

Efficient Parallel Runtime Bounds Checking with the TAU Performance System

John C. Linford, Sameer Shende, Allen D. Malony
{jlinford, sameer, malony}@paratools.com
ParaTools, Inc., Eugene, OR

Andrew Wissink
U.S. Army Aviation Development Directorate - AFDD
NASA Ames Research Center, Moffett Field, CA

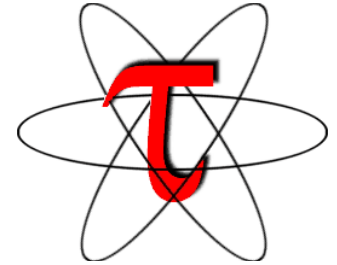
10-12 September 2013
Waltham, MA

Download slides from:

<http://www.paratools.com/hpec13>

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)
- **Integrates with 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 does TAU support?

C/C++

CUDA UPC

OpenCL

Fortran

OpenACC

Python
GPI

pthread

Intel MIC

Java MPI

OpenMP

Intel GNU

LLVM PGI Cray Sun

MinGW

Linux Windows AIX

Insert
yours
here

BlueGene Fujitsu ARM

NVIDIA Kepler OS X

Availability on New Systems

- Intel compilers with Intel MPI on Intel Xeon Phi™ (MIC)
- GPI with Intel Linux x86_64 Infiniband clusters
- IBM BG/Q and Power 7 Linux with IBM XL UPC compilers
- NVIDIA Kepler K20 with CUDA 5.0 with NVCC
- Fujitsu Fortran/C/C++ MPI compilers on the K computer
- PGI compilers with OpenACC support on NVIDIA systems
- Cray CX30 Sandybridge Linux systems with Intel compilers
- Cray CCE compilers with OpenACC support on Cray XK7
- AMD OpenCL libs with GNU on AMD Fusion cluster systems
- MPC compilers on TGCC Curie system (Bull, Linux x86_64)
- GNU compilers on ARM Linux clusters (MontBlanc, BSC)
- Cray CCE compilers with OpenACC on Cray XK6 (K20)
- Microsoft MPI with Mingw compilers under Windows Azure
- LLVM and GNU compilers under Mac OS X, IBM BGQ

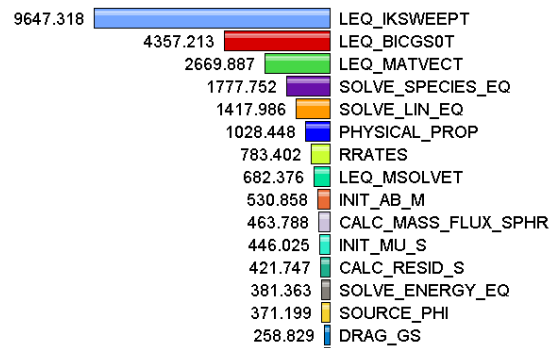
What Can TAU Do?

- **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

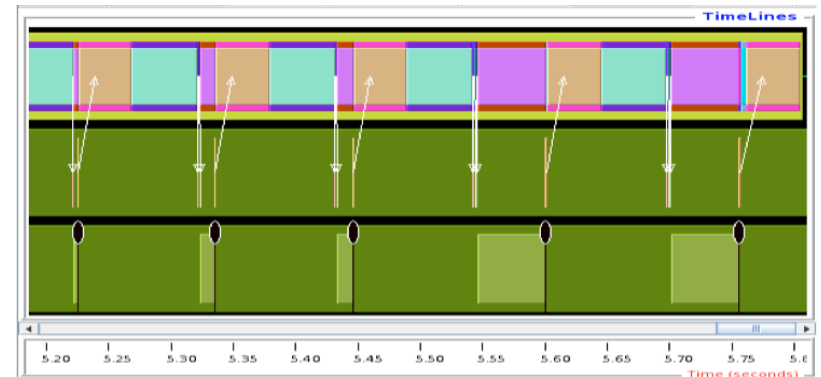
Profiling and Tracing

Profiling

Value: Exclusive
Units: seconds



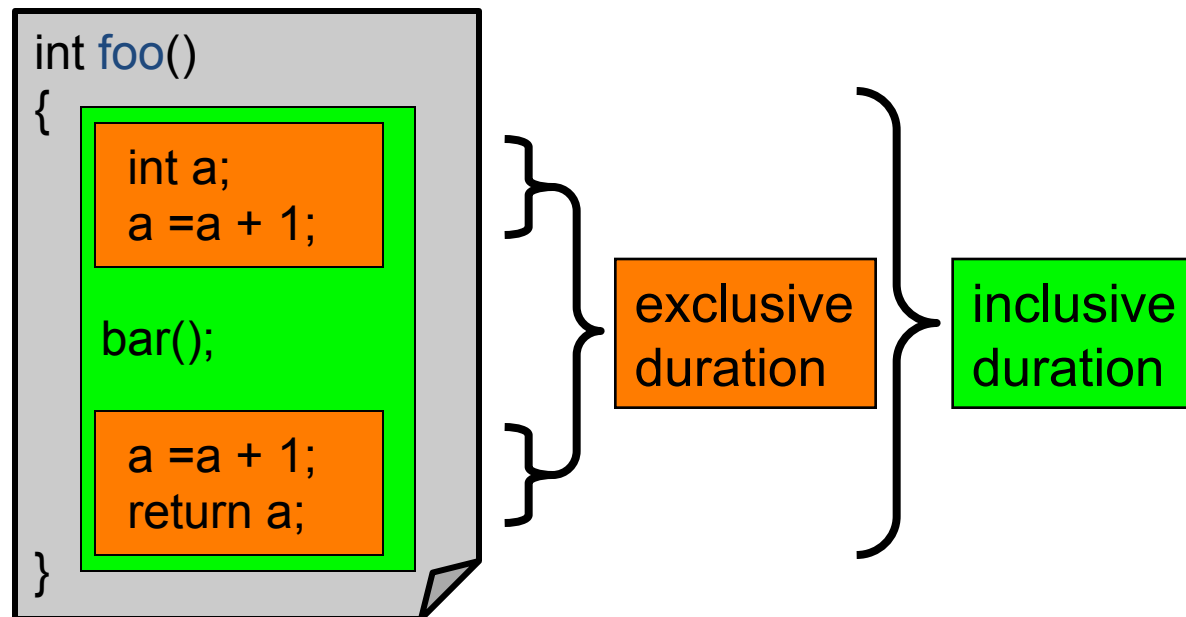
Tracing



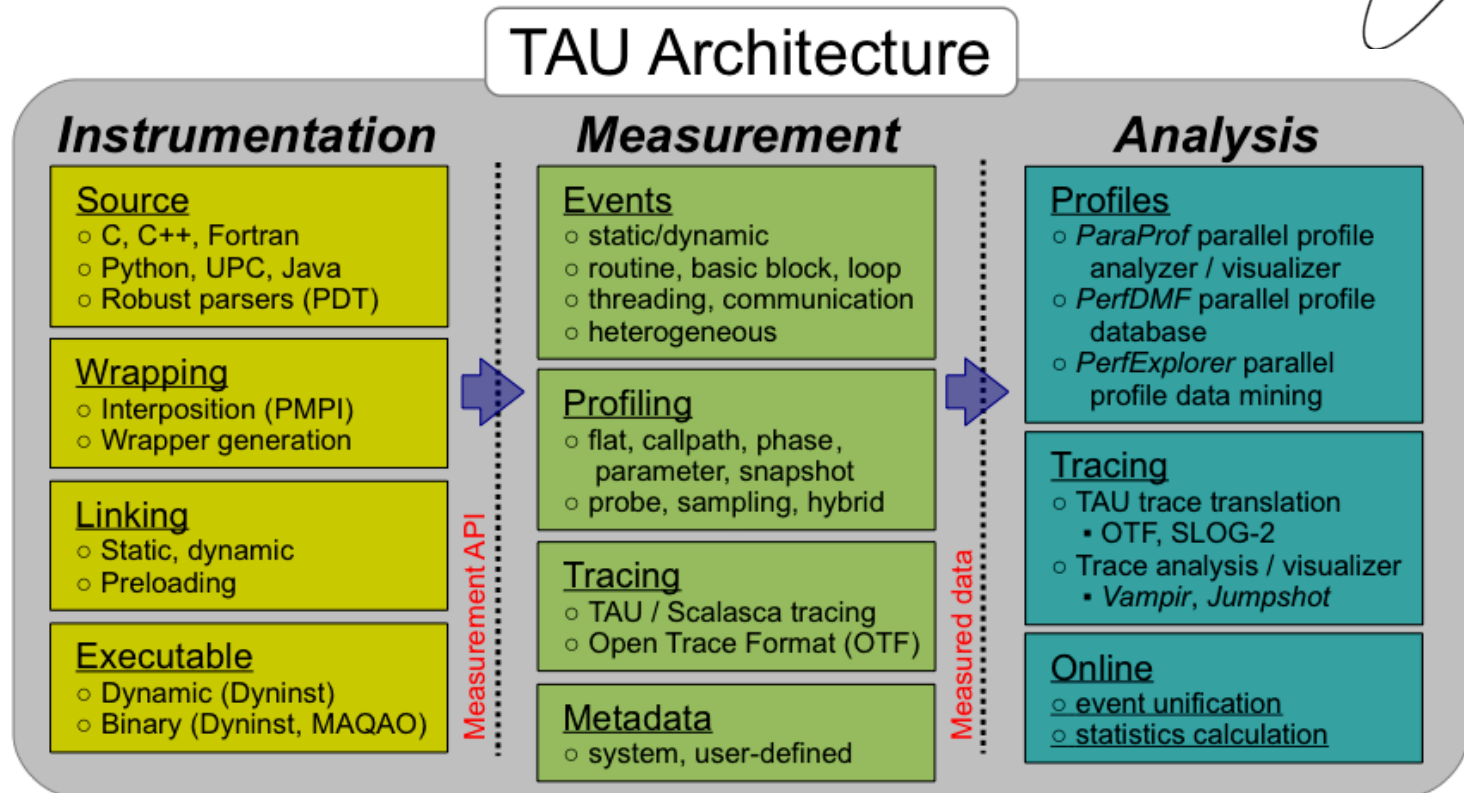
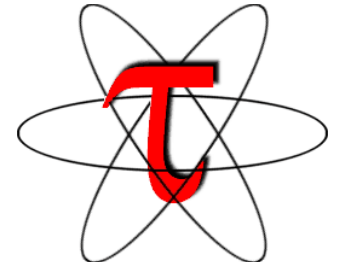
- **Profiling** shows you **how much** (total) time was spent in each routine
- **Tracing** shows you **when** the events take place on a timeline
- Metrics can be time or hardware performance counters (cache misses, instructions)
- TAU can automatically instrument your source code using a package called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

Inclusive vs. Exclusive Measurements

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



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

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

```

xterm
NODE 0;CONTEXT 0;THREAD 0:
-----
%Time      Exclusive    Inclusive    #Call    #Subrs    Inclusive Name
          msec      total msec
-----
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()
-----

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
-----
27.1
1%

```

Interval events show *duration*

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

```

% export TAU_CALLPATH_DEPTH=0
% export TAU_TRACK_HEAP=1

```

Atomic Events and Context Events

```
xterm
```

%Time	Exclusive msec	Inclusive total msec	#Call	#Subrs	Inclusive Name usec/call
100.0	0.253	1.106	1	44	1106701 int main(int, char **) C
93.2	1.031	1.031	1	0	1031311 MPI_Init()
6.0	1	66	40	320	1650 void func(int, int) C
5.7	63	63	40	0	1588 MPI_Barrier()
0.8	9	9	1	0	9119 MPI_Finalize()
0.1	1	1	120	0	10 MPI_Recv()
0.0	0.141	0.141	120	0	1 MPI_Send()
0.0	0.085	0.085	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()

Atomic events

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

NumSamples	MaxValue	MinValue	MeanValue	Std. Dev.	Event Name
40	40	40	40	0	Message size for broadcast
365	5.139E+04	44.39	3.091E+04	1.234E+04	Heap Memory Used (KB) : Entry
40	5.139E+04	3097	3.114E+04	1.227E+04	Heap Memory Used (KB) : Entry : MPI_Barrier()
40	5.139E+04	1.13E+04	3.134E+04	1.187E+04	Heap Memory Used (KB) : Entry : MPI_Bcast()
1	2067	2067	2067	0	Heap Memory Used (KB) : Entry : MPI_Comm_rank()
1	2066	2066	2066	0	Heap Memory Used (KB) : Entry : MPI_Comm_size()
1	5.139E+04	5.139E+04	5.139E+04	0.0006905	Heap Memory Used (KB) : Entry : MPI_Finalize()
1	57.56	57.56	57.56	0	Heap Memory Used (KB) : Entry : MPI_Init()
120	5.139E+04	1.13E+04	3.134E+04	1.187E+04	Heap Memory Used (KB) : Entry : MPI_Recv()
120	5.139E+04	1.129E+04	3.134E+04	1.187E+04	Heap Memory Used (KB) : Entry : MPI_Send()
1	44.39	44.39	44.39	0	Heap Memory Used (KB) : Entry : int main(int, char **) C
40	5.036E+04	2068	3.011E+04	1.227E+04	Heap Memory Used (KB) : Entry : void func(int, int) C

Context events are atomic events with executing context

Controls depth of executing context shown in profiles

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

Context Events with Callpath

```

xterm
NODE 0:CONTEXT 0:THREAD 0:
-----
%Time    Exclusive    Inclusive    #Call    #Subrs    Inclusive Name
         msec      total msec
-----
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
-----
3.7 1%

```

```

% 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

Debugging

Multi-language Application Debugging

```
% export TAU_MAKEFILE=$TAU/Makefile.tau-icpc-papi-mpi-pdt
% export TAU_OPTIONS='-optMemDbg -optVerbose'
% make F90=tau_f90.sh CC=tau_cc.sh CXX=tau_cxx.sh
% qsub -IX -l walltime=00:30:00 -l select=1:mpiprocs=16:ncpus=16 -l
place=scatter:excl -A <your_account> -q standard

% export TAU_MEMDBG_PROTECT_ABOVE=1
% export TAU_MEMDBG_PROTECT_BELOW=1
% export TAU_MEMDBG_PROTECT_FREE=1
% mpirun -np 16 ./matmult
% paraprof
```

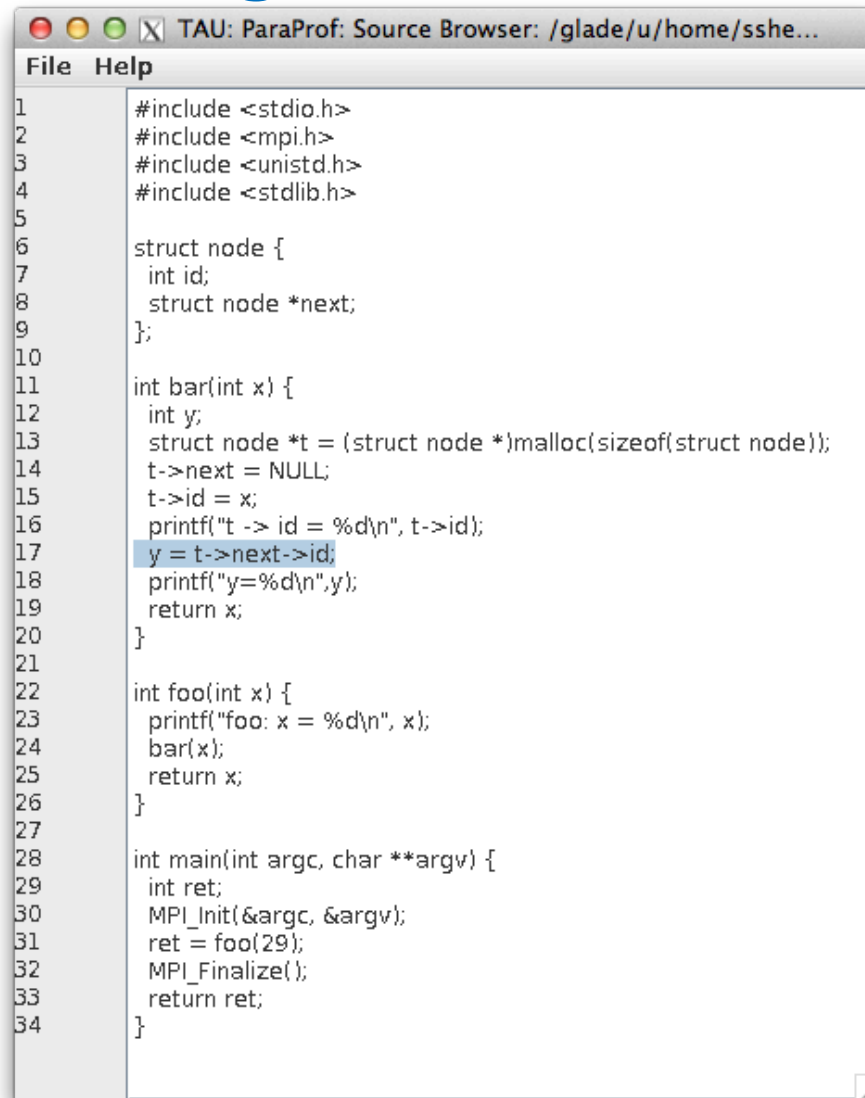
Multi-language Application Debugging

TAU: ParaProf Manager

- Applications
- Standard Applications
 - Default App
 - Default Exp
 - debug.ppk
 - TIME
- Default (jdbc:h2:/User:)
- alcf (jdbc:postgresql://)
- perexplorer_working (
- taudb (jdbc:postgresql

TrialField	Value
Name	debug.ppk
Application ID	0
Experiment ID	0
Trial ID	0
BACKTRACE(1) 1	[bar] [/glade/u/home/sshende/pkgs/workshop/debug/foo.c:17] [/glade/u/home/sshende/pkgs/workshop/debug/foo.tau_rewritten]
BACKTRACE(1) 2	[foo] [/glade/u/home/sshende/pkgs/workshop/debug/foo.c:24] [/glade/u/home/sshende/pkgs/workshop/debug/foo.tau_rewritten]
BACKTRACE(1) 3	[main] [/glade/u/home/sshende/pkgs/workshop/debug/foo.c:31] [/glade/u/home/sshende/pkgs/workshop/debug/foo.tau_rewritten]
BACKTRACE(1) 4	[_libc_start_main] [(unknown):0] [/lib64/libc-2.12.so]
BACKTRACE(1) 5	[_start] [(unknown):0] [/glade/u/home/sshende/pkgs/workshop/debug/foo.tau_rewritten]
CPU Cores	8
CPU MHz	2600.000
CPU Type	Intel(R) Xeon(R) CPU E5-2670 0 @ 2.60GHz
CPU Vendor	GenuineIntel
CWD	/glade/u/home/sshende/pkgs/workshop/debug
Cache Size	20480 KB
Command Line	/glade/u/home/sshende/pkgs/workshop/debug/foo.tau_rewritten
Executable	/glade/u/home/sshende/pkgs/workshop/debug/foo.tau_rewritten
File Type Index	0
File Type Name	ParaProf Packed Profile
Hostname	ys0101
Local Time	2013-03-31T17:13:40-06:00
MPI Processor Name	ys0101
Memory Size	32988340 kB
Node Name	ys0101
OS Machine	x86_64
OS Name	Linux
OS Release	2.6.32-220.13.1.el6.x86_64
OS Version	#1 SMP Thu Mar 29 11:46:40 EDT 2012
SIGNAL	Segmentation fault
Starting Timestamp	1364771620575991
TAU Architecture	x86_64
TAU Config	-iowrapper -mpiinc=/opt/ibmhpc/pe1209/mpich2/intel/include64 -mpilib=/opt/ibmhpc/pe1209/mpich2/intel/lib64 -cc=icc -c+ +...
TAU Makefile	/glade/u/home/sshende/pkgs/tau-2.22.2/x86_64/lib/Makefile.tau-icpc-mpi-pdt
TAU Version	2.22.2

Location of segmentation violation



The screenshot shows a source browser window titled "TAU: ParaProf: Source Browser: /glade/u/home/sshe...". The window contains C code with line numbers 1 through 34. The code defines a linked list structure and functions bar, foo, and main. A segmentation violation is highlighted on line 17: `y = t->next->id;`. The code includes headers for stdio, mpi, unistd, and stdlib. The bar function prints the id of a node and the id of its next node. The foo function prints the value of x and calls bar. The main function initializes MPI, calls foo with 29, and finalizes MPI.

```
1 #include <stdio.h>
2 #include <mpi.h>
3 #include <unistd.h>
4 #include <stdlib.h>
5
6 struct node {
7     int id;
8     struct node *next;
9 };
10
11 int bar(int x) {
12     int y;
13     struct node *t = (struct node *)malloc(sizeof(struct node));
14     t->next = NULL;
15     t->id = x;
16     printf("t -> id = %d\n", t->id);
17     y = t->next->id;
18     printf("y=%d\n",y);
19     return x;
20 }
21
22 int foo(int x) {
23     printf("foo: x = %d\n", x);
24     bar(x);
25     return x;
26 }
27
28 int main(int argc, char **argv) {
29     int ret;
30     MPI_Init(&argc, &argv);
31     ret = foo(29);
32     MPI_Finalize();
33     return ret;
34 }
```


Memory Leak Detection

```
% export TAU_MAKEFILE=$TAU/Makefile.tau-icpc-papi-mpi-pdt
% export PATH=$TAU/../../bin:$PATH
% export TAU_OPTIONS='-optMemDbg -optVerbose'
% make F90=tau_f90.sh CC=tau_cc.sh CXX=tau_cxx.sh
% qsub -IX -l walltime=00:30:00 -l
select=1:mpiprocs=16:ncpus=16 -l place=scatter:excl -A
<your_account> -q standard

% export TAU_TRACK_MEMORY_LEAKS=1
% mpirun -np 16 ./matmult
% paraprof
```

Multi-language Memory Leak Detection

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

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}]						
▼ MPL_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

Support Acknowledgments



UNIVERSITY
OF OREGON

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)

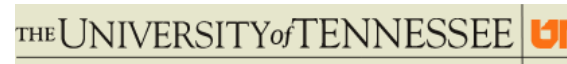
- Glassbox, SI-2



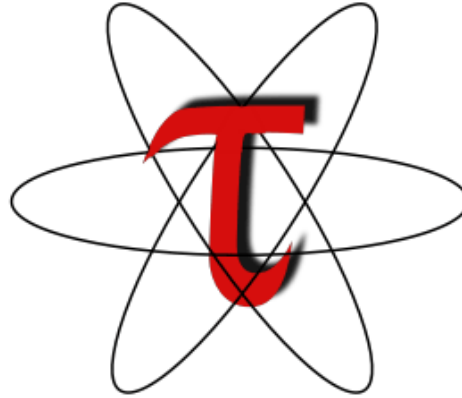
University of Tennessee, Knoxville

T.U. Dresden, GWT

Juelich Supercomputing Center



Download TAU from U. Oregon



<http://tau.uoregon.edu>

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

Free download, open source, BSD license