHeaderInst | Enable Instrumentation of headers |
| -optTrackUPCR | Track UPC runtime layer routines (used with tau_upc.sh) |
| -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) … |

**ParaTools**

# Compile-Time Options (contd.)

**Optional parameters for the TAU_OPTIONS environment variable:**
**% tau_compiler.sh**

| | |
|---|---|
| -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.

…

ParaTools

UNIVERSITY OF OREGON

# 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>) |

**ParaTools**

O | UNIVERSITY OF OREGON

# 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 | 1 | Setting to 1 is necessary for event based sampling to resolve addresses with OMPT. Setting to 0 allows the user to do offline address translation. |

ParaTools

UNIVERSITY OF OREGON

# 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_ALINGMENT | 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 |

**ParaTools**

UNIVERSITY OF OREGON