

Improving Application Performance Using the TAU Performance System

Sameer Shende, John C. Linford
{sameer, jlinford}@paratools.com

ParaTools, Inc and University of Oregon.

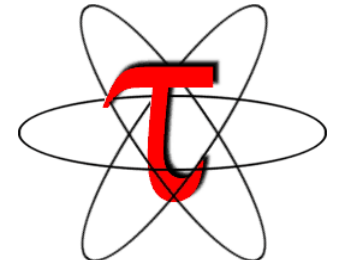
April 4-5, 2013, CG1, NCAR, UCAR

Download slides from:

<http://www.paratools.com/sea13>

TAU Performance System[®]

<http://tau.uoregon.edu/>



- **Tuning and Analysis Utilities (18+ 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)
- **Easy to integrate in application frameworks**

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?

What Can TAU Do?

- **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

- **Multi-language debugging**

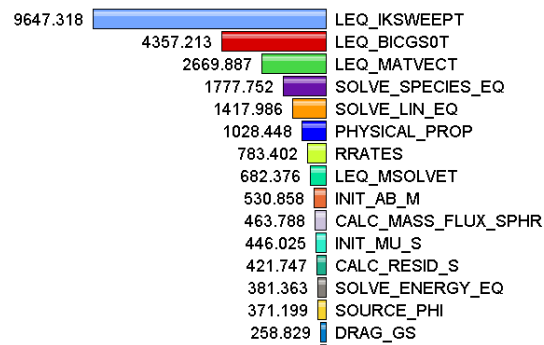
- Identify the source location of a crash by unwinding the system callstack
- Identify memory errors (off-by-one, etc.)

- Profiling and tracing can measure **time** as well as **hardware performance counters** (cache misses, instructions) from your CPU
- TAU can **automatically instrument** your source code using a package called PDT for routines, loops, I/O, memory, phases, etc.
- TAU runs on **all HPC platforms** and it is free (BSD style license)
- TAU includes instrumentation, measurement and analysis tools

Profiling and Tracing

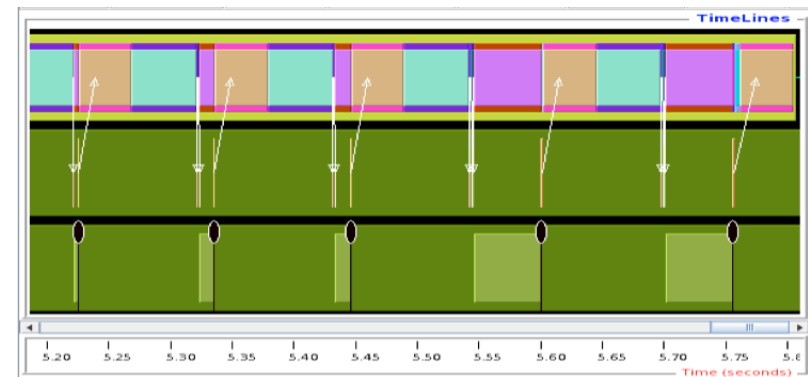
Profiling

Value: Exclusive
Units: seconds



- **Profiling** shows you how much (total) time was spent in each routine
- Metrics can be time or hardware performance counters (cache misses, instructions)
- TAU can automatically instrument your source code using a package called PDT for routines, loops, I/O, memory, phases, etc.

Tracing



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

What does TAU support?

C/C++

CUDA UPC

OpenCL

Python

Fortran

OpenACC

GPI

Java MPI

pthread

Intel MIC

OpenMP

Intel GNU

LLVM

PGI

Cray

Sun

MinGW

Linux

Windows

AIX

Insert
yours
here

BlueGene

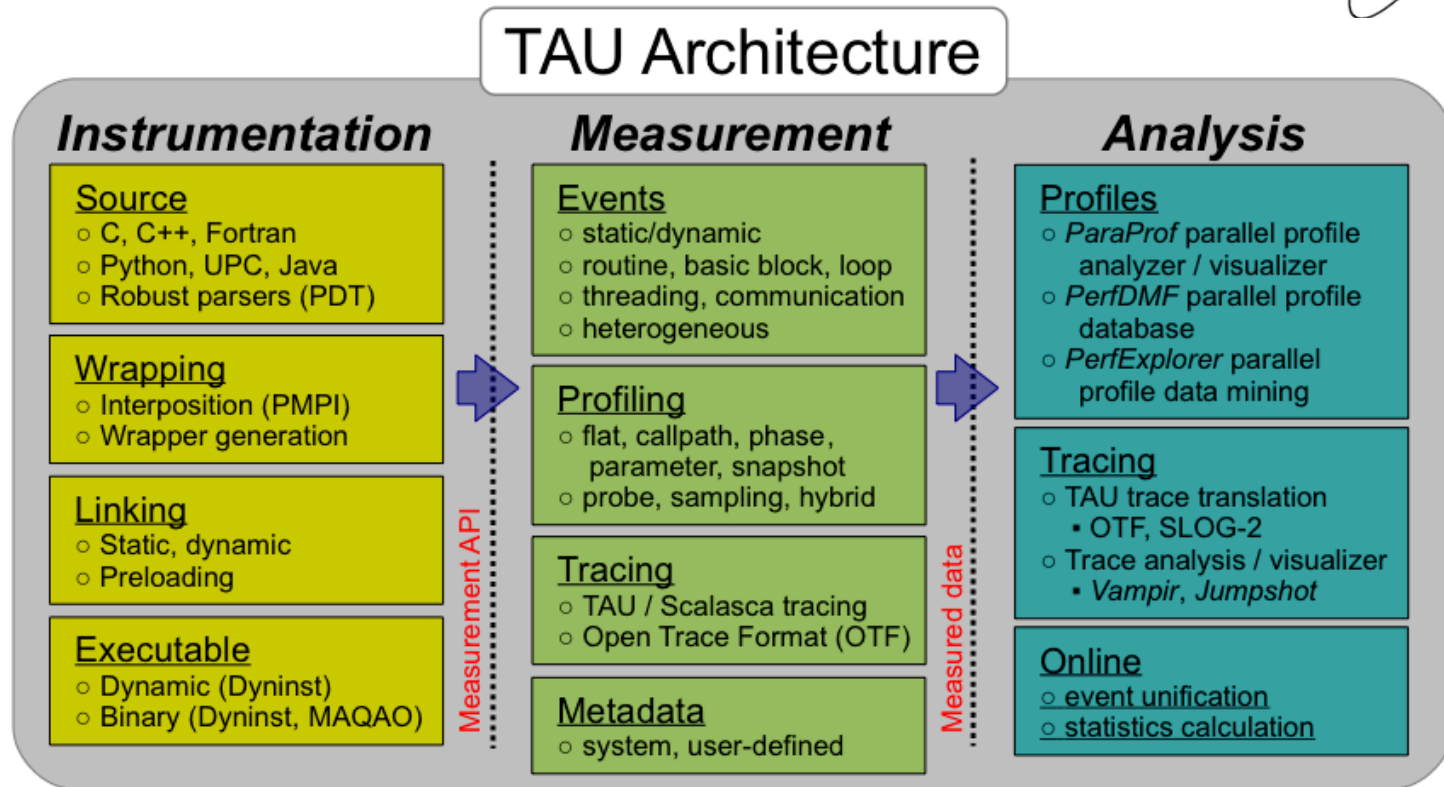
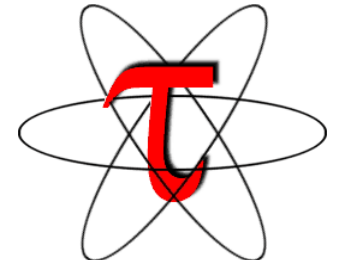
Fujitsu

ARM

NVIDIA Kepler

OS X

The TAU Architecture



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)

Direct Observation Events

Interval events (begin/end events)

- Measures exclusive & inclusive durations between events
- Metrics monotonically increase
- Example: Wall-clock timer

Atomic events (trigger with data value)

- Used to capture performance data state
- Shows extent of variation of triggered values (min/max/mean)
- Example: heap memory consumed at a particular point

Code events

- Routines, classes, templates
- Statement-level blocks, loops
- Example: for-loop begin/end

Interval and Atomic Events in TAU

NODE 0;CONTEXT 0;THREAD 0:

%Time	Exclusive msec	Inclusive total msec	#Call	#Subrs	Inclusive Name usec/call
100.0	0.187	1.105	1	44	1105659 int main(int, char **) C
93.2	1.030	1.030	1	0	1030654 MPI_Init()
5.9	0.879	65	40	320	1637 void func(int, int) C
4.6	51	51	40	0	1277 MPI_Barrier()
1.2	13	13	120	0	111 MPI_Recv()
0.8	9	9	1	0	9328 MPI_Finalize()
0.0	0.137	0.137	120	0	1 MPI_Send()
0.0	0.086	0.086	40	0	2 MPI_Bcast()
0.0	0.002	0.002	1	0	2 MPI_Comm_size()
0.0	0.001	0.001	1	0	1 MPI_Comm_rank()

Interval events show **duration**

USER EVENTS Profile :NODE 0, CONTEXT 0, THREAD 0

NumSamples	MaxValue	MinValue	MeanValue	Std. Dev.	Event Name
365	5.138E+04	44.39	3.09E+04	1.234E+04	Heap Memory Used (KB) : Entry
365	5.138E+04	2064	3.115E+04	1.21E+04	Heap Memory Used (KB) : Exit
40	40	40	40	0	Message size for broadcast

Atomic events (triggered with value) show **extent of variation** (min/max/mean)

27.1 1%

% export TAU_CALLPATH_DEPTH=0
 % export TAU_TRACK_HEAP=1

Context Events with Callpath

NODE 0:CONTEXT 0:THREAD 0:					
%Time	Exclusive msec	Inclusive total msec	#Call	#Subrs	Inclusive Name usec/call
100.0	0.357	1.114	1	44	1114040 int main(int, char **) C
92.6	1.031	1.031	1	0	1031066 MPI_Init()
6.7	72	74	40	320	1865 void func(int, int) C
0.7	8	8	1	0	8002 MPI_Finalize()
0.1	1	1	120	0	12 MPI_Recv()
0.1	0.608	0.608	40	0	15 MPI_Barrier()
0.0	0.136	0.136	120	0	1 MPI_Send()
0.0	0.095	0.095	40	0	2 MPI_Bcast()
0.0	0.001	0.001	1	0	1 MPI_Comm_size()
0.0	0	0	1	0	0 MPI_Comm_rank()

USER EVENTS Profile :NODE 0, CONTEXT 0, THREAD 0					
NumSamples	MaxValue	MinValue	MeanValue	Std. Dev.	Event Name
365	5.139E+04	44.39	3.091E+04	1.234E+04	Heap Memory Used (KB) : Entry
1	44.39	44.39	44.39	0	Heap Memory Used (KB) : Entry : int main(int, char **) C
1	2068	2068	2068	0	Heap Memory Used (KB) : Entry : int main(int, char **) C => MPI_Comm_rank()
1	2066	2066	2066	0	Heap Memory Used (KB) : Entry : int main(int, char **) C => MPI_Comm_size()
1	5.139E+04	5.139E+04	5.139E+04	0	Heap Memory Used (KB) : Entry : int main(int, char **) C => MPI_Finalize()
1	57.58	57.58	57.58	0	Heap Memory Used (KB) : Entry : int main(int, char **) C => MPI_Init()
40	5.036E+04	2069	3.011E+04	1.228E+04	Heap Memory Used (KB) : Entry : int main(int, char **) C => void func(int, int) C
40	5.139E+04	3098	3.114E+04	1.227E+04	Heap Memory Used (KB) : Entry : void func(int, int) C => MPI_Barrier()
40	5.139E+04	1.13E+04	3.134E+04	1.187E+04	Heap Memory Used (KB) : Entry : void func(int, int) C => MPI_Bcast()
120	5.139E+04	1.13E+04	3.134E+04	1.187E+04	Heap Memory Used (KB) : Entry : void func(int, int) C => MPI_Recv()
120	5.139E+04	1.13E+04	3.134E+04	1.187E+04	Heap Memory Used (KB) : Entry : void func(int, int) C => MPI_Send()
365	5.139E+04	2065	3.116E+04	1.21E+04	Heap Memory Used (KB) : Exit

```
% export TAU_CALLPATH_DEPTH=2
% export TAU_TRACK_HEAP=1
```

Callpath shown on context events

Direct Instrumentation Options in TAU

Source Code Instrumentation

- Automatic instrumentation using pre-processor based on static analysis of source code (PDT), creating an instrumented copy
- Compiler generates instrumented object code
- Manual instrumentation

Library Level Instrumentation

- Statically or dynamically linked wrapper libraries
 - MPI, I/O, memory, etc.
- Wrapping external libraries where source is not available

Runtime pre-loading and interception of library calls

Binary Code instrumentation

- Rewrite the binary, runtime instrumentation

Virtual Machine, Interpreter, OS level instrumentation

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 $<
```

Binary Rewriting Instrumentation

- Support for **Intel, PGI, and GNU** compilers
- Specify a list of routines to instrument
- Specify the TAU measurement library to be injected
- **DyninstAPI:**

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

- **MAQAO:**

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

- **Pebil:**

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

- Execute the application to get measurement data:

```
% mpiexec ./a.inst
```

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

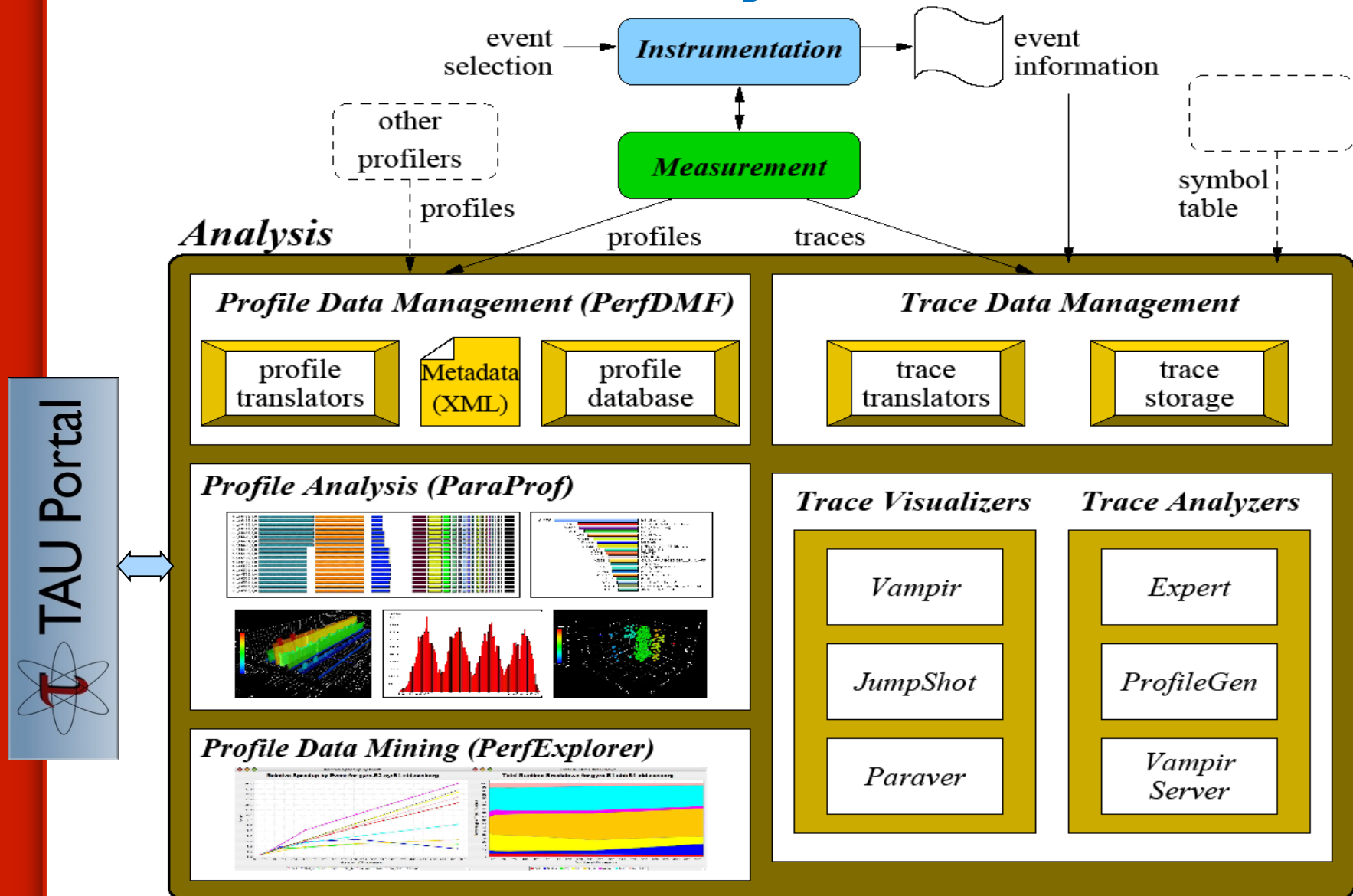
Preloading a library at runtime

- 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*

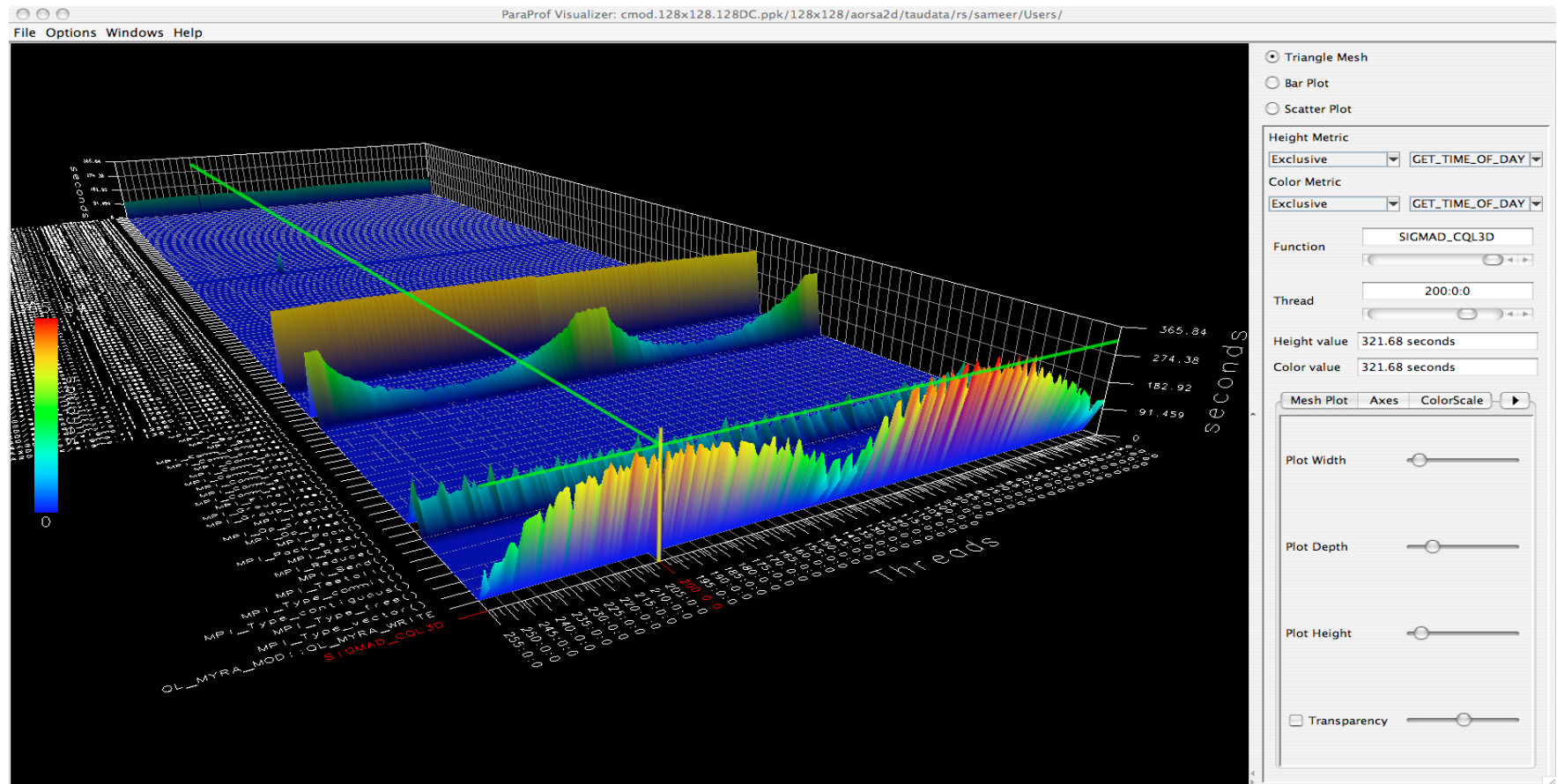
Linker based substitution

- 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

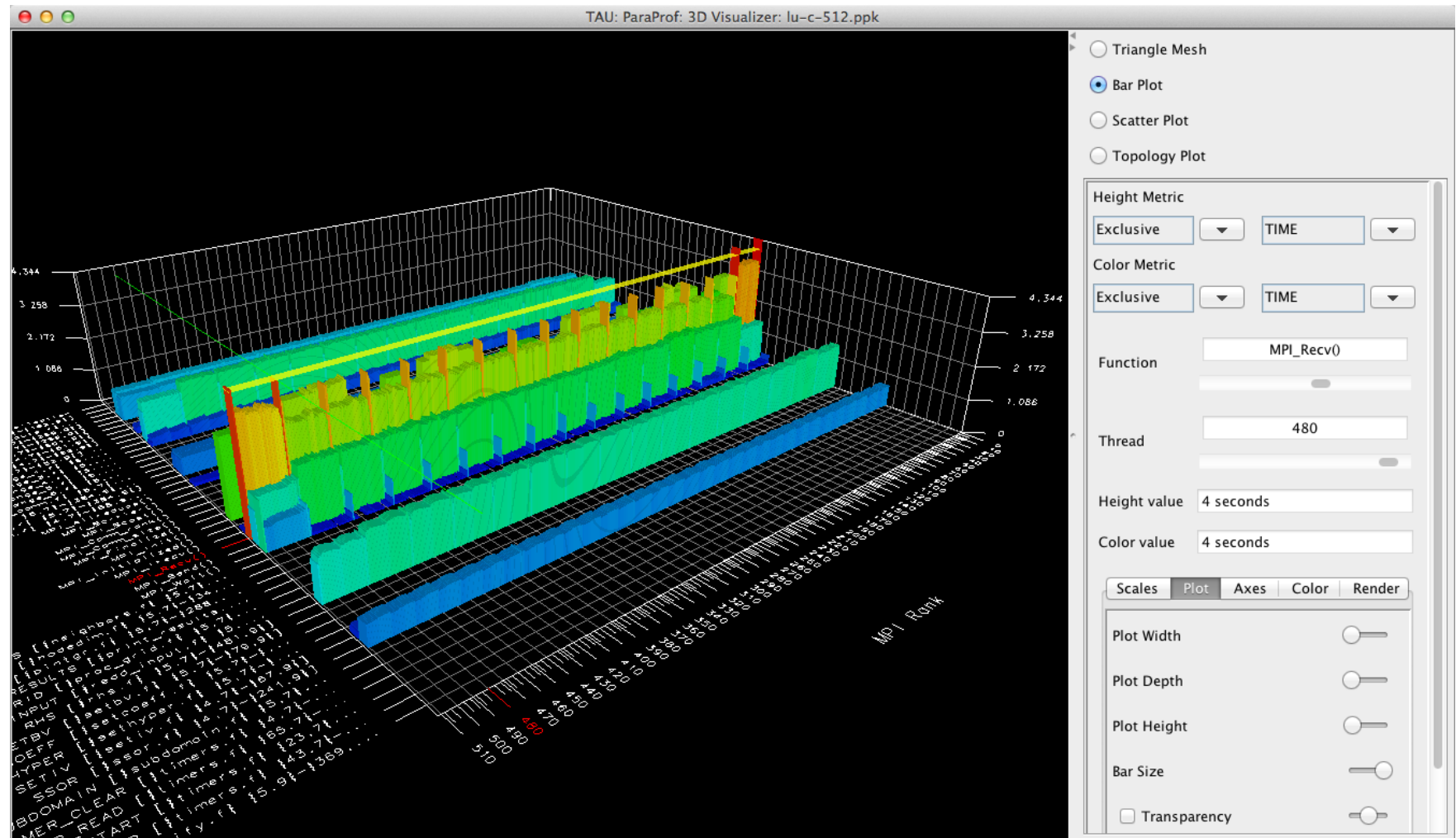
Performance Analysis



ParaProf 3D Profile Browser

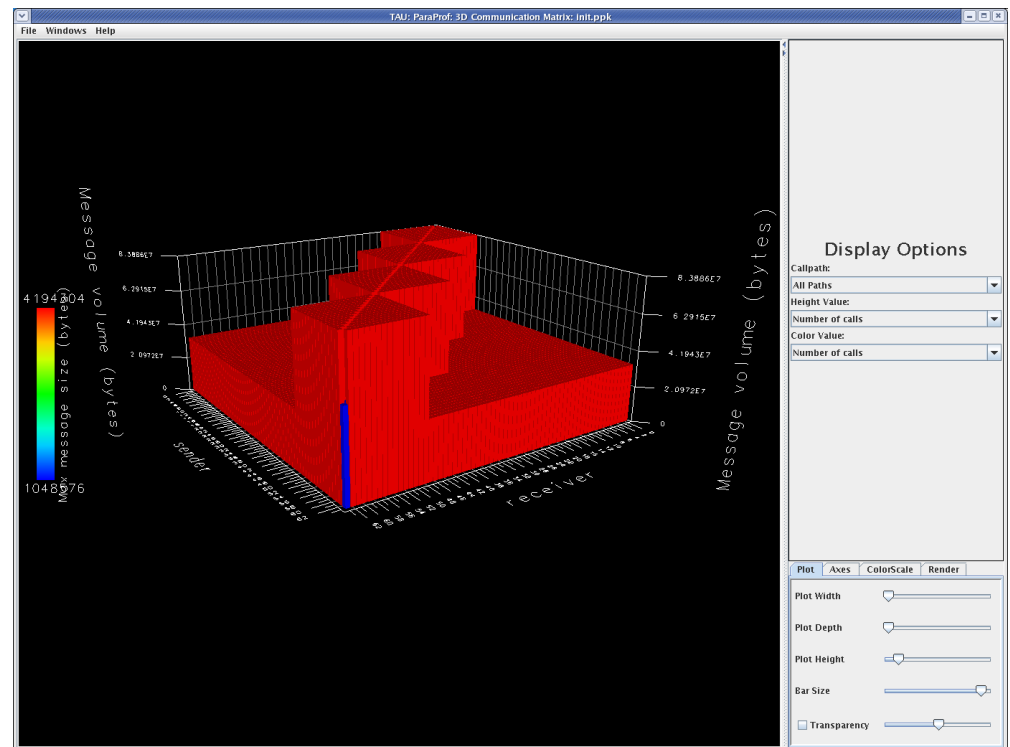
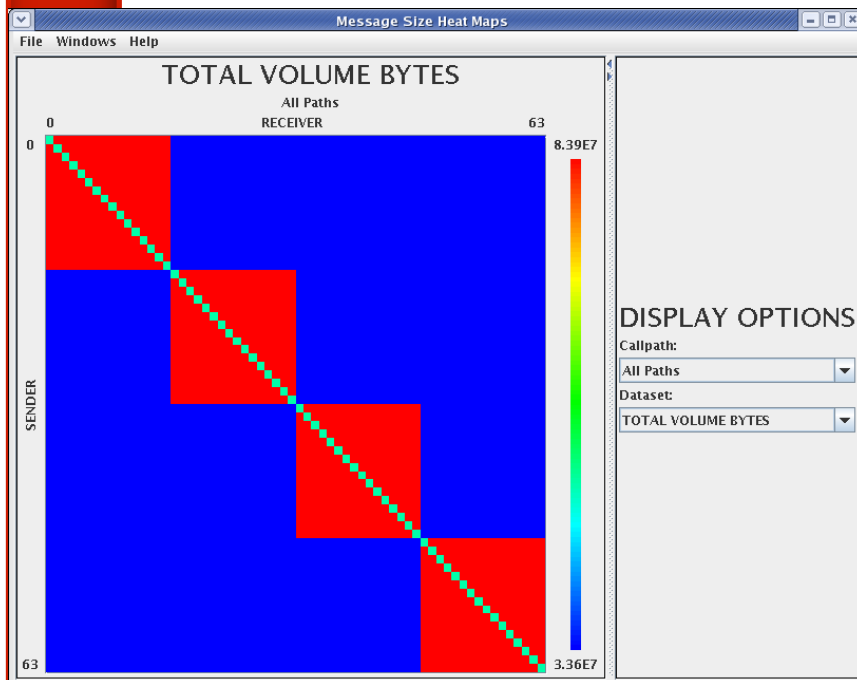


ParaProf 3D Profile Browser

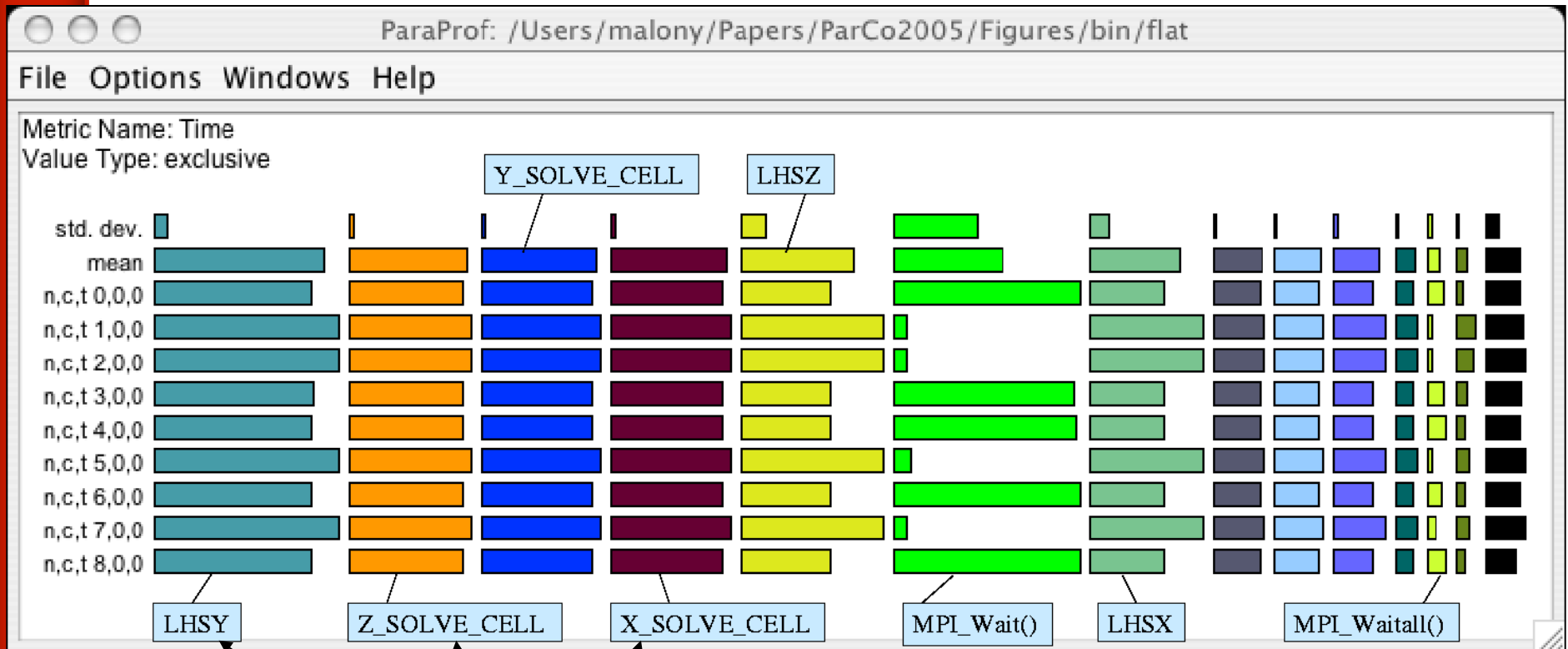


Communication Matrix Display

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



NAS BT – Flat Profile



Application routine names
reflect phase semantics

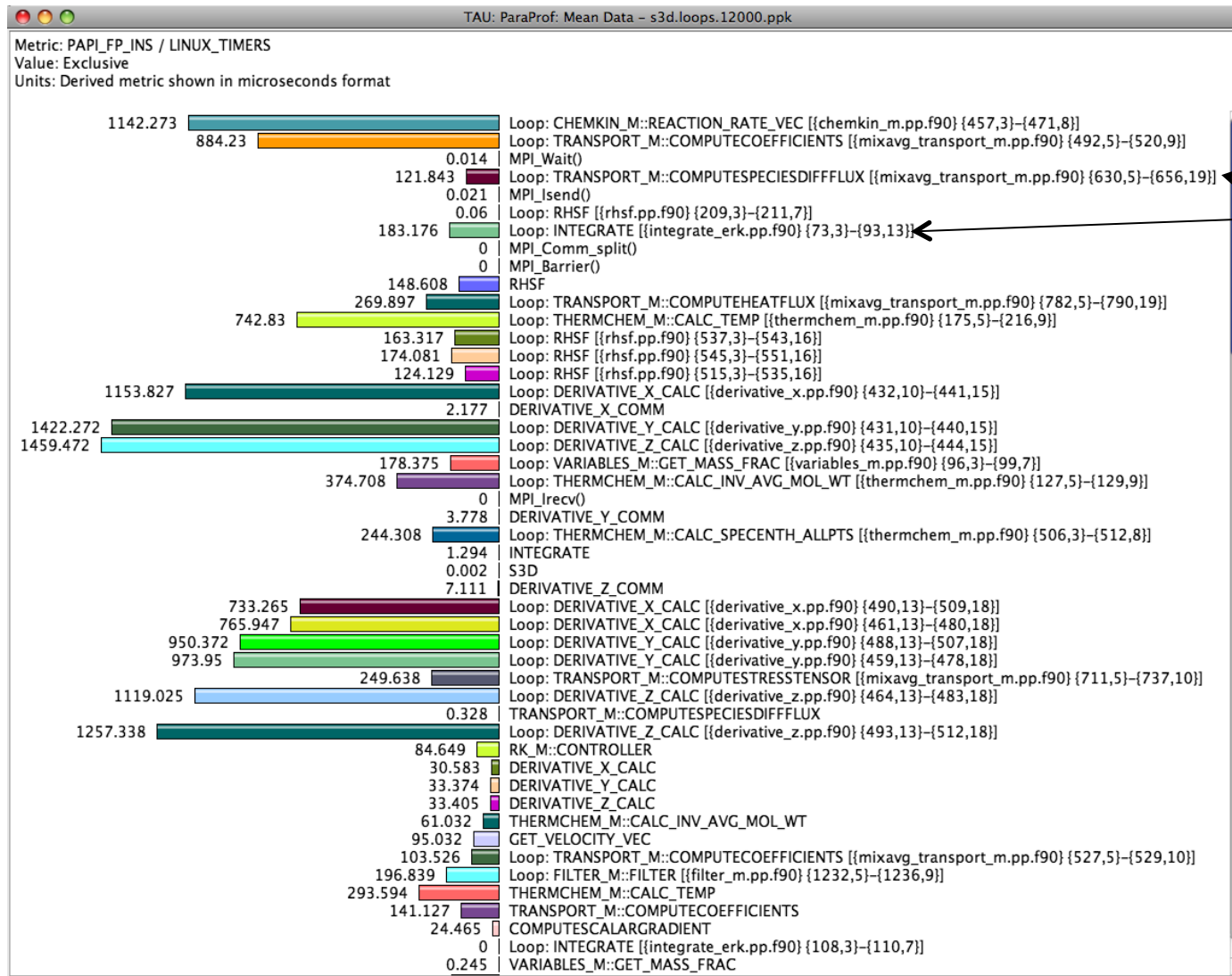
How is MPI_Wait()
distributed relative to
solver direction?

NAS BT – Phase Profile

Main phase shows nested phases and immediate events



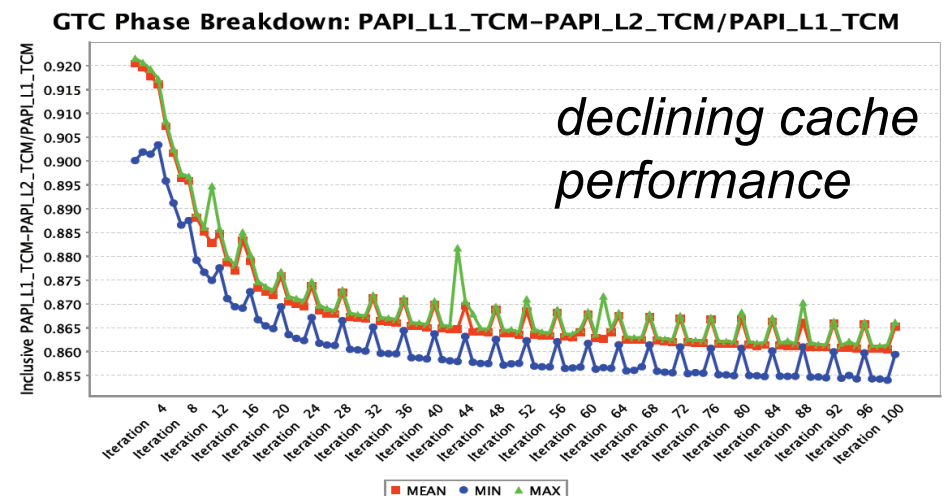
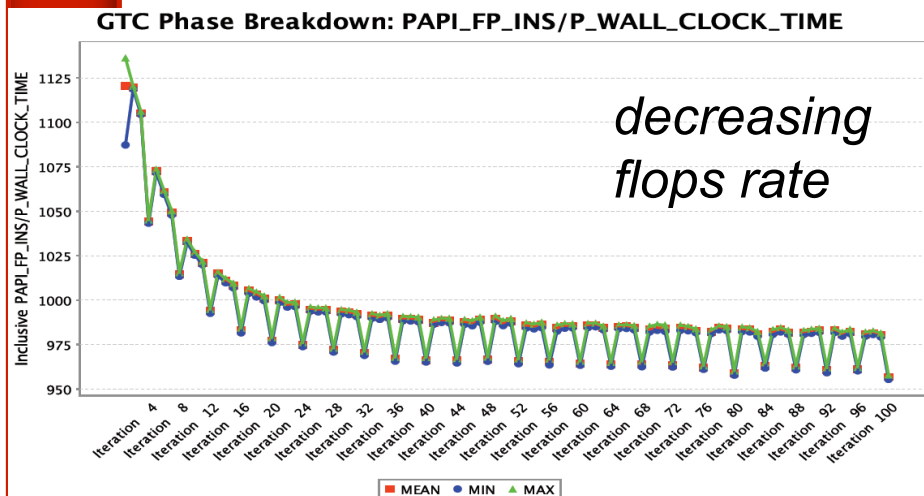
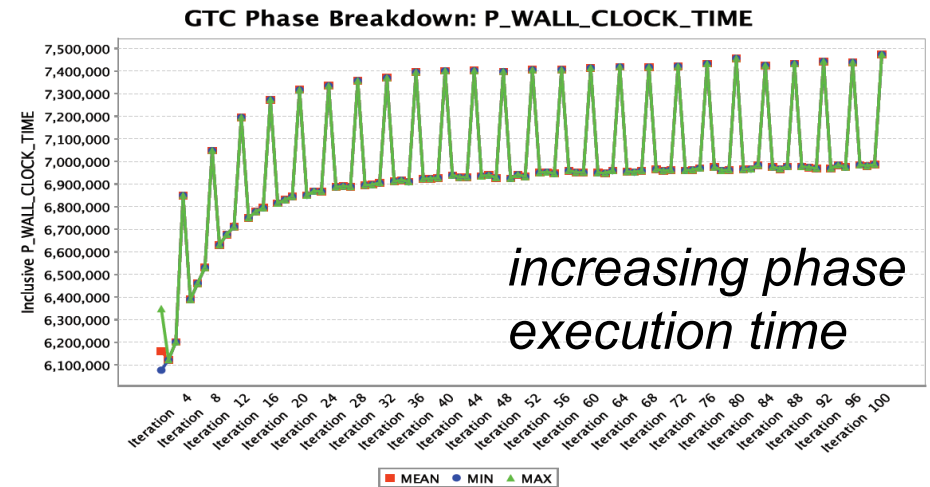
Derived Metrics Help Identify Potential Bottlenecks



Low MFLOPS
in loops?

Phase Profiling of HW Counters

- GTC particle-in-cell simulation of fusion turbulence
- Phases assigned to iterations
- Poor temporal locality for one important data
- Automatically generated by PE2 python script



Using TAU: Simplest Case

Uninstrumented code:

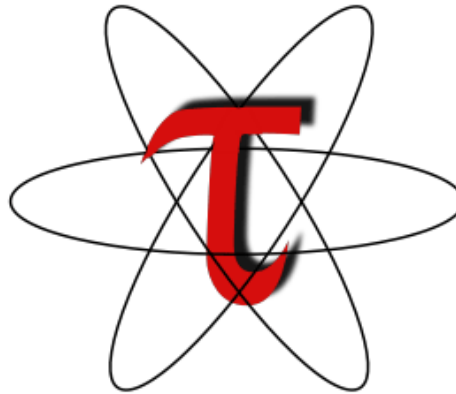
- % mpirun.lsf ./a.out

With TAU:

- % module load workshop tau
- % mpirun.lsf **tau_exec** ./a.out
- % paraprof

**Come to the tutorial
Thursday and Friday!**

Download TAU from U. Oregon



<http://tau.uoregon.edu>

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

Free download, open source, BSD license

Support Acknowledgments

US Department of Energy (DOE)

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

Department of Defense (DoD)

- PETTT, HPCMP

National Science Foundation (NSF)

- Glassbox, SI-2

University of Tennessee, Knoxville

T.U. Dresden, GWT

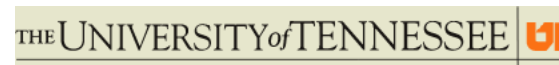
Juelich Supercomputing Center

**And a special
thanks to UCAR!**

ParaTools



UNIVERSITY
OF OREGON



UNIVERSITY OF OREGON