

Parallel Performance Evaluation Tools: TAU, PAPI, Scalasca and Vampir

Two day tutorial at LLNL
Building 453 R1001 (Armadillo), May 26-27, Livermore, CA
Sameer Shende
sameer@paratools.com
<http://www.paratools.com/lnl09>

ParaTools

Outline

	Slide #
• Outline and workshop goals	2
• Part I: TAU: A quick reference	8
• Part II: Introduction to performance engineering	77
• Part III: PAPI	105
• Part IV: TAU	119
• Part V: Vampir/VNG	266
• Part VI: Scalasca/KOJAK	307
• Lab Session: PAPI, TAU, Vampir and Scalasca examples	385

ParaTools

Outline

- Day 1:
 - Introduction to performance evaluation tools: TAU, PAPI, Scalasca, and Vampir
 - Hands-on:
 - TAU instrumentation at routine, loop level, PAPI hardware performance counter data collection, derived metrics, analyzing performance using TAU's paraprof profile browser, using Performance database (PerfDMF), memory evaluation, leak detection
- Day 2:
 - Scalasca, TAU PerfExplorer, VampirServer
 - Hands-on:
 - Scalasca bottleneck detection tools, PerfExplorer, trace visualization, workshop examples including the NAS Parallel Benchmarks 3.1

ParaTools

3

Workshop Goals

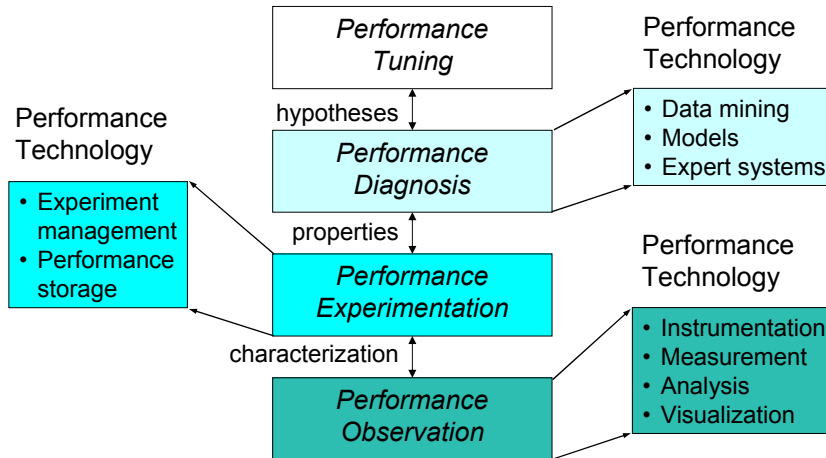
- This tutorial is an introduction to portable performance evaluation tools.
- You should leave here with a better understanding of...
 - Concepts and steps involved in performance evaluation
 - Understanding key concepts in improving and understanding code performance
 - How to collect and analyze data from hardware performance counters using PAPI
 - How to instrument your programs with TAU
 - Automatic instrumentation at the routine level and outer loop level
 - Manual instrumentation at the loop/statement level
 - Measurement options provided by TAU
 - Environment variables used for choosing metrics, generating performance data
 - How to use the TAU's profile browser, ParaProf
 - How to use TAU's database for storing and retrieving performance data
 - General familiarity with TAU's use for Fortran, Python, C++, C, MPI for mixed language programming
 - How to generate trace data in different formats
 - How to use Scalasca for detecting performance bottlenecks
 - How to analyze trace data using Vampir, and Jumpshot
 - Facilities provided by the Eclipse PTP integrated development environment for parallel programs

ParaTools

4

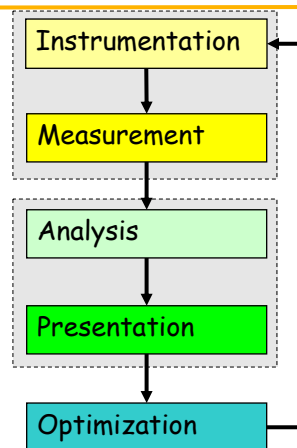
Performance Engineering

- Optimization process
- Effective use of performance technology



Performance Optimization Cycle

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

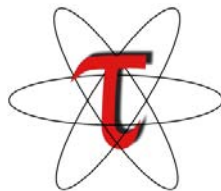


More Information

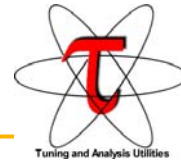
- PAPI References:
 - PAPI documentation page available from the PAPI website:
<http://icl.cs.utk.edu/papi/>
- TAU References:
 - TAU Users Guide and papers available from the TAU website:
<http://tau.uoregon.edu/>
- VAMPIR References
 - VAMPIR-NG website
<http://www.vampir.eu/>
- Scalasca/KOJAK References
 - Scalasca documentation page
<http://www.scalasca.org/>

TAU: A Quick Reference

Part I: TAU: A Tutorial



TAU Performance System



- <http://tau.uoregon.edu/>
- Multi-level performance instrumentation
 - Multi-language automatic source instrumentation
- Flexible and configurable performance measurement
- Widely-ported parallel performance profiling system
 - Computer system architectures and operating systems
 - Different programming languages and compilers
- Support for multiple parallel programming paradigms
 - Multi-threading, message passing, mixed-mode, hybrid
- Integration in complex software, systems, applications

ParaTools

9

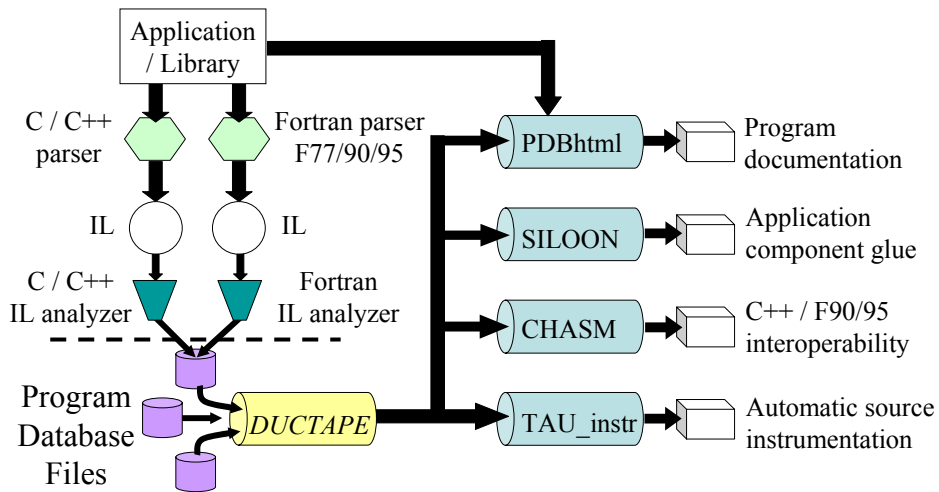
What is TAU?

- TAU is a performance evaluation tool
- It supports both parallel profiling and tracing
- Profiling shows you how much (total) time was spent in each routine (event)
- Tracing shows you *when* events take place in each process along a timeline
- Profiling and tracing can measure time as well as hardware performance counters
- TAU uses a package called *PDT* for automatic instrumentation of the source code
- With PDT, TAU can instrument routine, loop, phase, I/O, and memory
- TAU can also use your compiler to insert the instrumentation at routine boundaries
- TAU can *throttle* the insignificant lightweight routines at runtime to reduce perturbation. It can also *subtract* the timer overhead at runtime to compensate.
- TAU runs on all HPC platforms and it is free (BSD style license)
- TAU has instrumentation, measurement and analysis tools
 - ParaProf is TAU's 3D profile browser, PerfDMF is the TAU database tool
- **To use TAU, all you need to do is set a couple of environment variables and substitute the name of your compiler with a TAU shell script (e.g., tau_f90.sh)**

ParaTools

10

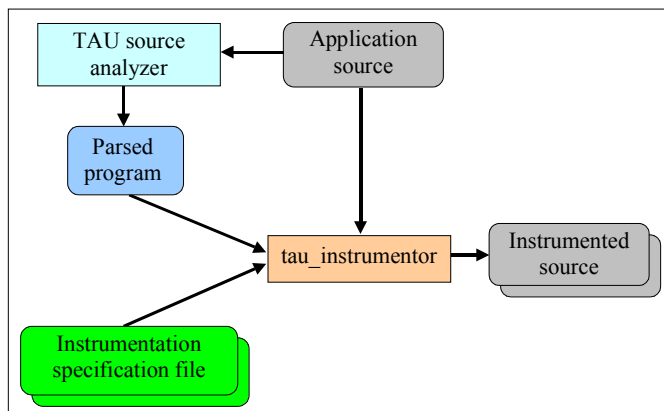
Program Database Toolkit (PDT)



ParaTools

11

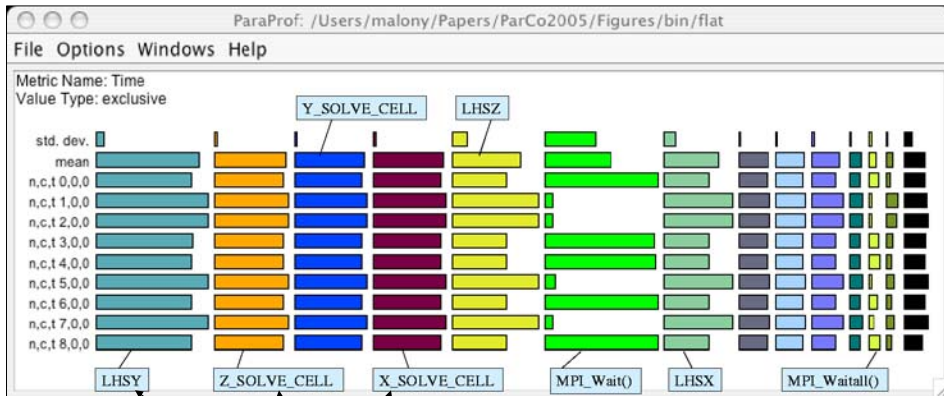
Automatic Source-Level Instrumentation in TAU using Program Database Toolkit (PDT)



ParaTools

12

ParaProf – Flat Profile (NAS BT)



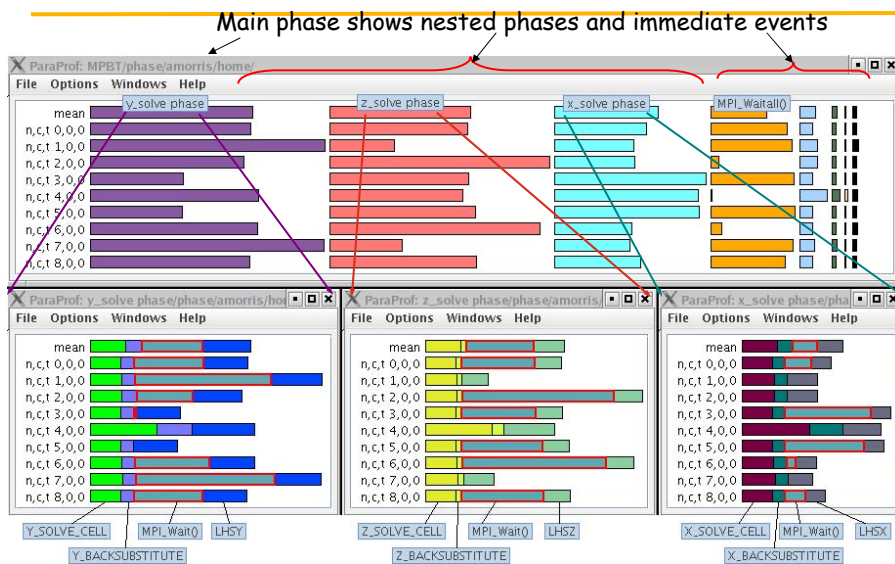
Application routine names reflect phase semantics

How is MPI_Wait() distributed relative to solver direction?

ParaTools

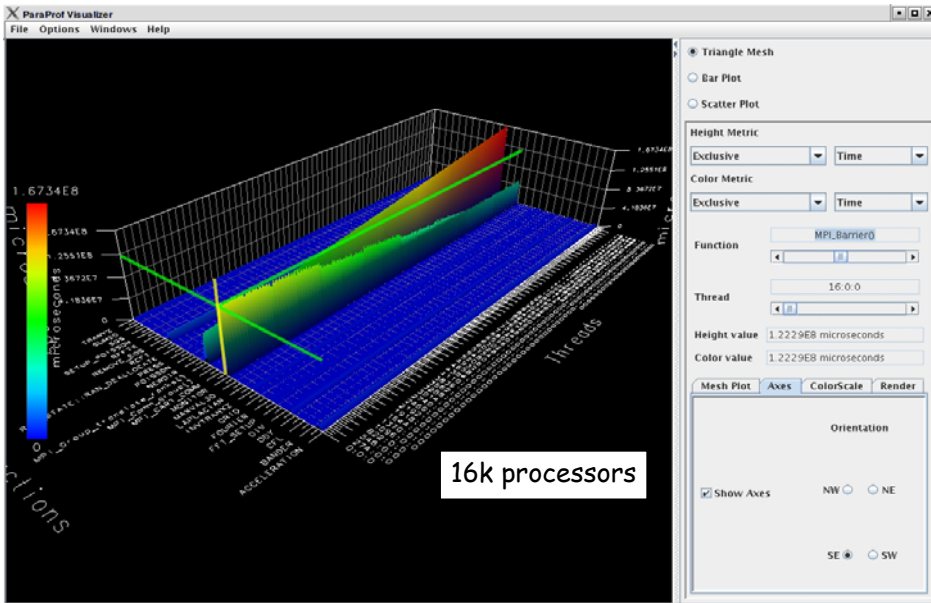
13

ParaProf – Phase Profile (NAS BT)

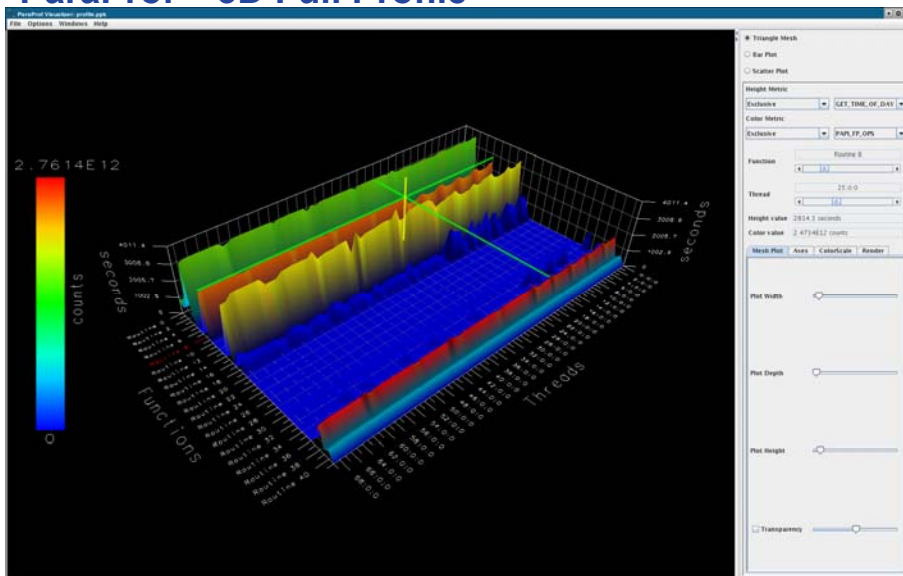


14

ParaProf – 3D Full Profile (Miranda, LLNL)



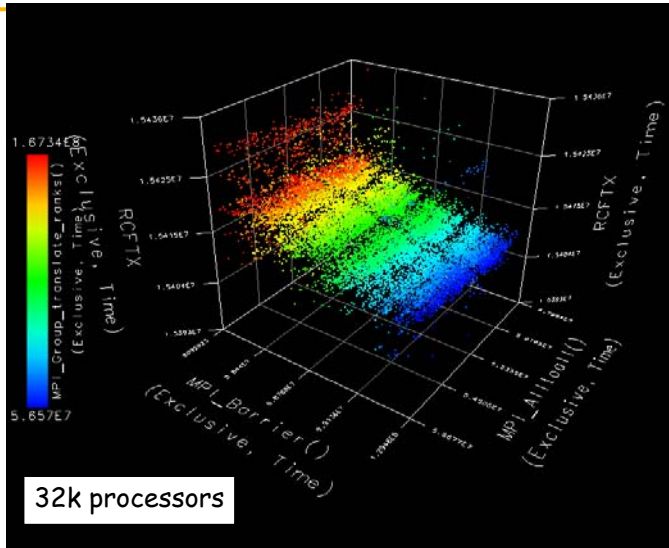
ParaProf – 3D Full Profile



ParaTools

ParaProf – 3D Scatterplot

- Each point is a “thread” of execution
- A total of four metrics shown in relation
- 3D profile visualization library – JOGL



ParaTools

ParaProf’s Source Browser: Loop Level Instrumentation

TAU: ParaProf: Function Data Window: s3d_callpath_papi.ppk

Name: Loop: TRANSPORT_MCOMPUTESPECIESDIFFFLUX [micavg_transport_m_ppf90] (630,5)-(656,19)

Metric Name: PAPI_PP_INS / GET_TIME_OF_DAY

Value: Exclusive

	1.088	std. dev.
114.979		mean
117.62		RCI 0.0,0
115.134		RCI 1.0,0
114.709		RCI 2.0,0
114.615		RCI 3.0,0
113.547		RCI 4.0,0
114.581		RCI 5.0,0
114.837		RCI 6.0,0
114.789		RCI 7.0,0

TAU: ParaProf: Function Data Window: s3d_callpath_papi.ppk

Name: Loop: TRANSPORT_MCOMPUTESPECIESDIFFFLUX [micavg_transport_m_ppf90] (630,5)-(656,19)

Metric Name: GET_TIME_OF_DAY

Value: Exclusive percent

	0.91%	std. dev.
12.786%		mean
11.931%		RCI 0.0,0
12.19%		RCI 1.0,0
12.248%		RCI 2.0,0
12.798%		RCI 3.0,0
12.33%		RCI 4.0,0
12.241%		RCI 5.0,0
12.223%		RCI 6.0,0
12.226%		RCI 7.0,0

TAU: ParaProf: Function Data Window: s3d_callpath_papi.ppk

Name: Loop: TRANSPORT_MCOMPUTESPECIESDIFFFLUX [micavg_transport_m_ppf90] (630,5)-(656,19)

Metric Name: PAPI_LL_DCM

Value: Exclusive

Units: counts

	816336.1	std. dev.
5.07019		mean
5.06929		RCI 0.0,0
5.0719		RCI 1.0,0
5.06919		RCI 2.0,0
5.07019		RCI 3.0,0
5.07089		RCI 4.0,0
5.07119		RCI 5.0,0
5.07129		RCI 6.0,0
5.06929		RCI 7.0,0

TAU: ParaProf: Source Browser: /mnt/rapion/Users/raameer/rts/taodata/s3d/harness/flat/papi8

```

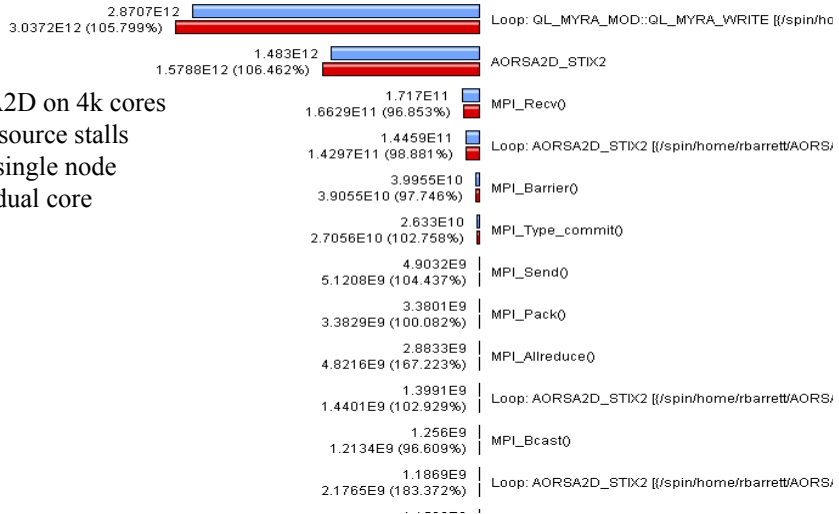
606  grad_piv%(...,i) = grad_piv%(...,i)*amole%(...)
607  end do
608
609  ! compute grad_P
610  if (daro_switch) then
611    allocate(grad_P(m,n,2))
612    grad_P = 0.0
613    if (vary_in_1 == 1) then
614      call derivative_1(m,n,nz, Press, grad_P(...,1), scale_in_1)
615    end if
616    if (vary_in_2 == 1) then
617      call derivative_2(m,n,nz, Press, grad_P(...,2), scale_in_2)
618    end if
619    if (vary_in_3 == 1) then
620      call derivative_3(m,n,nz, Press, grad_P(...,3), scale_in_3)
621    end if
622  end if
623
624  ! Changed by Rasmus - 01/26/95
625  ! Do scaling to new units
626
627  ! grad_P/press and amole*grad_Temp can be optimized by division before the loop.
628  ! compute diffusive flux for species n in direction x.
629  diffFlux(...,i,spec) = 0.0
630  DIRECTION = dir%1
631  SPECIES = dir%2,spec
632
633  IF (daro_switch) then
634    ! driving force includes gradient to mole fraction and bars-diffusion
635    diffFlux(...,i,spec) = - Ds_piv%(...,i) * ( grad_P(...,DIRECTION)
636      + %val(...,i) * grad_piv%(...,DIRECTION) )
637      + (1 - mole%(...,i)) * grad_P(...,DIRECTION)/P0000
638  else
639    ! driving force is just the gradient to mole fraction
640    diffFlux(...,i,spec) = - Ds_piv%(...,i) * ( grad_P(...,DIRECTION)
641      + %val(...,i) * grad_piv%(...,DIRECTION) )
642  end if
643
644  ! Add thermal diffusion
645  if (thermal_diffusion) then
646    diffFlux(...,i,spec) = diffFlux(...,i,spec) +
647      - Ds_piv%(...,i) * %val(...,i) * temp
648    amole% = grad_P(...,DIRECTION) / temp
649  end if
650
651  ! compute contribution to all species diffusive flux
652  ! this will ensure that the sum of the diffusive fluxes is zero
653  diffFlux(...,i,spec) = diffFlux(...,i,spec) - diffFlux(...,i,0)
654
655  enddo SPECIES
656  enddo DIRECTION
657
658  if (daro_switch) then
659    deallocate(grad_P)
660  end if
661
662  return
663 end subroutine computeSpeciesDiffFlux
664
665
666
667
668
669
670
671
672
673
674
        
```

Comparing Effects of MultiCore Processors

Metric: PAPI_RES_STL
Value: Exclusive
Units: counts

■ C:\iter.350x350.4096pes.sn.loops.BARRIER.ppk - Mean
■ C:\iter.350x350.2048pes.dc.loops.BARRIER.ppk - Mean

- AORSA2D on 4k cores
- PAPI resource stalls
- Blue is single node
- Red is dual core



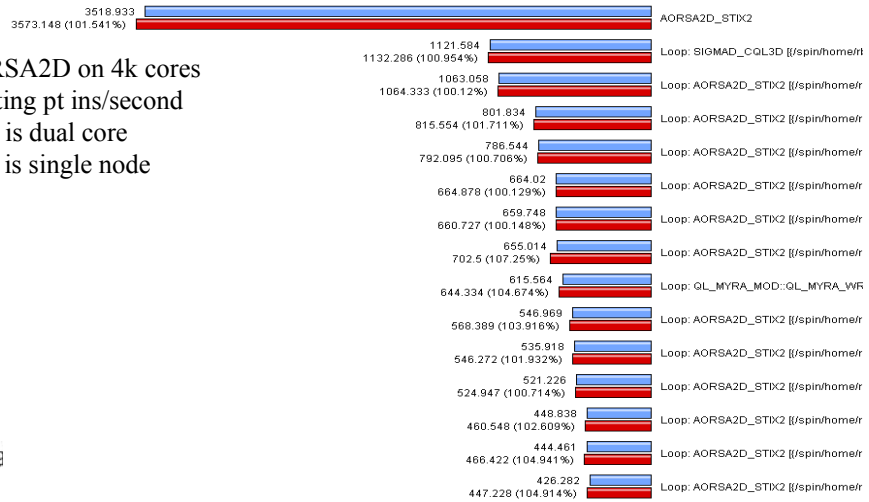
Para

Comparing FLOPS: MultiCore Processors

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

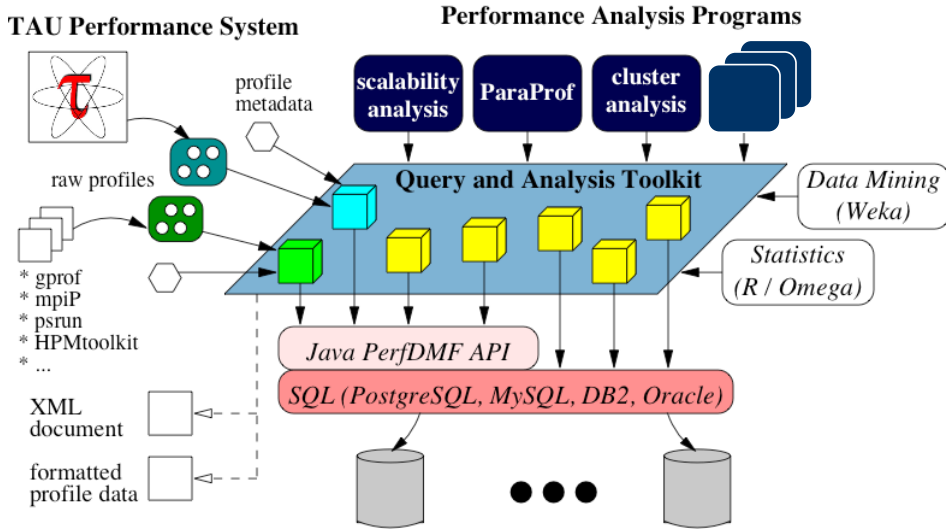
■ C:\iter.350x350.2048pes.dc.loops.BARRIER.ppk - Mean
■ C:\iter.350x350.4096pes.sn.loops.BARRIER.ppk - Mean

- AORSA2D on 4k cores
- Floating pt ins/second
- Blue is dual core
- Red is single node



Para

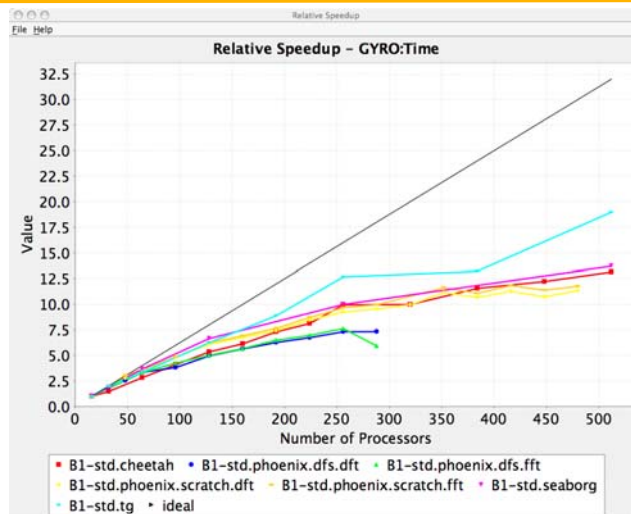
PerfDMF: Performance Data Mgmt. Framework



ParaTools

21

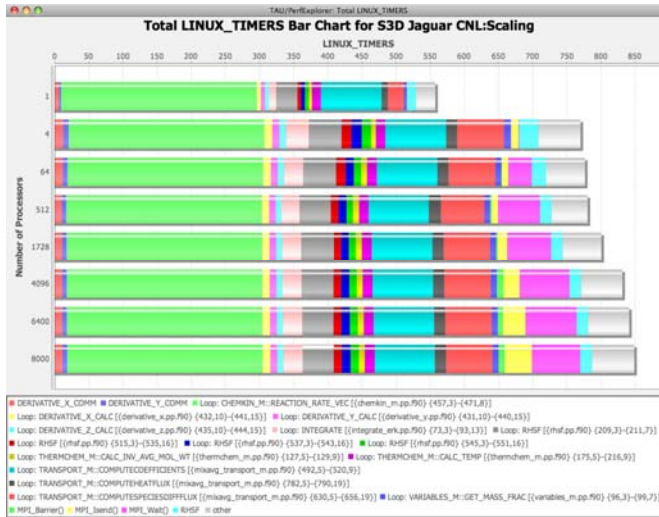
PerfExplorer: Comparing Relative Speedup on Different Architectures



ParaTools

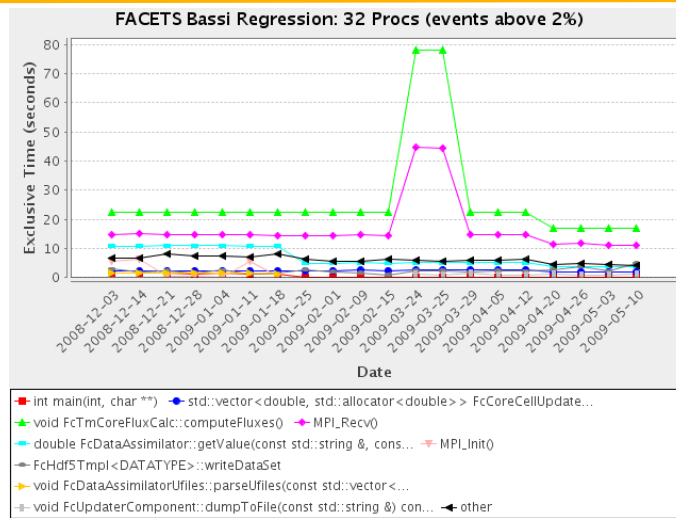
22

Usage Scenarios: Evaluate Scalability



ParaTools

Performance Regression Testing



ParaTools

Profiling

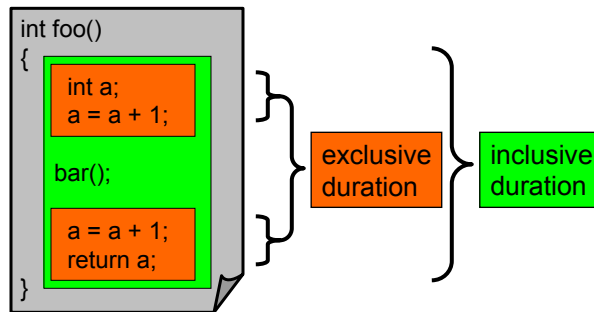
- Recording of aggregated information
 - Counts, time, ...
- ... about program and system entities
 - Functions, loops, basic blocks, ...
 - Processes, threads
- Methods
 - Event-based sampling (indirect, statistical)
 - Open|SpeedShop, PerfSuite, HPCToolkit, gprof,...
 - Direct measurement (deterministic)
 - TAU, VampirTrace, Scalasca,...

Direct Observation: Events

- Event types
 - Interval events (begin/end events)
 - measures performance between begin and end
 - metrics monotonically increase
 - Atomic events
 - used to capture performance data state
- Code events
 - Routines, classes, templates
 - Statement-level blocks, loops
- User-defined events
 - Specified by the user
- Abstract mapping events

Inclusive and Exclusive Profiles

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



ParaTools

27

Terminology – Example

- For routine "int main()":
- Exclusive time
 - 100-20-50-20=10 secs
- Inclusive time
 - 100 secs
- Calls
 - 1 call
- Subrs (no. of child routines called)
 - 3
- Inclusive time/call
 - 100secs

```
int main( )
{ /* takes 100 secs */

    f1(); /* takes 20 secs */
    f2(); /* takes 50 secs */
    f1(); /* takes 20 secs */

    /* other work */
}

/*
Time can be replaced by counts
from PAPI e.g., PAPI_FP_INS. */
```

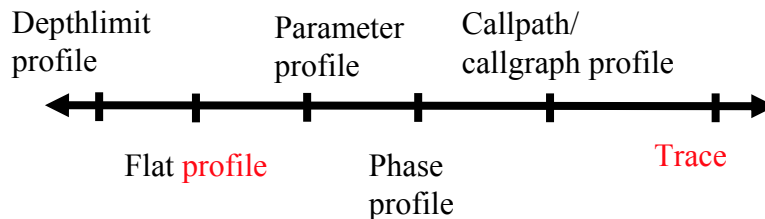
ParaTools

28

Flat and Callpath Profiles

- Static call graph
 - Shows all parent-child calling relationships in a program
- Dynamic call graph
 - Reflects actual execution time calling relationships
- Flat profile
 - Performance metrics for when event is active
 - Exclusive and inclusive
- Callpath profile
 - Performance metrics for calling path (event chain)
 - Differentiate performance with respect to program execution state
 - Exclusive and inclusive

Performance Evaluation Alternatives



Each alternative has:
- one metric/counter
- multiple counters

Volume of performance data

Tracing Measurement

Process A:

```
void master {
  trace(ENTER, 1);
  ...
  trace(SEND, B);
  send(B, tag, buf);
  ...
  trace(EXIT, 1);
}
```



1	master
2	worker
3	...

MONITOR

...				
58	A	ENTER	1	
60	B	ENTER	2	
62	A	SEND	B	
64	A	EXIT	1	
68	B	RECV	A	
69	B	EXIT	2	
...				

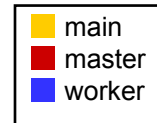
Process B:

```
void worker {
  trace(ENTER, 2);
  ...
  recv(A, tag, buf);
  trace(RECV, A);
  ...
  trace(EXIT, 2);
}
```

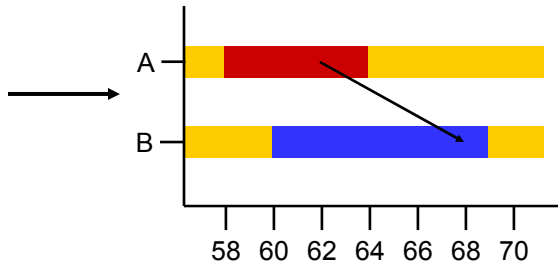
ParaTools

Tracing Analysis and Visualization

1	master
2	worker
3	...



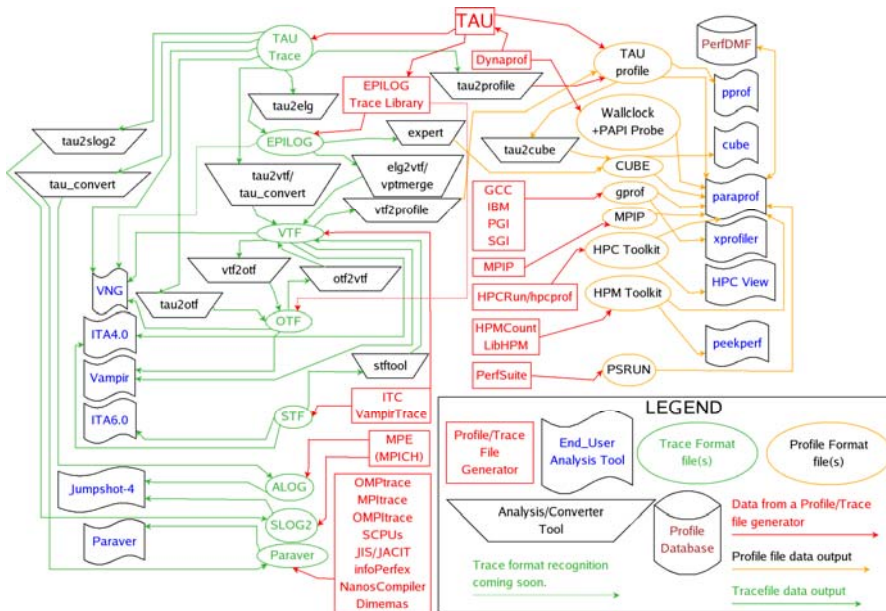
...				
58	A	ENTER	1	
60	B	ENTER	2	
62	A	SEND	B	
64	A	EXIT	1	
68	B	RECV	A	
69	B	EXIT	2	
...				



ParaTools



Building Bridges to Other Tools



Trace Formats

- Different tools produce different formats
 - Differ by event types supported
 - Differ by ASCII and binary representations
 - Vampir Trace Format (VTF)
 - KOJAK (EPILOG)
 - Jumpshot (SLOG-2)
 - Paraver
- Open Trace Format (OTF)
 - Supports interoperability between tracing tools

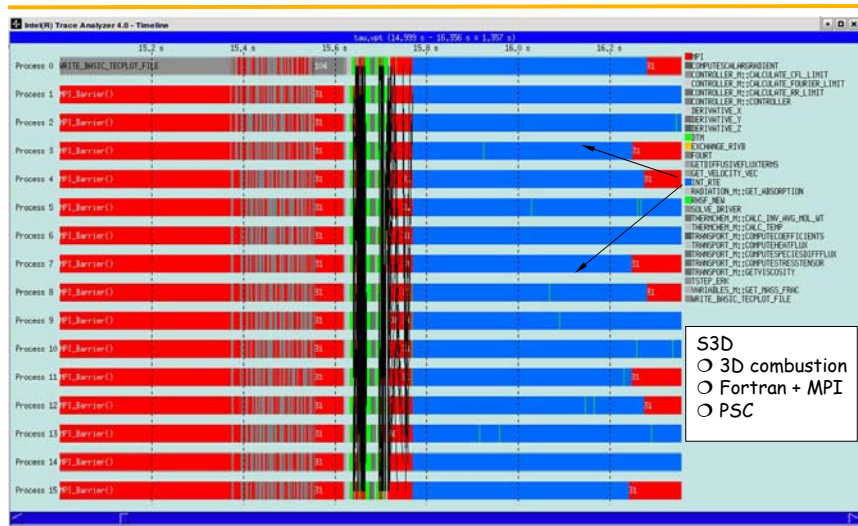
Profiling / Tracing Comparison

- Profiling
 - ☺ Finite, bounded performance data size
 - ☺ Applicable to both direct and indirect methods
 - ☹ Loses time dimension (not entirely)
 - ☹ Lacks ability to fully describe process interaction
- Tracing
 - ☺ Temporal and spatial dimension to performance data
 - ☺ Capture parallel dynamics and process interaction
 - ☹ Some inconsistencies with indirect methods
 - ☹ Unbounded performance data size (large)
 - ☹ Complex event buffering and clock synchronization

ParaTools

35

Vampir – Trace Analysis (TAU-to-VTF3) (S3D)



ParaTools

36

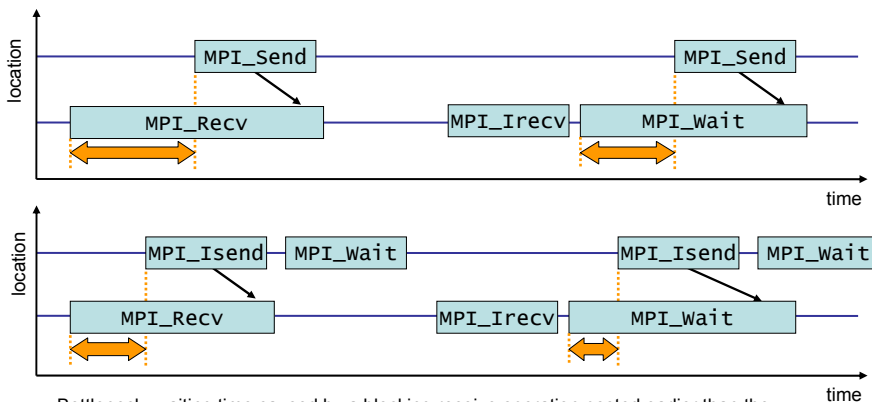
Vampir – Trace Zoomed (S3D)



ParaTools

37

Scalasca Bottleneck Detection Tool: Late Sender = Early Receive Bottleneck

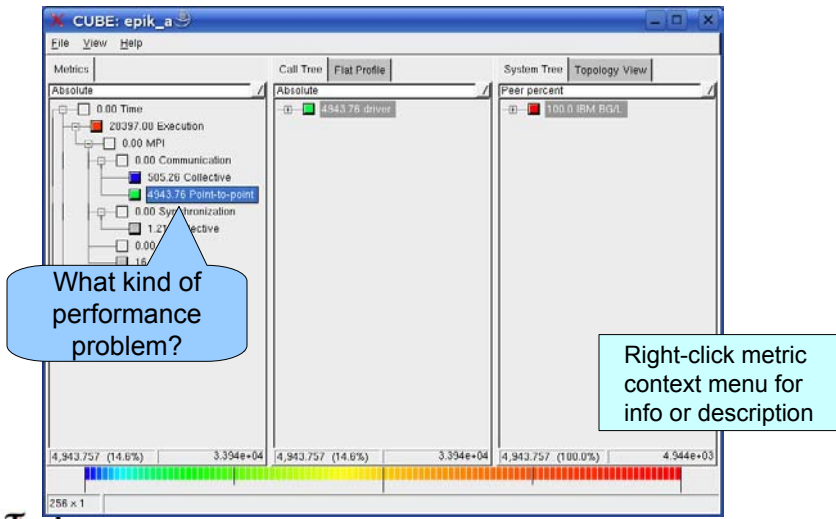


- Bottleneck: waiting time caused by a blocking receive operation posted earlier than the corresponding send operation. Scalasca has a library of different bottlenecks.
- Scalasca can measure the yellow arrows in the trace. It analyzes and generates a profile containing the bottlenecks as metrics (shown in yellow). The profiles may be viewed in TAU's paraprof browser or Scalasca's CUBE browser.

ParaTools

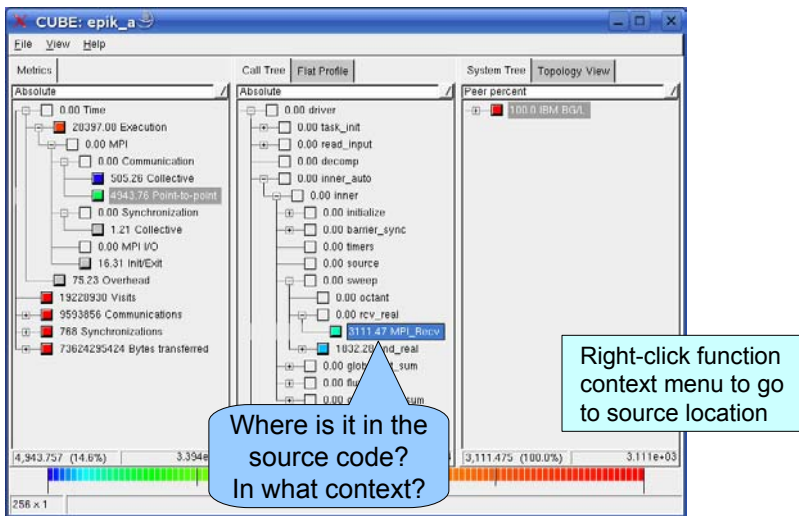
38

Scalasca: CUBE Profile Browser



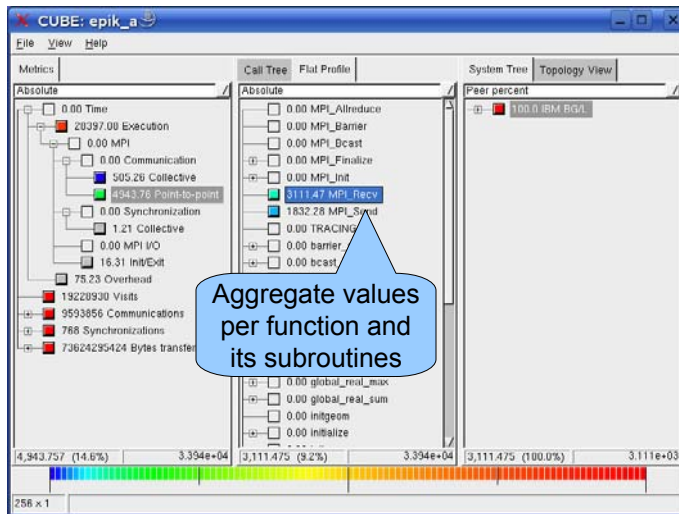
ParaTools

CUBE call tree dimension



ParaTools

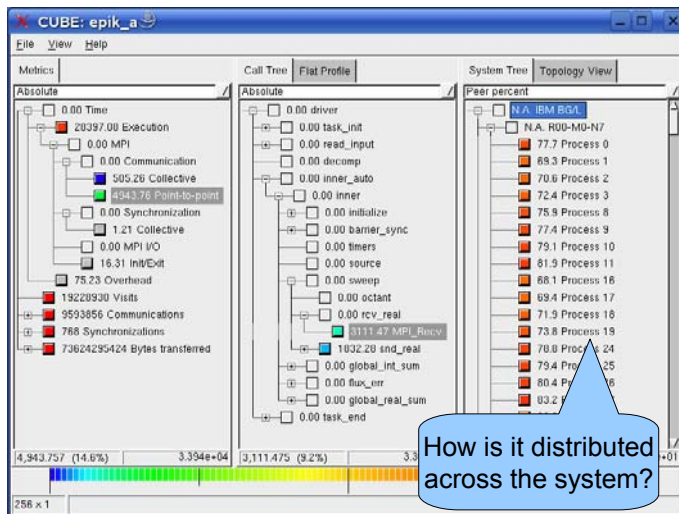
Alternative: Flat profile



ParaTools

41

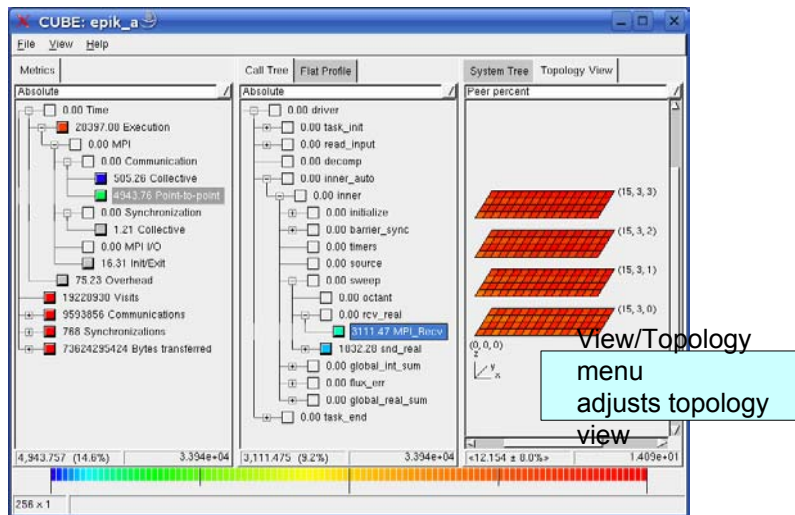
System tree dimension



ParaTools

42

Alternative: Topology display



ParaTools

43

Steps of Performance Evaluation

- Collect basic routine-level timing profile to determine where most time is being spent
- Collect routine-level hardware counter data to determine types of performance problems
- Collect callpath profiles to determine sequence of events causing performance problems
- Conduct finer-grained profiling and/or tracing to pinpoint performance bottlenecks
 - Loop-level profiling with hardware counters
 - Tracing of communication operations

ParaTools

44

Using TAU: A brief Introduction

- TAU supports several measurement options (profiling, tracing, profiling with hardware counters, etc.)
- Each measurement configuration of TAU corresponds to a unique stub makefile that is generated when you configure it
- To instrument source code using PDT
 - Choose an appropriate TAU stub makefile in <arch>/lib:
% setenv TAU_MAKEFILE /usr/global/tools/tau/training/tau-2.18.2/bgp/lib/Makefile.tau-mpi-pdt
% setenv TAU_OPTIONS '-optVerbose ...' (see tau_compiler.sh -help)
And use tau_f90.sh, tau_cxx.sh or tau_cc.sh as Fortran, C++ or C compilers:
% mpif90 foo.f90
changes to
% tau_f90.sh foo.f90
- Execute application and analyze performance data:
 - % pprof (for text based profile display)
 - % paraprof (for GUI)

ParaTools

45

TAU Measurement Configuration

```
% cd /usr/global/tools/tau/training/tau-2.18.2/bgp/lib; ls Makefile.*
Makefile.tau-pdt
Makefile.tau-mpi-pdt
Makefile.tau-opari-openmp-mpi-pdt
Makefile.tau-mpi-scalasca-epilog-pdt
Makefile.tau-mpi-vampirtrace-pdt
Makefile.tau-mpi-papi-pdt
Makefile.tau-papi-mpi-openmp-opari-pdt
Makefile.tau-pthread-pdt...
```

- For an MPI+F90 application, you may want to start with:

```
Makefile.tau-mpi-pdt
  - Supports MPI instrumentation & PDT for automatic source instrumentation
  - % setenv TAU_MAKEFILE
    /usr/global/tools/tau/training/tau-2.18.2/bgp/lib/Makefile.tau-
    mpi-pdt
  - % tau_f90.sh matrix.f90 -o matrix
```

ParaTools

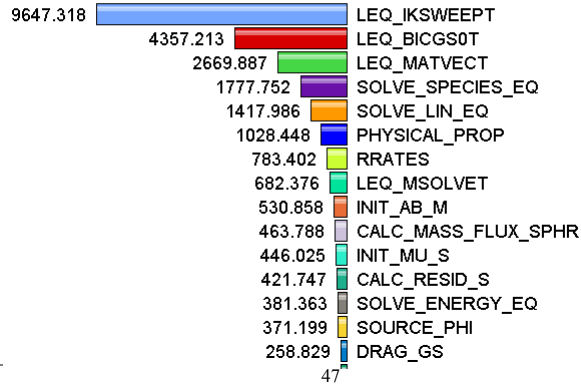
46

Usage Scenarios: Routine Level Profile

- Goal: What routines account for the most time? How much?

- Flat profile with wallclock time:

Metric: P_VIRTUAL_TIME
Value: Exclusive
Units: seconds



ParaTools_

47

Solution: Generating a flat profile with MPI

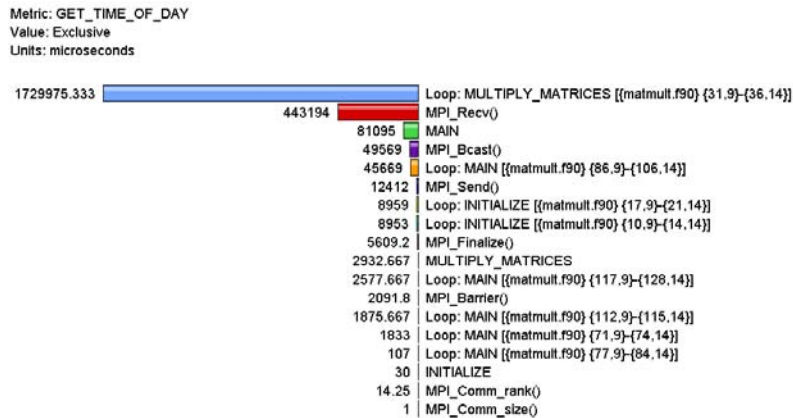
```
% setenv TAU_MAKEFILE /usr/global/tools/tau/training/tau-2.18.2/bgp  
/lib/Makefile.tau-mpi-pdt  
% set path=(/usr/global/tools/tau/training/tau-2.18.2/bgp/bin $path)  
OR  
% source /usr/global/tools/tau/training/src/tau.cshrc [ or  
tau.bashrc]  
% make F90=tau_f90.sh  
(Or edit Makefile and change F90=tau_f90.sh)  
  
% qsub run.job  
% paraprof --pack app.ppk  
Move the app.ppk file to your desktop.  
  
% paraprof app.ppk
```

ParaTools_

48

Usage Scenarios: Loop Level Instrumentation

- Goal: What loops account for the most time? How much?
- Flat profile with wallclock time with loop instrumentation:



Par

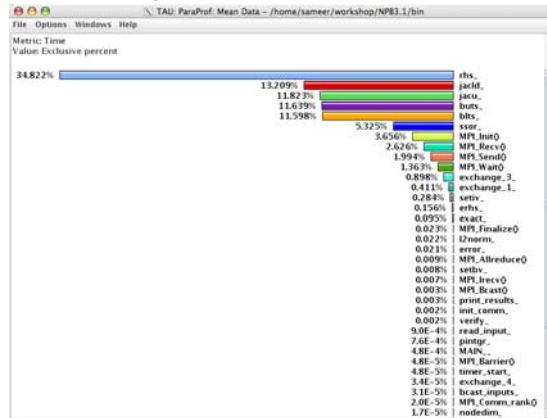
Solution: Generating a loop level profile

```
% setenv TAU_MAKEFILE /usr/global/tools/tau/training/tau-2.18.2/bgp  
/lib/Makefile.tau-mpi-pdt  
% setenv TAU_OPTIONS '-optTauSelectFile=select.tau -optVerbose'  
% cat select.tau  
BEGIN_INSTRUMENT_SECTION  
loops routine="#"  
END_INSTRUMENT_SECTION  
  
% set path=(/usr/global/tools/tau/training/tau-2.18.2/bgp/bin $path)  
% make F90=tau_f90.sh  
(Or edit Makefile and change F90=tau_f90.sh)  
% qsub run.job  
% paraprof --pack app.ppk  
Move the app.ppk file to your desktop.  
  
% paraprof app.ppk
```

ParaTools

Usage Scenarios: Compiler-based Instrumentation

- Goal: Easily generate routine level performance data using the compiler instead of PDT for parsing the source code



ParaTools

51

Use Compiler-Based Instrumentation

```
% setenv TAU_MAKEFILE /usr/global/tools/tau/training/tau-2.18.2/bgp
/lib/Makefile.tau-mpi
% setenv TAU_OPTIONS '-optCompInst -optVerbose'
% % set path=(/usr/global/tools/tau/training/tau-2.18.2/bgp/bin
$path)
% make F90=tau_f90.sh
(Or edit Makefile and change F90=tau_f90.sh)

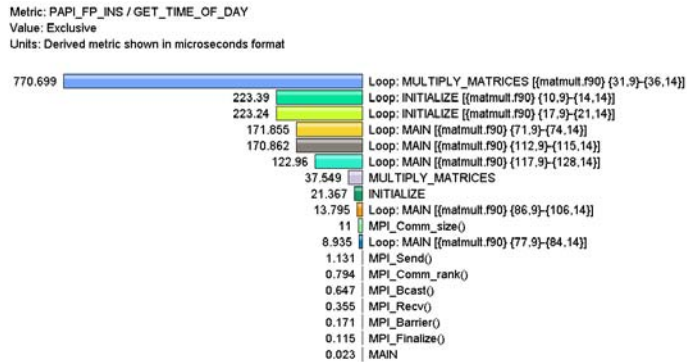
% qsub run.job
% paraprof --pack app.ppk
Move the app.ppk file to your desktop.
% paraprof app.ppk
```

ParaTools

52

Usage Scenarios: Calculate mflops in Loops

- Goal: What MFlops am I getting in all loops?
- Flat profile with PAPI_FP_INS/OPS and time (-multiplicounters) with loop instrumentation:



Par

53

Generate a PAPI profile with 2 or more counters

```
% setenv TAU_MAKEFILE /usr/global/tools/tau/training/tau-2.18.2/bgp
/lib/Makefile.tau-papi-mpi-pdt

% setenv TAU_OPTIONS '-optTauSelectFile=select.tau -optVerbose'

% cat select.tau
BEGIN_INSTRUMENT_SECTION
loops routine="#"
END_INSTRUMENT_SECTION

% set path=(usr/global/tools/tau/training/tau-2.18.2/bgp/bin $path)
% make F90=tau_f90.sh
(Or edit Makefile and change F90=tau_f90.sh)
% setenv COUNTER1 GET_TIME_OF_DAY
% setenv COUNTER2 PAPI_FP_INS
OR
% setenv TAU_METRICS TIME:PAPI_FP_INS
% qsub run.job
% paraprof --pack app.ppk
Move the app.ppk file to your desktop.
% paraprof app.ppk
Choose Options -> Show Derived Panel -> Arg 1 = PAPI_FP_INS,
Arg 2 = GET_TIME_OF_DAY, Operation = Divide -> Apply, choose.
```

ParaTools

54

Derived Metrics in ParaProf

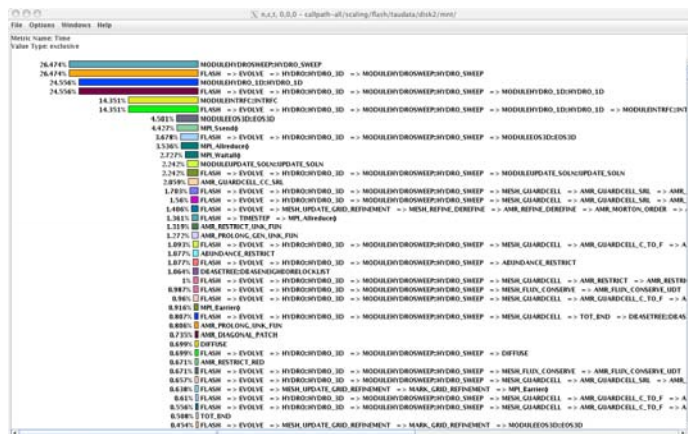
The image shows three screenshots of the ParaProf application. The top-left screenshot shows a list of applications with 'GET_TIME_OF_DAY' selected. The top-right screenshot shows the 'Derived Metrics' dialog box with 'PAR_PROF_GET_TIME_OF_DAY' selected. The bottom screenshot shows the 'Derived Metrics' dialog box with 'PAR_PROF_GET_TIME_OF_DAY' selected and the 'Apply operation' button highlighted.

ParaTools

55

Usage Scenarios: Generating Callpath Profile

- Goal: Who calls my MPI_Barrier()? Where?
- Callpath profile for a given callpath depth:

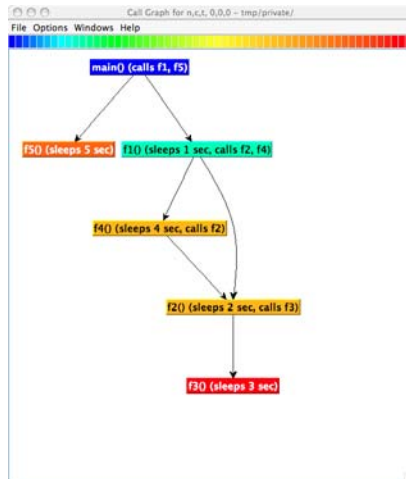


ParaTools

56

Callpath Profile

- Generates program callgraph



ParaTools

57

Generate a Callpath Profile

```
% setenv TAU_MAKEFILE /usr/global/tools/tau/training/tau-2.18.2/bgp
/lib/Makefile.tau-callpath-mpi-pdt
% set path=(/usr/global/tools/tau/training/tau-2.18.2/bgp/bin $path)
% make F90=tau_f90.sh
(Or edit Makefile and change F90=tau_f90.sh)
% setenv TAU_CALLPATH_DEPTH 100

% qsub run.job
% paraprof --pack app.ppk
  Move the app.ppk file to your desktop.
% paraprof app.ppk
(Windows -> Thread -> Call Graph)

NOTE: In TAU v2.18.1+, you may choose to just set:
% setenv TAU_CALLPATH 1
instead of recompiling your code with the above stub makefile.
Any TAU instrumented executable can generate callpath profiles.
```

ParaTools

58

Usage Scenario: Detect Memory Leaks

The screenshot shows two windows from the ParaTools interface. The top window, titled 'TAU: ParaProf: Mean Context Events - mem.ppk', displays a table of memory events. The bottom window, titled 'User Event Window: mem.ppk', shows a detailed view of a specific memory leak event with a bar chart.

Name	NumSamples	MaxValue	MinValue	MeanValue	Std. Dev.
MAIN [(matrix.f90) {141,7}-(146,22)]					
MATRICES:ALLOCATE_MATRICES [(matrix.f90) {10,7}-(13,38)]					
MEMORY LEAK! malloc size <file=matrix.f90, variable=C, line=11>	1	8,000,000	8,000,000	8,000,000	0
malloc size <file=matrix.f90, variable=A, line=11>	1	8,000,000	8,000,000	8,000,000	0
malloc size <file=matrix.f90, variable=B, line=11>	1	8,000,000	8,000,000	8,000,000	0
malloc size <file=matrix.f90, variable=C, line=11>	1	8,000,000	8,000,000	8,000,000	0
MATRICES:DEALLOCATE_MATRICES [(matrix.f90) {14,7}-(17,40)]					
free size <file=matrix.f90, variable=A, line=15>	1	8,000,000	8,000,000	8,000,000	0
free size <file=matrix.f90, variable=B, line=15>	1	8,000,000	8,000,000	8,000,000	0

The bottom window shows the details for the 'MEMORY LEAK! malloc size <file=matrix.f90, variable=C, line=11>' event. It includes a bar chart with the following data:

Value	Mean	Std. Dev.
8000000	n.c.t	0.0
8000000	n.c.t	1.0
8000000	n.c.t	2.0
8000000	n.c.t	3.0
0		

ParaTools

59

Detect Memory Leaks

```
% setenv TAU_MAKEFILE /usr/global/tools/tau/training/tau-2.18.2/bgpp
/lib/Makefile.tau-mpi-pdt
% setenv TAU_OPTIONS '-optDetectMemoryLeaks -optVerbose'
% set path=(/usr/global/tools/tau/training/tau-2.18.2/bgpp/bin $path)
% make F90=tau_f90.sh
(Or edit Makefile and change F90=tau_f90.sh)
% setenv TAU_CALLPATH_DEPTH 100

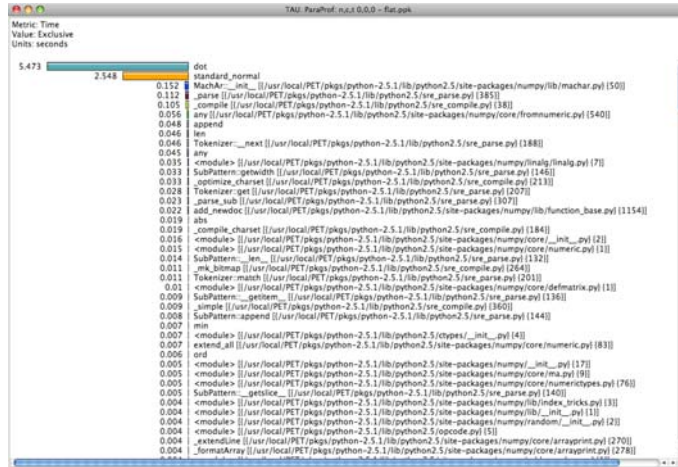
% qsub run.job
% paraprof --pack app.ppk
  Move the app.ppk file to your desktop.
% paraprof app.ppk
(Windows -> Thread -> Context Event Window -> Select thread -> select...
  expand tree)
(Windows -> Thread -> User Event Bar Chart -> right click LEAK
-> Show User Event Bar Chart)
```

ParaTools

60

Usage Scenarios: Instrument a Python program

- Goal: Generate a flat profile for a Python program



ParaTools

61

Usage Scenarios: Instrument a Python program

Original code:

```
% cat foo.py
#!/usr/bin/env python
import numpy
ra=numpy.random
la=numpy.linalg

size=2000
a=ra.standard_normal((size, size))
b=ra.standard_normal((size, size))
c=la.linalg.dot(a, b)
print c
```

Create a wrapper:

```
% cat wrapper.py
#!/usr/bin/env python

# setenv PYTHONPATH $PET_HOME/pkg/tau-2.17.3/ppc64/lib/bindings-gnu-python-pdt

import tau

def OurMain():
    import foo

tau.run('OurMain()')
```

ParaTools

62

Generate a Python Profile

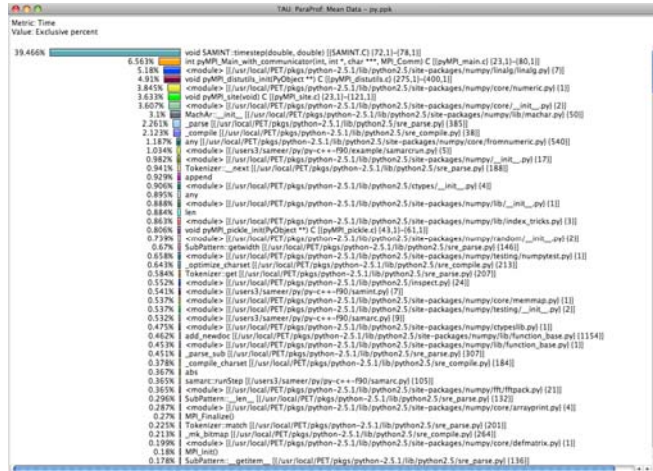
```
% setenv TAU_MAKEFILE /usr/global/tools/tau/ibm64
    /lib/Makefile.tau-python-pdt
% set path=(/usr/global/tools/tau/ibm64/bin $path)
% cat wrapper.py
import tau
def OurMain():
    import foo
    tau.run('OurMain()')
Uninstrumented:
% ./foo.py
Instrumented:
% setenv PYTHONPATH <taudir>/ibm64/<lib>/bindings-python-pdt
(same options string as TAU_MAKEFILE)
% setenv LD_LIBRARY_PATH <taudir>/x86_64/lib/bindings-python-pdt\:
$LD_LIBRARY_PATH
% ./wrapper.py

Wrapper invokes foo and generates performance data
% pprof/paraprof
```

ParaTools

Usage Scenarios: Mixed Python+F90+C+pyMPI

- Goal: Generate multi-level instrumentation for Python+MPI+C+F90+C++ ...



ParaTools

Generate a Multi-Language Profile w/ Python

```

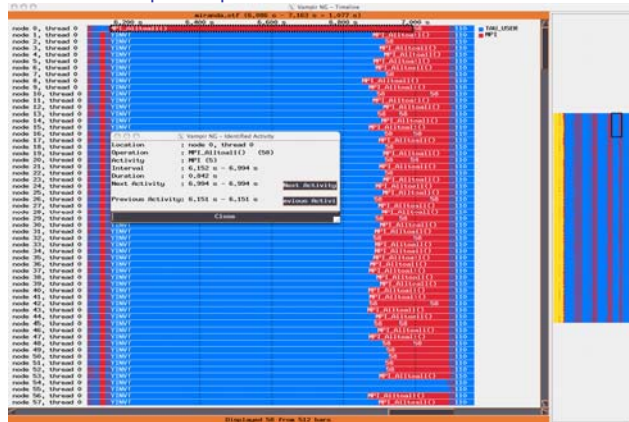
% setenv TAU_MAKEFILE /usr/global/tools/tau/ibm64
      /lib/Makefile.tau-python-mpi-pdt
% set path=(/usr/global/tools/tau/ibm64/bin $path)
% setenv TAU_OPTIONS '-optShared -optVerbose...'
(Python needs shared object based TAU library)
% make F90=tau_f90.sh CXX=tau_cxx.sh CC=tau_cc.sh (build libs, pyMPI w/TAU)
% cat wrapper.py
import tau
def OurMain():
    import App
    tau.run('OurMain()')
Uninstrumented:
% mpirun.lsf /usr/global/tools/.unsupported/pyMPI-2.5b0/bin/pyMPI ./App.py
Instrumented:
% setenv PYTHONPATH <taudir>/x86_64/<lib>/bindings-python-mpi-pdt
(same options string as TAU_MAKEFILE)
% setenv LD_LIBRARY_PATH <taudir>/x86_64/lib/bindings-python-mpi-pdt\:
$LD_LIBRARY_PATH
% mpirun -np 4 /usr/global/tools/.unsupported/pyMPI-2.5b0-TAU/bin/pyMPI
./wrapper.py (Instrumented pyMPI with wrapper.py)

```

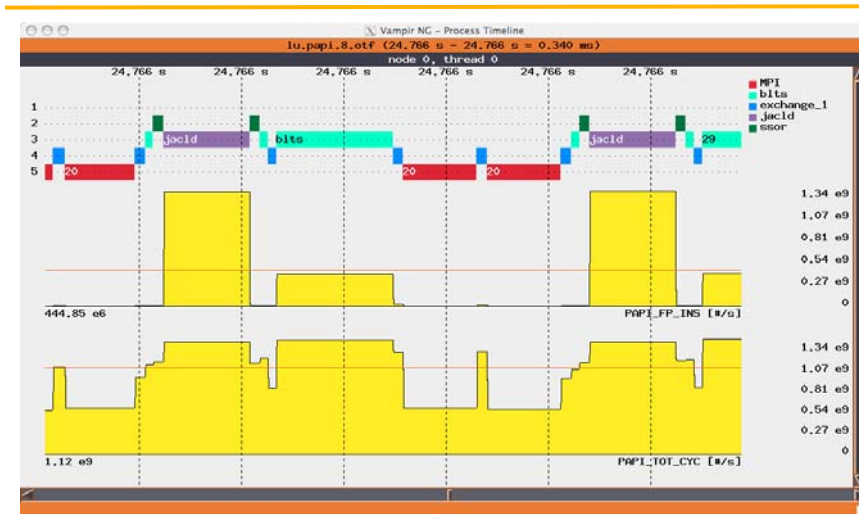
65

Usage Scenarios: Generating a Trace File

- Goal: Identify the temporal aspect of performance. What happens in my code at a given time? When?
- Event trace visualized in Vampir/Jumpshot



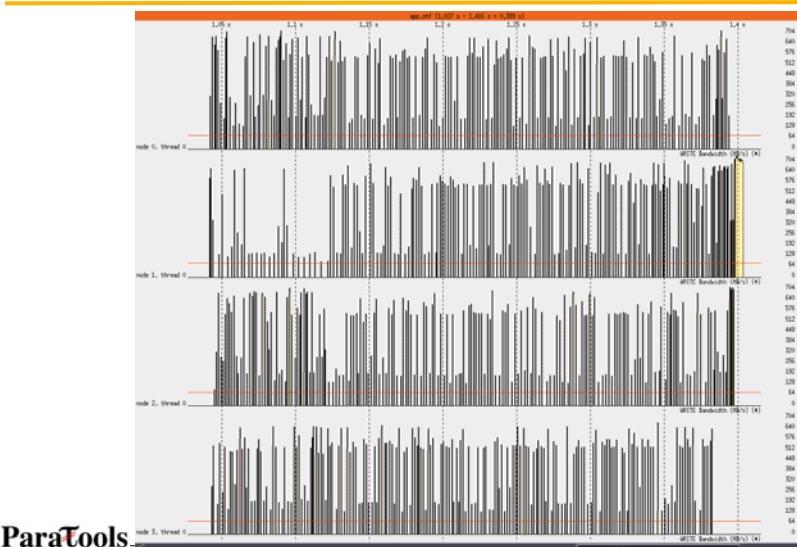
VNG Process Timeline with PAPI Counters



ParaTools

67

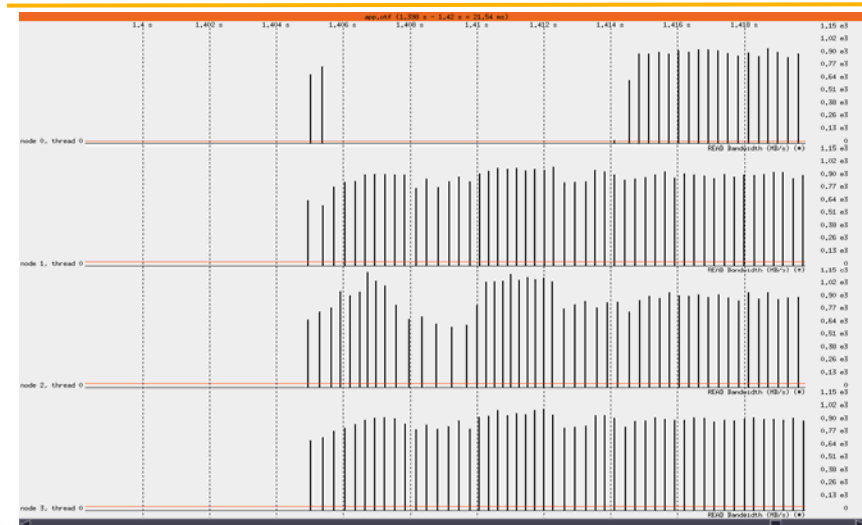
Vampir Counter Timeline Showing I/O BW



ParaTools

68

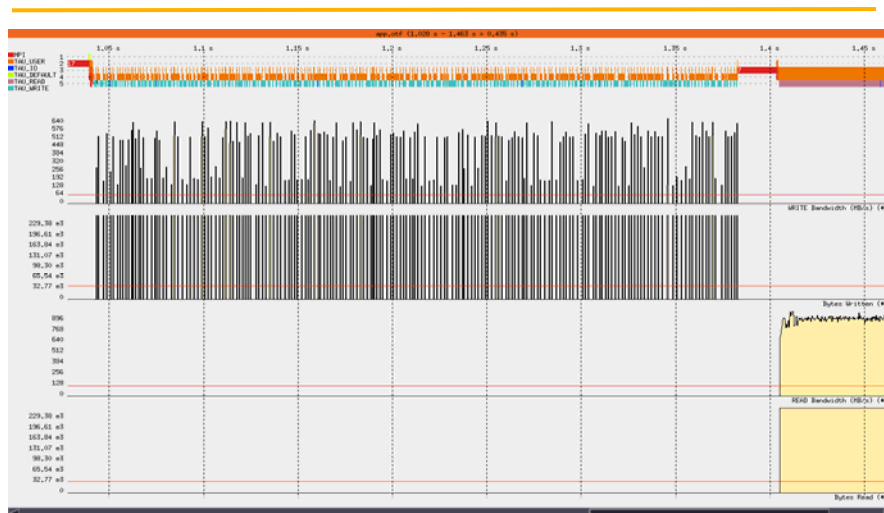
TAU: I/O Instrumentation for Read Bandwidth: Vampir



ParaTools

69

Vampir Process Timeline for Rank 0 (IOR, LLNL)



ParaTools

70

Generate a Trace File

```

% setenv TAU_MAKEFILE /usr/global/tools/tau/training/tau-
2.18.2/bgp
                               /lib/Makefile.tau-mpi-pdt-trace
or setenv TAU_TRACE 1 (in TAU v2.18.2+)
% set path=(/usr/global/tools/tau/training/tau-
2.18.2/bgp/bin $path)
% make F90=tau_f90.sh
(Or edit Makefile and change F90=tau_f90.sh)
% qsub run.job
% tau_treemerge.pl
(merges binary traces to create tau.trc and tau.edf files)
JUMPSHOT:
% tau2slog2 tau.trc tau.edf -o app.slog2
% jumpshot app.slog2
OR
VAMPIR:
% tau2otf tau.trc tau.edf app.otf -n 4 -z
(4 streams, compressed output trace)
% vampir app.otf
(or vng client with vngd server).

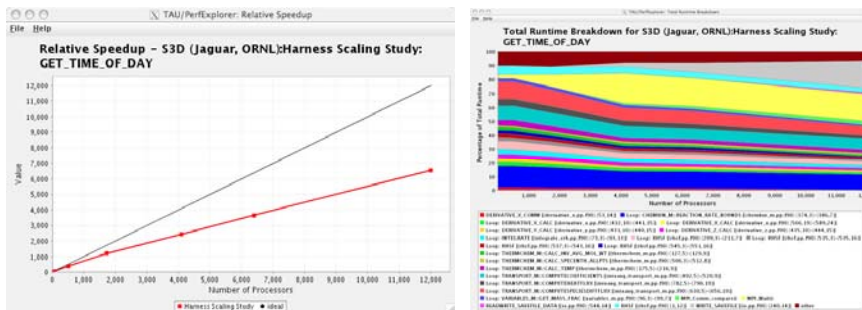
```

ParaTools

71

Usage Scenarios: Evaluate Scalability

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



ParaTools

72

Evaluate Scalability using PerfExplorer Charts

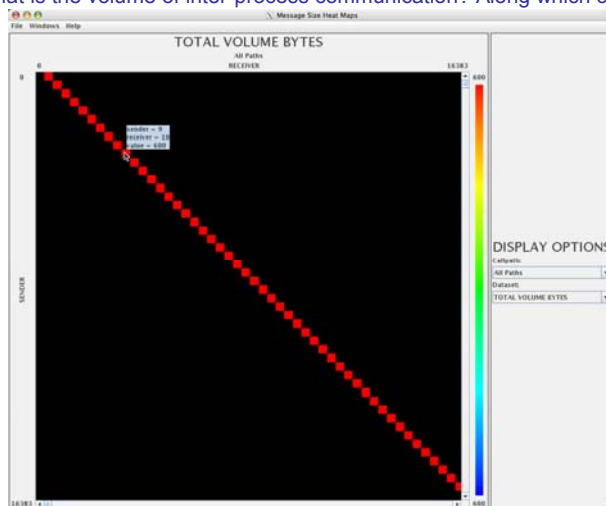
```
% setenv TAU_MAKEFILE /usr/global/tools/tau/training/tau-2.18.2/bgp
/lib/Makefile.tau-mpi-pdt
% set path=(/usr/global/tools/tau/training/tau-2.18.2/bgp/bin $path)
% make F90=tau_f90.sh
(Or edit Makefile and change F90=tau_f90.sh)
% qsub run1p.job
% paraprof --pack 1p.ppk
% qsub run2p.job ...
% paraprof --pack 2p.ppk ... and so on.
On your client:
% perfdmf_configure --create-default
(Chooses derby, blank user/passwd, yes to save passwd, defaults)
% perfexplorer_configure
(Yes to load schema, defaults)
% paraprof
(load each trial: DB -> Add Trial -> Type (Paraprof Packed Profile) -> OK) OR use
perfdmf_loadtrial
Then,
% perfexplorer
(Select experiment, Menu: Charts -> Speedup)
```

ParaTools

73

Communication Matrix Display

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



ParaTools

74

Evaluate Scalability using PerfExplorer Charts

```
% setenv TAU_MAKEFILE /usr/global/tools/tau/training/tau-2.18.2/bgp
/lib/Makefile.tau-mpi-pdt
% set path=(/usr/global/tools/tau/training/tau-2.18.2/bgp/bin $path)
% make F90=tau_f90.sh
(Or edit Makefile and change F90=tau_f90.sh)
% setenv TAU_COMM_MATRIX 1

% qsub run.job (setting the environment variables)

% paraprof
(Windows -> Communication Matrix)
```

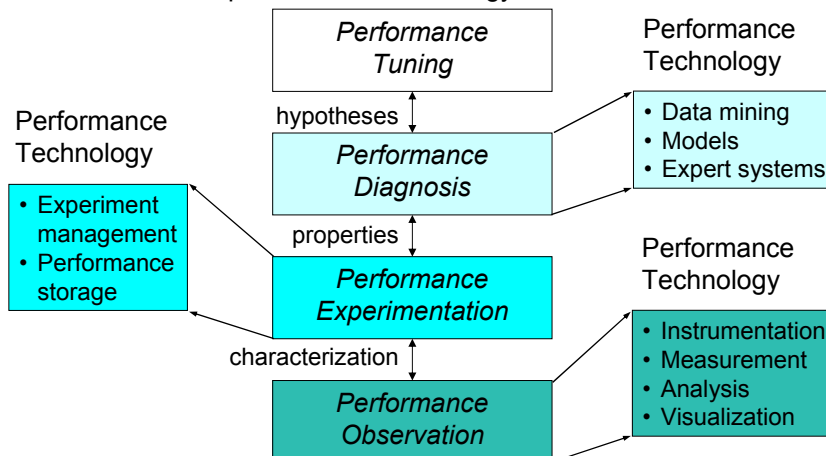
Labs

- Add one of
source /usr/global/tools/tau/training/src/tau.bashrc
or
source /usr/global/tools/tau/training/src/tau.cshrc
to the end of your **.login** file (for bash or csh/tcsh users respectively)
These files contain LLNL specific location information.
- wget <http://www.paratools.com/lnl09/workshop.tar.gz>
or
cp /usr/global/tools/tau/training/src/workshop.tar.gz .
and follow the README file.

Part II: Introduction to Performance Engineering

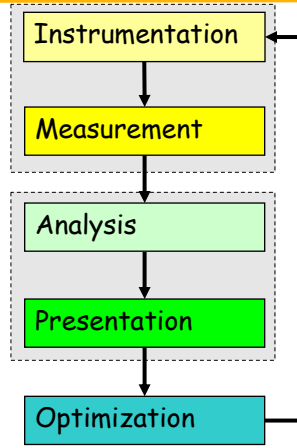
Performance Engineering

- Optimization process
- Effective use of performance technology



Performance Optimization Cycle

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



Parallel Performance Properties

- Parallel code performance is influenced by both sequential and parallel factors?
- Sequential factors
 - Computation and memory use
 - Input / output
- Parallel factors
 - Thread / process interactions
 - Communication and synchronization

Performance Observation

- Understanding performance requires observation of performance properties
- Performance tools and methodologies are primarily distinguished by what observations are made and how
 - What aspects of performance factors are seen
 - What performance data is obtained
- Tools and methods cover broad range

Metrics and Measurement

- Observability depends on measurement
- A metric represents a type of measured data
 - Count, time, hardware counters
- A measurement records performance data
 - Associates with program execution aspects
- Derived metrics are computed
 - Rates (e.g., flops)
- Metrics / measurements decided by need

Execution Time

- Wall-clock time
 - Based on realtime clock
- Virtual process time
 - Time when process is executing
 - serial time and system time
 - Does not include time when process is stalled
- Parallel execution time
 - Runs whenever any parallel part is executing
 - Global time basis

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* (ESB)
- Performance attribution is inferred
 - Determined by execution context (state)
 - Observation resolution determined by interrupt period
 - Performance data associated with context for period

Direct Observation: Events

- Event types
 - Interval events (begin/end events)
 - measures performance between begin and end
 - metrics monotonically increase
 - Atomic events
 - used to capture performance data state
- Code events
 - Routines, classes, templates
 - Statement-level blocks, loops
- User-defined events
 - Specified by the user
- Abstract mapping events

Direct Observation: Instrumentation

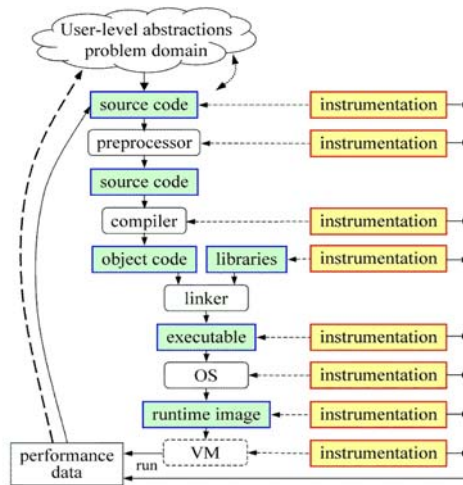
- Events defined by instrumentation access
- Instrumentation levels
 - Source code
 - Object code
 - Runtime system
 - Library code
 - Executable code
 - Operating system
- Different levels provide different information
- Different tools needed for each level
- Levels can have different granularity

Direct Observation: Techniques

- Static instrumentation
 - Program instrumented prior to execution
- Dynamic instrumentation
 - Program instrumented at runtime
- Manual and automatic mechanisms
- Tool required for automatic support
 - Source time: preprocessor, translator, compiler
 - Link time: wrapper library, preload
 - Execution time: binary rewrite, dynamic
- Advantages / disadvantages

Direct Observation: Mapping

- Associate performance data with high-level semantic abstractions
- Abstract events at user-level provide semantic context



ParaTools

89

Indirect Observation: Events/Triggers

- Events are actions external to program code
 - Timer countdown, HW counter overflow, ...
 - Consequence of program execution
 - Event frequency determined by:
 - Type, setup, number enabled (exposed)
- Triggers used to invoke measurement tool
 - Traps when events occur (interrupt)
 - Associated with events
 - May add differentiation to events

ParaTools

90

Indirect Observation: Context

- When events trigger, execution context determined at time of trap (interrupt)
 - Access to PC from interrupt frame
 - Access to information about process/thread
 - Possible access to call stack
 - requires call stack unwinder
- Assumption is that the context was the same during the preceding period
 - Between successive triggers
 - Statistical approximation valid for long running programs

Direct / Indirect Comparison

- Direct performance observation
 - ☺ Measures performance data exactly
 - ☺ Links performance data with application events
 - ☹ Requires instrumentation of code
 - ☹ Measurement overhead can cause execution intrusion and possibly performance perturbation
- Indirect performance observation
 - ☺ Argued to have less overhead and intrusion
 - ☺ Can observe finer granularity
 - ☺ No code modification required (may need symbols)
 - ☹ Inexact measurement and attribution

Measurement Techniques

- When is measurement triggered?
 - External agent (indirect, asynchronous)
 - interrupts, hardware counter overflow, ...
 - Internal agent (direct, synchronous)
 - through code modification
- How are measurements made?
 - Profiling
 - summarizes performance data during execution
 - per process / thread and organized with respect to context
 - Tracing
 - trace record with performance data and timestamp
 - per process / thread

Measured Performance

- Counts
- Durations
- Communication costs
- Synchronization costs
- Memory use
- Hardware counts
- System calls

Critical issues

- Accuracy
 - Timing and counting accuracy depends on resolution
 - Any performance measurement generates overhead
 - Execution on performance measurement code
 - Measurement overhead can lead to intrusion
 - Intrusion can cause perturbation
 - alters program behavior
- Granularity
 - How many measurements are made
 - How much overhead per measurement
- Tradeoff (general wisdom)
 - Accuracy is inversely correlated with granularity

Performance Problem Solving Goals

- Answer questions at multiple levels of interest
 - High-level performance data spanning dimensions
 - machine, applications, code revisions, data sets
 - examine broad performance trends
 - Data from low-level measurements
 - use to predict application performance
- Discover general correlations
 - performance and features of external environment
 - Identify primary performance factors
- Benchmarking analysis for application prediction
- Workload analysis for machine assessment

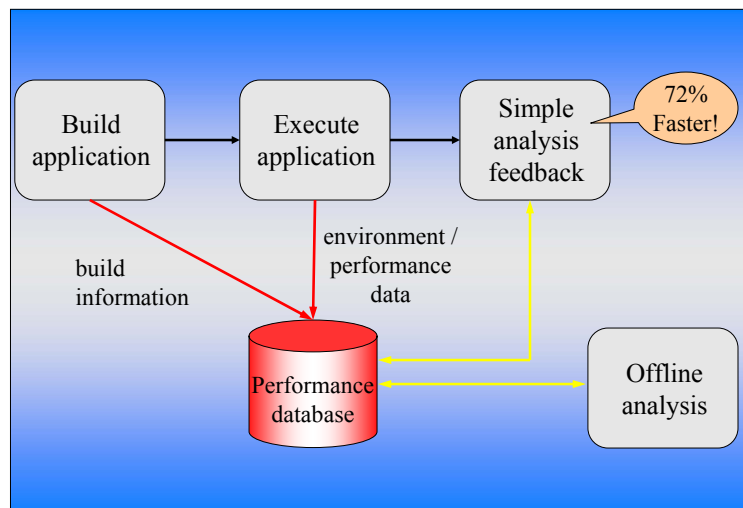
Performance Analysis Questions

- How does performance vary with different compilers?
- Is poor performance correlated with certain OS features?
- Has a recent change caused unanticipated performance?
- How does performance vary with MPI variants?
- Why is one application version faster than another?
- What is the reason for the observed scaling behavior?
- Did two runs exhibit similar performance?
- How are performance data related to application events?
- Which machines will run my code the fastest and why?
- Which benchmarks predict my code performance best?

ParaTools

97

Automatic Performance Analysis



ParaTools

98

Performance Data Management

- Performance diagnosis and optimization involves multiple performance experiments
- Support for common performance data management tasks augments tool use
 - Performance experiment data and metadata storage
 - Performance database and query
- What type of performance data should be stored?
 - Parallel profiles or parallel traces
 - Storage size will dictate
 - Experiment metadata helps in meta analysis tasks
- Serves tool integration objectives

Metadata Collection

- Integration of metadata with each parallel profile
 - Separate information from performance data
- Three ways to incorporate metadata
 - Measured hardware/system information
 - CPU speed, memory in GB, MPI node IDs, ...
 - Application instrumentation (application-specific)
 - Application parameters, input data, domain decomposition
 - Capture arbitrary name/value pair and save with experiment
 - Data management tools can read additional metadata
 - Compiler flags, submission scripts, input files, ...
 - Before or after execution
- Enhances analysis capabilities

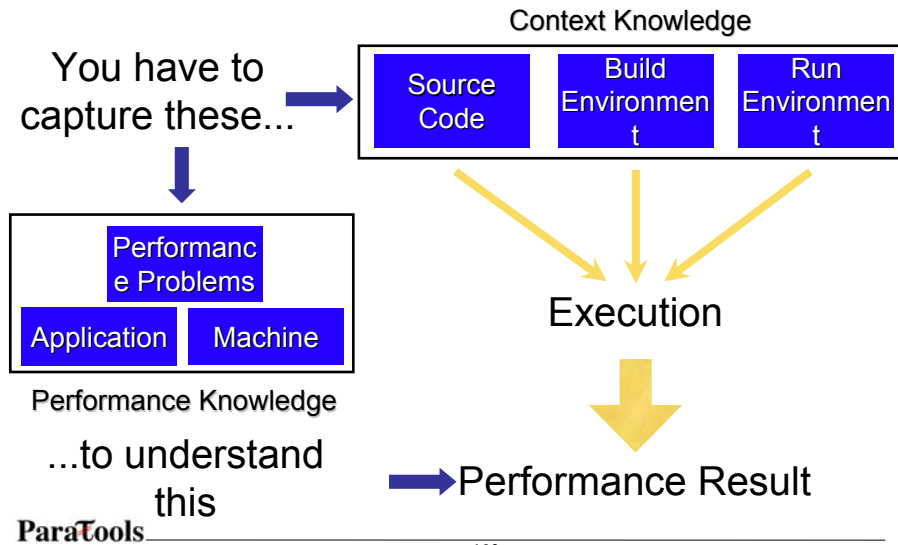
Performance Data Mining

- Conduct parallel performance analysis in a systematic, collaborative and reusable manner
 - Manage performance complexity and automate process
 - Discover performance relationship and properties
 - Multi-experiment performance analysis
- Data mining applied to parallel performance data
 - Comparative, clustering, correlation, characterization, ...
 - Large-scale performance data reduction
- Implement extensible analysis framework
 - Abstraction / automation of data mining operations
 - Interface to existing analysis and data mining tools

How to explain performance?

- Should not just redescrbed performance results
- Should explain performance phenomena
 - What are the causes for performance observed?
 - What are the factors and how do they interrelate?
 - Performance analytics, forensics, and decision support
- Add *knowledge* to do more intelligent things
 - Automated analysis needs good informed feedback
 - Performance model generation requires interpretation
- Performance knowledge discovery framework
 - Integrating meta-information
 - Knowledge-based performance problem solving

Metadata and Knowledge Role



Performance Optimization Process

- Performance characterization
 - Identify major performance contributors
 - Identify sources of performance inefficiency
 - Utilize timing and hardware measures
- Performance diagnosis (Performance Debugging)
 - Look for conditions of performance problems
 - Determine if conditions are met and their severity
 - What and where are the performance bottlenecks
- Performance tuning
 - Focus on dominant performance contributors
 - Eliminate main performance bottlenecks

Part III: PAPI

University of Tennessee, Knoxville

What's PAPI?



- Middleware to provide a consistent programming interface for the performance counter hardware found in most major micro-processors.
- Countable events are defined in two ways:
 - platform-neutral *preset* events
 - Platform-dependent native events
- Presets can be **derived** from multiple *native events*
- All events are referenced by name and collected in EventSets for sampling
- Events can be **multiplexed** if counters are limited
- Statistical sampling implemented by:
 - Hardware overflow if supported by the platform
 - Software overflow with timer driven sampling

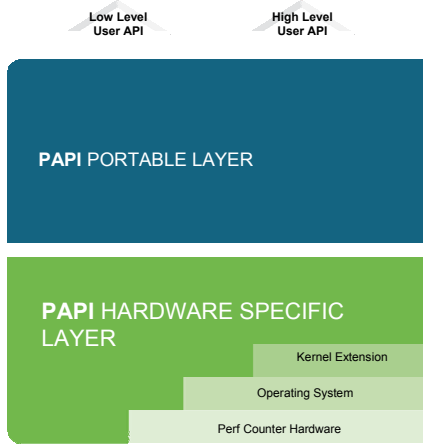
PAPI Counter Interfaces

PAPI provides 3 interfaces to the underlying counter hardware:

1. A Low Level API manages hardware events in user defined groups called EventSets, and provides access to advanced features.
2. A High Level API provides the ability to start, stop and read the counters for a specified list of events.
3. Graphical and end-user tools provide facile data collection and visualization.

ParaTools

3rd Party and GUI Tools



107

PAPI Preset Events

◆Preset Events

- Standard set of over 100 events for application performance tuning
- No standardization of the exact definition
- Mapped to either single or linear combinations of native events on each platform
- Use `papi_avail` utility to see what preset events are available on a given platform

Level 1 Cache

PAPIL1_DCH:	Level 1 data cache hits
PAPIL1_DCA:	Level 1 data cache accesses
PAPIL1_DCR:	Level 1 data cache reads
PAPIL1_DCW:	Level 1 data cache writes
PAPIL1_DCM:	Level 1 data cache misses
PAPIL1_ICH:	Level 1 instruction cache hits
PAPIL1_ICA:	Level 1 instruction cache accesses
PAPIL1_ICR:	Level 1 instruction cache reads
PAPIL1_ICW:	Level 1 instruction cache writes
PAPIL1_ICM:	Level 1 instruction cache misses
PAPIL1_TCH:	Level 1 total cache hits
PAPIL1_TCA:	Level 1 total cache accesses
PAPIL1_TCR:	Level 1 total cache reads
PAPIL1_TCW:	Level 1 total cache writes
PAPIL1_TCM:	Level 1 cache misses
PAPIL1_LDM:	Level 1 load misses
PAPIL1_STM:	Level 1 store misses

ParaTools

108

PAPI Native Events

- Native Events
 - Any event countable by the CPU
 - Same interface as for preset events
 - Use *papi_native_avail* utility to see all available native events
- Use *papi_event_chooser* utility to select a compatible set of events

```
PRESET,  
PAPI_L2_DCA,  
DERIVED_ADD,  
L2_LD:SELF:ANY:MESI,  
L2_ST:SELF:MESI
```

```
{  
  .papi_name = "L2_ST",  
  .papi_code = 0x20,  
  .papi_flags = PAPI_CACHE_ACCESS,  
  .papi_desc = "L2 store requests",  
  .papi_events = {  
    {  
      .papi_name = "MESI",  
      .papi_desc = "Any cacheLine access",  
      .papi_code = 0x1f  
    },  
    {  
      .papi_name = "I_SNOOP",  
      .papi_desc = "Invalid cacheLine",  
      .papi_code = 0x1e  
    },  
    {  
      .papi_name = "S_SNOOP",  
      .papi_desc = "Shared cacheLine",  
      .papi_code = 0x1d  
    },  
    {  
      .papi_name = "E_SNOOP",  
      .papi_desc = "Exclusive cacheLine",  
      .papi_code = 0x1c  
    },  
    {  
      .papi_name = "M_SNOOP",  
      .papi_desc = "Modified cacheLine",  
      .papi_code = 0x1b  
    }  
  },  
  {  
    .papi_name = "SELF",  
    .papi_desc = "This core",  
    .papi_code = 0x10  
  },  
  {  
    .papi_name = "BOTH_CORES",  
    .papi_desc = "Both cores",  
    .papi_code = 0x00  
  }  
},  
.papi_events = 7  
};
```

ParaTools

109

PAPI & Multicore

- Multicore is the (near term) future of Petascale computing
- Minimizing resource contention will be key
 - Memory bandwidth
 - Cache sharing
 - Bus and other resource contention

ParaTools

110

Multicore counter support

- AMD Barcelona
 - 4 L3 shared cache events:
 - READ_REQUEST_TO_L3_CACHE
 - L3_CACHE_MISSES
 - L3_FILLS_CAUSED_BY_L2_EVICTIONS
 - L3_EVICTIONS
 - First 3 are qualified per core
- Intel Core2 series:
 - SELF/ANY
 - L2 shared cache, bus, snoop
 - 39 events/~140 are core qualified
- Itanium Montecito:
 - SELF/ANY
 - 50 bus events (~1/6) are core qualified

Extending PAPI beyond the CPU

- PAPI has historically targeted on on-processor performance counters
- Several categories of off-processor counters exist
 - network interfaces: Myrinet, Infiniband, GigE
 - memory interfaces: Cray X1, SeaStar
 - thermal and power interfaces: ACPI, Im-sensors
 - accelerators?
- CHALLENGE:
 - Extend the PAPI interface to address multiple counter domains
 - Preserve the PAPI calling semantics, ease of use, and platform independence for existing applications

Motivation

- Performance counters also exist in off-cpu resources
- All information is valuable for performance optimization
- Increasing cpu counts & power demands place greater importance on:
 - Thermal health and management
 - Power consumption
- Multicore systems require careful resource balancing
- Higher processor & core counts make communications metrics more critical:
 - Bandwidth
 - Latency
 - Dropped packets
 - Bytes transferred

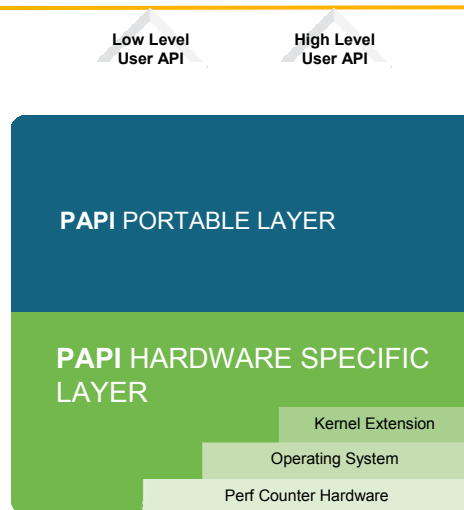
Limitations

- Interfaces are often obscure, unexposed or non-standard
- Performance data (accelerators) can be vastly different than cpus
- Measurements are usually system-wide and asynchronous
 - May not matter on dedicated single-task OS's like Cray Catamount and Blue Gene CNK
 - But matters more for Multicore
- Often very different time scales

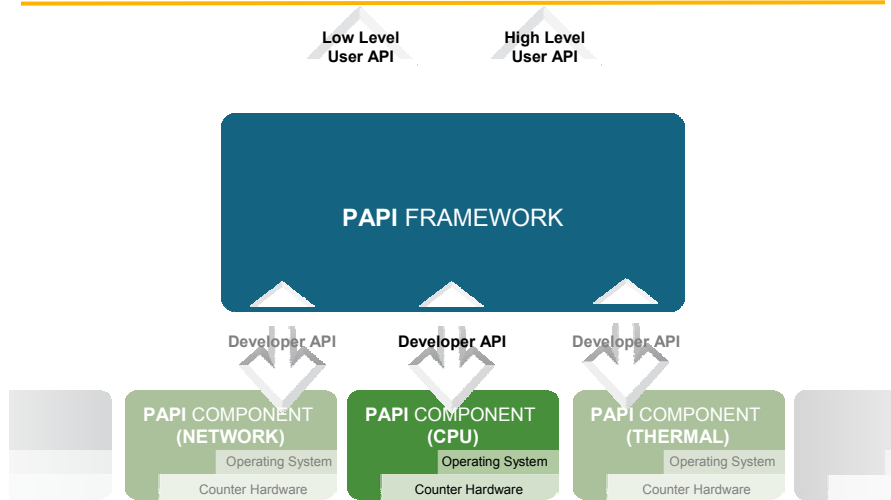
Component PAPI Goals

- Support simultaneous access to on- and off-processor counters
- Isolate hardware dependent code in separable 'component' modules
- Extend platform independent code to support multiple simultaneous components
- Add or modify API calls to support access to any of several components
- Modify build environment for easy selection and configuration of multiple available components

Monolithic 'PAPI Classic'



Component PAPI



ParaTools

117

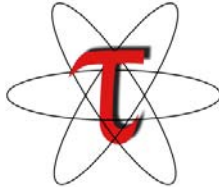
For more information

- PAPI Website: <http://icl.cs.utk.edu/papi/>
 - Software
 - Release notes
 - Documentation
 - Links to tools that use PAPI
 - Mailing/discussion lists

ParaTools

118

Part IV: TAU Internals



Performance Tools FAQ/Concerns

- Does it automatically instrument my code? At the routine level? At the outer-loop level?
- Can it show me where time is spent in my code? PAPI Flops? L1 data cache misses? Can I measure more than one quantity in a trial?
- Does the tool support profiling (runtime summarization) as well as tracing (time-line based displays)? What about profile snapshots? Callpath (parent-child) profiles? Can I use it to easily benchmark codes?
- Can I observe the performance data at runtime as the application executes?
- Can it show me memory utilization? Memory leaks? Mallocs/frees? When and where?
- What about I/O? Can I observe bandwidth of reads/writes? Volume of I/O? What about Kernel events? User space+Kernel?
- What is the typical overhead? Can I reduce it to < 5%? < 1%? Can it compensate and remove timer overhead from performance data? Can it throttle away instrumentation in lightweight routines at runtime to reduce overhead?
- I already have profile data from <XYZ> tool. Can it import my legacy data?
- I prefer <XYZ> performance tool for visualization. Can it hook up with this tool? Are there converters?

Performance Tools FAQ/Concerns (contd.)

- Can I use it for multi-core CPUs? Compare the performance of application running on a single vs. multi-core processor? Can I observe multi-core data snoops, invalidates?
- Can I share the performance data with my colleagues in a secure manner (web/database)? Can it automatically track progress of my application over time (~ 6 mos)? Can I use it for scalability studies? Over multiple platforms?
- Are the GUI client tools available under Linux? MS Windows? Apple?
- Does it run on all Cray, IBM, SGI, HP ... platforms? CNL? Catamount?
- Does it support MPI? MPI2? Threads? Hybrid MPI+Pthreads/MPI+OpenMP?
- Does it support Fortran? C++, C? Java? Python? Python+MPI+F90+C++...?
- Does it support Intel/PGI/PathScale/IBM/Cray/Sun compilers?
- Are tools available in command-line form & GUI? IDE GUI? Web-based? 3D?
- Is it already installed and supported on my HPC system? What about systems at NERSC? ANL? LLNL? LANL? NASA? DoD? NSF sites?...
- Is there support (phone/e-mail) available for the tool? Professional support? For instrumentation? Analysis?
- Will it work on the new <XYZ> HPC platform scheduled for release six months from now?
- Is it free? BSD license? ...

TAU Performance System® Project

- **Tuning and Analysis Utilities (15+ year project effort)**
- **Performance system framework for HPC systems**
 - Integrated, scalable, and flexible
 - Target parallel programming paradigms
- **Integrated toolkit for performance problem solving**
 - Instrumentation, measurement, analysis, and visualization
 - Portable performance profiling and tracing facility
 - Performance data management and data mining
- **Partners**
 - LLNL, ANL, LANL
 - Research Centre Jülich, TU Dresden

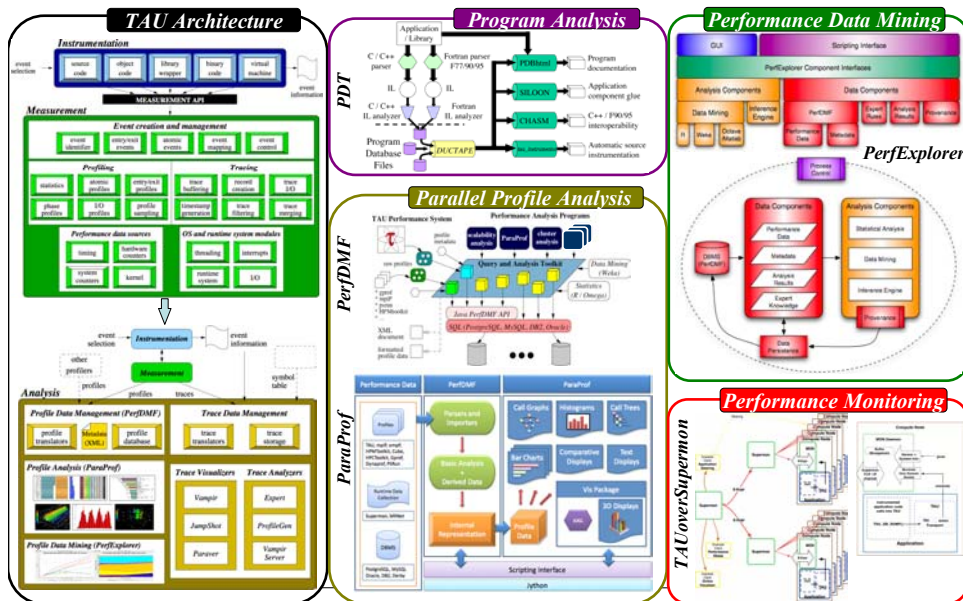
TAU Parallel Performance System Goals

- **Portable (open source) parallel performance system**
 - Computer system architectures and operating systems
 - Different programming languages and compilers
- Multi-level, multi-language performance instrumentation
- **Flexible and configurable performance measurement**
- Support for multiple parallel programming paradigms
 - Multi-threading, message passing, mixed-mode, hybrid, object oriented (generic), component-based
- Support for performance mapping
- Integration of leading performance technology
- **Scalable (very large) parallel performance analysis**

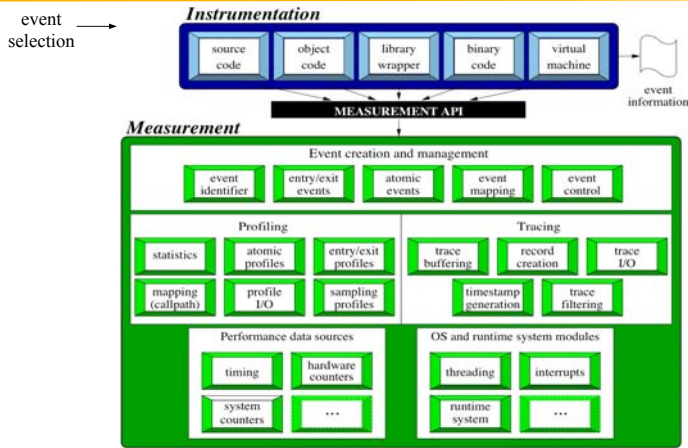
ParaTools

123

TAU Performance System Components



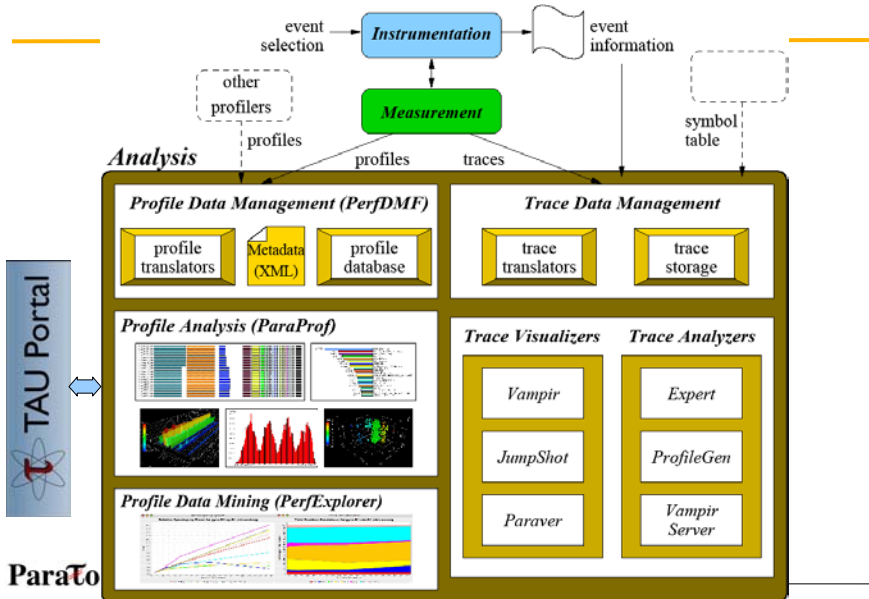
TAU Performance System Architecture



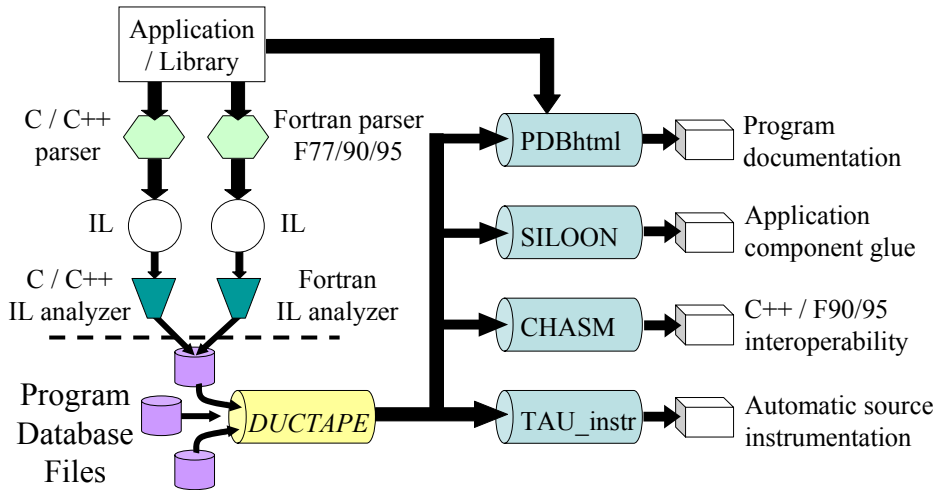
ParaTools

125

TAU Performance System Architecture



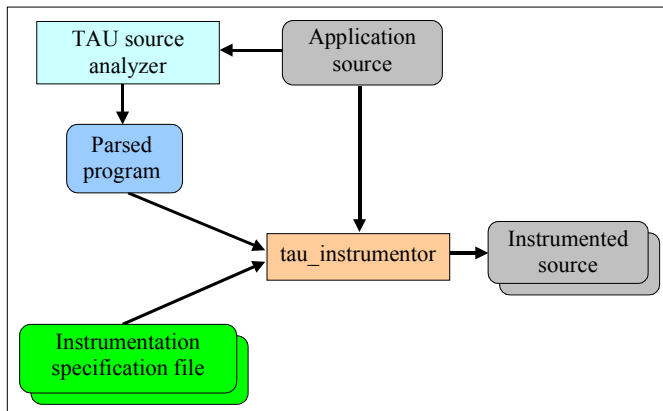
Program Database Toolkit (PDT)



ParaTools

127

Automatic Source-Level Instrumentation in TAU

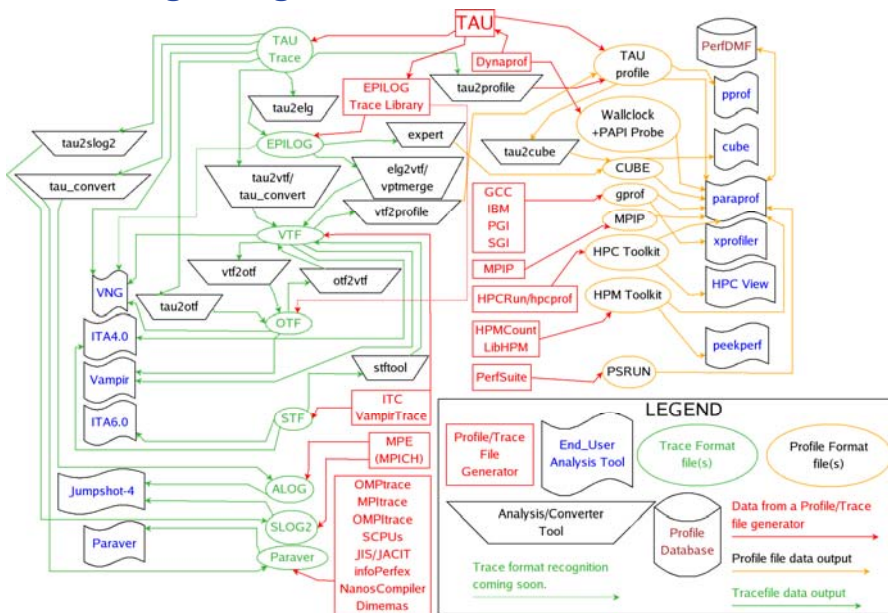


ParaTools

128



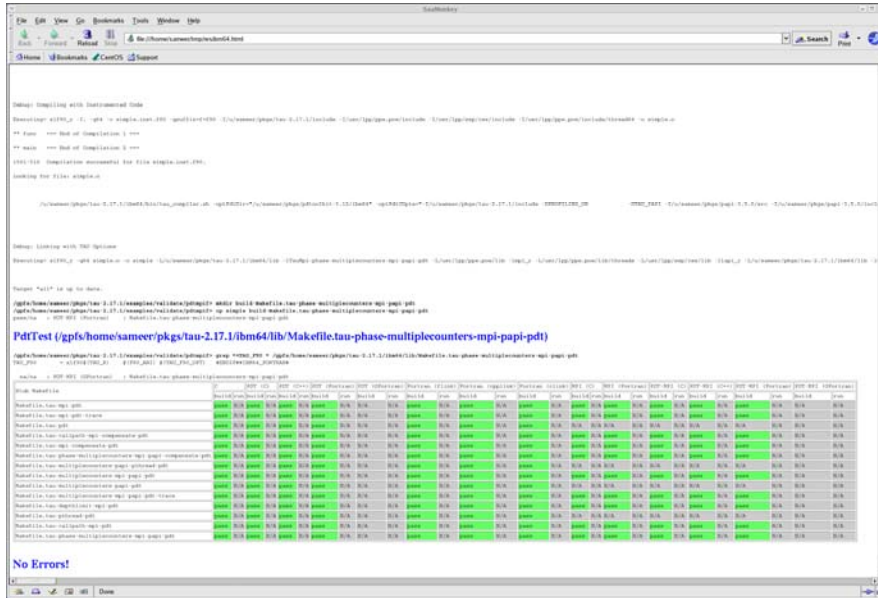
Building Bridges to Other Tools



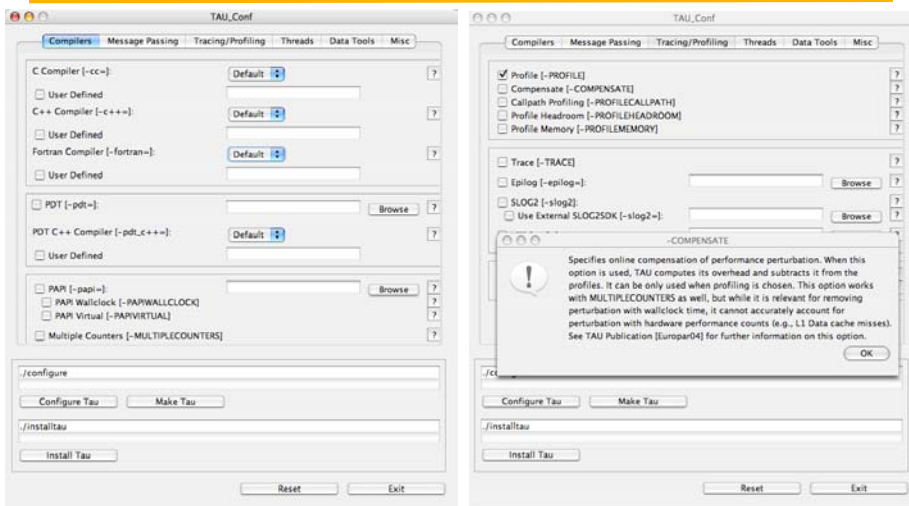
Installing TAU on 64bit AIX

- Install PAPI and PDT
 - PAPI:
 - ./configure –prefix=\$HOME/pkgs/papi-3.5.0;
 - Make ; make install
 - PDT:
 - ./configure –prefix=\$HOME/pkgs/pdt-3.13 –PGI
 - make; make install
- Install TAU:
 - ./installtau –pdt=\$HOME/pkgs/pdt-3.13 –papi=\$HOME/pkgs/papi-3.5.0 –c++=pgCC –cc=pgcc –fortran=pgi –mpiinc=<dir> –mpilib=<dir>
 - Configures multiple typically requested versions for you in ia64/lib/Makefile.tau-* configurations
 - tau_validate –html –build x86_64 >& results.html
 - mozilla results.html

Validating an Install



TAU_SETUP: A GUI for Installing TAU



Upgrading TAU v2.18 configurations to 2.18.2

- Upgrade TAU

- Previous installation in \$HOME/pkgs/tau-2.18
 - `cd tau-2.18.2`
 - `./upgradetau /usr/global/tools/pkgs/tau-2.18`
 - Builds all previous configurations in the current dir
 - You may also upgrade with a new package say PDT 3.14.1
 - `./upgradetau /usr/global/tools/pkgs/tau-2.18 -pdt=/usr/global/tools/pkgs/pdtoolkit-3.14.1`

- Validate your new installation

- `./tau_validate -html -build x86_64 >& results.html`
- `mozilla pwd`/results.html`

Using TAU

- Install TAU

`% ./configure [options]; make clean install`

- Replace the names of your compiler with `tau_f90.sh`, `tau_cxx.sh` and `tau_cc.sh` in your makefiles

- Set environment variables

- Choose the measurement option and compile your code:
 - `setenv TAU_MAKEFILE $TAU/Makefile.tau-mpi-pdt`
 - `setenv TAU_OPTIONS '-optVerbose -optKeepFiles -optPreProcess'`
- At runtime, if more than one metric is measured (-multiplecounters):
 - `setenv COUNTER1 GET_TIME_OF_DAY`
 - `setenv COUNTER2 PAPI_FP_INS`
 - `setenv COUNTER3 PAPI_NATIVE_<native_name>`
 - Use `papi_native_avail`, `papi_avail`, and `papi_event_chooser` to select these preset and native event names

- Build the application, run it, analyze performance data

Using TAU: A brief Introduction

- To instrument source code:
% **setenv TAU_MAKEFILE /usr/global/tools/tau/training/tau-2.18.2/bgp/lib/Makefile.tau-mpi-pdt**
And use tau_f90.sh, tau_cxx.sh or tau_cc.sh as Fortran, C++ or C compilers:
% **mpif90 foo.f90**
changes to
% **tau_f90.sh foo.f90**
- Execute application and then run:
% **pprof (for text based profile display)**
% **paraprof (for GUI)**
- LABS:
% **source /usr/global/tools/tau/training/src/tau.cshrc**
% **cp /usr/global/tools/tau/training/src/workshop.tar.gz .**
and follow instructions in README file

TAU Instrumentation Approach

- Support for standard program events
 - Routines
 - Classes and templates
 - Statement-level blocks
- Support for user-defined events
 - Begin/End events (“user-defined timers”)
 - Atomic events (e.g., size of memory allocated/freed)
 - Selection of event statistics
- Support definition of “semantic” entities for mapping
- Support for event groups
- Instrumentation optimization (eliminate instrumentation in lightweight routines)

TAU Instrumentation

- Flexible instrumentation mechanisms at multiple levels
 - Source code
 - manual (TAU API, TAU Component API)
 - automatic
 - C, C++, F77/90/95 (Program Database Toolkit (*PDT*))
 - OpenMP (directive rewriting (*Opari*), *POMP spec*)
 - Object code
 - pre-instrumented libraries (e.g., MPI using *PMPI*)
 - statically-linked and dynamically-linked
 - Executable code
 - dynamic instrumentation (pre-execution) (*DynInstAPI*)
 - virtual machine instrumentation (e.g., Java using *JVMPI*)
 - Python interpreter based instrumentation at runtime
 - Proxy Components

ParaTools

137

TAU Measurement Approach

- Portable and scalable parallel profiling solution
 - Multiple profiling types and options
 - Event selection and control (enabling/disabling, throttling)
 - Online profile access and sampling
 - Online performance profile overhead compensation
- Portable and scalable parallel tracing solution
 - Trace translation to Open Trace Format (OTF)
 - Trace streams and hierarchical trace merging
- Robust timing and hardware performance support
- Multiple counters (hardware, user-defined, system)
- Performance measurement for CCA component software

ParaTools

138

Using TAU

- Configuration
- Instrumentation
 - Manual
 - MPI – Wrapper interposition library
 - PDT- Source rewriting for C,C++, F77/90/95
 - Compiler-based instrumentation for C, C++, F90
 - OpenMP – Directive rewriting
 - Component based instrumentation – Proxy components
 - Binary Instrumentation
 - DynInstAPI – Runtime Instrumentation/Rewriting binary
 - Java – Runtime instrumentation
 - Python – Runtime instrumentation
- Measurement
- Performance Analysis

ParaTools

139

TAU Measurement System Configuration

- `configure [OPTIONS]`
 - `{-c++=<CC>, -cc=<cc>}` Specify C++ and C compilers
 - `-pdt=<dir>` Specify location of PDT
 - `-opari=<dir>` Specify location of Opari OpenMP tool
 - `-papi=<dir>` Specify location of PAPI
 - `-vampirtrace=<dir>` Specify location of VampirTrace
 - `-mpi[inc/lib]=<dir>` Specify MPI library instrumentation
 - `-dyninst=<dir>` Specify location of DynInst Package
 - `-shmem[inc/lib]=<dir>` Specify PSHMEM library instrumentation
 - `-python[inc/lib]=<dir>` Specify Python instrumentation
 - `-tag=<name>` Specify a unique configuration name
 - `-epilog=<dir>` Specify location of EPILOG
 - `-slog2` Build SLOG2/Jumpshot tracing package
 - `-otf=<dir>` Specify location of OTF trace package
 - `-arch=<architecture>` Specify architecture explicitly
(bgl, xt3,ibm64,ibm64linux...)
 - `{-pthread, -sproc}` Use pthread or SGI sproc threads
 - `-openmp` Use OpenMP threads
 - `-jdk=<dir>` Specify Java instrumentation (JDK)
 - `-fortran=[vendor]` Specify Fortran compiler

ParaTools

140

TAU Measurement System Configuration

- `configure [OPTIONS]`
 - TRACE Generate binary TAU traces
 - PROFILE (default) Generate profiles (summary)
 - PROFILECALLPATH Generate call path profiles
 - PROFILEPHASE Generate phase based profiles
 - PROFILEMEMORY Track heap memory for each routine
 - PROFILEHEADROOM Track memory headroom to grow
 - MULTIPLECOUNTERS Use hardware counters + time
 - COMPENSATE Compensate timer overhead
 - CPUTIME Use usertime+system time
 - PAPIWALLCLOCK Use PAPI's wallclock time
 - PAPIVIRTUAL Use PAPI's process virtual time
 - SGITIMERS Use fast IRIX timers
 - LINUXTIMERS Use fast x86 Linux timers

TAU Measurement Configuration – Examples

- `./configure -arch=x86_64 -pdt=/usr/global/tools/pkgs/pdtoolkit-3.14 -mpi` Configure using PDT and MPI
- `./configure -arch=x86_64 -papi=/usr/global/tools/pkgs/papi-3.6.2 -pdt=<dir> -mpi -MULTIPLECOUNTERS; make clean install`
 - Use PAPI counters (one or more) with C/C++/F90 automatic instrumentation. Also instrument the MPI library.
- Typically configure multiple measurement libraries
- Each configuration creates a unique `<arch>/lib/Makefile.tau<options>` stub makefile. It corresponds to the configuration options used. e.g.,
 - `$(PET_HOME)/tau/x86_64/lib/Makefile.tau-mpi-pdt`
 - `$(PET_HOME)/tau/x86_64/lib/Makefile.tau-multiplecounters-mpi-papi-pdt`

TAU Measurement Configuration – Examples

```
% cd $(PET_HOME)/tau/x86_64/lib; ls Makefile.*pgi
```

```
Makefile.tau-pdt
```

```
Makefile.tau-mpi-pdt
```

```
Makefile.tau-callpath-mpi-pdt
```

```
Makefile.tau-mpi-pdt-trace
```

```
Makefile.tau-mpi-compensate-pdt
```

```
Makefile.tau-multiplecounters-mpi-papi-pdt
```

```
Makefile.tau-multiplecounters-mpi-papi-pdt-trace
```

```
Makefile.tau-mpi-papi-pdt-epilog-scalasca-trace
```

```
Makefile.tau-pdt...
```

- For an MPI+F90 application, you may want to start with:

```
Makefile.tau-mpi-pdt
```

- Supports MPI instrumentation & PDT for automatic source instrumentation for PGI compilers

ParaTools

143

Configuration Parameters in Stub Makefiles

- Each TAU stub Makefile resides in <tau>/<arch>/lib directory
- Variables:
 - **TAU_CXX** Specify the C++ compiler used by TAU
 - **TAU_CC, TAU_F90** Specify the C, F90 compilers
 - **TAU_DEFS** Defines used by TAU. Add to CFLAGS
 - **TAU_LDFLAGS** Linker options. Add to LDFLAGS
 - **TAU_INCLUDE** Header files include path. Add to CFLAGS
 - **TAU_LIBS** Statically linked TAU library. Add to LIBS
 - **TAU_SHLIBS** Dynamically linked TAU library
 - **TAU_MPI_LIBS** TAU's MPI wrapper library for C/C++
 - **TAU_MPI_FLIBS** TAU's MPI wrapper library for F90
 - **TAU_FORTRANLIBS** Must be linked in with C++ linker for F90
 - **TAU_CXXLIBS** Must be linked in with F90 linker
 - **TAU_INCLUDE_MEMORY** Use TAU's malloc/free wrapper lib
 - **TAU_DISABLE** TAU's dummy F90 stub library
 - **TAU_COMPILER** Instrument using tau_compiler.sh script
- Each stub makefile encapsulates the parameters that TAU was configured with
- It represents a specific instance of the TAU libraries. TAU scripts use stub makefiles to identify what performance measurements are to be performed.

ParaTools

144

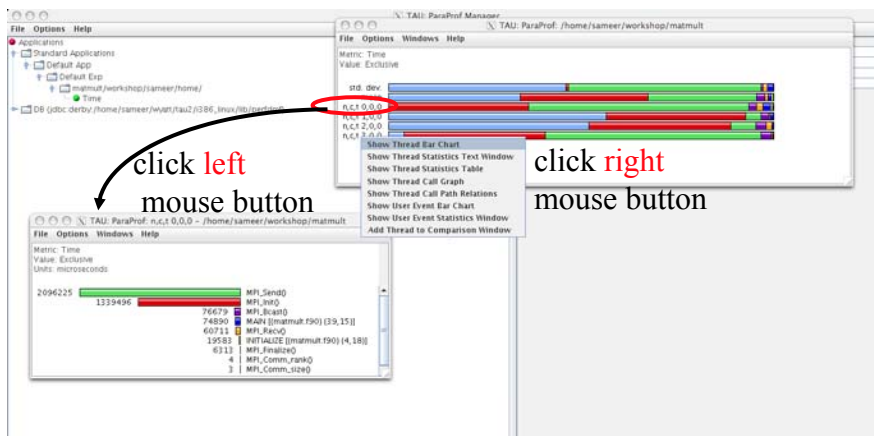
Using TAU

- **Install TAU**
% configure [options]; make clean install
- **Typically modify application makefile and choose TAU configuration**
 - Select TAU's stub makefile, change name of compiler in Makefile
- % setenv TAU_MAKEFILE /usr/global/tools/tau/training/tau-2.18.2/bgp/lib/Makefile.tau-mpi-pdt
- % setenv TAU_OPTIONS '-optVerbose -optKeepFiles ...'
- – F90 = tau_f90.sh CXX = tau_cxx.sh CC = tau_cc.sh
- **Set environment variables**
 - Directory where profiles/traces are to be stored/counter selection
- **Execute application**
% qsub run.cray.job
- **Analyze performance data**
 - paraprof, vampir, pprof, paraver ...

ParaTools

145

ParaProf Main Window



% paraprof mat mult.ppk

ParaTools

146

TAU's MPI Wrapper Interposition Library

- Uses standard MPI Profiling Interface
 - Provides name shifted interface
 - MPI_Send = PMPI_Send
 - Weak bindings
- Interpose TAU's MPI wrapper library between MPI and TAU
 - -Impi replaced by -ITauMpi -Ipmi -Impi
- No change to the source code!
 - Just re-link the application to generate performance data
 - `setenv TAU_MAKEFILE <dir>/<arch>/lib/Makefile.tau-mpi -[options]`
 - Use `tau_cxx.sh`, `tau_f90.sh` and `tau_cc.sh` as compilers

Runtime MPI Shared Library Instrumentation

- We can now interpose the MPI wrapper library for applications that have already been compiled
 - No re-compilation or re-linking necessary!
- Uses LD_PRELOAD for Linux
- On AIX, TAU uses MPI_EUILIB / MPI_EUILIBPATH
- Simply compile TAU with MPI support and prefix your MPI program with `tauex`
 - `% mpirun -np 4 tauex a.out`
- Requires shared library MPI - does not work on XT3
- Approach will work with other shared libraries

-PROFILE Configuration Option

- Generates flat profiles (one for each MPI process)
 - It is the default option.
- Uses wallclock time (gettimeofday() sys call)
- Calculates exclusive, inclusive time spent in each timer and number of calls

% pprof

%Time	Exclusive msec	Inclusive msec	#Call	#Subrs	Inclusive Name
100.0	1	3111.293	1	15	191293269 apputv
99.8	3,667	3110.483	3	37817	63427920 read_input
67.1	491	2108.326	37200	37200	3420 exchange_1
44.0	4,464	1125.159	9300	18600	9157 write
29.5	1181.436	1181.436	18600	0	4217 MPI_Recv()
16.2	6,778	56,407	9300	18600	6065 write
3.9	50,142	50,142	18204	0	2811 MPI_Send()
2.6	24,491	31,031	301	602	10398 rwa
2.6	7,501	7,501	9300	0	607 MPI_Init()
2.6	838	6,594	604	1812	10918 exchange_3
2.4	6,990	6,990	9300	0	709 MPI_Finalize()
0.2	4,989	4,989	608	0	8206 MPI_Init()
0.2	0.44	450	1	4	40081 MPI_Init()
0.1	398	399	1	39	399634 MPI_Init()
0.1	140	247	1	47616	247056 MPI_Init()
0.1	131	131	57252	0	2 MPI_Init()
0.1	89	103	1	2	103160 MPI_Init()
0.1	0.966	96	1	2	96488 MPI_Init()
0.0	95	95	9	0	10603 MPI_Init()
0.0	26	44	1	7937	44678 MPI_Init()
0.0	24	24	608	0	40 MPI_Init()
0.0	15	15	1	5	15630 MPI_Init()
0.0	7	12	3	1700	12330 MPI_Init()
0.0	3	3	3	0	2993 MPI_Init()
0.0	3	3	3	0	491 MPI_Init()
0.0	1	1	1	6	3674 MPI_Init()
0.0	0.116	0.837	1	0	1007 MPI_Init()
0.0	0.512	0.512	1	0	637 MPI_Init()
0.0	0.121	0.121	1	2	353 MPI_Init()
0.0	0.024	0.191	1	2	191 MPI_Init()
0.0	0.1003	0.1003	1	6	17 MPI_Init()

ParaTools

-MULTIPLECOUNTERS Configuration Option

- Instead of one metric, profile or trace with more than one metric
 - Set environment variables COUNTER[1-25] to specify the metric
 - % setenv COUNTER1 GET_TIME_OF_DAY
 - % setenv COUNTER2 PAPI_L2_DCM
 - % setenv COUNTER3 PAPI_FP_OPS
 - % setenv COUNTER4 PAPI_NATIVE_<native_event>
 - % setenv COUNTER5 P_WALL_CLOCK_TIME ...
 - OR
 - % setenv TAU_METRICS GET_TIME_OF_DAY:PAPI_L2_DCM:PAPI_FP_OPS...
- When used with -TRACE option, the first counter must be GET_TIME_OF_DAY
 - % setenv COUNTER1 GET_TIME_OF_DAY
 - Provides a globally synchronized real time clock for tracing
- multiplecounters appears in the name of the stub Makefile
- Often used with -papi=<dir> to measure hardware performance counters and time
- papi_native_avail and papi_avail are two useful tools

ParaTools

Papi_avail

```
cfel.sameer 66> ./papi_avail | more
Available events and hardware information.
-----
Vendor string and code   : GenuineIntel (1)
Model string and code   : Itanium 2 (1)
CPU Revision            : 5.000000
CPU Megahertz           : 1500.000000
CPU's in this Node     : 28
Nodes in this System    : 1
Total CPU's            : 28
Number Hardware Counters : 4
Max Multiplex Counters  : 32
-----
The following correspond to fields in the PAPI_event_info_t structure.

Name          Code          Avail  Deriv  Description (Note)
PAPI_L1_DCM   0x80000000   Yes    No     Level 1 data cache misses
PAPI_L1_ICM   0x80000001   Yes    No     Level 1 instruction cache misses
PAPI_L2_DCM   0x80000002   Yes    Yes    Level 2 data cache misses
PAPI_L2_ICM   0x80000003   Yes    No     Level 2 instruction cache misses
PAPI_L3_DCM   0x80000004   Yes    Yes    Level 3 data cache misses
PAPI_L3_ICM   0x80000005   Yes    No     Level 3 instruction cache misses
PAPI_L1_TCM   0x80000006   Yes    Yes    Level 1 cache misses
PAPI_L2_TCM   0x80000007   Yes    No     Level 2 cache misses

...

```

Papi_native_avail

```
cfel.sameer 67> ./papi_native_avail | more
Available native events and hardware information.
-----
Vendor string and code   : GenuineIntel (1)
Model string and code   : Itanium 2 (1)
CPU Revision            : 5.000000
CPU Megahertz           : 1500.000000
CPU's in this Node     : 28
Nodes in this System    : 1
Total CPU's            : 28
Number Hardware Counters : 4
Max Multiplex Counters  : 32
-----
The following correspond to fields in the PAPI_event_info_t structure.

Symbol          Event Code  Long Description
Register Name[n]
Register Value[n]
ALAT_CAPACITY_MISS_ALL      0x40000000  ALAT Entry Replaced -- both integer and floating point i
nstructions
ALAT_CAPACITY_MISS_FP      0x40000001  ALAT Entry Replaced -- only floating point instructions
ALAT_CAPACITY_MISS_INT     0x40000002  ALAT Entry Replaced -- only integer instructions

...

```

Papi_event_chooser on IA-64

```

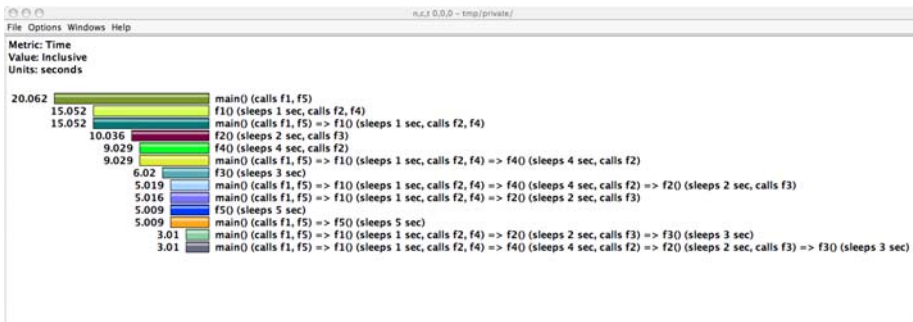
cfe1.sameer 68> ./papi_event_chooser
Usage: eventChooser NATIVE|PRESET evt1 evt2 ...

cfe1.sameer 72> ./papi_event_chooser PRESET PAPI_FP_OPS PAPI_L1_DCM PAPI_TOT_CYC
Test case eventChooser: Available events which can be added with given events.
-----
Vendor string and code   : GenuineIntel (1)
Model string and code    : Itanium 2 (1)
CPU Revision             : 5.000000
CPU Megahertz           : 1500.000000
CPU's in this Node      : 28
Nodes in this System    : 1
Total CPU's             : 28
Number Hardware Counters : 4
Max Multiplex Counters  : 32
-----
Name          Derived Description (Mgr. Note)
PAPI_L1_ICM   No      Level 1 instruction cache misses ()
PAPI_L2_ICM   No      Level 2 instruction cache misses ()
PAPI_L3_ICM   No      Level 3 instruction cache misses
    
```

PAPI_L1_ICM may be counted with these counters on IA64

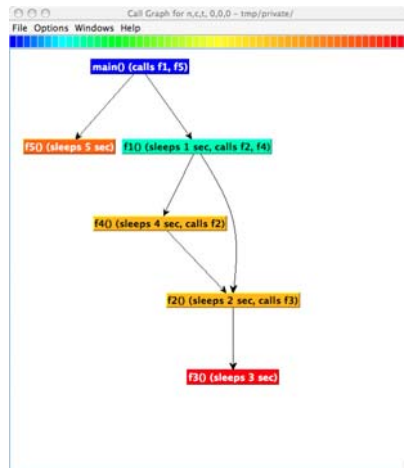
-PROFILECALLPATH Configuration Option

- Generates profiles that show the calling order (edges & nodes in callgraph)
 - A=>B=>C shows the time spent in C when it was called by B and B was called by A
 - Control the depth of callpath using `TAU_CALLPATH_DEPTH` env. Variable
 - `-callpath` in the name of the stub Makefile name
 - In TAU 2.18.2+, any executable can generate callpath profiles using
 - `% setenv TAU_CALLPATH 1`



-PROFILECALLPATH Configuration Option

- Generates program callgraph



ParaTools

155

Profile Measurement – Three Flavors

- **Flat profiles**
 - Time (or counts) spent in each routine (nodes in callgraph).
 - Exclusive/inclusive time, no. of calls, child calls
 - E.g.: MPI_Send, foo, ...
- **Callpath Profiles**
 - Flat profiles, **plus**
 - Sequence of actions that led to poor performance
 - Time spent along a calling path (edges in callgraph)
 - E.g., “main=> f1 => f2 => MPI_Send” shows the time spent in MPI_Send when called by f2, when f2 is called by f1, when it is called by main. Depth of this callpath = 4 (TAU_CALLPATH_DEPTH environment variable)
- **Phase based profiles**
 - Flat profiles, **plus**
 - Flat profiles under a phase (nested phases are allowed)
 - Default “main” phase has all phases and routines invoked outside phases
 - Supports static or dynamic (per-iteration) phases
 - E.g., “IO => MPI_Send” is time spent in MPI_Send in IO phase

ParaTools

156

-DEPTHLIMIT Configuration Option

- Allows users to enable instrumentation at runtime based on the depth of a calling routine on a callstack.
 - Disables instrumentation in all routines a certain depth away from the root in a callgraph
- TAU_DEPTH_LIMIT environment variable specifies depth
 - % setenv TAU_DEPTH_LIMIT 1
enables instrumentation in only “main”
 - % setenv TAU_DEPTH_LIMIT 2
enables instrumentation in main and routines that are directly called by main
- Stub makefile has -depthlimit in its name:
setenv TAU_MAKEFILE <taudir>/<arch>/lib/Makefile.tau-mpi-depthlimit-pdt

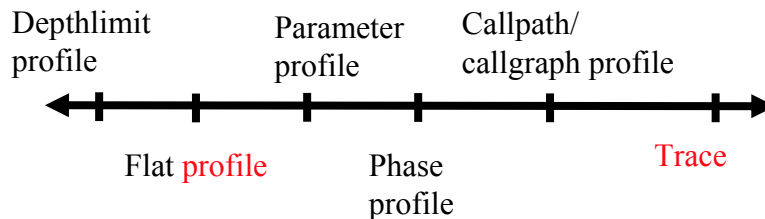
-COMPENSATE Configuration Option

- Specifies online compensation of performance perturbation
- TAU computes its timer overhead and subtracts it from the profiles
- Works well with time or instructions based metrics
- Does not work with level 1/2 data cache misses
- setenv TAU_COMPENSATE 1 (in TAU v2.18.2+)

-TRACE Configuration Option

- Generates event-trace logs, rather than summary profiles
- Traces show when and where an event occurred in terms of location and the process that executed it
- Traces from multiple processes are merged:
 - % tau_treemerge.pl
 - generates tau.trc and tau.edf as merged trace and event definition file
- TAU traces can be converted to Vampir's OTF/VTF3, Jumpshot SLOG2, Paraver trace formats:
 - % tau2otf tau.trc tau.edf app.otf
 - % tau2vtf tau.trc tau.edf app.vpt.gz
 - % tau2slog2 tau.trc tau.edf -o app.slog2
 - % tau_convert -paraver tau.trc tau.edf app.prv
- Stub Makefile has **-trace** in its name
 - % setenv TAU_MAKEFILE <taudir>/<arch>/lib/
 - Makefile.tau-mpi-pdt-**trace**-pdt
 - In TAU 2.18.2+ you may simply use with any configuration, at runtime:
 - % setenv TAU_TRACE 1

Performance Evaluation Alternatives



Each alternative has:
- one metric/counter
- multiple counters

→
Volume of performance data

-PROFILEPARAM Configuration Option

- Idea: partition performance data for individual functions based on runtime parameters
- Enable by configuring with **-PROFILEPARAM**
- TAU call: TAU_PROFILE_PARAM1L (value, "name")
- Simple example:

```
void foo(long input) {  
    TAU_PROFILE("foo", "", TAU_DEFAULT);  
    TAU_PROFILE_PARAM1L(input, "input");  
    ... }  
}
```

Workload Characterization

- 5 seconds spent in function "foo" becomes
 - 2 seconds for "foo [<input> = <25>]"
 - 1 seconds for "foo [<input> = <5>]"
 - ...
- Currently used in MPI wrapper library
 - Allows for partitioning of time spent in MPI routines based on parameters (message size, message tag, destination node)
 - Can be extrapolated to infer specifics about the MPI subsystem and system as a whole

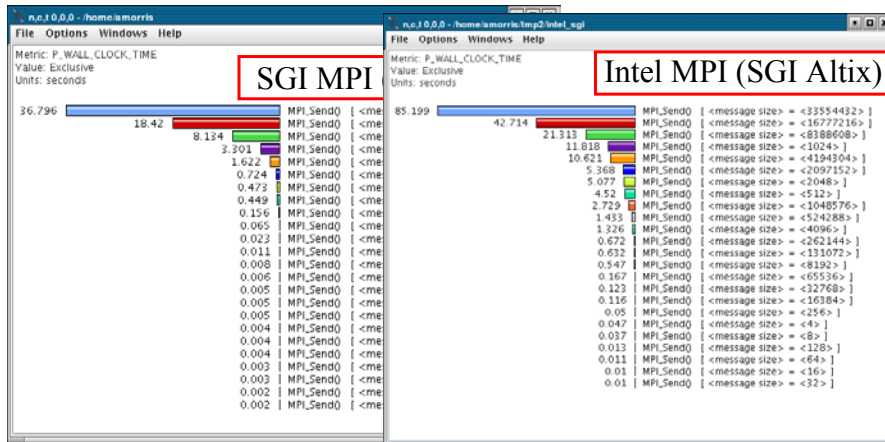
Workload Characterization

```
#include <stdio.h>
#include <mpi.h>
int buffer[8*1024*1024];

int main(int argc, char **argv) {
    int rank, size, i, j;
    MPI_Init(&argc, &argv);
    MPI_Comm_size( MPI_COMM_WORLD, &size );
    MPI_Comm_rank( MPI_COMM_WORLD, &rank );
    for (i=0;i<1000;i++)
        for (j=1;j<=8*1024*1024;j*=2) {
            if (rank == 0) {
                MPI_Send(buffer, j, MPI_INT, 1, 42, MPI_COMM_WORLD);
            } else {
                MPI_Status status;
                MPI_Recv(buffer, j, MPI_INT, 0, 42, MPI_COMM_WORLD, &status);
            }
        }
    MPI_Finalize();
}
```

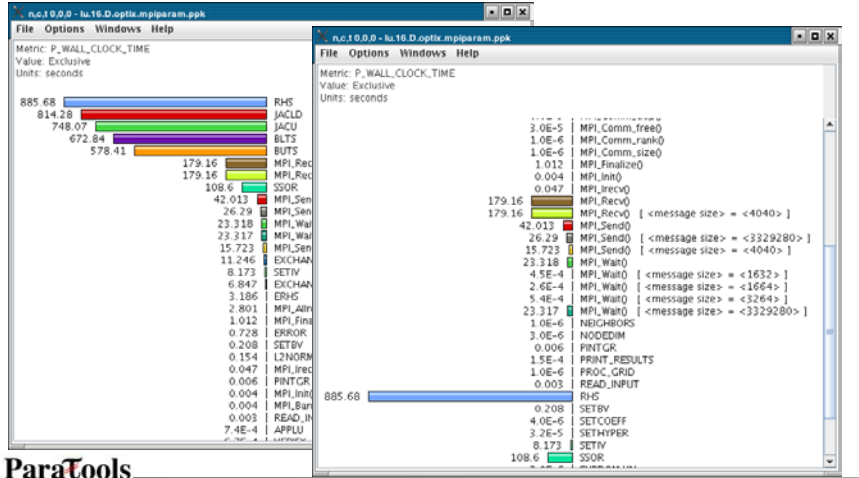
Workload Characterization

```
% icc mpi.c -lmpi
% mpirun -np 2 tauex a.out
```



Workload Characterization

- MPI Results (NAS Parallel Benchmark 3.1, LU class D on

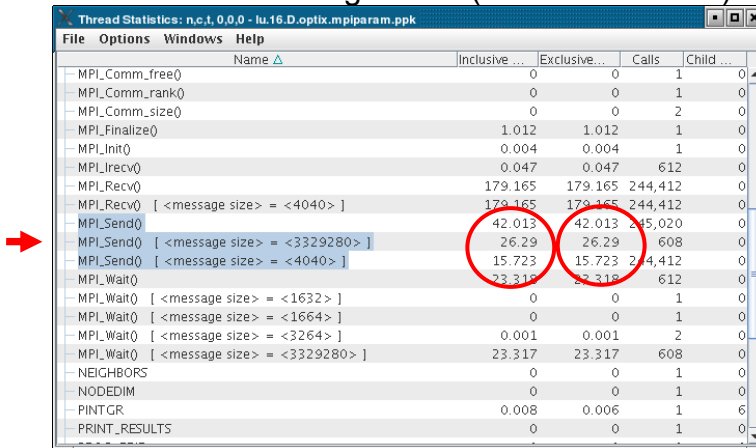


ParaTools

165

Workload Characterization

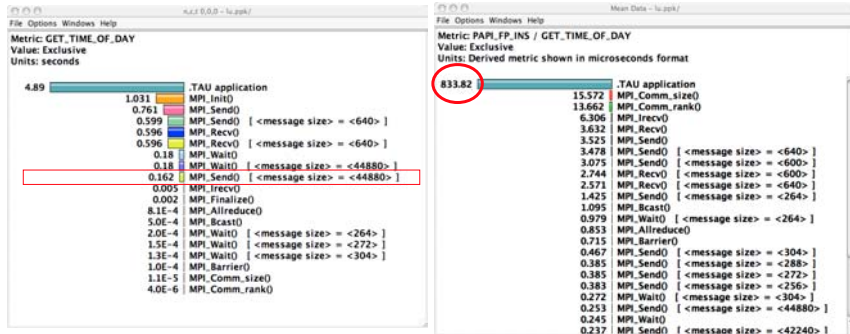
- Two different message sizes (~3.3MB and ~4K)



ParaTools

166

Job Tracking: ParaProf profile browser



LU spent 0.162 seconds sending messages of size 44880

It got 833.82 Mflops!

ParaTools

167

Memory Profiling in TAU

- Configuration option **-PROFILEMEMORY**
 - Records global heap memory utilization for each function
 - Takes one sample at beginning of each function and associates the sample with **function name**
- Configuration option **-PROFILEHEADROOM**
 - Records headroom (amount of free memory to grow) for each function
 - Takes one sample at beginning of each function and associates it with the **callstack** [TAU_CALLPATH_DEPTH env variable]
 - Useful for debugging memory usage on IBM BG/L.
- Independent of instrumentation/measurement options selected
- No need to insert macros/calls in the source code
- User defined atomic events appear in profiles/traces

ParaTools

168

Memory Profiling in TAU (Atomic events)

Sorted By: number of userEvents

NumSamples	Max	Min	Mean	Std. Dev	Name
252032	2022.7	1181.2	1534.3	410.04	MODULEHYDRO_ID::HYDRO_ID - Heap Memory (KB)
252032	2022.8	1181.7	1534.3	410.04	MODULEINTRFC::INTRFC - Heap Memory (KB)
104559	2023.2	331.13	1526.6	409.54	MODULEEOS3D::EOS3D - Heap Memory (KB)
63008	2022.7	1182	1534.3	410.01	MODULEUPDATE_SOLN::UPDATE_SOLN - Heap Memory (KB)
55545	2023.3	333.07	1514.2	408.31	DBASETREE::DBASENEIGHBORBLOCKLIST - Heap Memory (KB)
51374	2023	1179.4	1497.7	402.53	AMR_PROLONG_GEN_UNK_FUN - Heap Memory (KB)
42120	2022.7	1187.5	1533.5	409.83	ABUNDANCE_RESTRICT - Heap Memory (KB)
41958	2023	346.12	1514.9	408.39	AMR_RESTRICT_UNK_FUN - Heap Memory (KB)
31832	2022.8	1187.4	1534.1	409.91	AMR_RESTRICT_RED - Heap Memory (KB)
31504	2022.7	1181.8	1534.3	410.04	DIFFUSE - Heap Memory (KB)
26042	2023	1179.2	1501.9	403.61	AMR_PROLONG_UNK_FUN - Heap Memory (KB)

Flash2 code profile (-PROFILEMEMORY) on IBM BlueGene/L [MPI rank 0]

ParaTools

169

Memory Profiling in TAU

- Instrumentation based observation of global heap memory (not per function)
 - call TAU_TRACK_MEMORY()
 - call TAU_TRACK_MEMORY_HEADROOM()
 - Triggers one sample every 10 secs
 - call TAU_TRACK_MEMORY_HERE()
 - call TAU_TRACK_MEMORY_HEADROOM_HERE()
 - Triggers sample at a specific location in source code
 - call TAU_SET_INTERRUPT_INTERVAL(seconds)
 - To set inter-interrupt interval for sampling
 - call TAU_DISABLE_TRACKING_MEMORY()
 - call TAU_DISABLE_TRACKING_MEMORY_HEADROOM()
 - To turn off recording memory utilization
 - call TAU_ENABLE_TRACKING_MEMORY()
 - call TAU_ENABLE_TRACKING_MEMORY_HEADROOM()
 - To re-enable tracking memory utilization

ParaTools

170

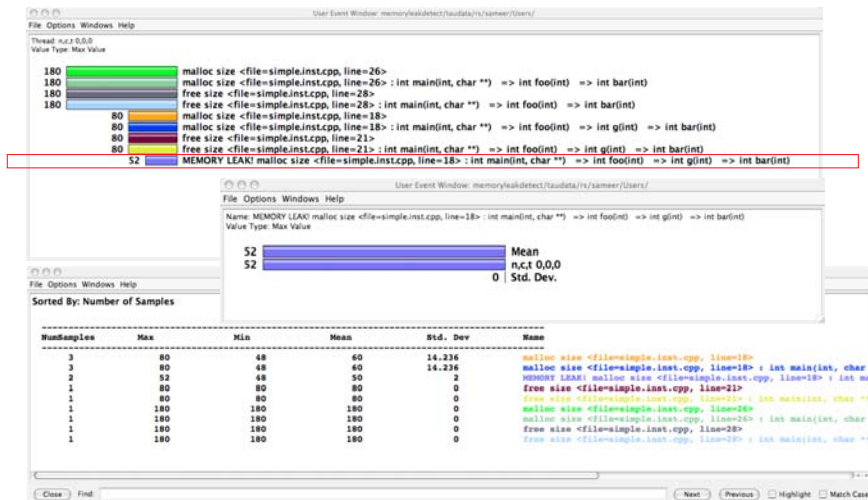
Detecting Memory Leaks in C/C++

- TAU wrapper library for malloc/realloc/free
- During instrumentation, specify
 - optDetectMemoryLeaks option to TAU_COMPILER
 - % setenv TAU_OPTIONS '-optVerbose -optDetectMemoryLeaks'
 - % setenv TAU_MAKEFILE <taudir>/<arch>/lib/Makefile.tau-mpi-pdt...
 - % tau_cxx.sh foo.cpp ...
- Tracks each memory allocation/de-allocation in parsed files
- Correlates each memory event with the executing callstack
- At the end of execution, TAU detects memory leaks
- TAU reports leaks based on allocations and the executing callstack
- Set **TAU_CALLPATH_DEPTH** environment variable to limit callpath data
 - default is 2
- Future work
 - Support for C++ new/delete planned
 - Support for Fortran 90/95 allocate/deallocate planned

ParaTools

171

Memory Leak Detection



ParaTools

172

Detecting Memory Leaks in Fortran

```
subroutine foo(x)
  integer:: x
  integer, allocatable :: A(:), B(:), C(:)

  print *, "inside foo"
  allocate(A(x), B(x), C(x))
  deallocate(A, C)
  print *, "exiting foo"

end subroutine foo

program main
  call foo(5)
end program main
```

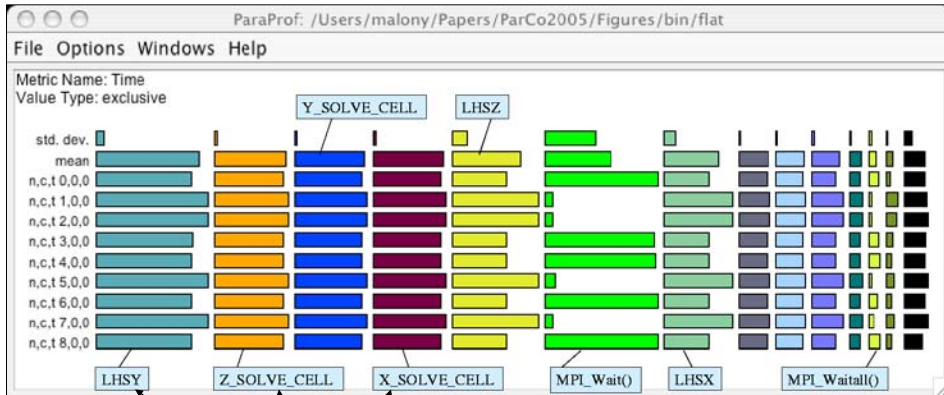
Detecting Memory Leaks in Fortran

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

NumSamples MaxValue MinValue MeanValue Std. Dev. Event Name

1	5	5	5	0	MEMORY LEAK! malloc size <file=simple.f, variable=B, line=6> : MAIN => FOO
1	5	5	5	0	free size <file=simple.f, variable=A, line=7>
1	5	5	5	0	free size <file=simple.f, variable=A, line=7> : MAIN => FOO
1	5	5	5	0	free size <file=simple.f, variable=C, line=7>
1	5	5	5	0	free size <file=simple.f, variable=C, line=7> : MAIN => FOO
1	5	5	5	0	malloc size <file=simple.f, variable=A, line=6>
1	5	5	5	0	malloc size <file=simple.f, variable=A, line=6> : MAIN => FOO
1	5	5	5	0	malloc size <file=simple.f, variable=B, line=6>
1	5	5	5	0	malloc size <file=simple.f, variable=B, line=6> : MAIN => FOO
1	5	5	5	0	malloc size <file=simple.f, variable=C, line=6>
1	5	5	5	0	malloc size <file=simple.f, variable=C, line=6> : MAIN => FOO

Phase Profiling (NAS BT, Flat Profile)



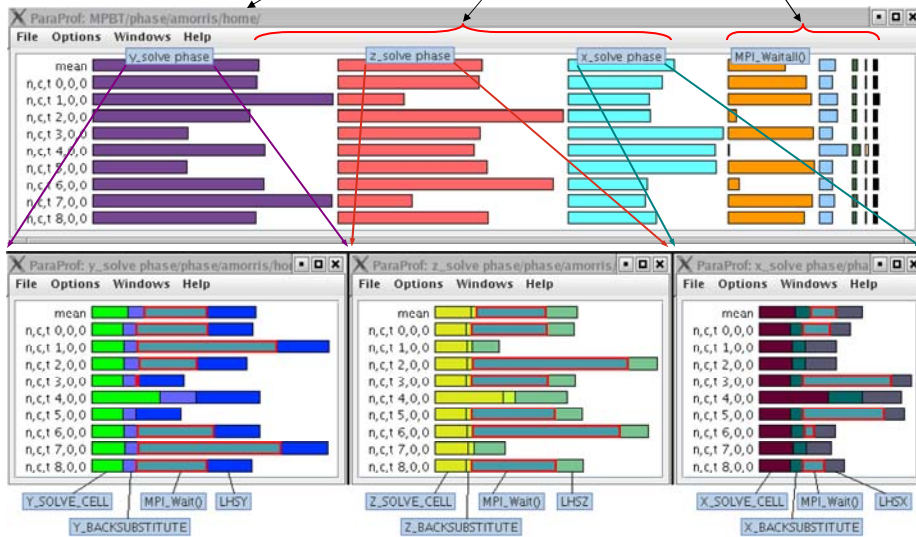
Application routine names reflect phase semantics

How is MPI_Wait() distributed relative to solver direction?

ParaTools

NAS BT – Phase Profile (Main and X, Y, Z)

Main phase shows nested phases and immediate events



TAU Timers and Phases

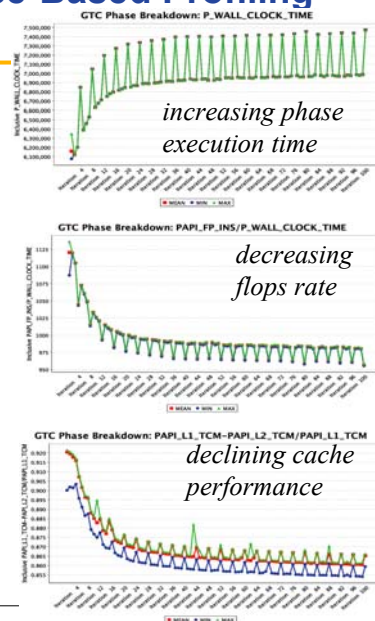
- **Static timer**
 - Shows time spent in all invocations of a routine (foo)
 - E.g., “foo()” 100 secs, 100 calls
- **Dynamic timer**
 - Shows time spent in each invocation of a routine
 - E.g., “foo() 3” 4.5 secs, “foo 10” 2 secs (invocations 3 and 10 respectively)
- **Static phase**
 - Shows time spent in all routines called (directly/indirectly) by a given routine (foo)
 - E.g., “foo() => MPI_Send()” 100 secs, 10 calls shows that a total of 100 secs were spent in MPI_Send() when it was called by foo.
- **Dynamic phase**
 - Shows time spent in all routines called by a given invocation of a routine.
 - E.g., “foo() 4 => MPI_Send()” 12 secs, shows that 12 secs were spent in MPI_Send when it was called by the 4th invocation of foo.

ParaTools

177

Performance Dynamics: Phase-Based Profiling

- Profile phases capture performance with respect to application-defined ‘phases’ of execution
 - Separate full profile produce for each phase
- GTC particle-in-cell simulation of fusion turbulence
- Phases assigned to iterations
- Data change affects cache



ParaTools

178

TAU's MPI Wrapper Interposition Library

- Uses standard MPI Profiling Interface
 - Provides name shifted interface
 - MPI_Send = PMPI_Send
 - Weak bindings
- Interpose TAU's MPI wrapper library between MPI and TAU
 - `-Impi` replaced by `-ITauMpi -Ipmi -Impi`
- No change to the source code! Just **re-link** the application to generate performance data
 - `setenv TAU_MAKEFILE`
`<dir>/<arch>/lib/Makefile.tau-mpi-[options]`
 - Use `tau_cxx.sh`, `tau_f90.sh` and `tau_cc.sh` as compilers

ParaTools

179

Using TAU

- Install TAU
 - Configuration
 - Measurement library creation
- Instrument application
 - Manual or automatic source instrumentation
 - Instrumented library (e.g., MPI – wrapper interposition library)
 - Binary instrumentation
- ➡ • **Create performance experiments**
 - **Integrate with application build environment**
 - **Set experiment variables**
- Execute application
- Analyze performance

ParaTools

180

Integration with Application Build Environment

- Try to minimize impact on user's application build procedures
- Handle process of parsing, instrumentation, compilation, linking
- Dealing with Makefiles
 - Minimal change to application Makefile
 - Avoid changing compilation rules in application Makefile
 - No explicit inclusion of rules for process stages
- Some applications do not use Makefiles
 - Facilitate integration in whatever procedures used
- Two techniques:
 - TAU shell scripts (tau_<compiler>.sh)
 - Invokes all PDT parser, TAU instrumenter, and compiler
 - TAU_COMPILER

ParaTools

181

Using Program Database Toolkit (PDT)

- 1. Parse the Program to create foo.pdb:**

```
% cxxparse foo.cpp -I/usr/local/mydir -DMYFLAGS ...  
or  
% cparse foo.c -I/usr/local/mydir -DMYFLAGS ...  
or  
% f95parse foo.f90 -I/usr/local/mydir ...  
% f95parse *.f -omerged.pdb -I/usr/local/mydir -R free
```
- 2. Instrument the program:**

```
% tau_instrumentor foo.pdb foo.f90 -o foo.inst.f90  
-f select.tau
```
- 3. Compile the instrumented program:**

```
% ifort foo.inst.f90 -c -I/usr/local/mpi/include -o foo.o
```

182

Tau_[cxx,cc,f90].sh – Improves Integration in Makefiles

```
# set TAU_MAKEFILE and TAU_OPTIONS env vars
CC = tau_cc.sh
F90 = tau_f90.sh
CFLAGS =
LIBS = -lm
OBJS = f1.o f2.o f3.o ... fn.o

app: $(OBJS)
    $(F90) $(LDFLAGS) $(OBJS) -o $@ $(LIBS)
.c.o:
    $(CC) $(CFLAGS) -c $<
.f90.o:
    $(F90) $(FFLAGS) -c $<
```

ParaTools

183

Automatic Instrumentation

- We now provide compiler wrapper scripts
 - Simply replace `mpxf90` with `tau_f90.sh`
 - Automatically instruments Fortran source code, links with TAU MPI Wrapper libraries.
- Use `tau_cc.sh` and `tau_cxx.sh` for C/C++

Before

```
CXX = mpCC
F90 = mpxf90_r
CFLAGS =
LIBS = -lm
OBJS = f1.o f2.o f3.o ... fn.o

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

After

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

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

ParaTools

184

TAU_COMPILER Commandline Options

- See `<taudir>/<arch>/bin/tau_compiler.sh -help`
- Compilation:

```
% mpxlf90 -c foo.f90
```

Changes to

```
% f95parse foo.f90 $(OPT1)
% tau_instrumentor foo.pdb foo.f90 -o foo.inst.f90 $(OPT2)
% ftn -c foo.f90 $(OPT3)
```
- Linking:

```
% ftn foo.o bar.o -o app
```

Changes to

```
% ftn foo.o bar.o -o app $(OPT4)
```
- Where options `OPT[1-4]` default values may be overridden by the user:
`F90 = tau_f90.sh`

TAU_COMPILER Options

- Optional parameters for `$(TAU_COMPILER)`: [`tau_compiler.sh -help`]
 - optVerbose Turn on verbose debugging messages
 - optComplnst Use compiler based instrumentation
 - optDetectMemoryLeaks Turn on debugging memory allocations/de-allocations to track leaks
 - optKeepFiles Does not remove intermediate .pdb and .inst.* files
 - optPreProcess Preprocess Fortran sources before instrumentation
 - optTauSelectFile="" Specify selective instrumentation file for tau_instrumentor
 - optLinking="" Options passed to the linker. Typically `$(TAU_MPI_FLIBS) $(TAU_LIBS) $(TAU_CXXLIBS)`
 - optCompile="" Options passed to the compiler. Typically `$(TAU_MPI_INCLUDE) $(TAU_INCLUDE) $(TAU_DEFS)`
 - optPdtF95Opts="" Add options for Fortran parser in PDT (f95parse/gfparse)
 - optPdtF95Reset="" Reset options for Fortran parser in PDT (f95parse/gfparse)
 - optPdtCOpts="" Options for C parser in PDT (cparse). Typically `$(TAU_MPI_INCLUDE) $(TAU_INCLUDE) $(TAU_DEFS)`
 - optPdtCxxOpts="" Options for C++ parser in PDT (cxxparse). Typically `$(TAU_MPI_INCLUDE) $(TAU_INCLUDE) $(TAU_DEFS)`
 - ...

Compiling Fortran Codes with TAU

- If your Fortran code uses free format in .f files (fixed is default for .f), you may use:
% setenv TAU_OPTIONS '-optPdtF95Opts="-R free" -optVerbose'
- To use the compiler based instrumentation instead of PDT (source-based):
% setenv TAU_OPTIONS '-optComplInst -optVerbose'
- If your Fortran code uses C preprocessor directives (#include, #ifdef, #endif):
% setenv TAU_OPTIONS '-optPreProcess -optVerbose -optDetectMemoryLeaks'
- To use an instrumentation specification file:
% setenv TAU_OPTIONS '-optTauSelectFile=mycmd.tau -optVerbose -optPreProcess'
% cat mycmd.tau
BEGIN_INSTRUMENT_SECTION
memory file="foo.f90" routine="#"
instruments all allocate/deallocate statements in all routines in foo.f90
loops file="*" routine="#"
io file="abc.f90" routine="FOO"
END_INSTRUMENT_SECTION

ParaTools

187

Overriding Default Options:TAU_COMPILER

```
% cat Makefile
F90 = tau_f90.sh
OBJS = f1.o f2.o f3.o ...
LIBS = -Lappdir -lapplib1 -lapplib2 ...

app: $(OBJS)
    $(F90) $(OBJS) -o app $(LIBS)
.f90.o:
    $(F90) -c $<
% setenv TAU_OPTIONS '-optVerbose
    -optTauSelectFile=select.tau -optKeepFiles'
% setenv TAU_MAKEFILE <taudir>/x86_64/lib/Makefile.tau-mpi-pdt
```

ParaTools

188

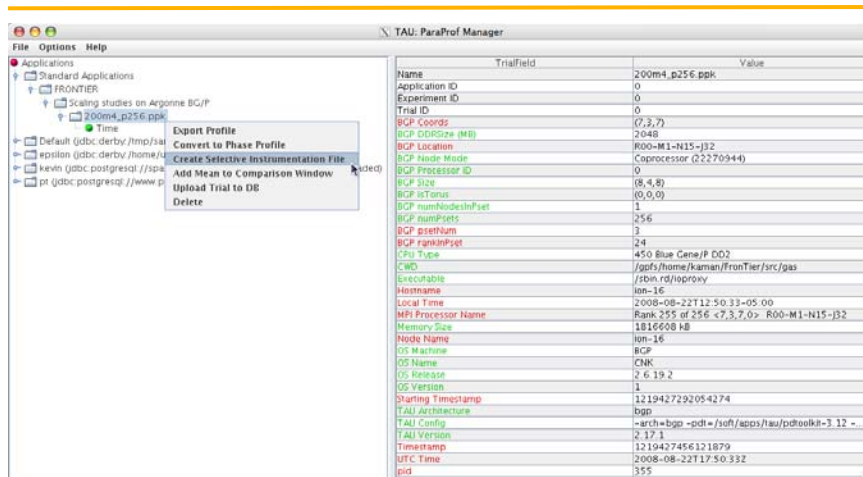
Optimization of Program Instrumentation

- Need to eliminate instrumentation in frequently executing lightweight routines
- Throttling of events at runtime (default in tau-2.17.2+):
% `setenv TAU_THROTTLE 1`
Turns off instrumentation in routines that execute over 100000 times (TAU_THROTTLE_NUMCALLS) and take less than 10 microseconds of inclusive time per call (TAU_THROTTLE_PERCALL). Use TAU_THROTTLE=0 to disable.
- Selective instrumentation file to filter events
% `tau_instrumentor [options] -f <file> OR`
% `setenv TAU_OPTIONS '-optTauSelectFile=tau.txt'`
- Compensation of local instrumentation overhead
% `configure --COMPENSATE`
`or`
% `setenv TAU_COMPENSATE 1 (in tau-2.18.2+)`

ParaTools

189

ParaProf: Creating Selective Instrumentation File



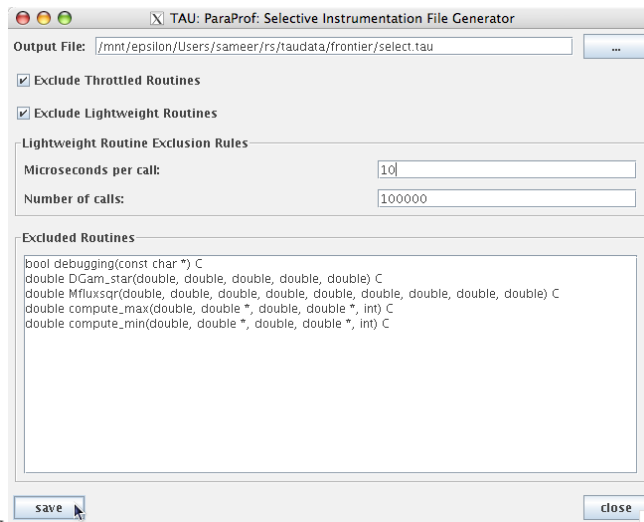
The screenshot shows the TAU ParaProf Manager interface. On the left, a file tree displays the directory structure: Applications > Standard Applications > FRONTIER > Scaling studies on Argonne BC/P > 200m4_p256.ppk. A context menu is open over the '200m4_p256.ppk' file, with options: Export Profile, Convert to Phase Profile, Create Selective Instrumentation File (highlighted), Add Mean to Comparison Window, Upload Trial to DB, and Delete. On the right, a table lists trial fields and their values for the selected trial.

Name	TrialField	Value
Application ID	0	200m4_p256.ppk
Experiment ID	0	
Trial ID	0	
BCP_Coords	(7,3,7)	
BCP_DDRSize (MB)	2048	
BCP_Location	R00-M1-N15-J32	
BCP_Node_Headr	Coprocessor (22270944)	
BCP_Processor ID	0	
BCP_Size	(8,4,8)	
BCP_nTorus	(0,0,0)	
BCP_numNodesInPset	1	
BCP_numPsets	256	
BCP_psetNum	3	
BCP_rankInPset	24	
Cmd_Line	450 Blue Gene/P D02	
CWD	/gafs/home/kaman/FrontTier/src/gas	
Executable	/sbin,rd/oproxy	
Hostname	ion-16	
Local Time	2008-08-22T12:50:33-05:00	
MPI_Processor Name	Rank 255 of 256 <7,3,7,0>: R00-M1-N15-J32	
Memory Size	1816608 kB	
Node Name	ion-16	
OS Machine	BCP	
OS Name	CRK	
OS Release	2.6.19.2	
OS Version	1	
Starting Timestamp	1219427292054274	
TAU Architecture	bgp	
TAU Config	-arch=bgp -pdt=/soft/apps/tau/pdtoolkit-3.12 -	
TAU Version	2.17.1	
Timestamp	1219427456121879	
UTC Time	2008-08-22T17:50:33Z	
pid	355	

ParaTools

190

Choosing Rules for Excluding Routines



ParaTools

191

Selective Instrumentation File

- Specify a list of routines to exclude or include (case sensitive)
- # is a wildcard in a routine name. It cannot appear in the first column.

```
BEGIN_EXCLUDE_LIST
Foo
Bar
D#EMM
END_EXCLUDE_LIST
```

- Specify a list of routines to include for instrumentation

```
BEGIN_INCLUDE_LIST
int main(int, char **)
F1
F3
END_INCLUDE_LIST
```

- Specify either an include list or an exclude list!

ParaTools

192

Selective Instrumentation File

- Optionally specify a list of files to exclude or include (case sensitive)

- * and ? may be used as wildcard characters in a file name

```
BEGIN_FILE_EXCLUDE_LIST
f*.f90
Foo?.cpp
END_FILE_EXCLUDE_LIST
```

- Specify a list of routines to include for instrumentation

```
BEGIN_FILE_INCLUDE_LIST
main.cpp
foo.f90
END_FILE_INCLUDE_LIST
```

Selective Instrumentation File

- User instrumentation commands are placed in INSTRUMENT section
- ? and * used as wildcard characters for file name, # for routine name
- \ as escape character for quotes
- Routine entry/exit, arbitrary code insertion
- Outer-loop level instrumentation

```
BEGIN_INSTRUMENT_SECTION
loops file="foo.f90" routine="matrix#"
memory file="foo.f90" routine="#"
io routine="matrix#"
[static/dynamic] phase routine="MULTIPLY"
dynamic [phase/timer] name="foo" file="foo.cpp" line=22 to line=35
file="foo.f90" line = 123 code = " print *, \" Inside foo\""
exit routine = "int foo()" code = "cout <<\"exiting foo\"<<endl;"
END_INSTRUMENT_SECTION
```

Instrumentation Specification

```
% tau_instrumentor
Usage : tau_instrumentor <pdbfile> <sourcefile> [-o <outputfile>] [-noinline]
[-g groupname] [-i headerfile] [-c|-c++|-fortran] [-f <instr_req_file> ]
For selective instrumentation, use -f option
% tau_instrumentor foo.pdb foo.cpp -o foo.inst.cpp -f selective.dat
% cat selective.dat
# Selective instrumentation: Specify an exclude/include list of routines/files.
BEGIN_EXCLUDE_LIST
void quicksort(int *, int, int)
void sort_5elements(int *)
void interchange(int *, int *)
END_EXCLUDE_LIST

BEGIN_FILE_INCLUDE_LIST
Main.cpp
Foo?.c
*.C
END_FILE_INCLUDE_LIST
# Instruments routines in Main.cpp, Foo?.c and *.C files only
# Use BEGIN_[FILE]_INCLUDE_LIST with END_[FILE]_INCLUDE_LIST
```

195

Instrumentation of OpenMP Constructs

- **OpenMP Pragma And Region Instrumentor** [UTK, FZJ]
- Source-to-Source translator to insert **POMP** calls around OpenMP constructs and API functions
- **Done:** Supports
 - **Fortran77** and **Fortran90**, OpenMP 2.0
 - **C** and **C++**, OpenMP 1.0
 - **POMP** Extensions
 - EPILOG and TAU POMP implementations
 - Preserves source code information (**#line line file**)
- **tau_ompcheck**
 - Balances OpenMP constructs (DO/END DO) and detects errors
 - Invoked by tau_compiler.sh prior to invoking Opari
- KOJAK Project website <http://icl.cs.utk.edu/kojak>



OpenMP API Instrumentation

- Transform
 - `omp_#_lock()` → `pomp_#_lock()`
 - `omp_#_nest_lock()` → `pomp_#_nest_lock()`

- [# = `init` | `destroy` | `set` | `unset` | `test`]

- POMP version
 - Calls omp version internally
 - Can do extra stuff before and after call

Example: !\$OMP PARALLEL DO Instrumentation

```
call pomp_parallel_fork(d)
!$OMP PARALLEL other-clauses...
  call pomp_parallel_begin(d)
  call pomp_do_enter(d)
  !$OMP DO schedule-clauses, ordered-clauses,
           lastprivate-clauses
    do loop
  !$OMP END DO NOWAIT
  call pomp_barrier_enter(d)
  !$OMP BARRIER
  call pomp_barrier_exit(d)
  call pomp_do_exit(d)
  call pomp_parallel_end(d)
!$OMP END PARALLEL DO
call pomp_parallel_join(d)
```

Opari Instrumentation: Example

```
pomp_for_enter(&omp_rd_2);
#line 252 "stommel.c"
#pragma omp for schedule(static) reduction(+: diff) private(j)
    firstprivate (a1,a2,a3,a4,a5) nowait
for( i=i1;i<=i2;i++) {
    for(j=j1;j<=j2;j++){
        new_psi[i][j]=a1*psi[i+1][j] + a2*psi[i-1][j] + a3*psi[i][j+1]
            + a4*psi[i][j-1] - a5*the_for[i][j];
        diff=diff+fabs(new_psi[i][j]-psi[i][j]);
    }
}
pomp_barrier_enter(&omp_rd_2);
#pragma omp barrier
pomp_barrier_exit(&omp_rd_2);
pomp_for_exit(&omp_rd_2);
```

ParaTools

199

Using Opari with TAU

Configure TAU with Opari (used here with MPI and PDT)

```
% configure -opari -arch=x86_64 -mpi -
pdt=/usr/contrib/TAU/pdtoolkit-3.14.1
% make clean; make install
% setenv TAU_MAKEFILE /tau/<arch>/lib/Makefile.tau-...opari-...
% tau_cxx.sh -c foo.cpp
% tau_cxx.sh -c bar.f90
% tau_cxx.sh *.o -o app
```

ParaTools

200

Dynamic Instrumentation

- TAU uses DyninstAPI for runtime code patching
- Developed by U. Wisconsin and U. Maryland
- <http://www.dyninst.org>
- ***tau_run*** (mutator) loads measurement library
- Instruments mutatee
- MPI issues:
 - one mutator per executable image [TAU, DynaProf]
 - one mutator for several executables [Paradyn, DPCL]

Using DyninstAPI with TAU

```
Step I: Install DyninstAPI[Download from http://www.dyninst.org]  
% cd dyninstAPI-6/core; make  
Set DyninstAPI environment variables (including LD_LIBRARY_PATH)  
Step II: Configure TAU with Dyninst  
% configure -dyninst=/usr/local/dyninstAPI-6  
% make clean; make install  
Builds <taudir>/<arch>/bin/tau_run  
% tau_run [<-o outfile>] [-Xrun<libname>] [-f <select_inst_file>] [-v] <infile>  
% tau_run -o a.inst.out a.out  
Rewrites a.out  
% tau_run klargest  
Instruments klargest with TAU calls and executes it  
% tau_run -XrunTAUsh-papi a.out  
Loads libTAUsh-papi.so instead of libTAU.so for measurements
```

Virtual Machine Performance Instrumentation

- Integrate performance system with VM
 - Captures robust performance data (e.g., thread events)
 - Maintain features of environment
 - portability, concurrency, extensibility, interoperation
 - Allow use in optimization methods
- JVM Profiling Interface (JVMPi)
 - Generation of JVM events and hooks into JVM
 - Profiler agent (TAU) loaded as shared object
 - registers events of interest and address of callback routine
 - Access to information on dynamically loaded classes
 - No need to modify Java source, bytecode, or JVM

Using TAU with Java Applications

Step I: Sun JDK 1.4+ [download from www.javasoft.com]

Step II: Configure TAU with JDK (v 1.2 or better)

```
% configure -jdk=/usr/java2 -TRACE -PROFILE
```

```
% make clean; make install
```

Builds <taudir>/<arch>/lib/libTAU.so

For Java (without instrumentation):

```
% java application
```

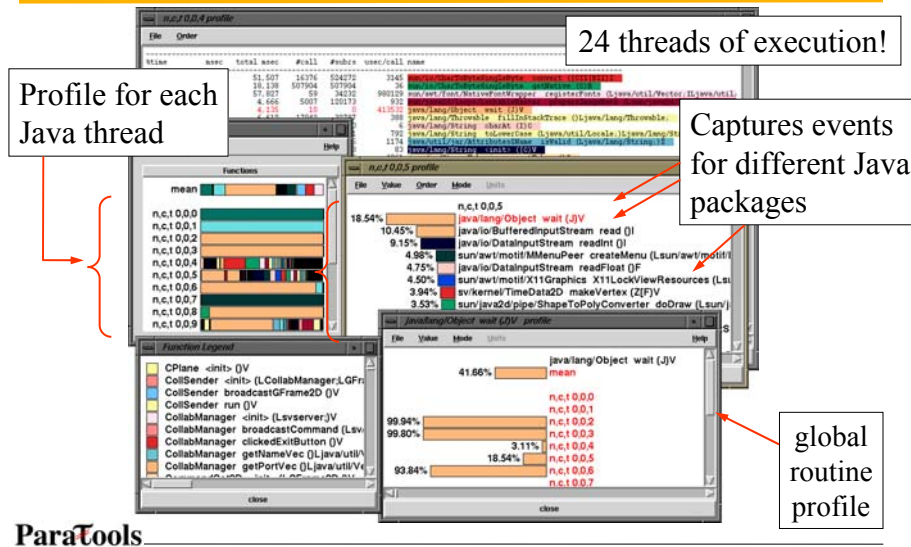
With instrumentation:

```
% java -XrunTAU application
```

```
% java -XrunTAU:exclude=sun/io,java application
```

Excludes sun/io/* and java/* classes

TAU Profiling of Java Application (SciVis)



205

Using TAU with Python Applications

Step I: Configure TAU with Python

```
% configure -pythoninc=/usr/include/python2.5/include
% make clean; make install
```

Builds <taudir>/<arch>/lib/<bindings>/pytau.py and tau.py packages for manual and automatic instrumentation respectively

```
% setenv PYTHONPATH $PYTHONPATH:<taudir>/<arch>/lib/[[dir]]
```

Python Automatic Instrumentation Example

```
#!/usr/bin/env/python
```

```
import tau
```

```
from time import sleep
```

```
def f2():
```

```
    print " In f2: Sleeping for 2 seconds "
```

```
    sleep(2)
```

```
def f1():
```

```
    print " In f1: Sleeping for 3 seconds "
```

```
    sleep(3)
```

```
def OurMain():
```

```
    f1()
```

```
tau.run('OurMain()')
```

Running:

```
% setenv PYTHONPATH  
<tau>/<arch>/lib/bindings-  
python
```

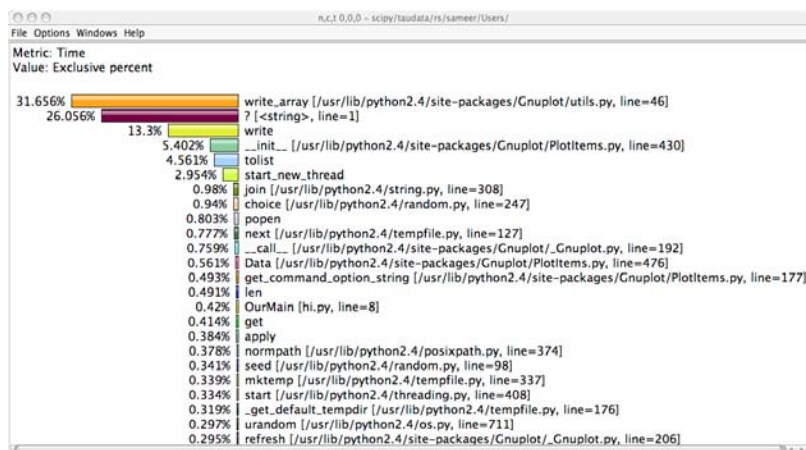
```
% ./auto.py
```

```
Instruments OurMain, f1, f2,  
print...
```

ParaTools

207

Python Instrumentation: SciPy



ParaTools

208

Performance Analysis

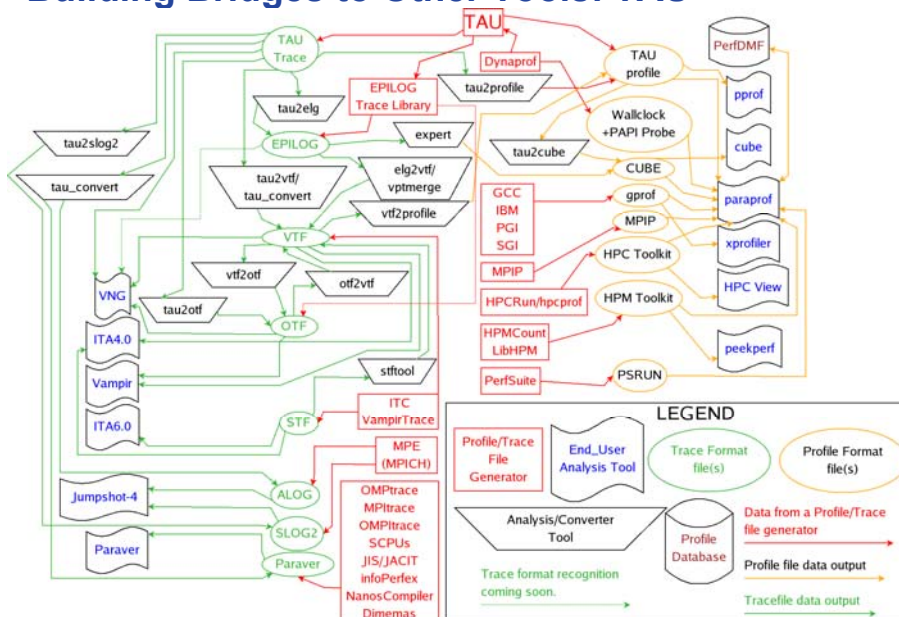
- paraprof profile browser (GUI)
- pprof (text based profile browser)
- TAU traces can be exported to many different tools
 - Vampir/VNG [T.U. Dresden] (formerly Intel (R) Trace Analyzer)
 - EXPERT [FZJ]
 - Jumpshot (bundled with TAU) [Argonne National Lab] ...

ParaTools

209



Building Bridges to Other Tools: TAU



TAU Performance System Interfaces

- PDT [U. Oregon, LANL, FZJ] for instrumentation of C++, C99, F95 source code
- PAPI [UTK] for accessing hardware performance counters data
- DyninstAPI [U. Maryland, U. Wisconsin] for runtime instrumentation
- KOJAK [FZJ, UTK]
 - Epilog trace generation library
 - CUBE callgraph visualizer
 - Opari OpenMP directive rewriting tool
- Vampir/VNG Trace Analyzer [TU Dresden]
- VTF3/OTF trace generation library [TU Dresden] (available from TAU website)
- Paraver trace visualizer [CEPBA]
- Jumpshot-4 trace visualizer [MPICH, ANL]
- JVMPI from JDK for Java program instrumentation [Sun]
- Paraprof profile browser/PerfDMF database supports:
 - TAU format
 - Gprof [GNU]
 - HPM Toolkit [IBM]
 - MpiP [ORNL, LLNL]
 - Dynaprof [UTK]
 - PSRun [NCSA]

ParaTools

211

ParaProf – Manager Window



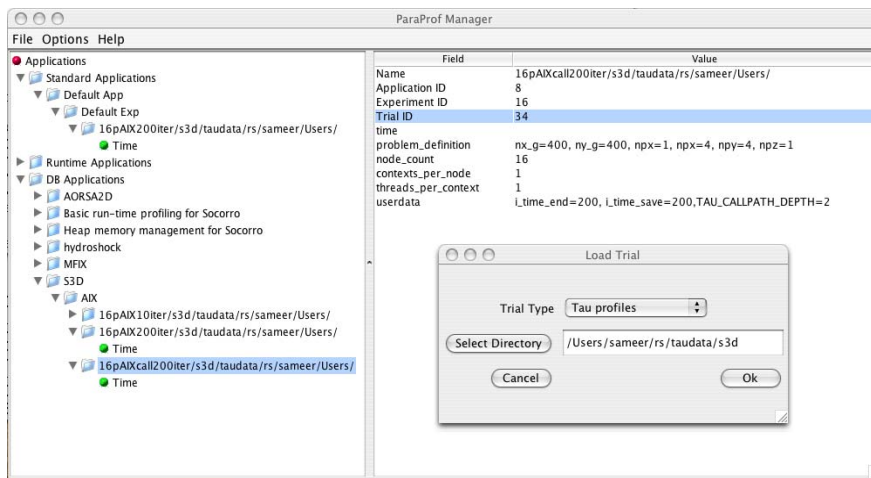
The screenshot shows the ParaProf Manager interface. On the left is a tree view of applications, with a label 'performance database' pointing to it. On the right is a table of metadata with columns 'Name', 'Field', '64 CPU', and 'Value'. A label 'metadata' points to this table. A 'Load Trial' dialog box is open in the foreground, showing a list of trial types including 'Tau profiles', 'Tau gprof.dat', 'Dynaprof', 'MpiP', 'jvmpiToolKit', 'Gprof', 'PSRun', 'ParaProf Packed Profile', 'Cube', and 'HPCToolKit'.

Name	Field	64 CPU	Value
Application ID		4	
Experiment ID		18	
Trial ID		85	
DATE			
COLLECTOR		64	
NOISE_COUNT		1	
CONTEXTS_PER_NODE		1	
THREADS_PER_CONTEXT		1	

ParaTools

212

Performance Database: Storage of MetaData



ParaTools

213

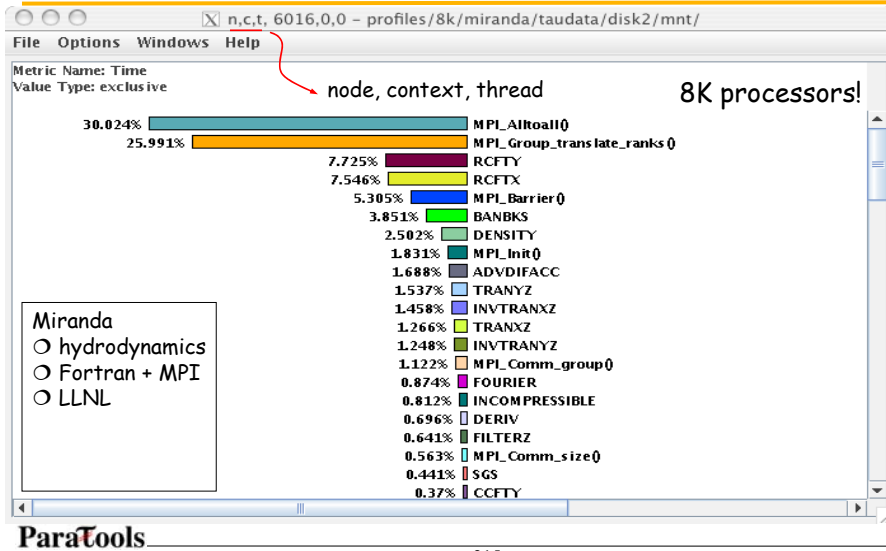
ParaProf Main Window (Lammps)



ParaTools

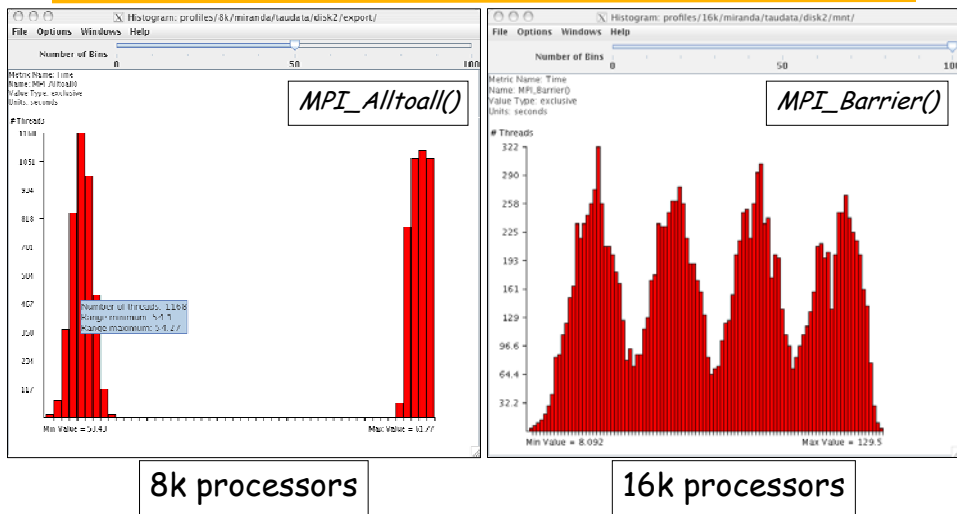
214

ParaProf – Flat Profile (Miranda)



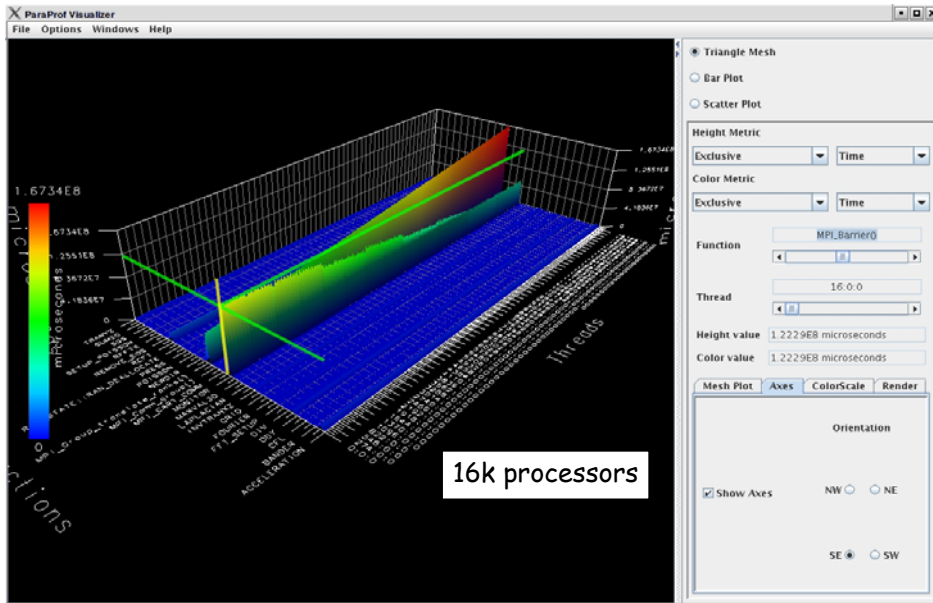
215

ParaProf – Histogram View (Miranda)



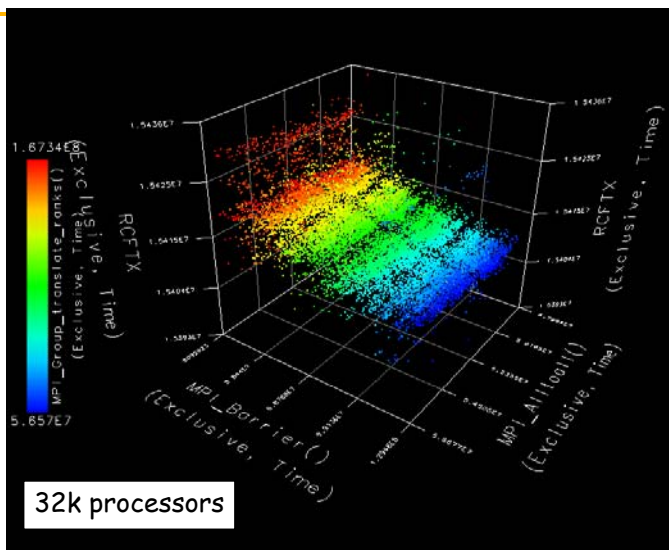
216

ParaProf – 3D Full Profile (Miranda)

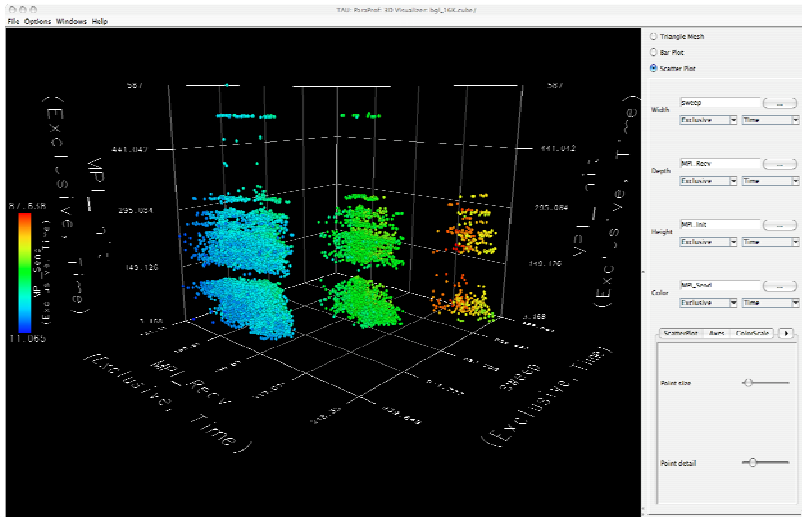


ParaProf – 3D Scatterplot (Miranda)

- Each point is a “thread” of execution
- A total of four metrics shown in relation
- ParaVis 3D profile visualization library
 - JOGL

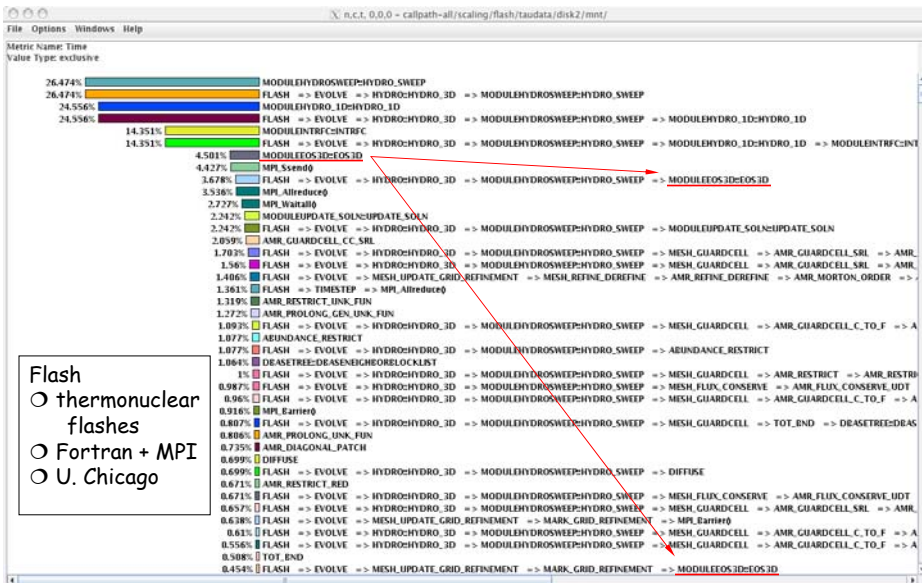


ParaProf – 3D Scatterplot (SWEEP3D CUBE)

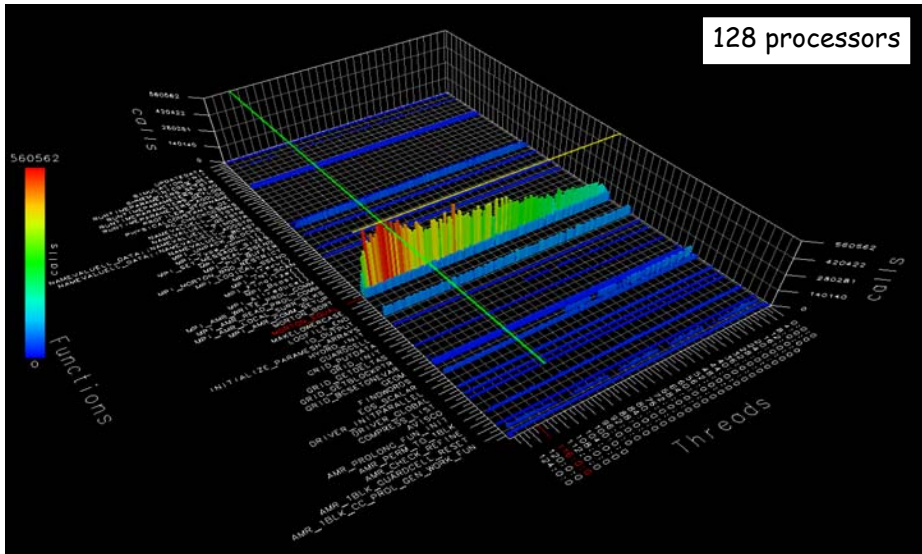


ParaTools

ParaProf – Callpath Profile (Flash)

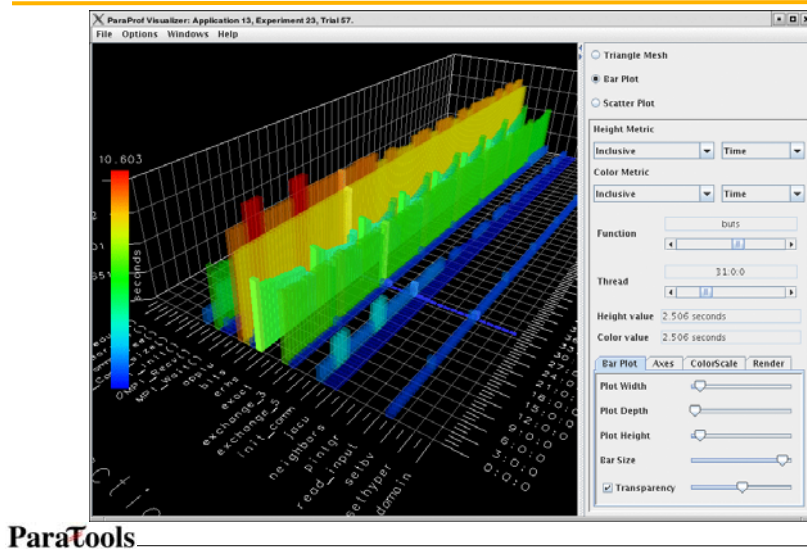


ParaProf – 3D Full Profile Bar Plot (Flash)



221

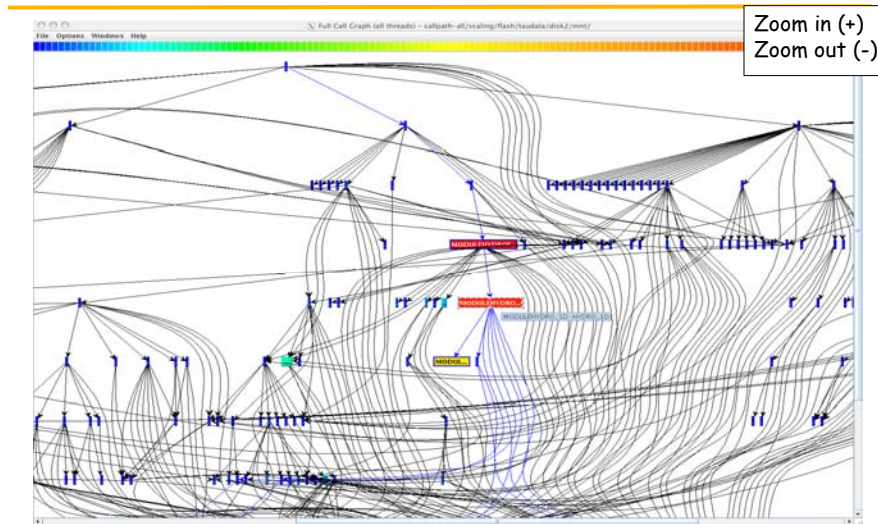
ParaProf Bar Plot (Zoom in/out +/-)



ParaTools

222

ParaProf – Callgraph Zoomed (Flash)



ParaTools

223

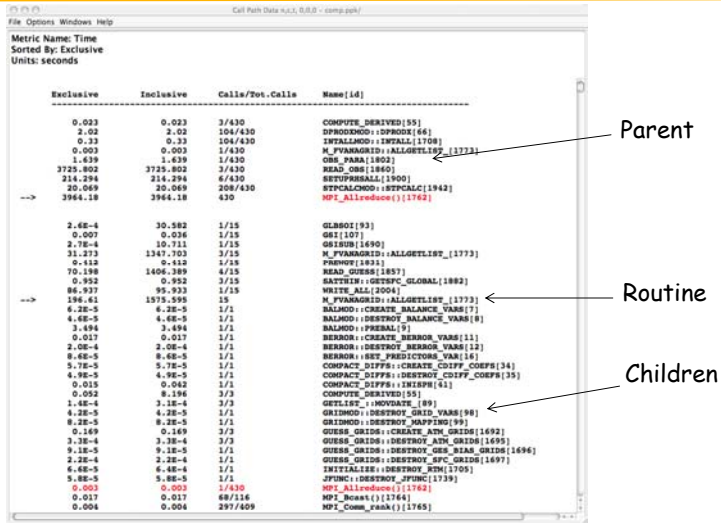
ParaProf - Thread Statistics Table (GSI)

Name	Inclusive Time	Exclusive Time	Calls	Child Calls
GSI	5,223.564	0.098	1	30
SPECMOD::INIT_SPEC_VARS	0.26	0.26	1	0
MPL_Init()	0.056	0.054	1	1
CSISUB	5,223.094	0.012	1	13
RADINFO::RADINFO_READ	0.103	0.101	1	1,196
PCPINFO::PCPINFO_READ	0.042	0.042	1	0
GLBSOI	5,212.171	0.024	1	12
MPL_Finalize()	1.004	1.004	1	0
OBS_PARA	3.635	0.181	1	56
JFUNC::CREATE_JFUNC	0.142	0.142	1	0
GUESS_GRIDS::CREATE_GES_BIAS_GRIDS	0.059	0.059	1	0
READ_GUESS	1,406.412	0.023	1	8
READ_OBS	3,770.188	0.016	1	6
MPL_Allreduce()	3,725.802	3,725.802	3	0
READ_BUFRTOS	44.369	0.254	1	871,535
SATTHIN::MAKEGVALS	0	0	1	0
W3FSZ1	0	0	1	1
BINARY_FILE_UTILITY::OPEN_BINARY_FILE	0.025	0.012	1	3
INITIALIZE::INITIALIZE_RTM	0.099	0.001	1	2
GUESS_GRIDS::CREATE_SFC_GRIDS	0	0	1	0
M_FVANAGRID::ALLGETLIST_	30.582	0	1	10
ERROR_HANDLER::DISPLAY_MESSAGE	0	0	1	0
JFUNC::SET_POINTER	0	0	1	0
OZINFO::OZINFO_READ	0.016	0.016	1	0
DETER_SUBDOMAIN	0.008	0.008	1	0
GRIDMOD::CREATE_MAPPING	0.005	0.005	1	0
INIT_COMMVARS	0.004	0.004	1	0
M_FVANAGRID::ALLGETLIST_	10.711	0	1	1
GRIDMOD::CREATE_GRID_VARS	0	0	1	0

ParaTools

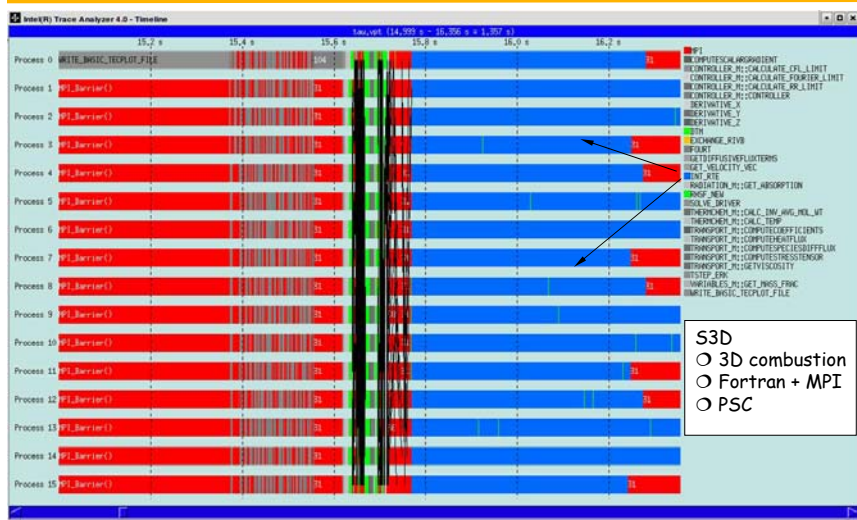
224

ParaProf - Callpath Thread Relations Window



ParaTools

Vampir – Trace Analysis (TAU-to-VTF3) (S3D)



ParaTools

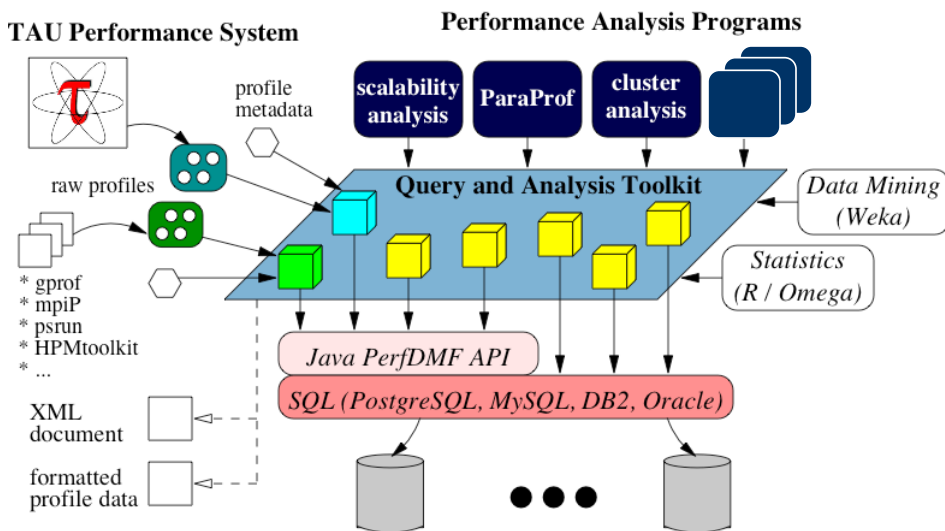
Vampir – Trace Zoomed (S3D)



ParaTools

227

PerfDMF: Performance Data Mgmt. Framework



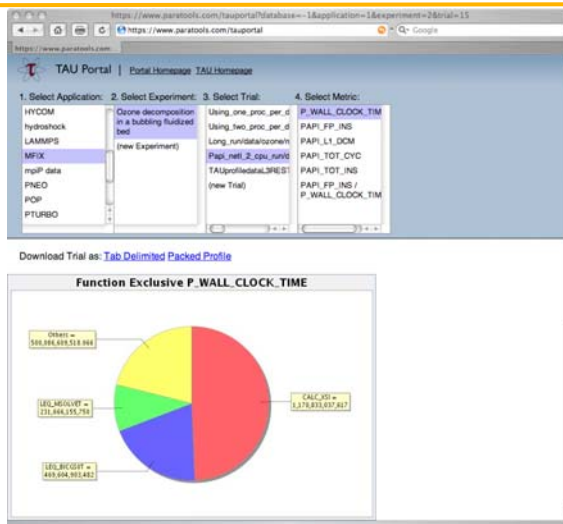
ParaTools

228

TAU Portal - www.paratools.com/tauportal



TAU Portal

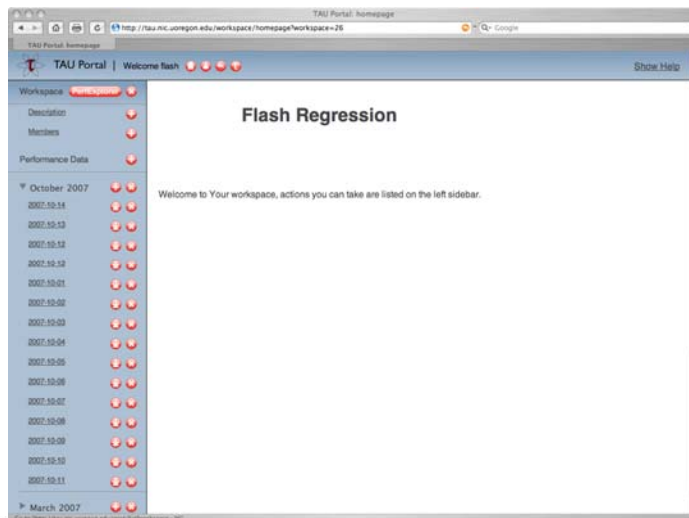


- Web-based access to TAU
- Support collaborative performance study
 - Secure performance data sharing
 - Does not require TAU installation
 - Launch TAU performance tools with Java WebStart
 - ParaProf, PerfExplorer
- FLASH regression testing
 - Nightly regression testcases
 - Uploaded to the database automatically
 - Interactive review of performance through TAU portal
 - Multi-experiment analysis

ParaTools

231

Portal: Nightly Performance Regression Testing



ParaTools

232

TAU Portal: Launch ParaProf/PerfExplorer

The screenshot shows the TAU Portal web interface. On the left, there is a sidebar with a tree view of performance data, including a 'Performance Data' section with a tree view of dates from 2007 to 2008. The main content area displays a detailed view of a specific profile for '2007-10-14'. The profile details include:

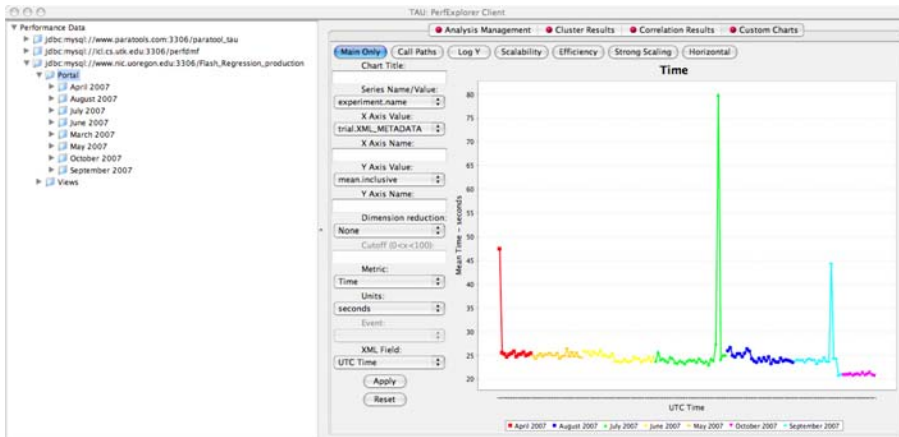
- Name: 2007-10-14
- Time: 2007-10-14 10:14:11 (S) (M) (Y)
- Number of Nodes: 4
- Clusters per Node: 1
- Threads per Cluster: 1
- Contexts: 1
- CPU Class: 7
- CPU MHz: 2962.554
- CPU Type: Intel(R) Xeon(R) CPU E5305
- CPU Model: XeonE5305
- OS: /opt/centos/binaries/linux
- Cache Size: 4096 KB
- Executable: /opt/centos/binaries/linux
- Hostname: alpha.us.oregon.edu
- Local Time: 2007-10-14 10:14:11
- URL Process Name: alpha.us.oregon.edu
- Memory Size: 4096736 KB
- Node Name: alpha.us.oregon.edu
- OS Release: x86_64
- OS Name: Linux
- OS Release: 2.6.9-5.0.0.0.1.el6
- OS Version: #1 SMP Thu Sep 27 18:28:28
- TAU Architecture: x86_64

On the right side of the interface, there is a 'Trial Run' section with instructions and a 'Header' section with a list of actions that can be performed on the profile, such as 'Click here to launch TAU's Performance Explorer and view the details of all the performance data in the workspace.'

ParaTools

233

PerfExplorer: Regression Testing



ParaTools

234

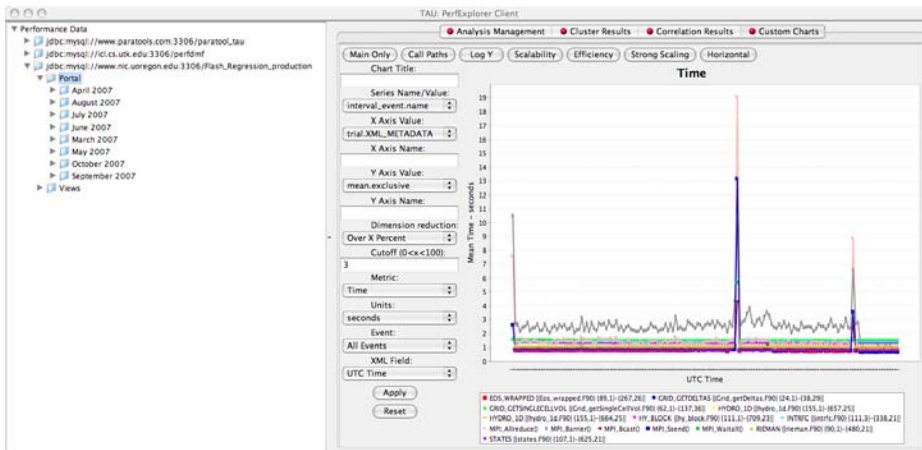
PerExplorer: Limiting Events (> 3%), Oct 2007



ParaTools

235

PerExplorer: Exclusive Time for Events (2007)



ParaTools

236

Using Performance Database (PerfDMF)

- **Configure PerfDMF (Done by each user)**
 - % perfdmf_configure --create-defaults
 - Choose derby, PostgreSQL, MySQL, Oracle or DB2
 - Hostname
 - Username
 - Password
 - Say yes to downloading required drivers (we are not allowed to distribute these)
 - Stores parameters in your ~/.ParaProf/perfdmf.cfg file
- **Configure PerfExplorer (Done by each user)**
 - % perfexplorer_configure
- **Execute PerfExplorer**
 - % perfexplorer

PerfDMF and the TAU Portal

- **Development of the TAU portal**
 - Common repository for collaborative data sharing
 - Profile uploading, downloading, user management
 - Paraprof, PerfExplorer can be launched from the portal using Java Web Start (no TAU installation required)
- **Portal URL**
 - <http://tau.nic.uoregon.edu>

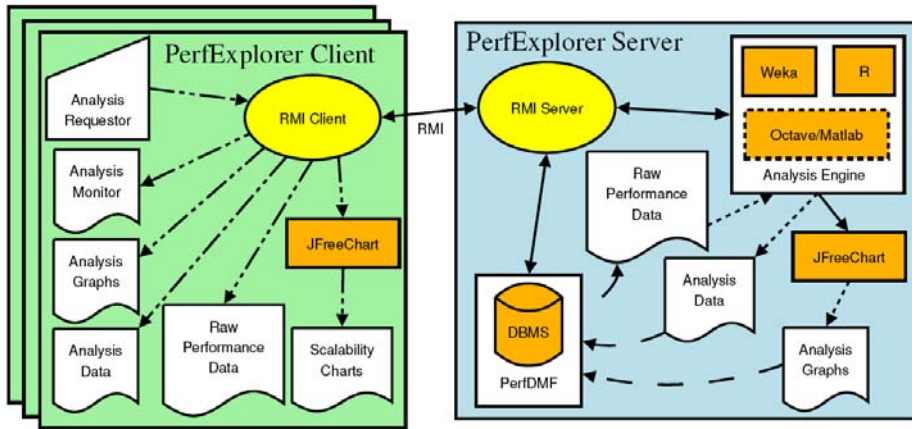
Performance Data Mining (Objectives)

- Conduct parallel performance analysis process
 - In a systematic, collaborative and reusable manner
 - Manage performance complexity
 - Discover performance relationship and properties
 - Automate process
- Multi-experiment performance analysis
- Large-scale performance data reduction
 - Summarize characteristics of large processor runs
- Implement extensible analysis framework
 - Abstraction / automation of data mining operations
 - Interface to existing analysis and data mining tools

Performance Data Mining (PerfExplorer)

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

Performance Data Mining (PerfExplorer)



PerfExplorer - Analysis Methods

- Data summaries, distributions, scatter plots
- Clustering
 - *k*-means
 - Hierarchical
- Correlation analysis
- Dimension reduction
 - PCA
 - Random linear projection
 - Thresholds
- Comparative analysis
- Data management views

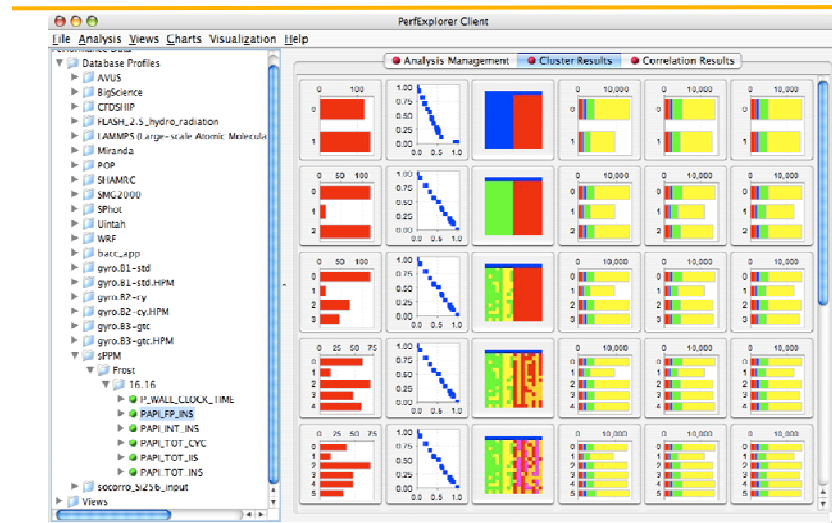
PerfExplorer - Cluster Analysis

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

ParaTools

243

PerfExplorer - Cluster Analysis (sPPM)

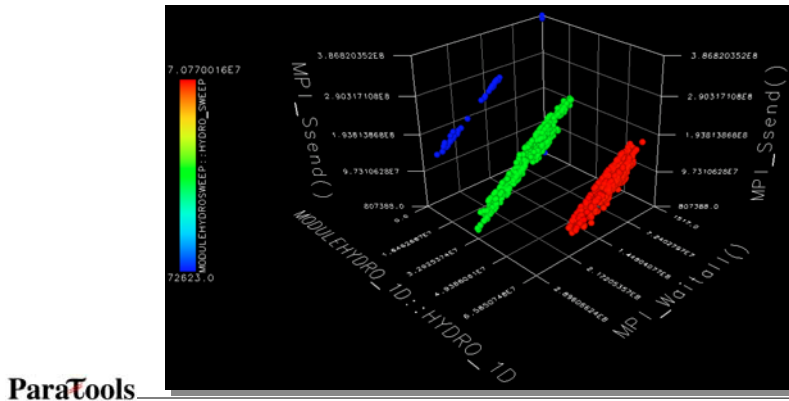


ParaTools

244

PerfExplorer - Cluster Analysis

- Four significant events automatically selected (from 16K processors)
- Clusters and correlations are visible

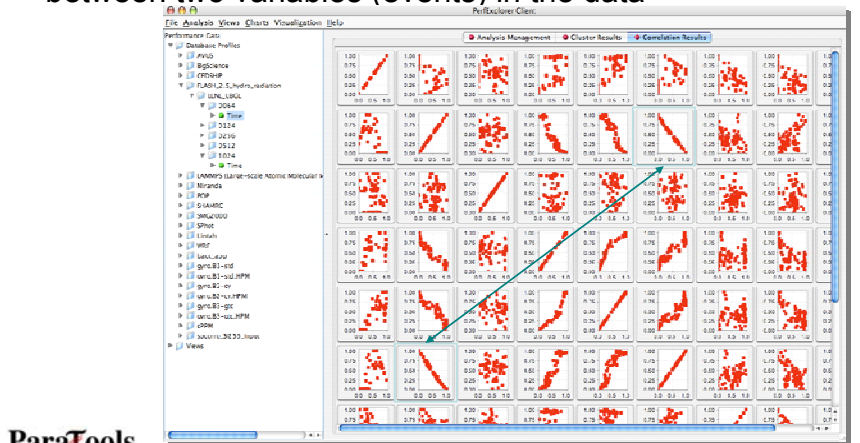


ParaTools

245

PerfExplorer - Correlation Analysis (Flash)

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

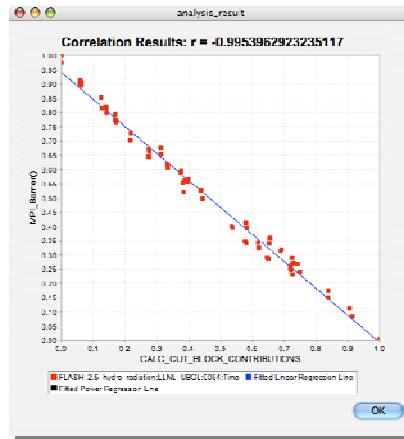


ParaTools

246

PerfExplorer - Correlation Analysis (Flash)

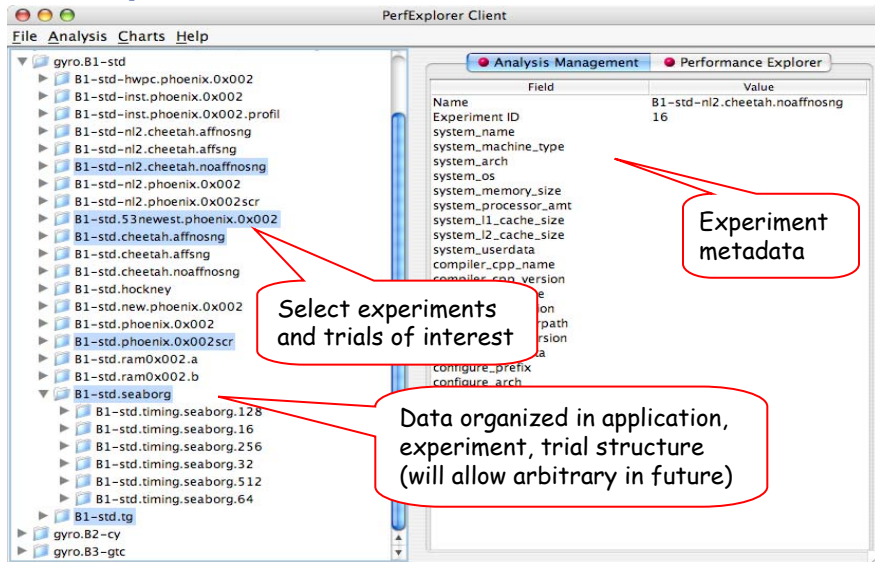
- -0.995 indicates strong, negative relationship
- As CALC_CUT_BLOCK_CONTRIBUTIONS() increases in execution time, MPI_Barrier() decreases



PerfExplorer - Comparative Analysis

- Relative speedup, efficiency
 - total runtime, by event, one event, by phase
- Breakdown of total runtime
- Group fraction of total runtime
- Correlating events to total runtime
- Timesteps per second
- Performance Evaluation Research Center (PERC)
 - PERC tools study (led by ORNL, Pat Worley)
 - In-depth performance analysis of select applications
 - Evaluation performance analysis requirements
 - Test tool functionality and ease of use

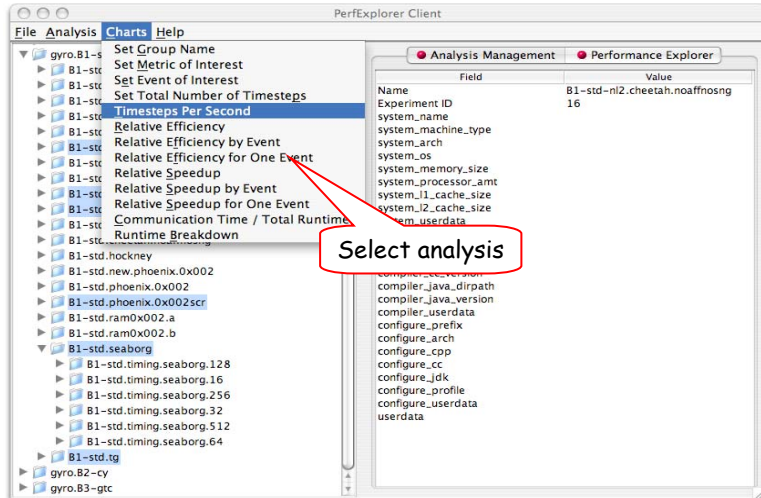
PerfExplorer - Interface



ParaTools

249

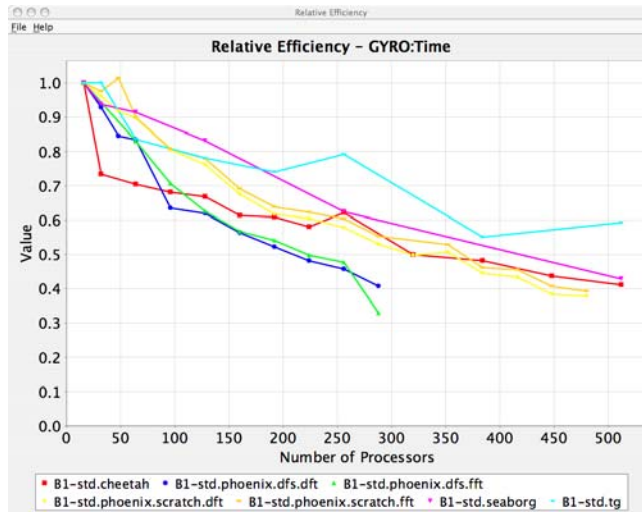
PerfExplorer - Interface



ParaTools

250

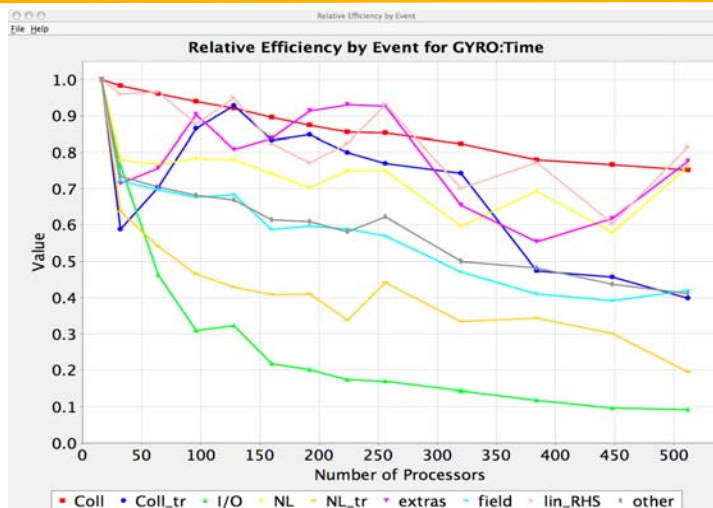
PerfExplorer - Relative Efficiency Plots



ParaTools

251

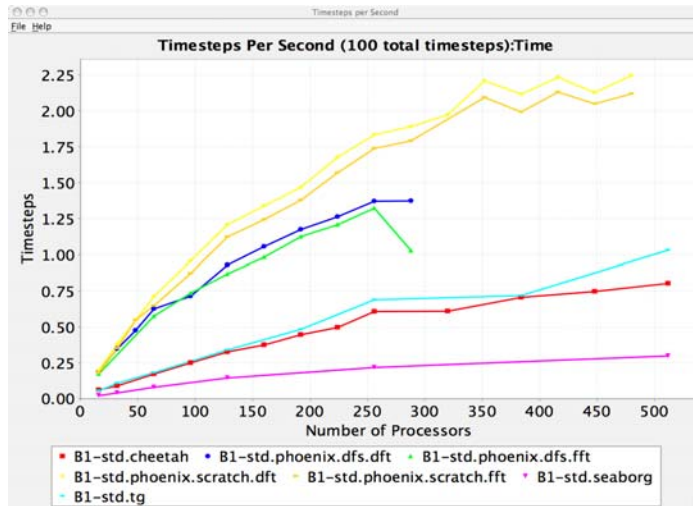
PerfExplorer - Relative Efficiency by Routine



ParaTools

252

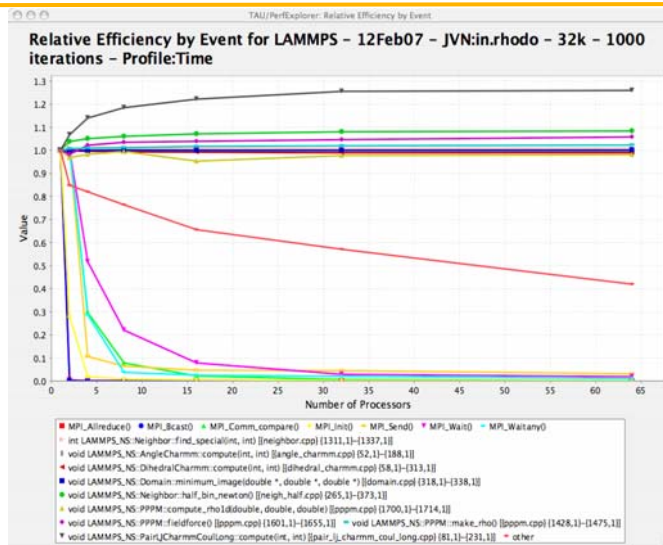
PerfExplorer - Timesteps Per Second



ParaTools

253

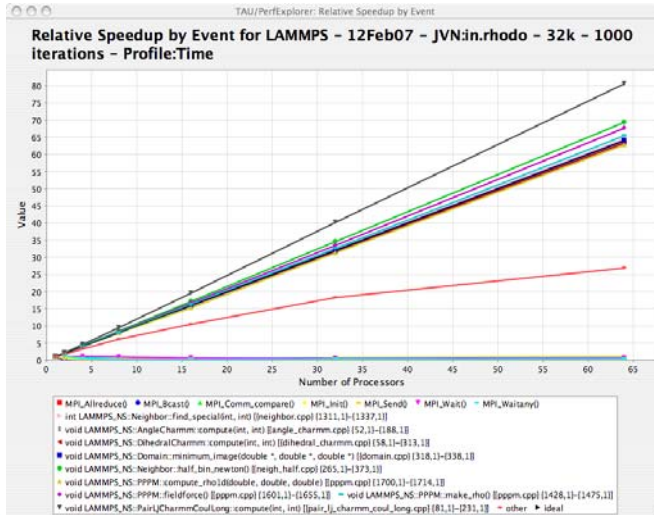
PerfExplorer - Relative Efficiency



ParaTools

254

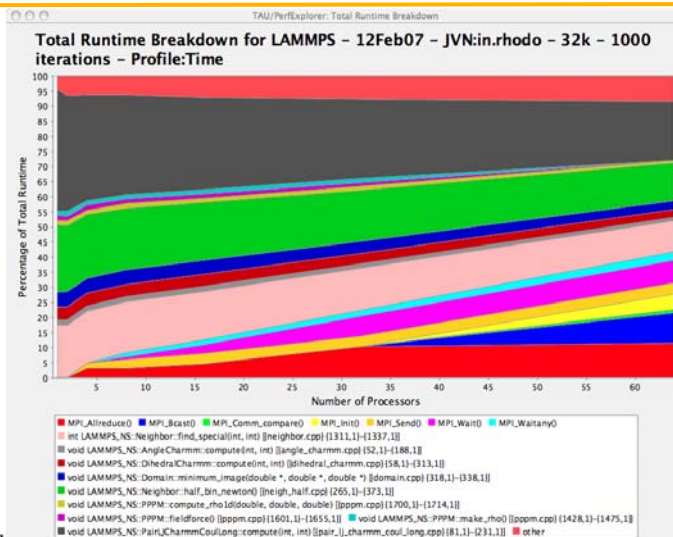
PerfExplorer - Relative Speedup by Event



ParaTools

255

PerfExplorer - Runtime Breakdown



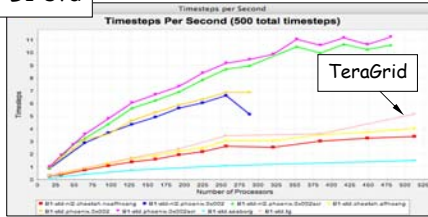
ParaTools

256

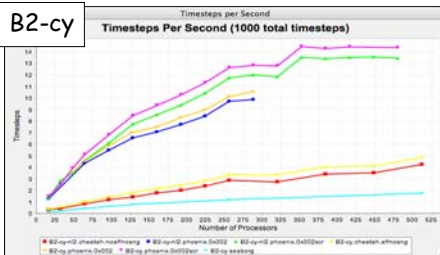
PerfExplorer - Timesteps per Second for GYRO

- Cray X1 is the fastest to solution
 - In all 3 tests
- FFT (nl2) improves time
 - B3-gtc only
- TeraGrid faster than p690
 - For B1-std?
- All plots generated automatically

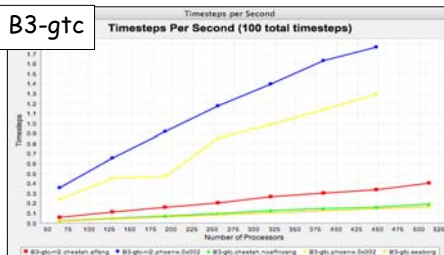
B1-std



B2-cy



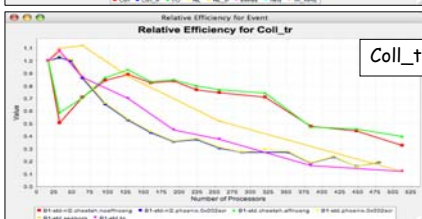
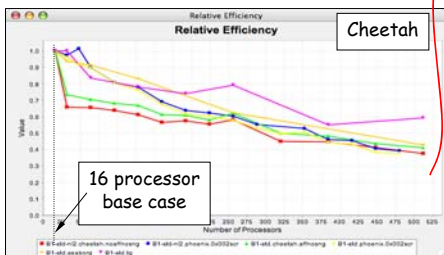
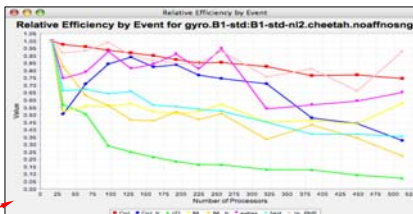
B3-gtc



ParaTools

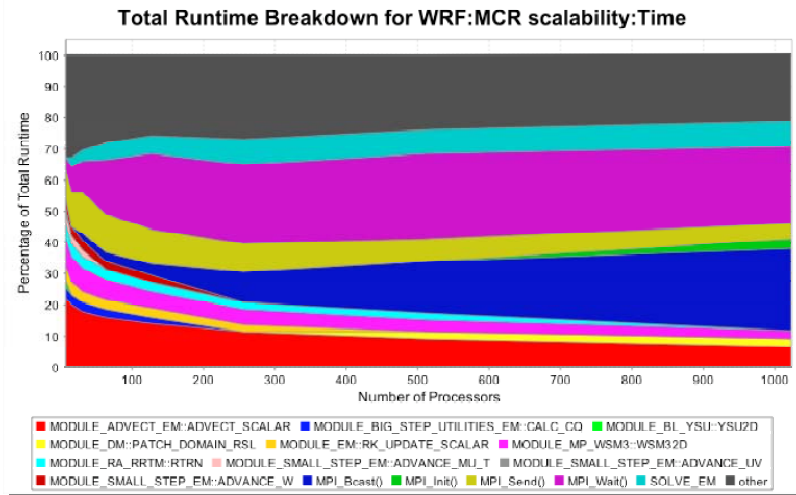
PerfExplorer - Relative Efficiency (B1-std)

- By experiment (B1-std)
 - Total runtime (Cheetah (red))
- By event for one experiment
 - Coll_tr (blue) is significant
- By experiment for one event
 - Shows how Coll_tr behaves for all experiments



ParaTools

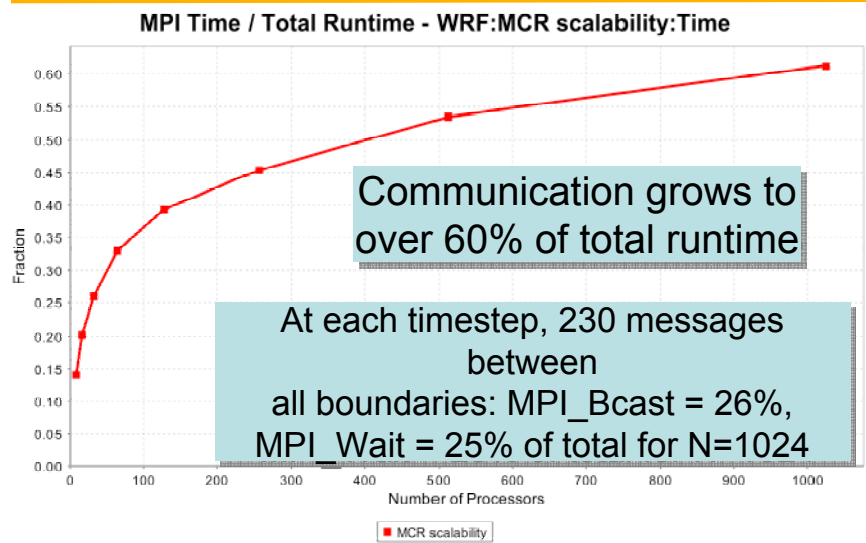
PerfExplorer - Runtime Breakdown



ParaTools

259

Group % of Total



260

TAU Integration with IDEs

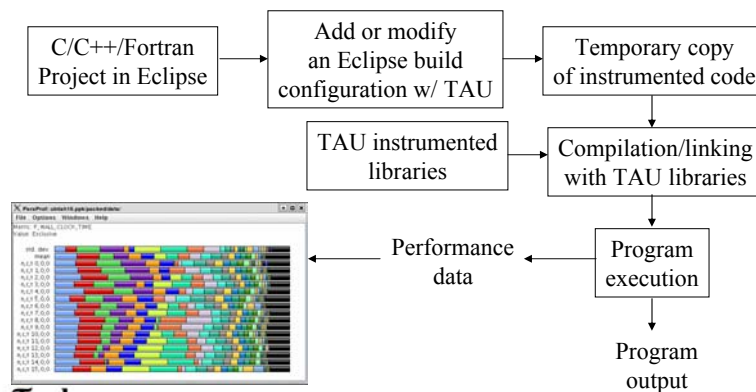
- High performance software development environments
 - Tools may be complicated to use
 - Interfaces and mechanisms differ between platforms / OS
- Integrated development environments
 - Consistent development environment
 - Numerous enhancements to development process
 - Standard in industrial software development
- Integrated performance analysis
 - Tools limited to single platform or programming language
 - Rarely compatible with 3rd party analysis tools
 - Little or no support for parallel projects

ParaTools

261

TAU and Eclipse

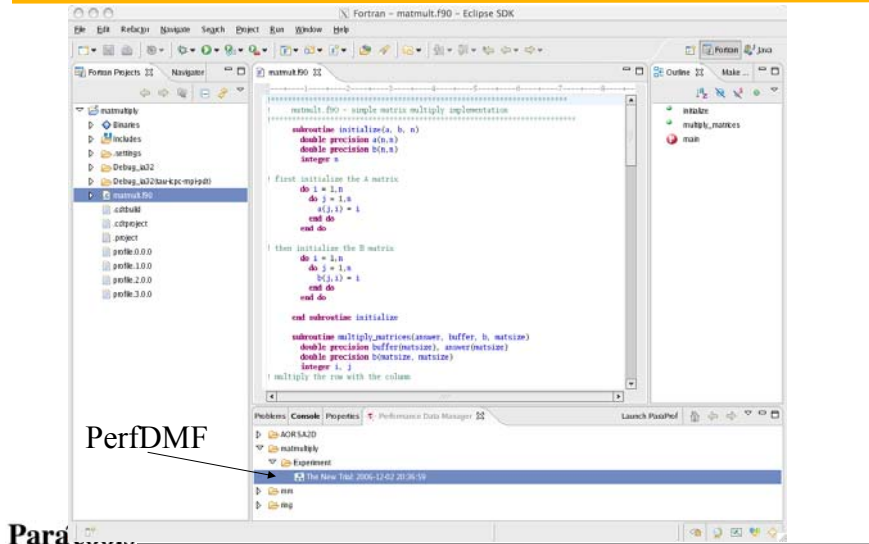
- Provide an interface for configuring TAU's automatic instrumentation within Eclipse's build system
- Manage runtime configuration settings and environment variables for execution of TAU instrumented programs



ParaTools

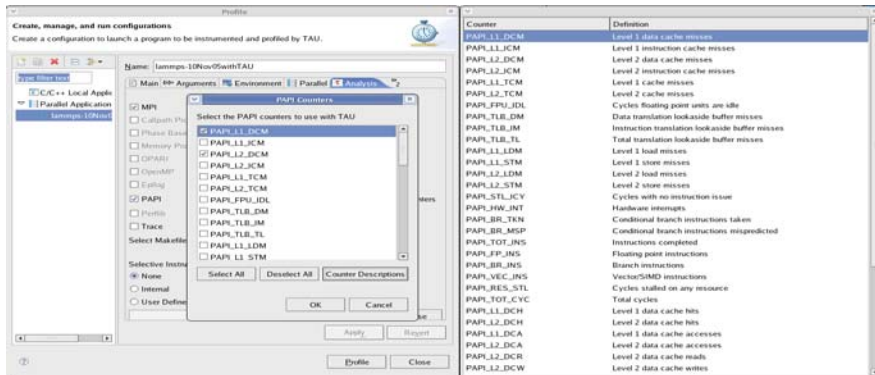
262

TAU and Eclipse



ParaTools

Choosing PAPI Counters with TAU in Eclipse



% /usr/global/tools/pkg/eclipse/eclipse

ParaTools

TAU Performance System Status

- Computing platforms (selected)
 - IBM SP/pSeries/BGL/Cell PPE, SGI Altix/Origin, Cray T3E/SV-1/X1/XT3, HP (Compaq) SC (Tru64), Sun, Linux clusters (IA-32/64, Alpha, PPC, PA-RISC, Power, Opteron), Apple (G4/5, OS X), Hitachi SR8000, NEC SX Series, Windows ...
- Programming languages
 - C, C++, Fortran 77/90/95, HPF, Java, Python
- Thread libraries (selected)
 - pthreads, OpenMP, SGI sproc, Java, Windows, Charm++
- Compilers (selected)
 - Intel, PGI, GNU, Fujitsu, Sun, PathScale, SGI, Cray, IBM, HP, NEC, Absoft, Lahey, Nagware, ...

Part V: VAMPIRTRACE & VAMPIR INTRODUCTION AND OVERVIEW

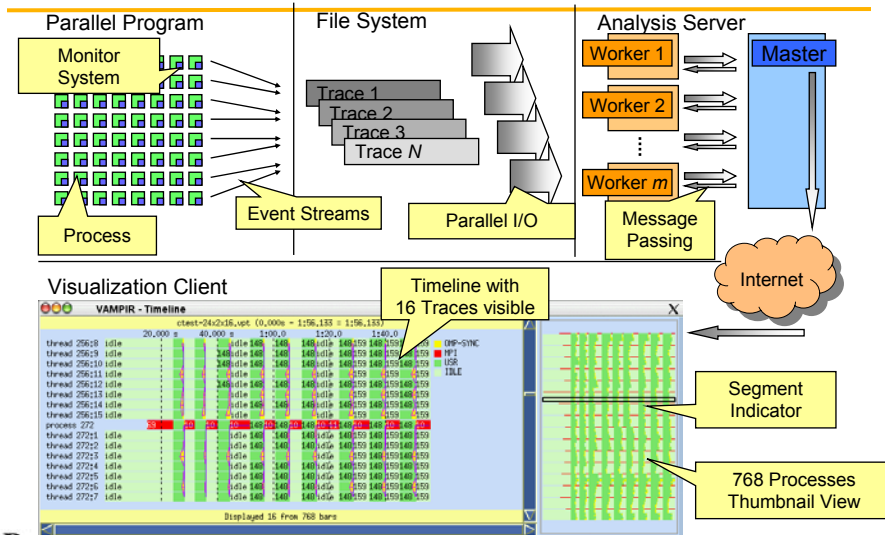
Overview

- Introduction
- Event Trace Visualization
- Vampir & VampirServer
- The Vampir Displays
 - Timeline
 - Process Timeline with Performance Counters
 - Summary Display
 - Message Statistics
- VampirTrace
 - Instrumentation & Run-Time Measurement
- Conclusions

ParaTools

267

VampirServer Architecture



ParaTools

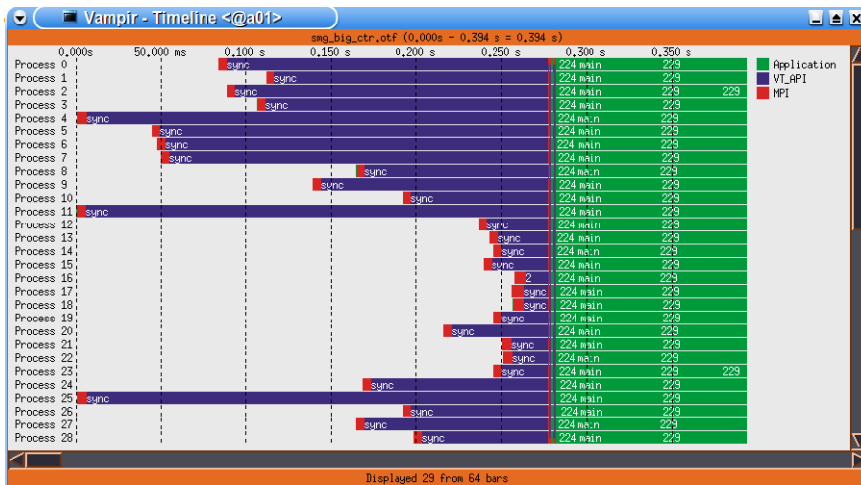
268

Vampir Displays

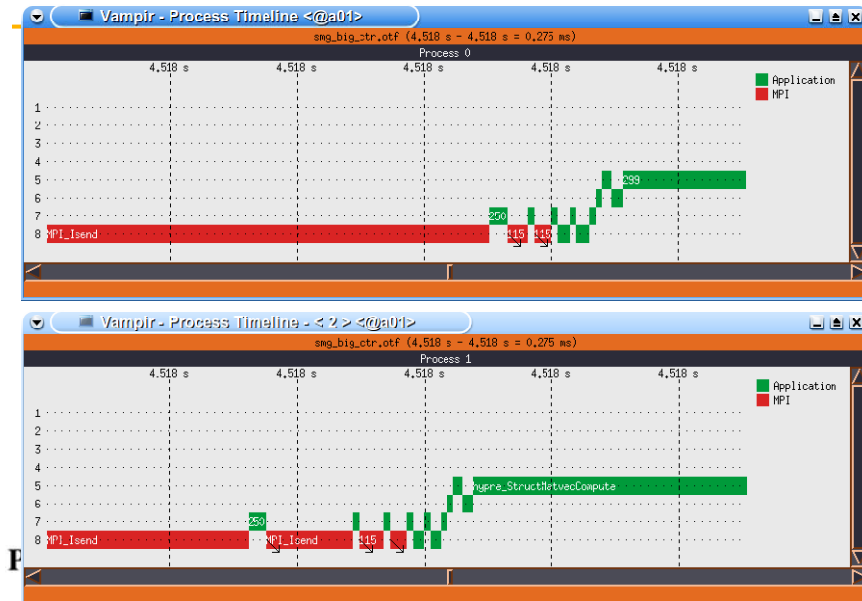
The main displays of Vampir:

- Global Timeline
- Process Timeline w/o Counters
- Statistic Summary
- Summary Timeline
- Message Statistics
- Collective Operation Statistics
- Counter Timeline
- Call Tree

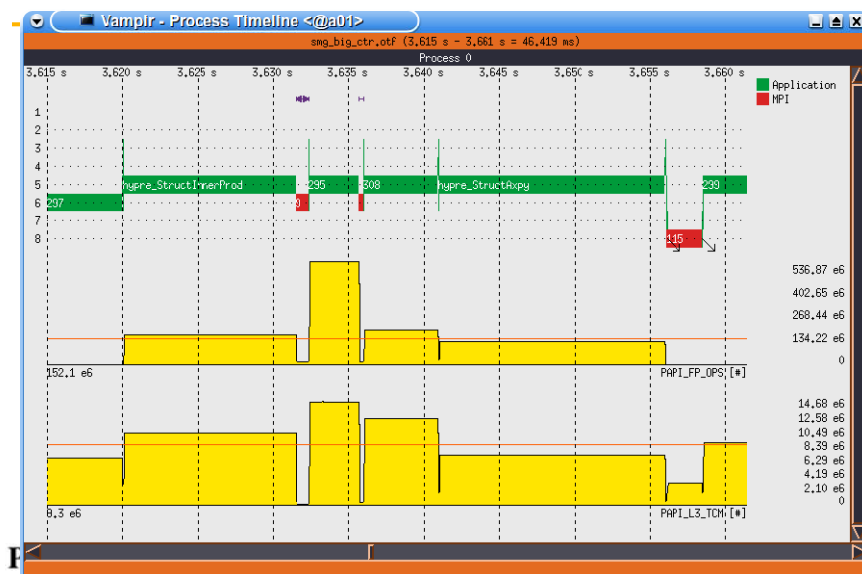
Vampir Global Timeline Display



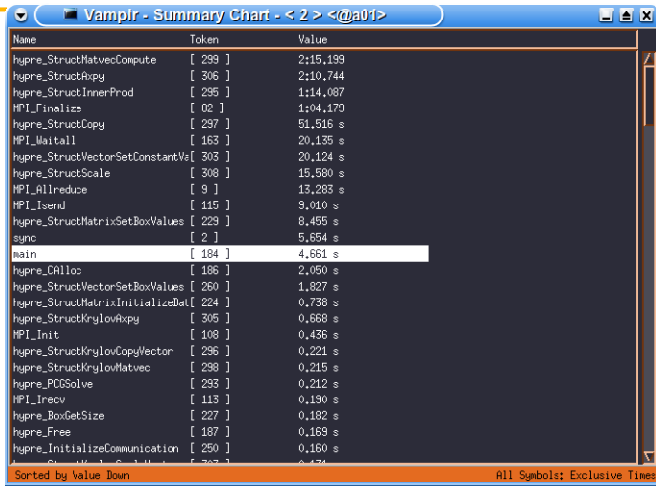
Process Timeline Display



Process Timeline with Counters



Statistic Summary Display



Vampir - Summary Chart - < 2 > <@a01>

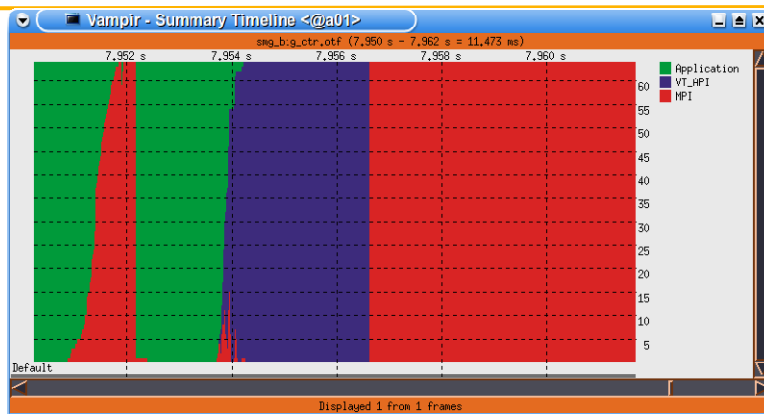
Name	Token	Value
hypre_StructMatvecCompute	[299]	2,15,199
hypre_StructAxpv	[306]	2,10,744
hypre_StructInnerProd	[295]	1,14,087
MPI_Finalize	[02]	1,04,170
hypre_StructCopy	[297]	51,516 s
MPI_Waitall	[163]	20,135 s
hypre_StructVectorSetConstantV	[303]	20,124 s
hypre_StructScale	[308]	15,580 s
MPI_Allreduce	[9]	13,283 s
MPI_Isend	[115]	8,010 s
hypre_StructMatrixSetBoxValues	[223]	8,456 s
sync	[2]	5,654 s
main	[184]	4,661 s
hypre_Calloc	[186]	2,050 s
hypre_StructVectorSetBoxValues	[260]	1,827 s
hypre_StructMatrixInitializeDat	[224]	0,738 s
hypre_StructKrylovAxpv	[305]	0,668 s
MPI_Init	[108]	0,436 s
hypre_StructKrylovCopyVector	[296]	0,221 s
hypre_StructKrylovMatvec	[298]	0,215 s
hypre_PCGSolve	[293]	0,212 s
MPI_Irecv	[113]	0,190 s
hypre_BoxGetSize	[227]	0,182 s
hypre_Free	[187]	0,169 s
hypre_InitializeCommunication	[250]	0,160 s

Sorted by value Down All Symbols: Exclusive Times

ParaTools

273

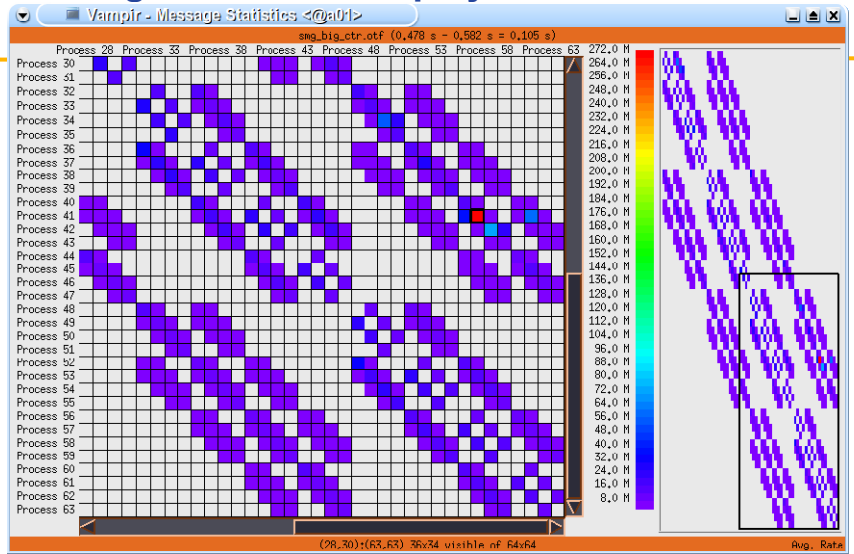
Summary Timeline Display



ParaTools

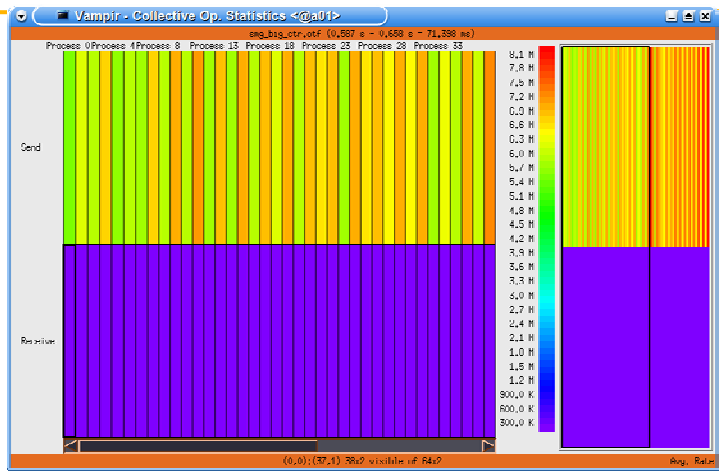
274

Message Statistics Display



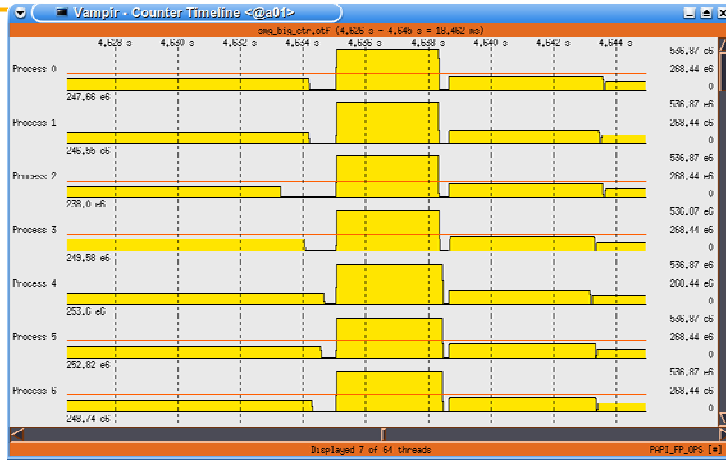
ParaTools

Collective Operation Statistics



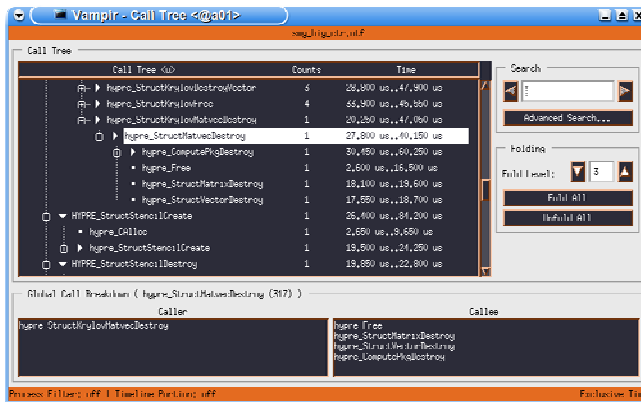
ParaTools

Counter Timeline Display



ParaTools

Call Tree Display



ParaTools

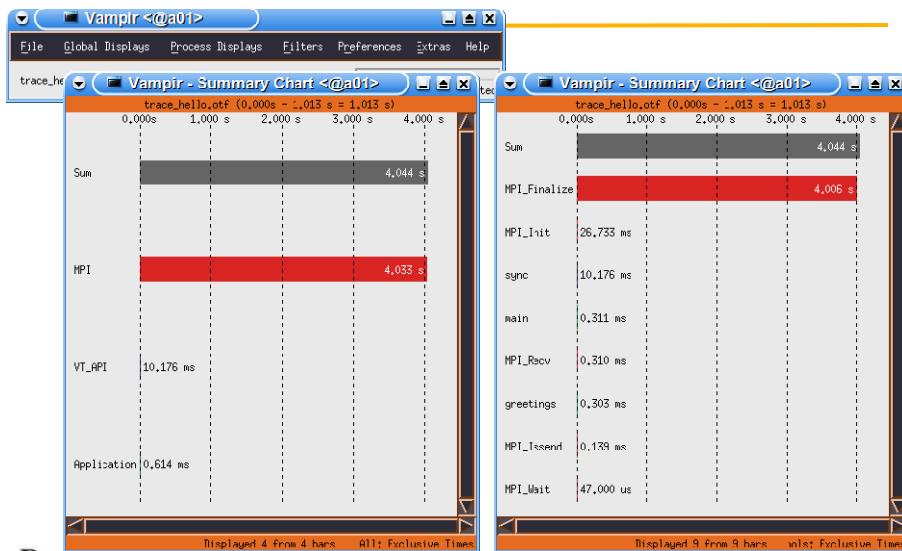
Open Trace Format (OTF)

- Open source trace file format
- Available at <http://www.tu-dresden.de/zih/otf/>
- Includes powerful libotf for reading/parsing/writing in custom applications
- multi-level API:
 - High level interface for analysis tools
 - Low level interface for trace libraries
- Actively developed in cooperation with the University of Oregon and the Lawrence Livermore National Laboratory

ParaTools

279

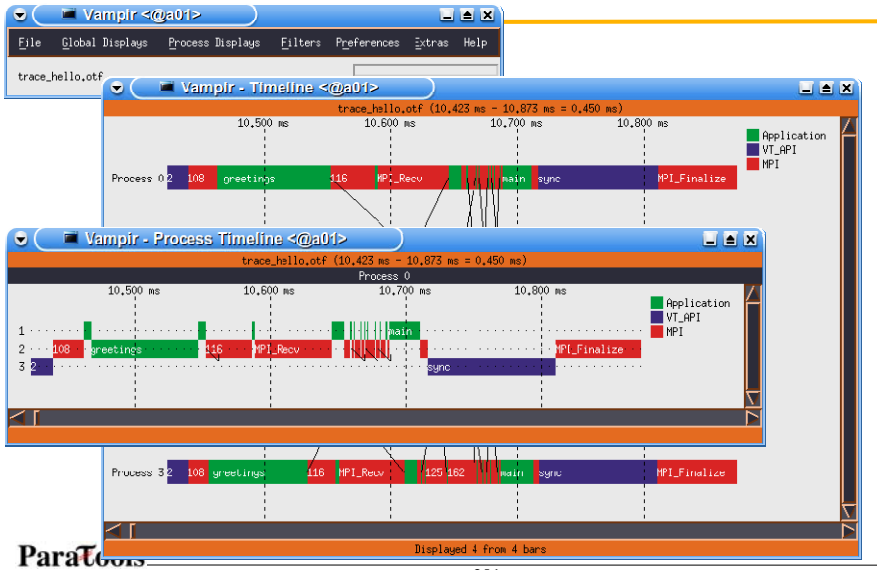
Hands-on: VampirServer



ParaTools

280

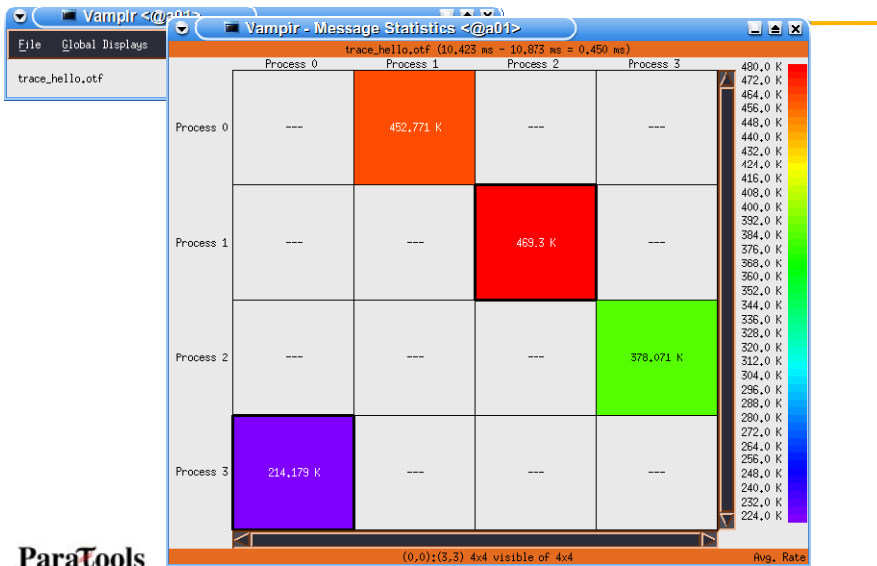
Hands-on: More Displays



ParaTools

281

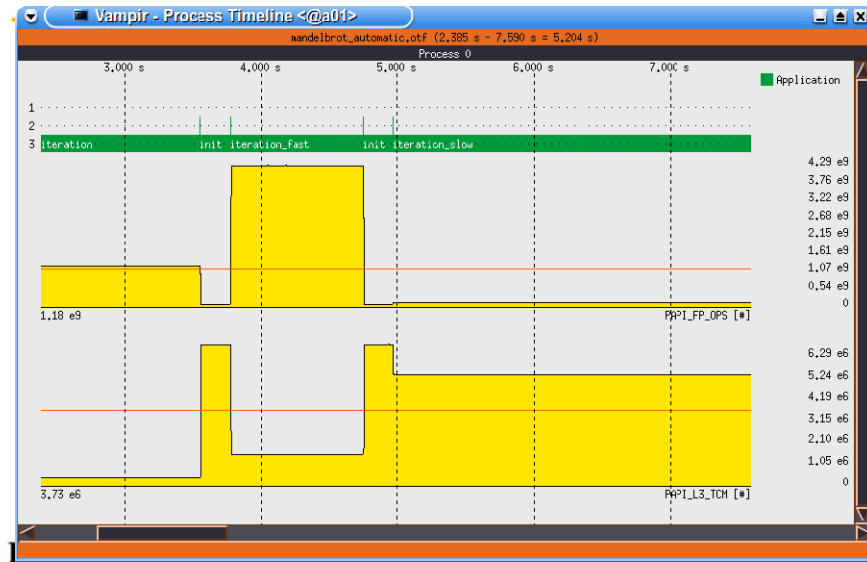
Hands-on: More Displays



ParaTools

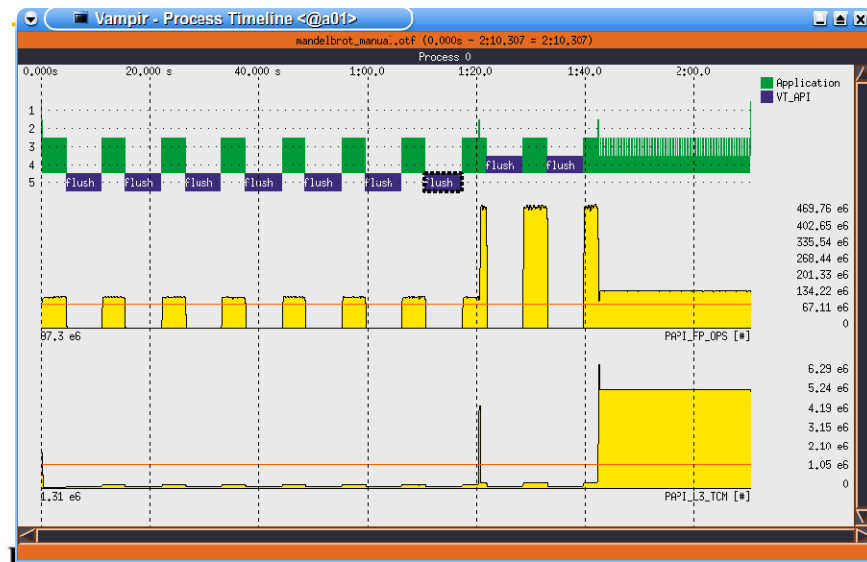
282

Hands-on: Performance Counters



283

Extra Manual Instrumentation



284

Finding Performance Bottlenecks

Finding Bottlenecks

- Trace Visualization
 - Vampir provides a number of display types
 - each allows many different options
- Advice
 - identify essential parts of an application (initialization, main iteration, I/O, finalization)
 - identify important components of the code (serial computation, MPI P2P, collective MPI, OpenMP)
 - make a hypothesis about performance problems
 - consider application's internal workings if known
 - select the appropriate displays
 - use statistic displays in conjunction with timelines

Finding Bottlenecks

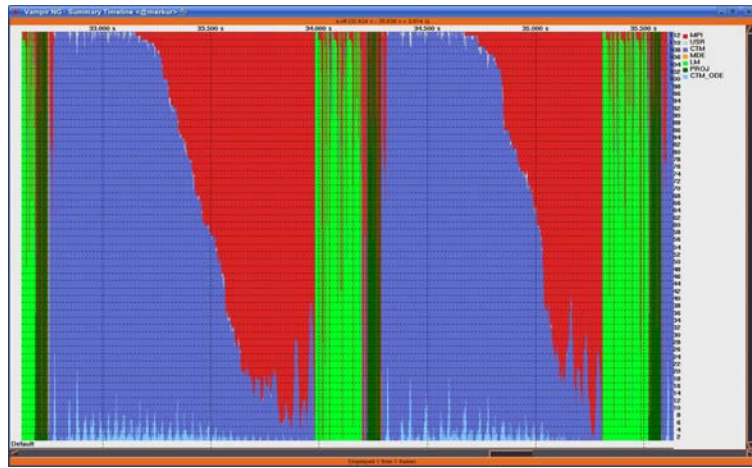
- Communication
- Computation
- Memory, I/O, etc
- Tracing itself

Bottlenecks in Communication

- communication as such (dominating over computation)
- late sender, late receiver
- point-to-point messages instead of collective communication
- unmatched messages
- overcharge of MPI's buffers
- bursts of large messages (bandwidth)
- frequent short messages (latency)
- unnecessary synchronization (barrier)

all of the above usually result in high MPI time share

Bottlenecks in Communication

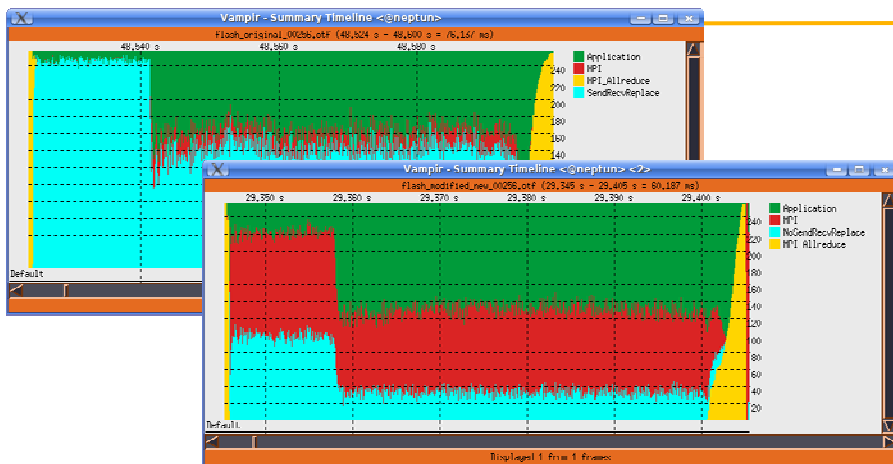


Example: prevalent communication

ParaTools

289

Bottlenecks in Communication

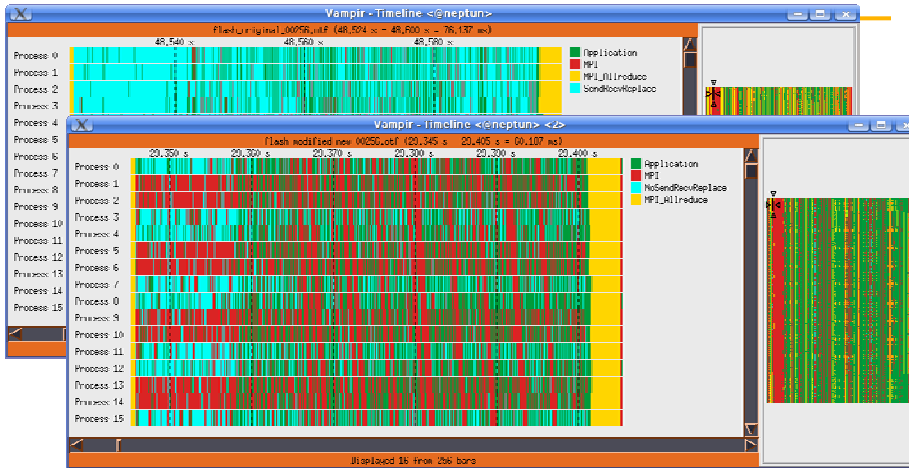


prevalent communication: MPI_Allreduce

ParaTools

290

Bottlenecks in Communication

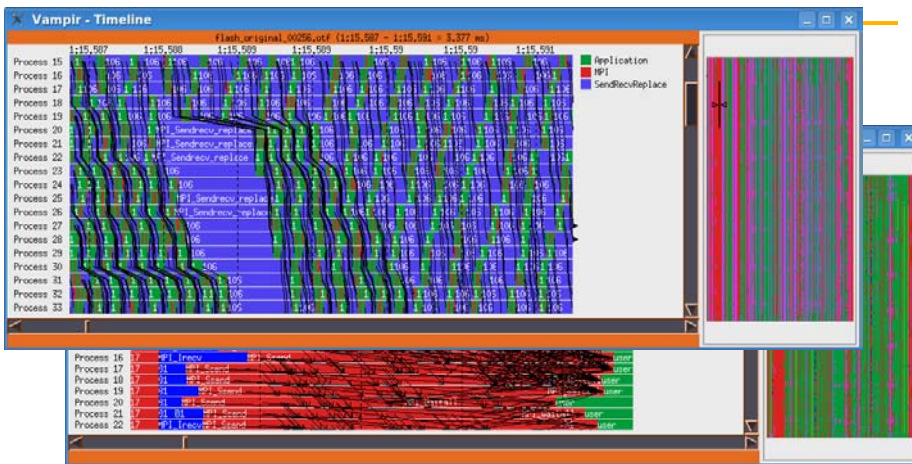


prevalent communication: timeline view

ParaTools

291

Bottlenecks in Communication

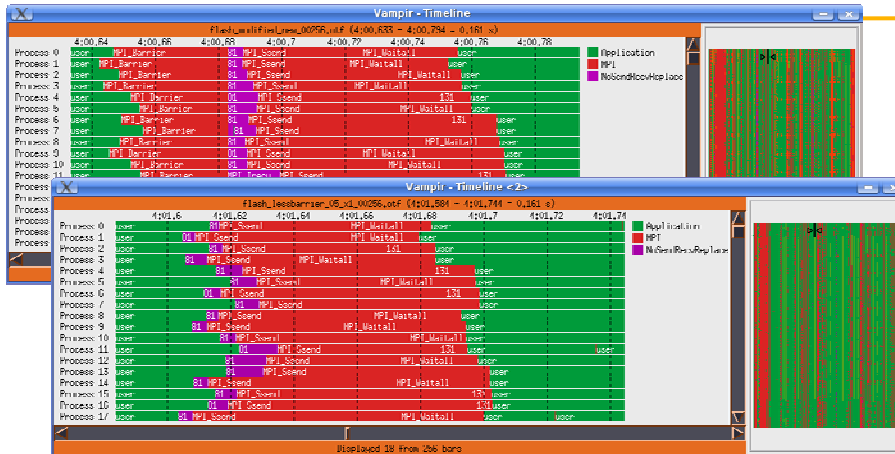


Propagated Delays in MPI_SendReceiveReplace

ParaTools

292

Bottlenecks in Communication

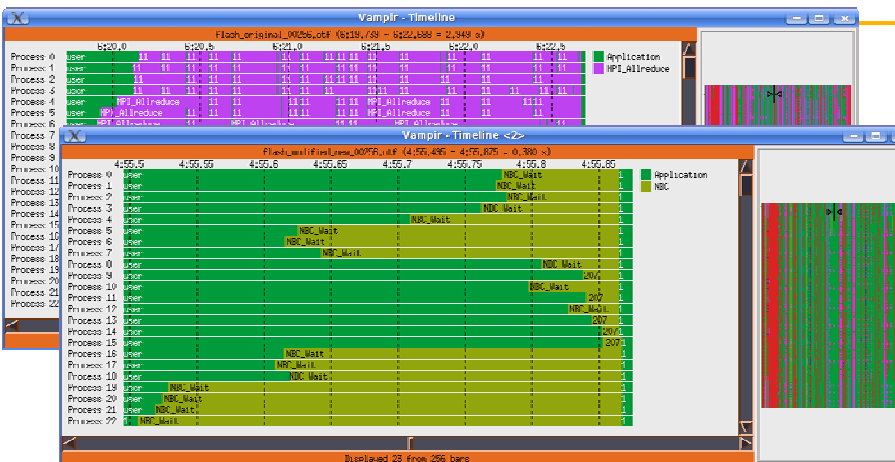


unnecessary MPI_Barriers

ParaTools

293

Bottlenecks in Communication



Patterns of Successive MPI_Allreduce Calls

ParaTools

294

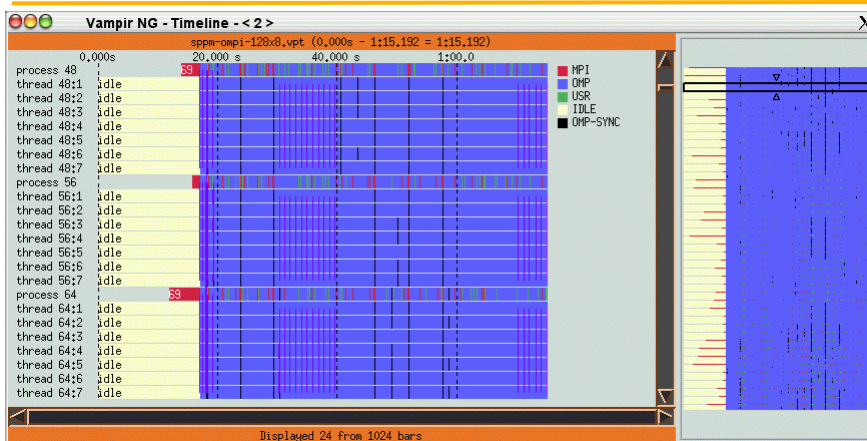
Further Bottlenecks

- unbalanced computation
 - single late comer
- strictly serial parts of program
 - idle processes/threads
- very frequent tiny function calls
- sparse loops

ParaTools

295

Further Bottlenecks



Example: Idle OpenMP threads

ParaTools

296

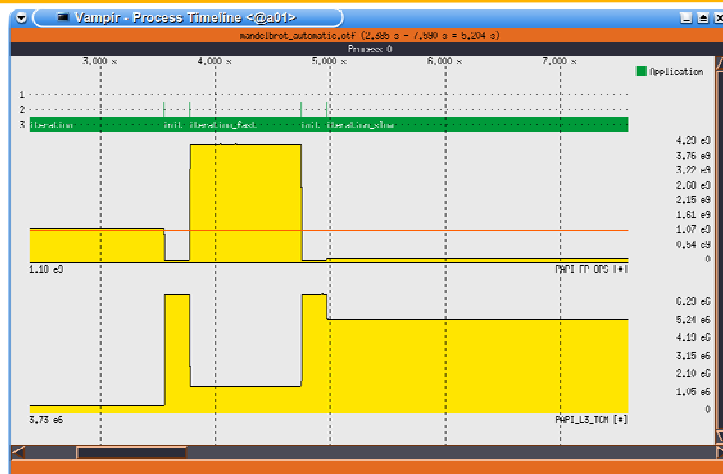
Bottlenecks in Computation

- memory bound computation
 - inefficient L1/L2/L3 cache usage
 - TLB misses
 - detectable via HW performance counters
- I/O bound computation
 - slow input/output
 - sequential I/O on single process
 - I/O load imbalance
- exception handling

ParaTools

297

Bottlenecks in Communication

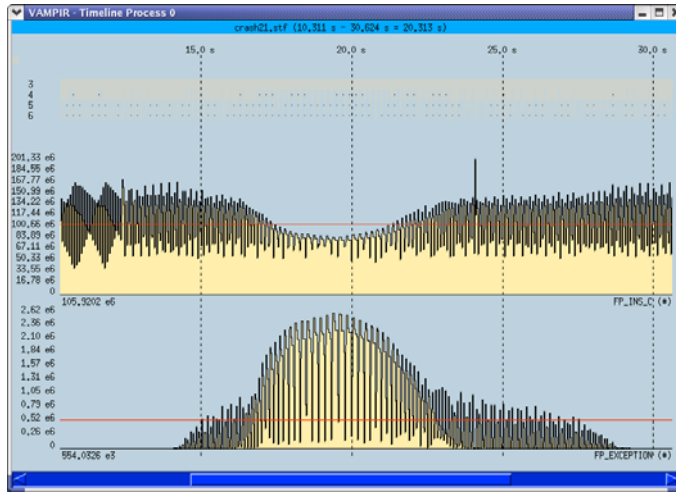


ParaTools

low FP rate due to heavy cache misses

298

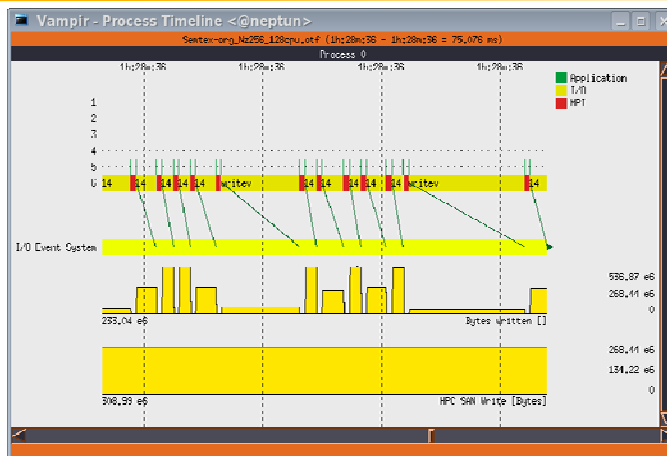
Bottlenecks in Communication



ParaTools low FP rate due to heavy FP exceptions

299

Bottlenecks in Communication



irregular slow I/O operations

ParaTools

300

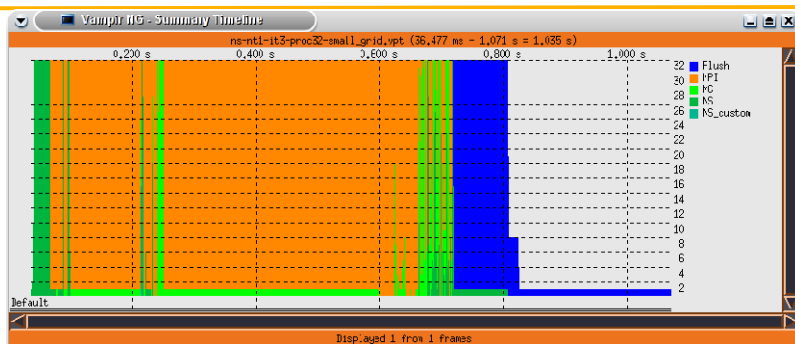
Effects due to Tracing Itself

- measurement overhead
 - esp. grave for tiny function calls
 - solve with selective instrumentation
- long/frequent/asynchronous trace buffer flushes
- too many concurrent counters
- heisenbugs

ParaTools

301

Effects due to Tracing Itself



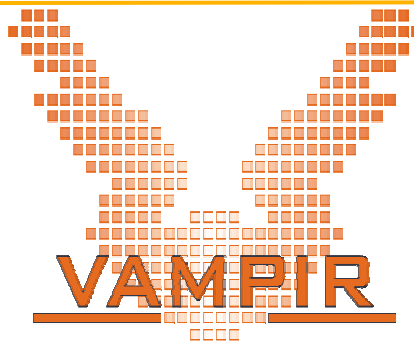
Trace buffer flushes are explicitly marked in the trace.
It is rather harmless at the end of a trace as shown here.

ParaTools

302

Conclusion and Outlook

- performance analysis very important in HPC
- use performance analysis tools for profiling and tracing
- do not spend effort in DIY solutions, e.g. like printf-debugging
- use tracing tools with some precautions
 - overhead
 - data volume
- let us know about problems and about feature wishes
- vampirsupport@zih.tu-dresden.de



Vampir and VampirTraces are available at <http://www.vampir.eu> and <http://www.tu-dresden.de/zih/vampirtrace/> ,

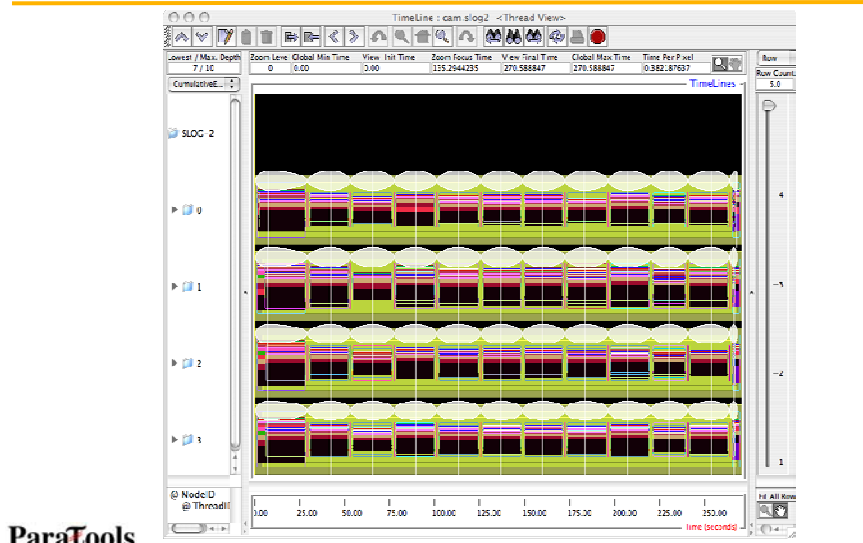
Jumpshot

- <http://www-unix.mcs.anl.gov/perfvis/software/viewers/index.htm>
- Developed at Argonne National Laboratory as part of the MPICH project
 - Also works with other MPI implementations
 - Installed on NAVO IBM and ERDC XT3/4
 - Jumpshot is bundled with the TAU package
- Java-based tracefile visualization tool for postmortem performance analysis of MPI programs
- Latest version is Jumpshot-4 for SLOG-2 format
 - Scalable level of detail support
 - Timeline and histogram views
 - Scrolling and zooming
 - Search/scan facility

ParaTools

305

Jumpshot



306

Part VI: KOJAK/Scalasca



Overview

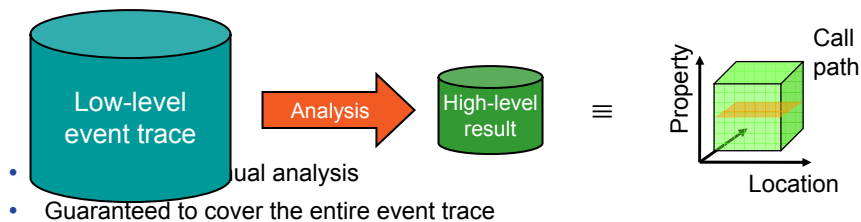
- Introduction
 - Motivation for automatic trace analysis
- Scalasca components and usage
 - instrumentation
 - measurement collection & automated analysis
 - analysis report exploration
- Demonstration
- Summary

Motivation

- Tracing offers critical insight into temporal behaviour of parallel execution unavailable from summarization
 - Inefficiencies manifest as wait states and imbalance
- Trace sizes proportional to number of processes/threads
 - as well as length of measurement and depth of detail
- Large-scale parallel traces must be carefully managed
 - minimization/elimination of disruptive file I/O
 - efficient parallel analysis of traces
 - effective hierarchical/graphical analysis presentation
- Simplification and ease-of-use
 - Automation of search for and classification of event patterns
 - Integration with trace visualizers to examine key instances

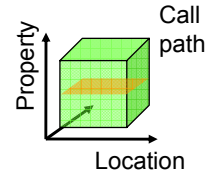
Automatic Trace Analysis

- Idea:
 - Automatic search for *patterns* of inefficient behaviour
 - Classification of behaviour
 - Quantification of significance



CUBE Result Browser

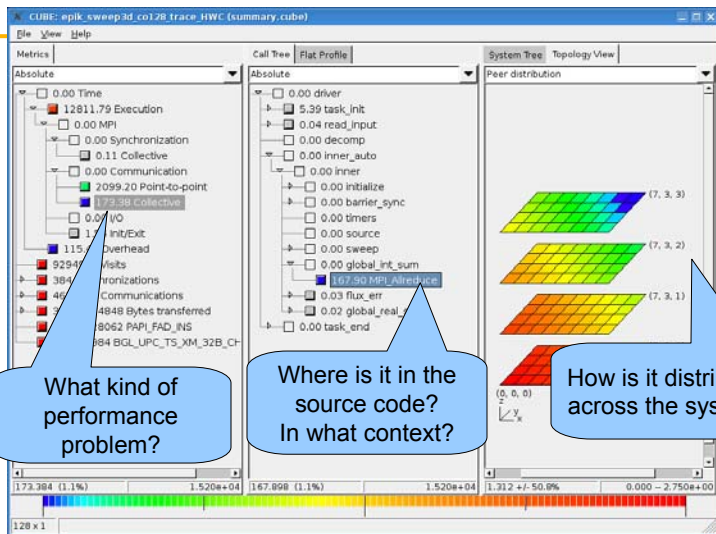
- Representation of results (severity matrix) along three hierarchical axes
 - Performance property
 - Call tree path
 - System location
- Three coupled tree browsers
- Each node displays severity
 - As colour: for easy identification of hotspots
 - As value: for precise comparison
 - Inclusive value when closed or exclusive when expanded
 - Customizable via display mode



ParaTools

311

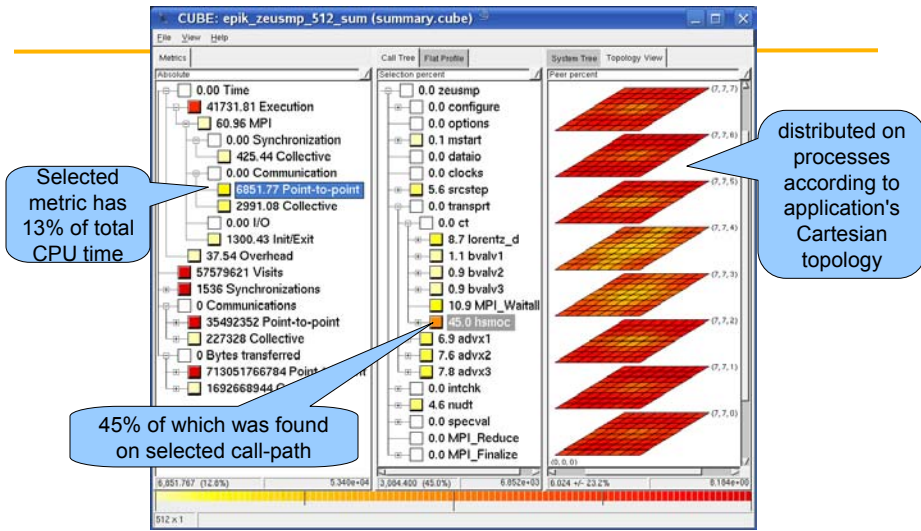
Basic Analysis Presentation



ParaTools

312

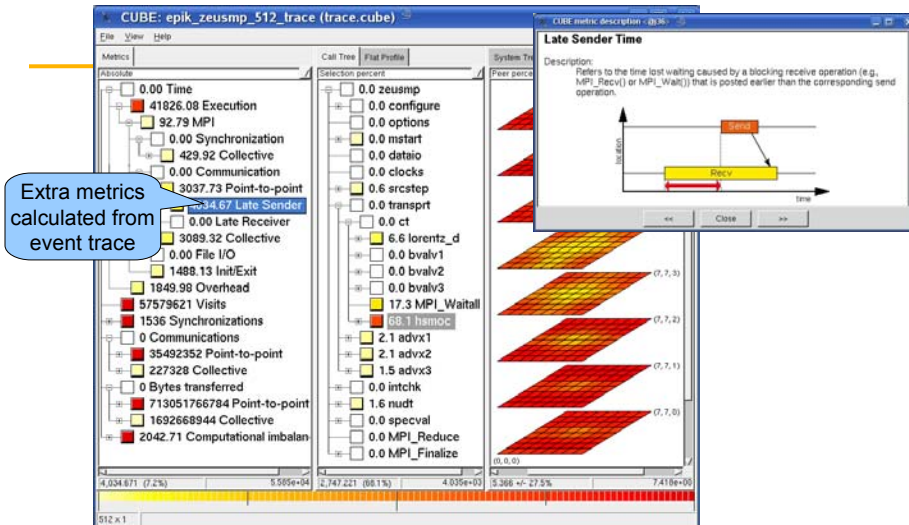
Summary Profile Analysis



ParaTools

313

Trace Pattern Analysis



ParaTools

314

Analysis Methodology

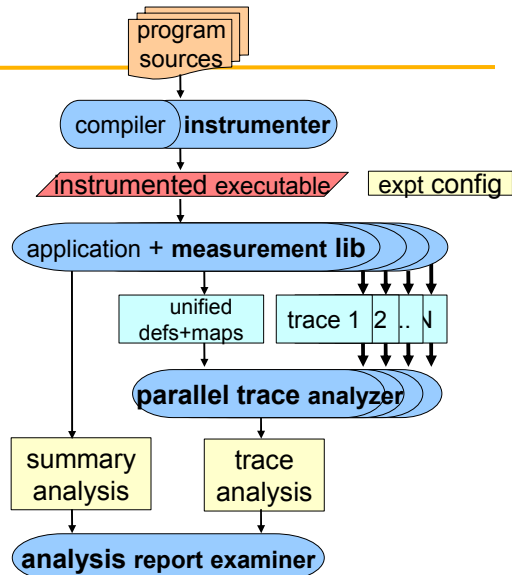
- Instrumentation of application executable and libraries
 - automatic MPI, OpenMP and function instrumentation
 - complementary manual region and phase instrumentation
- Execution of instrumented executable under control of configurable measurement collection & analysis nexus
 - commence from scalable runtime summary
 - identify excess instrumentation and trace buffer requirements
 - target tracing where it is most productive (and practical)
 - analyze traces using same resources as measurement
- Interactive analysis report exploration and algebra
 - examine severities and their locations
 - combine, compare and process reports
- Refine and repeat as necessary

ParaTools

315

Scalasca Components

- Scalasca instrumenter
= SKIN
- Scalasca measurement collector & analyzer
= SCAN
- Scalasca analysis report examiner
= SQUARE



ParaTools

316

Scalasca unified command: scalasca

- Run without action argument for basic usage info

```
% scalasca
usage: scalasca [-v][-n] {action}
1. prepare application objects and executable for measurement:
   scalasca -instrument <compile-or-link-command> # skin
2. run application under control of measurement system:
   scalasca -analyze <application-launch-command> # scan
3. interactively explore measurement analysis report:
   scalasca -examine <experiment-archive|report> # square
```
- Simply a convenience wrapper for action commands

Scalasca instrumenter: skin

- Usage: scalasca -instrument [opts] \$CC ...
 - **scalasca -instrument -user** mpicc -fast -c bar.c
 - **skin** mpif90 -Openmp -o foobar -fast foo.c bar.o -lm
- Processes source modules during compile & augments link with measurement library
 - Configures automatic function instrumentation capability of native compiler (if available)
 - All functions in source module(s) are instrumented
 - **[-pomp]** option enables processing of POMP directives
 - Optional manual source annotation of functions & regions
 - Replaces automatic function instrumentation
 - **[-user]** activates EPIK user-annotation API

Scalasca collector & analyzer: scan

- Usage: scalasca -analyze [opts] <launch command>
 - **scan** [opts] [launcher [args]] [target [target-args]]
- Prepares & runs measurement collection, with follow-on trace analysis (if appropriate)
 - [-n] preview without executing launches
 - [-s] enables runtime summarization [default]
 - [-t] enables trace collection & automatic pattern analysis
 - determines NP and/or NT (number of processes & threads) and MODE=vn|co|dual|smp (where appropriate)
 - names default measurement experiment archive `epik_${TARGET}_${MODE}${NP}x${NT}_[sum|trace]`
 - [-f filter] specifies file listing functions not to be measured
 - [-m metric1:metric2:...] includes hardware counter metrics

Scalasca analysis report explorer: square

- Usage: scalasca -examine <epik_archive | cubefile >
 - **scalasca -examine** epik_sweep3d_co32_trace
 - **square** epik_sweep3d_co32_trace/summary.cube
- Prepares & presents final analysis report
 - Checks EPIK archive directory for cubefiles
 - Remaps primitive initial analysis report(s) into refined formal report(s) with enriched metrics & metric hierarchies
 - epitome.cube -> summary.cube
 - scout.cube -> trace.cube
 - Presents refined report in CUBE3 browser
 - Trace analysis shown in preference to summary analysis
 - Additional reports can be loaded via File/Open menu

EPIK experiment archive

- Directory created by measurement library
 - Measurement aborts if archive already exists!
- Contains all files related to measurement
 - Measurement & analysis logs (epik.log, scout.log, etc.)
 - Primitive analysis reports (epitome.cube, scout.cube)
 - Refined analysis reports (summary.cube, trace.cube)
 - Process trace datafiles (ELG/*)
 - Unified definitions & map data (epik.esd, epik.map)
 - Miscellaneous (epik.conf, epik.filt, epik.path)

EPIK measurement configuration

- **epik_conf** reports current configuration
 - logged in measurement archive as epik.conf
- Read from EPIK.CONF file(s)
 - System default: \$SCALASCA_DIR/doc
 - Directory specified with EPIK_CONF environment variable [defaults to “.”]
- Over-ridden by environment variables
 - with same names as configuration file variables
- Over-ridden by **scan** command-line settings

Trace collection & analysis issues

- Process rank trace too large for trace collection buffers
 - Results in intermediate trace buffer flushes (with remainder flushed at measurement finalization)
 - Serious measurement perturbation!
- Irrelevant functions encumber analysis
 - Undesirable complexity and processing slowdown
 - Parallel trace analyzer requires memory more than twice largest rank (uncompressed) trace size to load entire trace
- Options
 - enlarge trace buffer size: `ELG_BUFFER_SIZE`
 - **cube3_score** utility provides estimate from summary
 - remove selected function instrumentation
 - specific function measurement filter (if supported!)

Selective instrumentation/measurement

- Unimportant functions can be determined from summary analysis report
 - form leaves of callpath-tree (w/o MPI)
 - negligible proportion of (exclusive) execution time
 - high proportion of (exclusive) visit count
 - **cube3_score -r** provides region breakdown & classification
 - MPI, USR (no MPI), COM (combined/intermediate)
- Eliminating pure user (USR) regions reduces overheads
 - runtime processing, storage & analysis
- Makes them “invisible” in the analysis
 - logically become part of their calling functions (as if they were in-lined by an optimizing compiler!)

Scalasca runtime summarization

- Event measurements accumulated and summarized for each call-path during runtime execution
- Summary report produced at finalization
- Provides overview of measured execution
 - contains call-path Visit frequency, Time, and MPI message statistics
 - *plus* optional hardware counter metrics
 - size independent of length of execution
- Scales to long execution measurements

Scalasca trace analysis

- Trace analysis based on parallel replay
 - enables scalability to thousands of processes
 - however, only suited to relatively brief measurements!
- Extends summary metric analysis
 - Summary can help configure selective tracing
- Allows execution performance properties to be more accurately determined and refined
- Can be combined with complementary runtime summary analysis
 - avoiding storage/processing overhead of hardware counter metrics in traces via direct summarization

Measurement support

- OpenMP compilers
 - GCC
 - IBM XL
 - Intel
 - Pathscale
 - PGI
 - Sun Studio
 - ...
- Supported functionality varies by language, version & system
- MPI libraries
 - MPICH 1 & 2
 - OpenMPI
 - Intel-MPI
 - IBM POE & BlueGene
 - Cray XT
 - Sun HPC ClusterTools
 - SGI MPToolkit
 - SiCortex MPI
 - Scali-MPI
 - HP-MPI
 - LAM

Basic use of Scalasca

- Automatic function instrumentation
 - Supported by most but not all compilers!
- Summary measurement experiment
- Summary analysis report exploration

- Trace collection & analysis experiment
- Trace pattern analysis report exploration

CUBE metrics dimension

What kind of performance problem?

Right-click metric context menu for info or description

ParaTools

329

CUBE call tree dimension

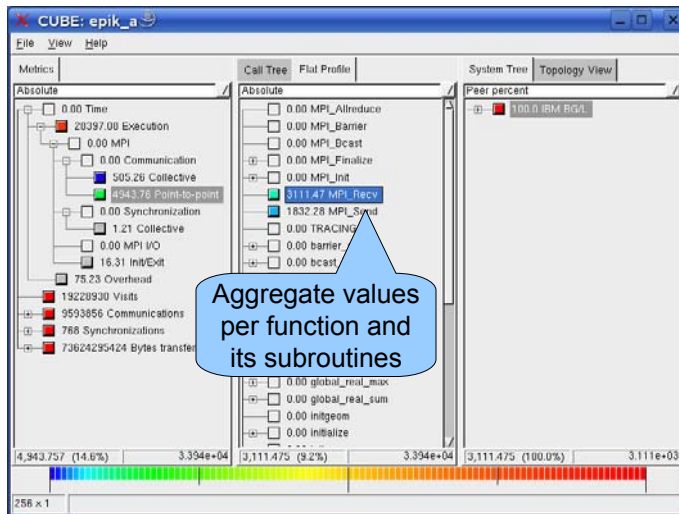
Where is it in the source code?
In what context?

Right-click function context menu to go to source location

ParaTools

330

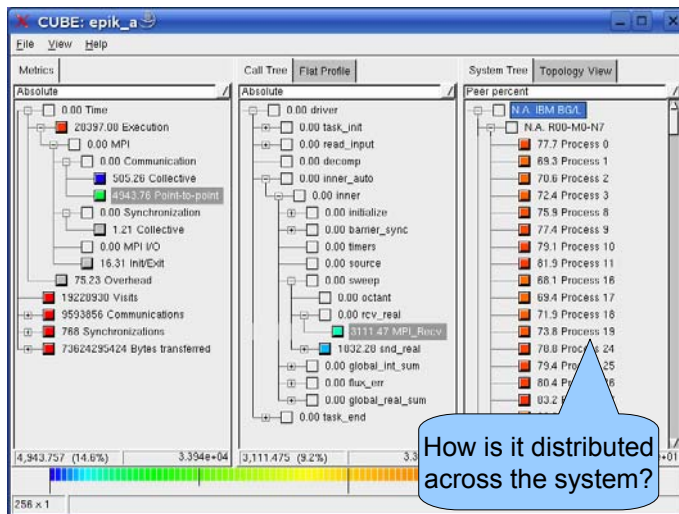
Alternative: Flat profile



ParaTools

331

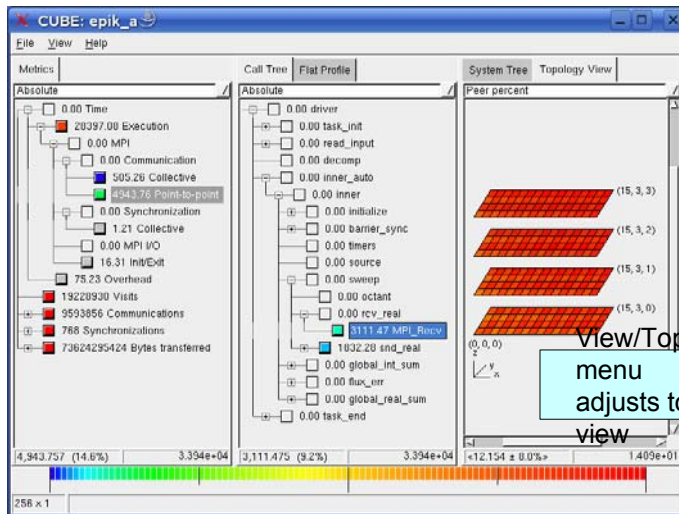
System tree dimension



ParaTools

332

Alternative: Topology display



ParaTools

333

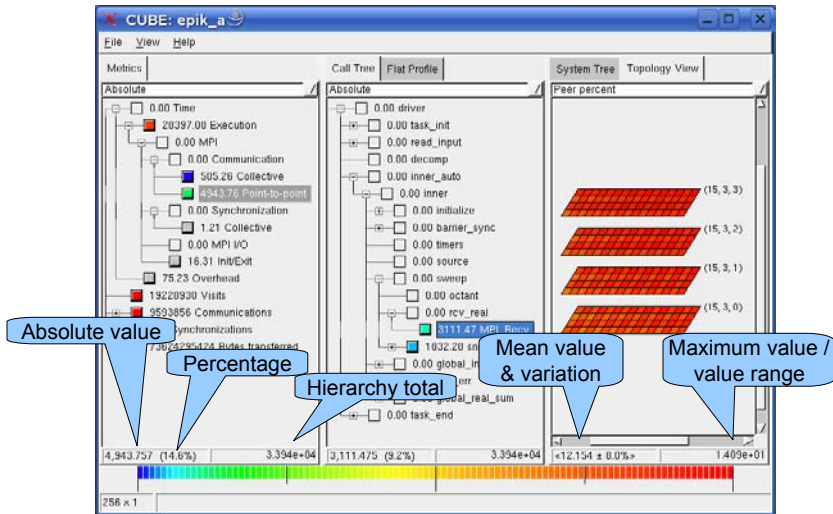
Topology display

- Topology information is recorded for
 - the hardware (supported on some systems)
 - MPI topologies (e.g., `MPI_Cart_create()`)
 - user-defined virtual topologies (under construction)
- Advantage
 - Better scalability than text-based system tree
- Restriction
 - Currently supports only 1D, 2D and 3D Cartesian topologies

ParaTools

334

Status fields



ParaTools

335

Display modes

- Absolute
 - Absolute values in seconds/number of occurrences
- Root percent
 - Percentage relative to the root node of the hierarchy
- External percent
 - Similar to “Root percent”, but relative to another data set
- Selection percent
 - Percentage relative to the node selected in the neighbouring column on the left

ParaTools

336

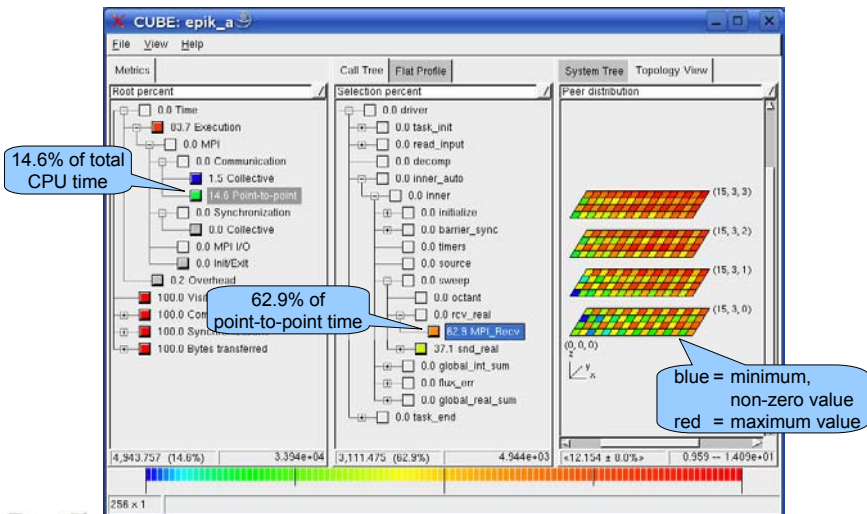
Display modes (system tree/topology only)

- Peer percent
 - Percentage relative to maximum of peer values (all values of the current leaf level)
- Peer distribution
 - Percentage relative to maximum and non-zero minimum of peer values

ParaTools

337

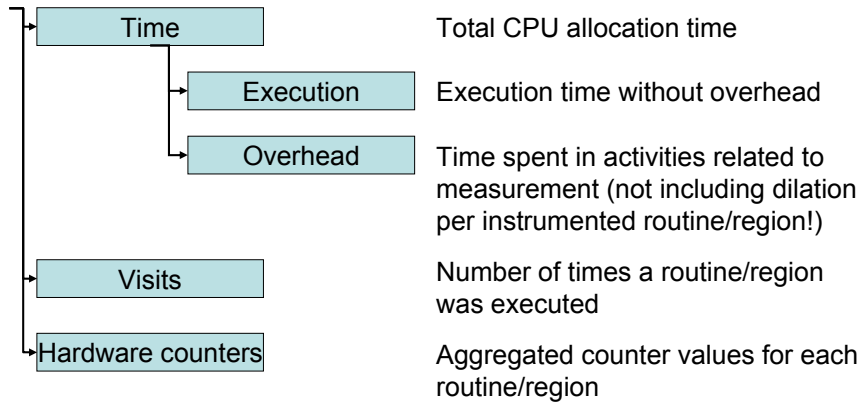
Display mode example



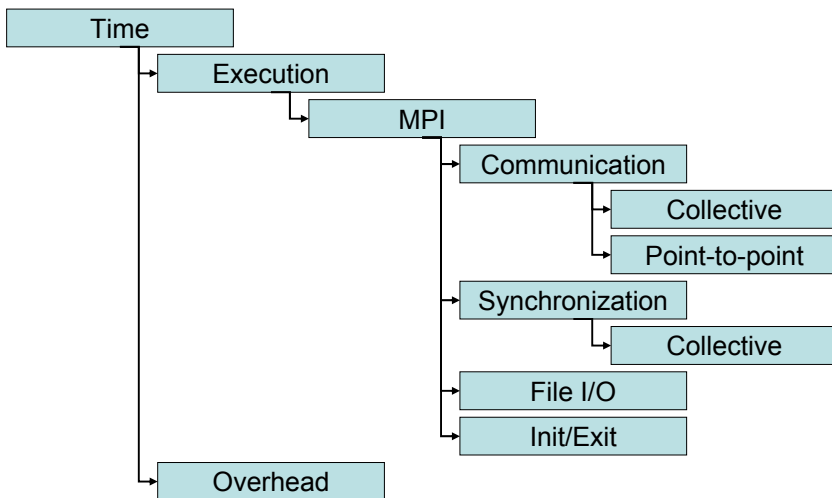
ParaTools

338

Generic metrics



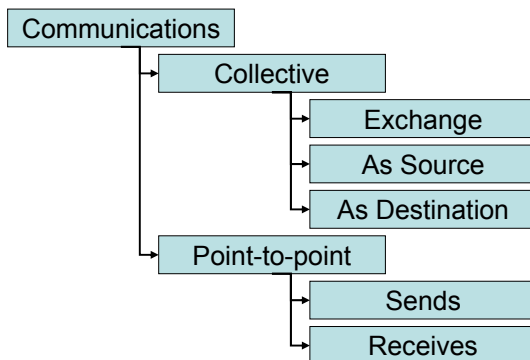
MPI Time hierarchy



MPI Time hierarchy (cont.)

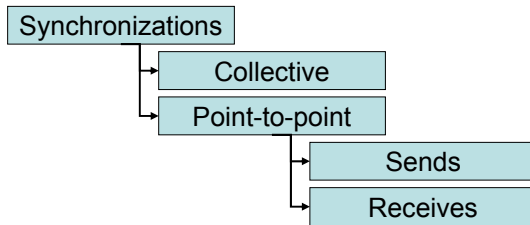
Time	Total CPU allocation time
Execution	Execution time without overhead
Overhead	Time spent in tasks related to measurement (not including dilation from instrumentation!)
MPI	Time spent in pre-instrumented MPI functions
Communication	Time spent in MPI communication calls, subdivided into collective and point-to-point
Synchronization	Time spent in MPI synchronization calls
File I/O	Time spent in MPI file I/O functions
Init/Exit	Time spent in MPI_Init() and MPI_Finalize()

MPI Communications hierarchy



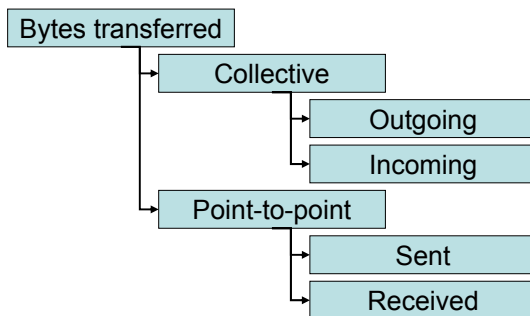
- Provides the number of calls to an MPI communication function of the corresponding class
- Zero-sized message transfers are considered *synchronization!*

MPI Synchronizations hierarchy



- Provides the number of calls to an MPI synchronization function of the corresponding class
- MPI synchronizations include zero-sized message transfers!

MPI Bytes transferred hierarchy



- Provides the number of bytes transferred by an MPI communication function of the corresponding class

Combined trace collection & analysis

- Modify jobscript
 - Use “scan -t” (or set EPK_TRACE=1)
 - Trace experiment EPK_TITLE set to
\$(TARGET)_\$(MODE)_\$(NP)_trace
 - Creates new experiment archive directory ./epik_\$(EPK_TITLE)
 - Trace unified & buffers flushed at measurement finalization
 - Automatic trace pattern analysis immediately follows
- Explore trace pattern analysis report using CUBE

Trace analysis output example

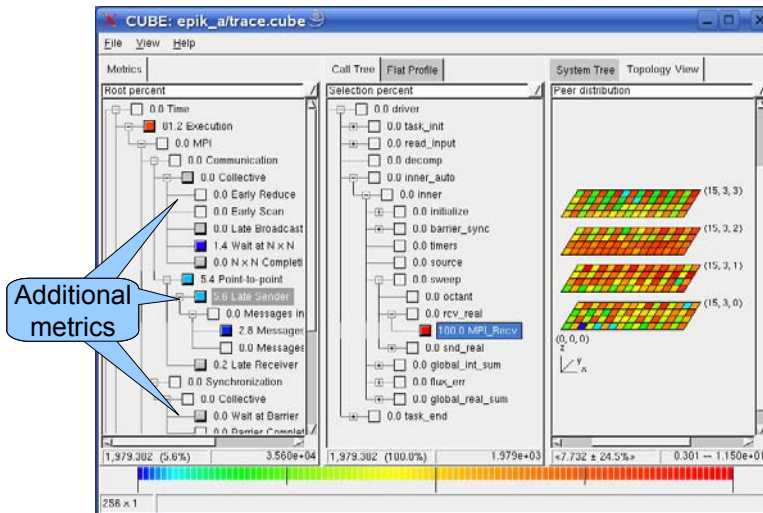
SCOUT

Analyzing experiment archive ./epik_sweep3d_co32_trace

```
Reading definition files ... done
Reading event trace files ... done
Preprocessing ... done
Analyzing event traces ... done
Writing report ... done
```

```
Total processing time: 4.083s
Total number of events: 5206596
Max. memory usage: 15.453 MB
```

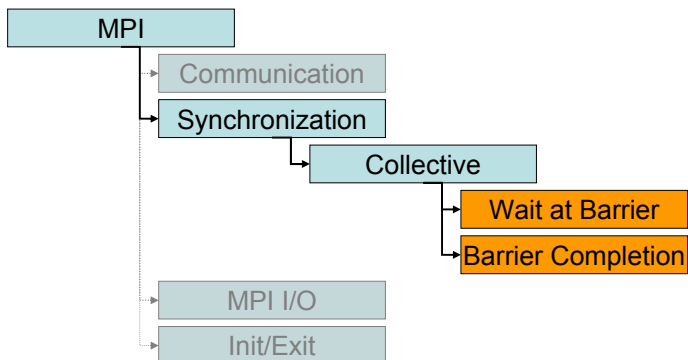
Trace analysis result



ParaTools

347

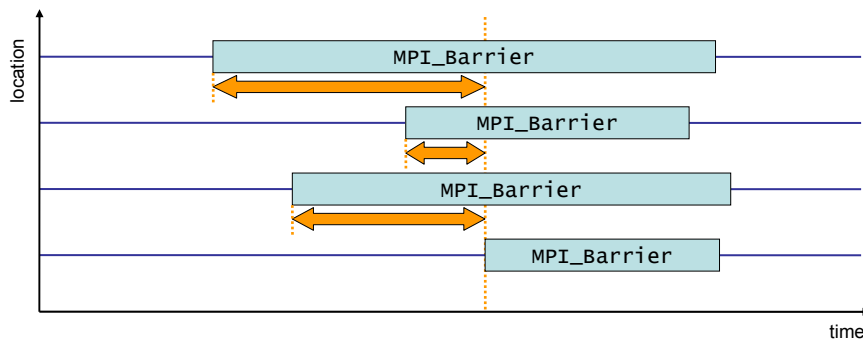
MPI collective synchronization time



ParaTools

348

Wait at Barrier = Early Barrier

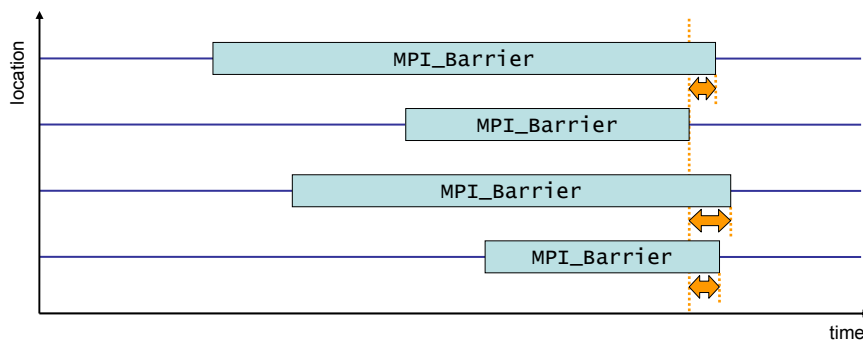


- Time spent waiting in front of a barrier call until the last process reaches the barrier operation
- Applies to: MPI_Barrier()

ParaTools

349

Barrier Completion

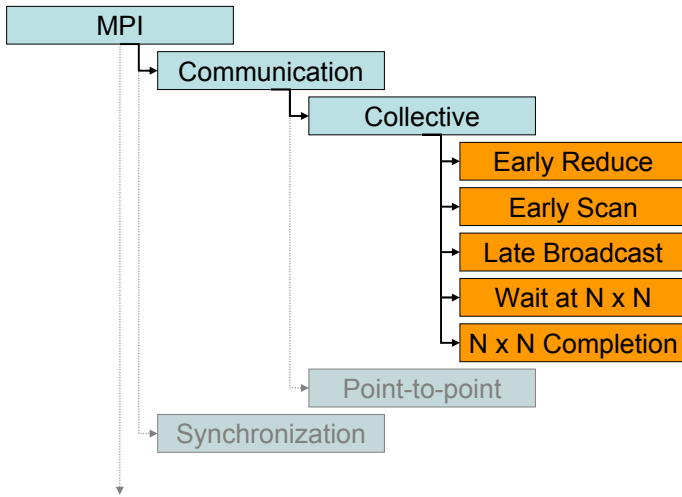


- Time spent in barrier after the first process has left the operation
- Applies to: MPI_Barrier()

ParaTools

350

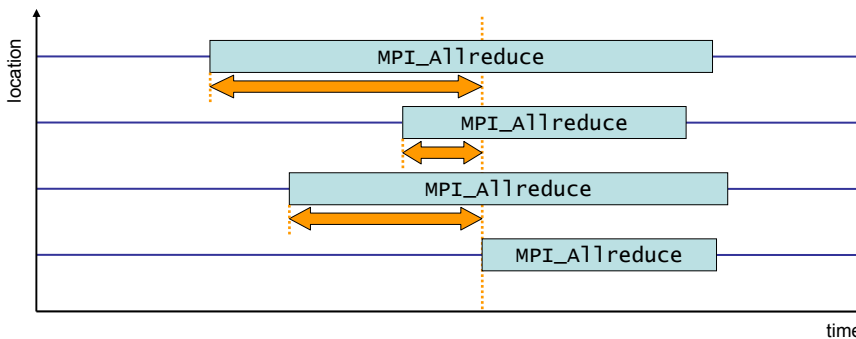
MPI collective communication time



ParaTools

351

Wait at N x N = Early N x N

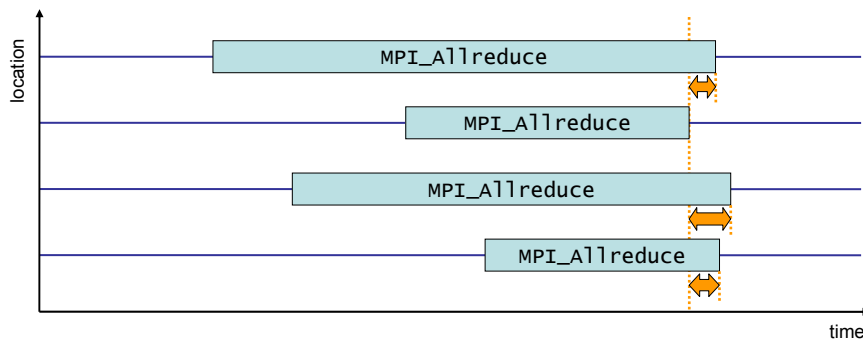


- Time spent waiting in front of a synchronizing collective operation call until the last process reaches the operation
- Applies to: `MPI_Allreduce()`, `MPI_Alltoall()`, `MPI_Alltoallv()`, `MPI_Allgather()`, `MPI_Allgatherv()`, `MPI_Reduce_scatter()`

ParaTools

352

N x N Completion

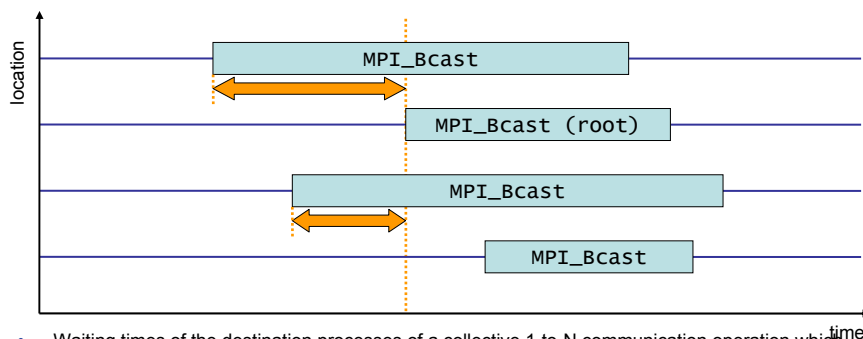


- Time spent in synchronizing collective operations after the first process has left the operation
- Applies to: `MPI_Allreduce()`, `MPI_Alltoall()`, `MPI_Alltoallv()`, `MPI_Allgather()`, `MPI_Allgatherv()`, `MPI_Reduce_scatter()`

ParaTools

353

Late Broadcast = Early Broadcast

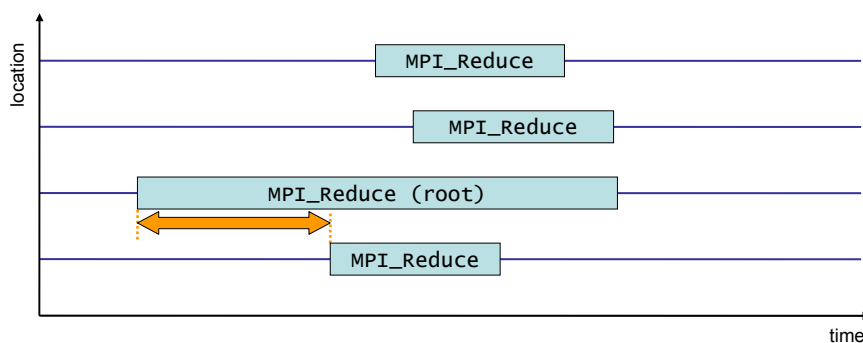


- Waiting times of the destination processes of a collective 1-to-N communication operation which enter the operation earlier than the source process (root)
 - Late Broadcast by source = Early Broadcast by destinations
- Applies to: `MPI_Bcast()`, `MPI_Scatter()`, `MPI_Scatterv()`

ParaTools

354

Early Reduce

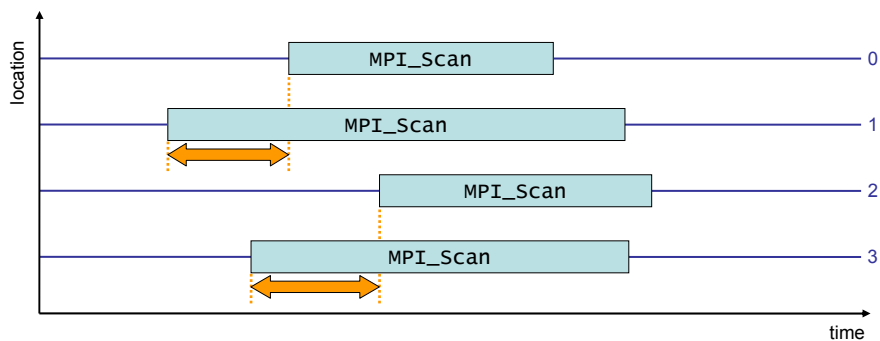


- Waiting time if the destination process (root) of a collective N-to-1 communication operation enters the operation earlier than its sending counterparts
- Applies to: MPI_Reduce(), MPI_Gather(), MPI_Gatherv()

ParaTools

355

Early Scan

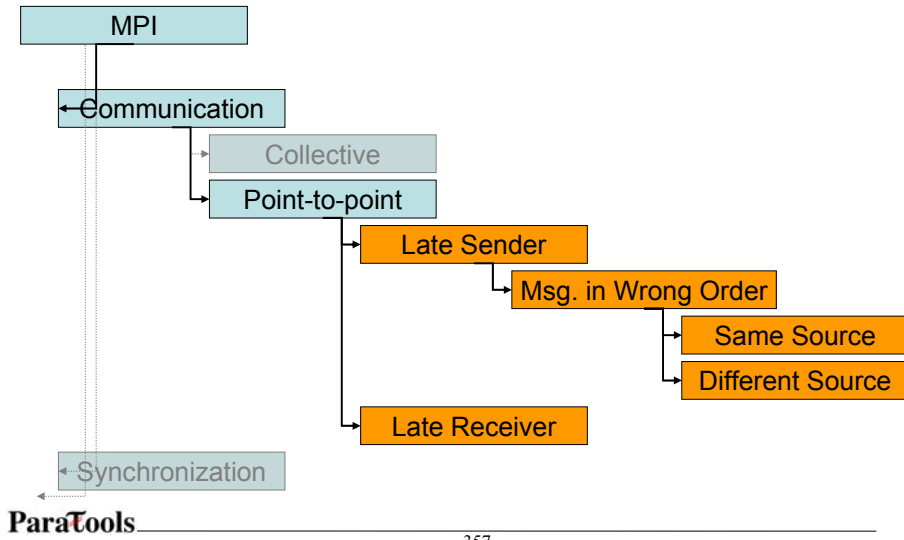


- Waiting time if process n enters a prefix reduction operation earlier than its sending counterparts (i.e., ranks $0..n-1$)
- Applies to: MPI_Scan()

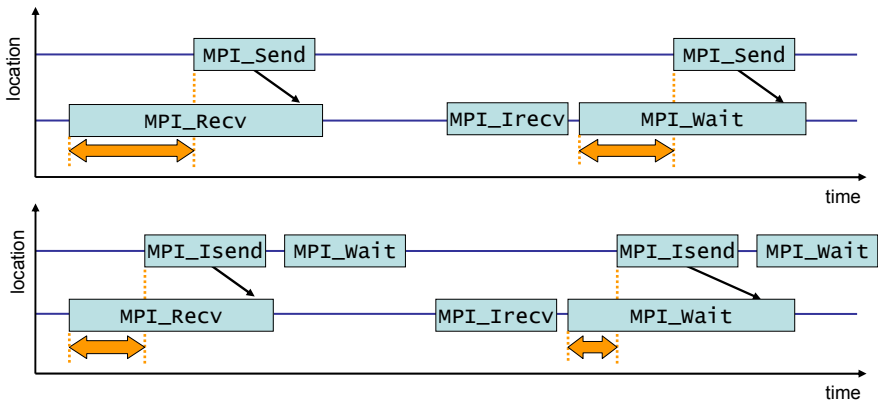
ParaTools

356

MPI point-to-point communication time

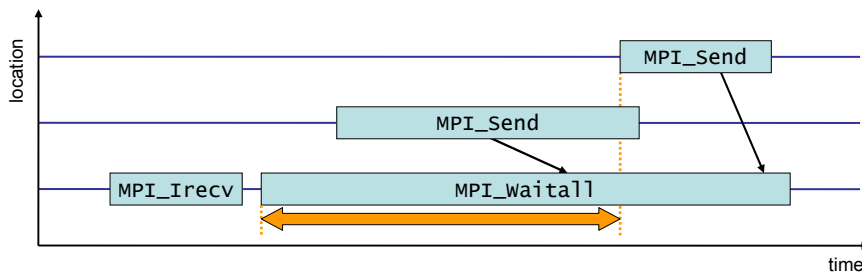


Late Sender = Early Receive



- Waiting time caused by a blocking receive operation posted earlier than the corresponding send operation
- Applies to blocking as well as non-blocking communication

Late Sender = Early Receive (cont.)

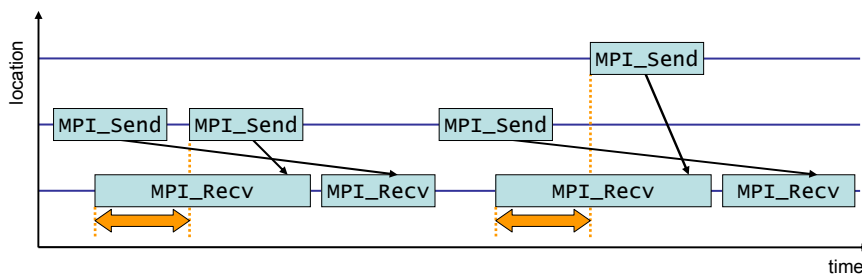


- While waiting for several messages, the maximum waiting time is accounted
- Applies to: MPI_waitall(), MPI_waitsome()

ParaTools

359

Late Sender, Messages in Wrong Order

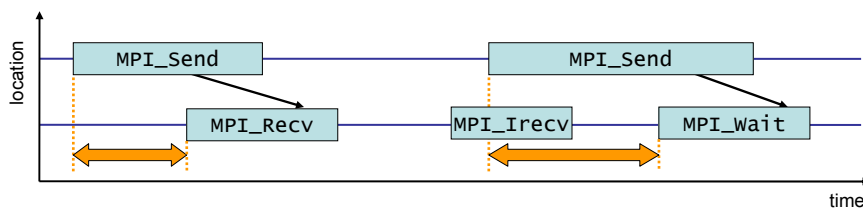


- Refers to Late Sender situations which are caused by messages received in wrong order
 - Early receive of message out of order
- Comes in two flavours:
 - Messages sent from same source location
 - Messages sent from different source locations

ParaTools

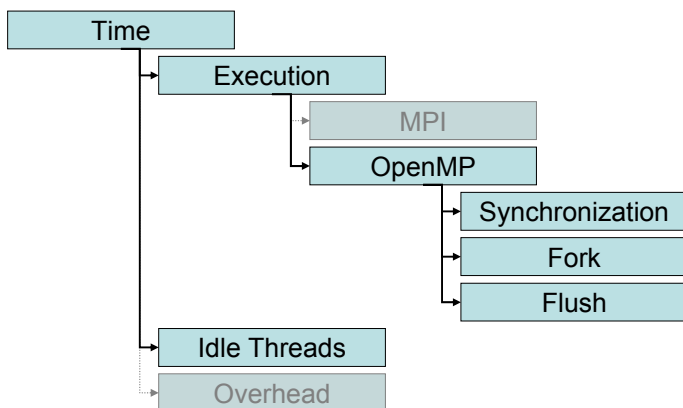
360

Late Receiver = Early Send



- Waiting time caused by a blocking send operation posted earlier than the corresponding receive operation
- Does not apply to non-blocking sends

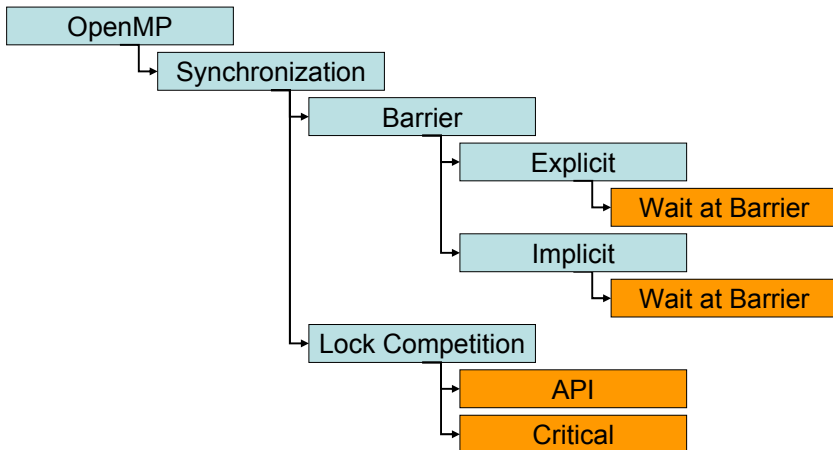
OpenMP Time hierarchy



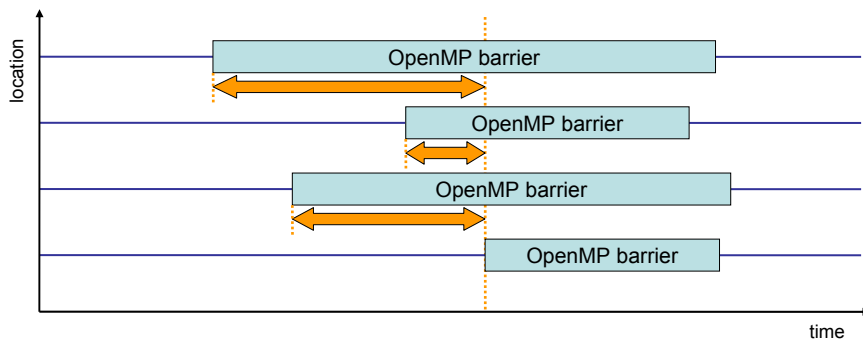
OpenMP Time hierarchy details

OpenMP	Time spent for all OpenMP-related tasks
Synchronization	Time spent synchronizing OpenMP threads
Fork	Time spent by master thread to create thread teams
Flush	Time spent in OpenMP flush directives
Idle Threads	Time spent idle on CPUs reserved for slave threads

OpenMP synchronization time



Wait at Barrier = Early Barrier

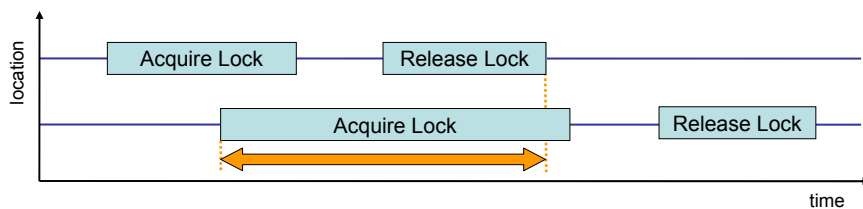


- Time threads spend waiting in front of a barrier call until the last thread reaches the barrier operation
- Applies to: Implicit/explicit barriers

ParaTools

365

Lock competition



- Time a thread spends waiting for a lock that is held by other threads until it is released and can be acquired by this thread
- Applies to: critical sections, OpenMP lock API

ParaTools

366

Other metrics

- LateReceivers/LateSenders
 - counts shown in hierarchies of Synchronizations & Communications below Sends & Receives respectively
- Computational Imbalance
 - load imbalance heuristic calculated as absolute difference from average exclusive execution time
- HWC metrics
 - shown as separate root metrics for each counter
 - only provided in summary reports

Intermediate use of Scalasca

- User-defined region instrumentation
 - EPIK annotation macros API
 - POMP annotation directives
 - Selective instrumentation
- Summary collections & analysis experiment
- Trace collection & analysis experiment
- Analysis report effectiveness score
- Customisation of measurement collection
 - Sizing of measurement data structures (e.g., trace buffers)
 - Function filter configuration
 - Optional HWC metrics
- Analysis report algebra

Instrumentation/measurement configuration

- Selective instrumentation
 - Adjust build not to (auto-)instrument particular modules
 - Separate/preprocess sources for functions in same module
 - Entirely avoids instrumentation & overhead
- Selective measurement via function filtering
 - Supported for GCC, IBM & Intel compilers
 - Specify text file listing names of functions (one per line, shell wildcarding) to ignore with EPK_FILTER
 - Use linker/decorated function names [Fortran/C++]

ParaTools

369

cube3_score with (sorted) region breakdown

```
% cube3_score -r epik_smg2000_mano_64/summary.cube | sort ...
flt type  max_tbc    time      % region
...
  MPI    2061936  293.30    23.89  MPI_Waitall
  COM    2346840   14.79     1.20  hypre_FinalizeCommunication
  COM    2346840   23.49     1.91  hypre_InitializeCommunication
  MPI    7495240   11.38     0.93  MPI_Irecv
  MPI    8149850   41.43     3.37  MPI_Isend
  USR    9426048   10.47     0.85  hypre_StructStencilElementRank
  USR    9426048   21.69     1.77  hypre_StructMatrixExtractPointerByIdx
  USR    11063016  16.80     1.37  hypre_MAlloc$AF10_5
  USR    11454432  25.82     2.10  hypre_MAlloc
  USR    11763336  26.90     2.19  hypre_CAlloc
  USR    23496576  38.16     3.11  hypre_Free

  ANY    162589938 1227.61  100.00  ALL (254 regions)
  MPI    17649090  456.64   37.20   + MPI ( 13 regions) pure MPI
  COM    9905832   321.80   26.21   + COM ( 32 regions) combined
MPI&USR
  USR    135034968 311.13  25.34   + USR (207 regions) pure User
```

max_tbc = est. maximum trace buffer capacity requirement (bytes/process)
to store all events that would be generated in an equivalent trace

ParaTools

370

cube3_score with trial region filter

```
% cube3_score -r -f smg2000.filt epik_smg2000_mano_64/summary.cube | sort
filt type max_tbc time % region
...
- MPI 2061936 293.30 23.89 MPI_Waitall
- COM 2346840 14.79 1.20 hypre_FinalizeCommunication
- COM 2346840 23.49 1.91 hypre_InitializeCommunication
- MPI 7495240 11.38 0.93 MPI_Irecv
- MPI 8149850 41.43 3.37 MPI_Isend
+ USR 9426048 10.47 0.85 hypre_StructStencilElementRank
+ USR 9426048 21.69 1.77 hypre_StructMatrixExtractPointerByIdx
+ USR 11063016 16.80 1.37 hypre_MAlloc$AF10_5
+ USR 11454432 25.82 2.10 hypre_MAlloc
+ USR 11763336 26.90 2.19 hypre_CAlloc
+ USR 23496576 38.16 3.11 hypre_Free

- ANY 162589938 1227.61 100.00 ALL (253 regions)
- MPI 17649090 456.64 37.20 + MPI (13 regions) pure MPI
- COM 9905832 321.80 26.21 + COM (32 regions) combined
MPI&USR
- USR 135034968 311.13 25.34 + USR (207 regions) pure User

+ FLT 103570824 182.11 14.83 FLT (9 regions) filtered
- FLT 59019114 1045.50 85.17 ALL-FLT (244 regions) remainder
```

ParaTools

371

Preparation of instrumented executable

- Auto-instrumentation of functions
 - Capability of most (but not all) compilers
 - Currently need separate Scalasca installations for each desired combination of MPI library & compiler suite
 - \$(PREP) \$(MPIFC) ...
 - \$(PREP) \$(MPICC) ...
 - \$(PREP) \$(MPICXX) ...
 - PREP="skin \$(SKIN_OPTS)" for instrumented build
 - PREP="" for uninstrumented build for production
- Auto-instrumentation plus API for user-defined regions
 - #include "epik_user.inc" or "epik_user.h"
 - % skin -user \$(MPIC) ...

Manual instrumentation options

- No instrumentation
 - \$(MPIC) [`kconfig -cflags`]
- MPI library instrumentation
 - \$(MPIC) [`kconfig -cflags`] `kconfig -libs`
- MPI library & EPIK user instrumentation
 - \$(MPIC) `kconfig -cflags` `kconfig -libs` **-DEPIK**
- `kconfig -cflags` is optional for source modules without explicit EPIK API
#include

EPIK instrumentation API dummy macros

- To use unmodified compile commands (without EPIK API include path) for sources with EPIK API calls, define dummy macros

```
#ifdef EPIK
#include "epik_user.inc" or "epik_user.h"
#else
#define EPIK_FUNC_REG(str)      /* undefined */
#define EPIK_FUNC_START()     /* undefined */
#define EPIK_FUNC_END()       /* undefined */
#define EPIK_USER_REG(id, str) /* undefined */
#define EPIK_USER_START(id)   /* undefined */
#define EPIK_USER_END(id)     /* undefined */
#endif
```

EPIK instrumentation API

- Manual phase annotation
 - EPIK_FUNC_REG("Fortran function/subroutine")
 - EPIK_FUNC_START()
 - EPIK_USER_REG(tloop, "<<time step>>")
 - EPIK_USER_START(tloop)
 - EPIK_USER_END(tloop)
 - EPIK_FUNC_END()
- Note matching of enter/start annotations
 - all possible exits must be annotated
 - regions must be correctly nested
 - C/C++ function names are automatically registered
 - Fortran function/routine names undefined if not preregistered

POMP instrumentation

- Uses pragma/comment directives to annotate regions

C/C++:	Fortran:
<code>#pragma pomp inst init</code>	<code>!POMP\$ INST INIT</code>
<code>#pragma pomp inst begin(tloop)</code>	<code>!POMP\$ INST BEGIN(tloop)</code>
<code> #pragma pomp inst altend(tloop)</code>	<code>!POMP\$ INST</code>
<code> ALTEND(tloop)</code>	
<code>#pragma pomp inst end(tloop)</code>	<code>!POMP\$ INST END(tloop)</code>
- Directives ignored unless activated with ***skin -pomp***
 - ***all directives in module instrumented***
- Current limitations
 - instrumentation inactive until “inst init”
 - no distinction of functions from other regions
 - last region exit must be marked “end”, all others as “altend”
 - doesn't support C99 `_Pragma` operator

Measurement configuration

- Example configuration
 - EPK_GDIR=/work/\$USER # archive location
 - EPK_TITLE=app_\$NP # experiment archive title
 - EPK_SUMMARY=1 # runtime summarisation
 - EPK_TRACE=0 # event trace collection
- New archive directory for each experiment
 - \$EPK_GDIR/epik_\$EPK_TITLE
 - contains intermediate data (e.g., trace files), log/config files and processed analyses
- Configured automatically (overridden) by *scan* args

Default EPIK.CONF configuration file extract

```
# E P I K configuration
- EPK_TITLE=a      # experiment archive title [scan -e]
- EPK_SUMMARY=1 # runtime summarization [scan -s]
- EPK_TRACE=0     # event trace collection [scan -t]
- EPK_FILTER=     # file listing functions to skip
- EPK_METRICS=    # colon-separated list of metrics [-m]

# E P I S O D E configuration
- ESD_PATHS=1024 # max. recorded call-paths
- ESD_FRAMES=32 # max. call-stack frames
- ESD_BUFFER_SIZE=100000 # definitions bytes

# E P I L O G configuration
- ELG_BUFFER_SIZE=10000000 # trace bytes
```

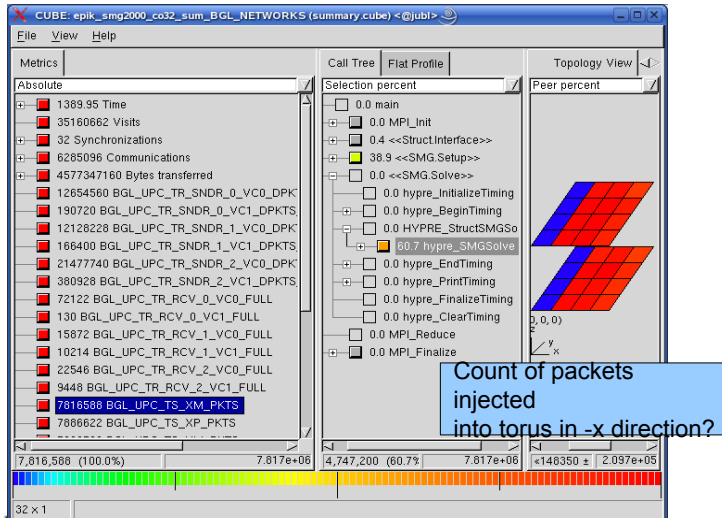
Hardware counter metrics

- Available counters (and their interpretation) are platform/processor-specific
 - considered separate root metrics in analyses
- Platform metrics specification
 - defines convenient groups of metrics
 - EPK_METRICS_SPEC=./METRICS.SPEC
- Group/list of counters to measure in experiment
 - EPK_METRICS=POWER4_DC # data-cache
 - EPK_METRICS=BGL_NETWORKS # torus & tree
 - EPK_METRICS=PM_CYC:PM_INST_CMPL
or PAPI_TOT_CYC:PAPI_TOT_INS

ParaTools

379

Scalasca summary experiment with HWC metrics



ParaTools

380

CUBE algebra tools

- CUBE files can be compared/combined with some useful command line tools
- Note that these work directly on CUBE files and not on archive directories
 - Reads CUBE2 & CUBE3 files, but only writes CUBE3 files
- General usage:
 - `cube3_tool [-o <output file>] <input file>`
- If no output file is specified, `tool.cube` is generated

CUBE algebra tools (2)

- `cube3_merge`
 - combines multiple analysis reports into integrated report
 - merges metric, call-path & system trees
 - takes metric severities from first available report
 - e.g., combine measurements of sets of HWC metrics in summary report(s) with a (non-HWC) trace analysis report into a “holistic” analysis report
- `% cube3_merge trace.cube summary_HWC[1234].cube`
 - Metrics listed in order of appearance in input reports
 - User-defined hierarchies of measured & derived HWC metrics not yet supported by CUBE3!

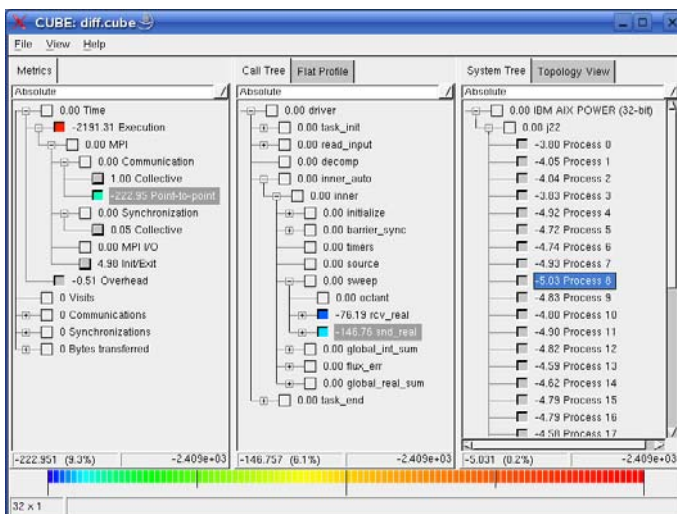
CUBE algebra tools (3)

- `cube3_mean`
 - Can eliminate “measurement noise” by averaging the results of several experiments
- `cube3_cut [-p prune] [-r root]`
 - Creates a new CUBE file without pruned subtrees and/or containing only the specified call tree node as new root(s)
- `cube3_diff`
 - Calculates the difference of two experiments
 - Useful to measure improvement/degradation due to a modification

ParaTools

383

Difference experiment: JUMP – JUBL (different architectures)



ParaTools

384

Labs!



Lab: PAPI, TAU, Vampir, and Scalasca/KOJAK

ParaTools

385

Lab Instructions (for LLNL systems)

Get `workshop.tar.gz` using:

```
% wget
```

```
  http://www.paratools.com/llnl09/workshop.tar.gz
```

Or

```
% cp
```

```
  /usr/global/tools/tau/training/src/workshop.tar.gz  
  .
```

```
% tar xzf workshop.tar.gz
```

```
source /usr/global/tools/tau/training/src/tau.cshrc
```

OR

```
source /usr/global/tools/tau/training/src/tau.bashrc
```

in your `.login` file and then follow the instructions
in the `README` file.

ParaTools

386

Lab Instructions

To profile a code using TAU:

1. Change the compiler name to `tau_cxx.sh`, `tau_f90.sh`, `tau_cc.sh`:
`F90 = tau_f90.sh`
2. Choose TAU stub makefile
`% setenv TAU_MAKEFILE`
`/usr/global/tools/tau/training/tau-2.18.2/bgp/lib/Makefile.tau-[options]`
3. If stub makefile has `-papi` in its name, set `COUNTER[1-<n>]` environment variables:
`% setenv COUNTER1 GET TIME OF DAY`
`% setenv COUNTER2 PAPI_L2_DCM`
`% setenv COUNTER3 PAPI_TOT_CYC ... OR`
`% setenv TAU_METRICS TIME:PAPI_L2_DCM:PAPI_TOT_CYC`
4. Build and run workshop examples, then run `pprof/paraprof`

ParaTools

387

Support Acknowledgements

- Department of Energy (DOE)
 - Office of Science contracts
 - LLNL-LANL-SNL ASC/NNSA Level 3 contract
 - LLNL ParaTools contracts
- Department of Defense (DoD)
 - PET
- National Science Foundation (NSF)
 - POINT
- University of Oregon
 - A. Malony, A. Morris, K. Huck, W. Spear, S. Biersdorff, A. Nataraj
- University of Tennessee, Knoxville
 - Dr. David Cronk and Dr. Shirley Moore
- T.U. Dresden, GWT
 - Dr. Wolfgang Nagel and Dr. Holger Brunst
- Research Centre Juelich
 - Dr. Bernd Mohr, Dr. Felix Wolf



ParaTools

388