TAU PERFORMANCE SYSTEM

Sameer Shende Alan Morris, Wyatt Spear, Scott Biersdorff Performance Research Lab

POINT

Allen D. Malony, Kevin Huck, Aroon Nataraj Department of Computer and Information Science University of Oregon





TAU Performance System

- <u>T</u>uning and <u>A</u>nalysis <u>U</u>tilities (16+ year project)
- Performance problem solving framework for HPC
 - Integrated, scalable, flexible, portable
 - Target all parallel programming / execution paradigms
- Integrated performance toolkit (open source)
 - Instrumentation, measurement, analysis, visualization
 - Widely-ported performance profiling / tracing

SC '09: Productive Performance Engineering of Petascale Applications with POINT and VI-HPS

Performance data management and data mining

TAU Performance System



Building Bridges to Other Tools





Direct Performance Observation

- Execution actions of interest 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 performance system to observe
- Observation is direct
 - Direct instrumentation of program (system) code (probes)
 - Instrumentation invokes performance measurement
 - Event measurement: performance data, meta-data, context
- Performance experiment
 - Actual events + performance measurements
- Contrast with (indirect) event-based sampling





TAU Instrumentation Approach

- Support for standard program events
 - Routines, classes and templates
 - Statement-level blocks
 - Begin/End events (Interval events)
- Support for user-defined events
 - Begin/End events specified by user
 - Atomic events (e.g., size of memory allocated/freed)
 - Flexible selection of event statistics
- Provides static events and dynamic events
- Enables "semantic" mapping
- Specification of event groups (aggregation, selection)
- Instrumentation optimization





TAU Event Interface

- Events have a type, a group association, and a name
- TAU events names are character strings
 - Powerful way to encode event information
 - Inefficient way to communicate each event occurrence
- TAU maps a new event name to an event ID
 - Done when event is first encountered (get event handle)
 - Event ID is used for subsequent event occurrences
 - Assigning a uniform event ID a priori is problematic
- A new event is identified by a new event name in TAU
 - Can create new event names at runtime

Allows for dynamic events (TAU renames events)
 PO Allows for context-base, parameter-based, phase

TAU Instrumentation Mechanisms

- Source code
 - Manual (TAU API, TAU component API)
 - Automatic (robust)
 - C, C++, F77/90/95 (Program Database Toolkit (PDT))
 - OpenMP (directive rewriting (Opari), POMP2 spec)
 - Library header wrapping
- Object code

- Pre-instrumented libraries (e.g., MPI using PMPI)
- Statically- and dynamically-linked (with LD_PRELOAD)
- Executable code
 - Binary and dynamic instrumentation (Dyninst)
 - Virtual machine instrumentation (e.g., Java using JVMPI)
- TAU_COMPILER to automate instrumentation process



Automatic Source-level Instrumentation



SC '09: Productive Performance Engineering of Petascale Applications with POINT and VI-HPS





MPI Wrapper Interposition Library

- Uses standard MPI Profiling Interface
 - Provides name shifted interface
 - MPI_Send = PMPI_Send
 - Weak bindings
- Create TAU instrumented MPI library
 - Interpose between MPI and TAU
 - Done during program link
 - Impi replaced by –ITauMpi –Ipmpi –Impi
 - No change to the source code!
 - Just re-link application to generate performance data





MPI Shared Library Instrumentation

- Interpose the MPI wrapper library for applications that have already been compiled – Avoid re-compilation or re-linking
- Requires shared library MPI
 - Uses LD_PRELOAD for Linux
 - On AIX use MPI_EUILIB / MPI_EUILIBPATH
 - Does not work on XT3
- Approach will work with other shared libraries
- Use TAU tauex

POINT

- % mpirun -np 4 tauex a.out



Selective Instrumentation File

- Specify a list of events to exclude or include
- # is a wildcard in a routine name **BEGIN EXCLUDE LIST** Foo Bar D#EMM END EXCLUDE LIST **BEGIN INCLUDE LIST** int main(int, char $*\overline{}$) **F1 F**3 END INCLUDE LIST



Selective Instrumentation File

- Optionally specify a list of files
- * and ? may be used as wildcard characters BEGIN_FILE_EXCLUDE_LIST f*.f90
 - Foo?.cpp
 - END_FILE_EXCLUDE_LIST BEGIN FILE INCLUDE LIST
 - main.cpp
 - foo.f90

END_FILE_INCLUDE_LIST



SC '09: Productive Performance Engineering of Petascale Applications with POINT and VI-HPS



Selective Instrumentation File

- User instrumentation commands
 - Placed in INSTRUMENT section
 - 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
```





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 OTF, EPILOG, Paraver, and SLOG2
 - Trace streams (OTF) and hierarchical trace merging
- Robust timing and hardware performance support
- Multiple counters (hardware, user-defined, system)

• Performance measurement of I/O and Linux kernel SC '09: Productive Performance Engineering of Petascale Applications with POINT and VI-HPS 17

TAU Measurement Mechanisms

- Parallel profiling
 - Function-level, block-level, statement-level
 - Supports user-defined events and mapping events
 - Support for flat, callgraph/callpath, phase profiling
 - Support for parameter and context profiling
 - Support for tracking I/O and memory (library wrappers)
 - Parallel profile stored (dumped, shapshot) during execution
- Tracing
 - All profile-level events
 - Inter-process communication events
 - Inclusion of multiple counter data in traced events



Types of Parallel Performance Profiling

- Flat profiles
 - Metric (e.g., time) spent in an event (callgraph nodes)
 - Exclusive/inclusive, # of calls, child calls
- Callpath profiles (Calldepth profiles)
 - Time spent along a calling path (edges in callgraph)
 - "main=> f1 => f2 => MPI Send" (event name)
 - TAU CALLPATH DEPTH environment variable
- Phase profiles
 - Flat profiles under a phase (nested phases are allowed)
 - Default "main" phase

- Supports static or dynamic (per-iteration) p



Performance Analysis

- Analysis of parallel profile and trace measurement
- Parallel profile analysis (ParaProf)
 - Java-based analysis and visualization tool
 - Support for large-scale parallel profiles
- Performance data management framework (PerfDMF)
- Parallel trace analysis
 - Translation to VTF (V3.0), EPILOG, OTF formats
 - Integration with Vampir / Vampir Server (TU Dresden)
 - Profile generation from trace data
- Online parallel analysis and visualization





ParaProf Profile Analysis Framework



Performance Data Management

- Provide an open, flexible framework to support common data management tasks
 - Foster multi-experiment performance evaluation
- Extensible toolkit to promote integration and reuse across available performance tools (PerfDMF)
 - Originally designed to address critical TAU requirements
 - Supported profile formats: TAU, CUBE (Scalasca), HPC Toolkit (Rice), HPM Toolkit (IBM), gprof, mpiP, psrun (PerfSuite), Open|SpeedShop, ...
 - Supported DBMS: PostgreSQL, MySQL, Oracle, DB2, Derby/Cloudscape
 - Profile query and analysis API
- Reference implementation for PERI-DB project



SC '09: Productive Performance Engineering of Petascale Applications with POINT and VI-HPS



PerfDMF Architecture



Metadata Collection

- Integration of XML metadata for each parallel profile
- Three ways to incorporate metadata
 - Measured hardware/system information (TAU, PERI-DB)
 - CPU speed, memory in GB, MPI node IDs, ...
 - Application instrumentation (application-specific)
 - TAU_METADATA() used to insert any name/value pair
 - Application parameters, input data, domain decomposition
 - PerfDMF data management tools can incorporate an XML file of additional metadata
 - Compiler flags, submission scripts, input files, ...
- Metadata can be imported from / exported to SC '09: Productive Performance Engineering of Petascale Applications with POINT and VI-HPS 25

Performance Data Mining / Analytics

- Conduct systematic and scalable analysis process
 - Multi-experiment performance analysis
 - Support automation, collaboration, and reuse
- Performance knowledge discovery framework
 - Data mining analysis applied to parallel performance data
 - comparative, clustering, correlation, dimension reduction, ...
 - Use the existing TAU infrastructure
- PerfExplorer v1 performance data mining framework
 - Multiple experiments and parametric studies
 - Integrate available statistics and data mining packages

• Weka, R. Matlab / Octave PO Apply data mining operations in interactive envi



How to explain performance?

- Should not just redescribe the 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
- Need to add knowledge to do more intelligent things
 - Automated analysis needs good informed feedback
 - iterative tuning, performance regression testing
 - Performance model generation requires interpretation
- · We need better methods and tools for
 - Integrating meta-information

POKowledge-based performance problems solving



PerfExplorer v2 – Requirements

- Component-based analysis process
 - Analysis operations implemented as modules
 - Linked together in analysis process and workflow
- Scripting
 - Provides process/workflow development and automation
- Metadata input, management, and access
- Inference engine
 - Reasoning about causes of performance phenomena
 - Analysis knowledge captured in expert rules
- Persistence of intermediate analysis results
- Provenance

POINT

- Provides historical record of analysis results





PerfExplorer v2 Architecture







Parallel Profile Analysis – pprof

📥 emacs@net	utron.cs.uor	egon.edu					•
Buffers Fil	es Tools.	Edit Search	Mule Help				
Reading Pro	ofile file	s in profile.*	:				
NODE O;CONT	TEXT O;THR	EAD O:					
%Time E×	clusive msec	Inclusive total msec	#Call	#Subrs	Inclusive usec/call	Name	
$\begin{bmatrix}$	$\begin{array}{c} 1\\ 3,667\\ 491\\ 6,461\\ .:18.436\\ 6,778\\ 50,142\\ 24,451\\ 7,501\\ 838\\ 6,590\\ 4,989\\ 0.44\\ 398\\ 140\\ 131\\ 89\\ 0.966\\ 24\\ 15\\ 4\\ 7\\ 3\\ 1\\ 0.966\\ 24\\ 15\\ 4\\ 7\\ 3\\ 1\\ 0.116\\ 0.512\\ 0.121\\ 0.024\\ 0.0$	3:11.293 3:10.463 2:08.326 1:25.159 1:18.436 56,407 50,142 31,031 7,501 6,594 6,590 4,989 400 399 247 131 103 96 95 44 15 12 8 3 0.837 0.512 0.353 0.191 0.433	$\begin{array}{c} 1\\ 3\\ 37200\\ 9300\\ 18600\\ 9300\\ 19204\\ 301\\ 9300\\ 604\\ 9300\\ 608\\ 1\\ 1\\ 57252\\ 1\\ 1\\ 57252\\ 1\\ 1\\ 57252\\ 1\\ 1\\ 57252\\ 1\\ 1\\ 1\\ 57252\\ 1\\ 1\\ 1\\ 608\\ 1\\ 1\\ 1\\ 1\\ 1\\ 1\\ 1\\ 1\\ 1\\ 1\\ 1\\ 1\\ 1\\$	15 37517 37200 18600 0 18600 0 602 0 1812 0 4 39 47616 0 2 2 0 7937 0 5 1700 5 1700 3 0 6 0 4 0 2 2 0 7937 0 5 1700 3 0 6 0 4 0 2 2 0 7937 0 5 1700 1800 0 1812 0 0 4 0 1812 0 0 0 1812 0 0 1812 0 0 1812 0 0 1812 0 0 0 1812 0 0 0 1812 0 0 0 1812 0 0 0 0 1812 0 0 0 0 1812 0 0 0 0 0 0 0 0 0 1812 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0	191293269 63487925 3450 9157 4217 6065 2611 103096 807 10918 709 8206 400081 399634 247086 2 103168 96458 10603 44878 400 15630 12335 2893 491 3874 1007 837 512 353 191	applu bcast_inputs exchange_1 buts MPI_Recv() blts MPI_Send() rhs jacld exchange_3 jacu MPI_Wait() init_comm MPI_Init() setiv exact erhs read_input MPI_Ecast() error MPI_Finalize() setbv 12norm MPI_Allreduce() pintgr MPI_Barrier() exchange_4 MPI_Keyval_create(exchange_6 MPI_Ima_continuou)

SC '09: Productive Performance Engineering of Petascale Applications with POINT and VI-HPS





Metadata for Each Experiment

🚺 TAU: ParaProf Manager

File Options Help

- Applications
- 🔶 🗐 Standard Applications
 - 🔶 🗐 Default App
 - ← □ Default Exp
 - f90/pdt_mpi/examples/tau2/amorris/home/
 - PAPI_FP_OPS

Ш

- GET_TIME_OF_DAY
- 👇 🔚 Default (jdbc:postgresql://spaceghost.cs.uoregon.edu:5
- 🔚 utonium (jdbc:postgresql://utonium.cs.uoregon.edu:543 🔶 🗂 spaceghost2 (jdbc:postgresql://spaceghost.cs.uoregon.
- 🗠 🔚 proton_mysql (jdbc:mysql://192.168.1.1:3306/perfdm
- 🗠 🔚 spaceghost_peri_milc (jdbc:postgresql://spaceghost.cs.u
- 🖕 🗂 proton_postgresgl (jdbc:postgresgl://192.168.1.1:543)
- 🔶 🗂 utonium_oracle (jdbc:oracle:thin:@//utonium.cs.uoregon. Image: period (idbc:postgresql://spaceghost.cs.uoregon.edu:5)

Multiple PerfDMF DBs

	TrialField	Value		
	Name	f90/pdt_mpi/examples/tau2/amorris/home/		
	Application ID	0		
	Experiment ID	0		
amples /tou? / amerris /home/	Trial ID	0		
amples/tauz/amorns/nome/	CPU Cores	2		
DAY	CPU MHz	2992.505		
spacedbast cs upredon edu:5	CPU Type	Intel(R) Xeon(R) CPU 5160 @ 3.00GHz		
spacegnost.cs.doregon.edu.5	CPU Vendor	GenuineIntel		
/utonium.cs.uoregon.edu.543	CWD	/home/amorris/tau2/examples/pdt_mpi/f90		
esql://spaceghost.cs.uoregon.	Cache Size	4096 KB		
//192.168.1.1:3306/perfdm	Executable	/home/amorris/tau2/examples/pdt_mpi/f		
c:postgresql://spaceghost.cs.(Hostname	demon.nic.uoregon.edu		
ostgresql://192.168.1.1:543	Local Time	2007-07-04T04:21:14-07:00		
e:thin:@//utonium.cs.uoregon	MPI Processor Name	demon.nic.uoregon.edu		
spacedhost os upredon edu:5	Memory Size	8161240 kB		
spacegnost.es.aoregon.eaa.s	Node Name	demon.nic.uoregon.edu		
	OS Machine	x86_64		
	OS Name	Linux		
	OS Release	2.6.9-42.0.3.EL.perfctrsmp		
	OS Version	#1 SMP Fri Nov 3 07:34:13 PST 2006		
	Starting Timestamp	1183548072220996		
	TAU Architecture	x86_64		
	TAU Config	-papi=/usr/local/packages/papi-3.5.0 -M		
	Timestamp	1183548074317538		
	UTC Time	2007-07-04T11:21:14Z		
	pid	11395		
00 100 Decidentics Decision 5	username	amorris		
SC US: Productive Performance Engineer	ring of Petascale Applications with	POINT and VI-HPS		
	33			

• 🗆 🗙



ParaProf – Stacked View

O O File Options Windows Help

Metric Name: Time Value Type: exclusive n,c,t 6488,U,U n,c,t 6489,0,0 n,c,t 6490,0,0 n,c,t 6491,0,0 n,c,t 6492,0,0 n,c,t 6493,0,0 n,c,t 6494,0,0 n,c,t 6495,0,0 n,c,t 6496,0,0 n,c,t 6497,0,0 n,c,t 6498,0,0 n,c,t 6499,0,0 n,c,t 6500,0,0 n,c,t 6501,0,0 n,c,t 6502,0,0 n,c,t 6503,0,0 n.c.t 6504.0.0 n.c.t 6505.0.0 n.c.t 6506.0.0 n,c,t 6507,0,0 MPI_Alltoall() n,c,t 6508,0,0 n,c,t 6509,0,0 n,c,t 6510,0,0 n,c,t 6511,0,0 n,c,t 6512,0,0 n,c,t 6513,0,0 n,c,t 6514,0,0 n,c,t 6515,0,0 n,c,t 6516,0,0 n,c,t 6517,0,0 n,c,t 6518,0,0 n,c,t 6519,0,0 n,c,t 6520,0,0 n,c,t 6521,0,0 n,c,t 6522,0,0 n,c,t 6523,0,0 n,c,t 6524,0,0 n,c,t 6525,0,0 n,c,t 6526,0,0 SC '09 Productive Performance Engineering of Pel n,c,t 6527,0,0 Applications n,c,t 6528,0,0 n,c,t 6529,0,0 351

ParaProf – Callpath Profile

000

File Options Windows Help

X n,c,t, 0,0,0 - callpath-all/scaling/flash/taudata/disk2/mnt/

Metric Name: Time Value Type: exclusive	
75.474%	
26.474%	
24.556%	MODULEHYDRO 1D:HYDRO 1D
24.556%	FLASH => EVOLVE => HYDRO::HYDRO_3D => MODULEHYDROSWEEP::HYDRO_SWEEP => MODULEHYDRO_1D::HYDRO_1D
14.351%	MODULEINTRFC::INTRFC
14.351%	FLASH => EVOLVE => HYDRO::HYDRO_3D => MODULEHYDROSWEEP::HYDRO_SWEEP => MODULEHYDRO_1D::HYDRO_1D => MODULEINTRFC::INT
4.	501% MODULEEOS3D::EOS3D
4	427% MPL_Ssendo
	3.678% FLASH => EVOLVE => HYDRO::HYDRO_3D => MODULEHYDROSWEEP::HYDRO_SWEEP => MODULEEOS3D::EOS3D
	3.536% MPLAIreduce0
	2.727% MODIFIED ATE SOLNEIDD ATE SOLN
	2.242% MODULEUPDATE_SOLIN.JEPDATE_SOLIN. 2.242% FLASE = > KONIVE = > HYDRO 3D => MODULEHYDROSWEEPHYDRO SWEEP => MODULEUPDATE SOLNHIPDATE SOLN
	1.703% IFLASH => EVOLVE => HYDRO::HYDRO: 3D => MODULEHYDROSWEEP::HYDRO_SWEEP => MESH_GUARDCELL => AMR_GUARDCELL_SRL => AMR_
	1.56% 📕 FLASH => EVOLVE => HYDRO::HYDRO_SD => MODULEHYDROSWEEP::HYDRO_SWEEP => MESH_GUARDCELL => AMR_GUARDCELL_SRL => AMR_
	1.406% 🗖 FLASH => EVOLVE => MESH_UPDATE_GRIQ_REFINEMENT => MESH_REFINE_DEREFINE => AMR_REFINE_DEREFINE => AMR_MORTON_ORDER => A
	1.361% FLASH => TIMESTEP => MPL_Alireduce¢
	1.319% AMR_RESTRICT_UNK_FUN
	1.272% AMR_PROLONG_GEN_UNK_FUN
	1.093% LASH => EVOLVE => HTDRO::HTDRO:3D => MUDULEHTDROSWEEP::HTDRO:SWEEP => MESH_GUARDUELL => AMR_GUARDUELL_UIU_F => A
<u> </u>	1.077% ABUNDANCE_RESTRICT
Flach	1.064% DRASTRFE: DRASTNFIGHRORRI OCKLIST
1 10311	1% FLASH => EVOLVE => HYDRO::HYDRO: 3D => MODULEHYDROSWEEP::HYDRO.SWEEP => MESH_GUARDCELL => AMR_RESTRICT => AMR_RESTRICT
O thermonuclear	0.987% FLASH => EVOLVE => HYDRO::HYDRO_3D => MODULEHYDROSWEEP::HYDRO_SWEEP => MESH_FLUX_CONSERVE => AMR_FLUX_CONSERVE_UDT
	0.96% 🗌 FLASH => EVOLVE => HYDRO::HYDRO_3D => MODULEHYDROSWEEP::HYDRO_SWEEP => MESH_GUARDCELL => AMR_GUARDCELL_C_TO_F => A
fleebee	0.916% MPI_Barrier¢
nasnes	0.807% FLASH => EVOLVE => HYDRO::HYDRO_3D => MODULEHYDROSWEEP::HYDRO_SWEEP => MESH_GUARDCELL => TOT_BND => DBASETREE::DBAS
\bigcirc Fortran + MPI	U.735% AMR_DIAGONAL_PATCH
Araonne	L671% AMR RESTRICT RED
<u> </u>	0.671% IFLASH => EVOLVE => HYDRO::HYDRO.3D => MODULEHYDROSWEEP::HYDROSWEEP => MESH_FLUX_CONSERVE => AMR_FLUX_CONSERVE_UDT
	0.657% IFLASH => EVOLVE => HYDRO::HYDRO_3D => MODULEHYDROSWEEP::HYDRO_SWEEP => MESH_GUARDCELL => AMR_GUARDCELL_SRL => AMR_
	0.638% [FLASH => EVOLVE => MESH_UPDATE_GRID_REFINEMENT => MARK_GRID_REFINEMENT => MPI_Barrier¢
	0.61% IFLASH => EVOLVE => HYDRO::HYDRO_3D => MODULEHYDROSWEEP::HYDRO_SWEEP => MESH_GUARDCELL => AMR_GUARDCELL_C_TO_F => A
	0.556% FLASH => EVOLVE => HYDRO::HYDRO_3D => MODULEHYDROSWEEP::HYDRO_SWEEP => MESH_GUARDCELL => AMR_GUARDCELL_C_TO_F => A
	U.508% [] TOT_BND 0.454% [] FLACU DODUC MECH URDATE OND DEFINIEMENT MARK OND DEFINIEMENT MODULEFOCOD-FOCOD
1	U+34% [[FLASH => EVOLVE => MESH_UPDATE_UKID_KEFINEMENT => MAKK_UKID_KEFINEMENT =¥ MUDULEEUS3U#EUS3D
ParaProf – Scalable Histogram





ParaProf – 3D View (Full Profile)



ParaProf – 3D View (Full Profile)

• 🗆 🗙

File Options Windows Help

X ParaProf Visualizer



ParaProf – 3D Scatterplot

- Each point is a "thread" of execution
- A total of four metrics shown in relation
- ParaProf's visualization library
 JOGL
- Miranda







Performance Mapping

• Example: Particles distributed on cube

```
Particle* P[MAX]; /* Array of particles */
int GenerateParticles() {
  /* distribute particles over all faces of the cube */
  for (int face=0, last=0; face < 6; face++) {</pre>
    /* particles on this face */
    int particles on this face = num(face);
    for (int i=last; i < particles on this face; i++) {</pre>
       /* particle properties are a function of face */
    P[i] = ... f(face);
                                                                   00
    last+= particles on this face;
                  SC '09: Productive Performance Engineering of Petascale Applications with POIND and VOHPS
                                     41
```

Performance Mapping



How much time (flops) spent processing face i particles?

POINT

• What is the distribution of performance among faces?

SC '09: Productive Performance Engineering of Petascale Applications with POINT and VI-HPS



No Mapping versus Mapping

- Typical performance tools report performance with respect to routines
- Does not provide support for manning

 TAU's performance mapping can observe performance with respect to scientist's programming and problem abstractions



NAS BT – Flat Profile



NAS BT – Phase Profile

Main phase shows nested phases and immediate events



Phase Profiling of HW Counters

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





Profile Snapshots in ParaProf

- Profile snapshots are parallel profiles recorded at runtime
- Show's porformance profile dynamics (all types TAU measurement



parallel profile snapshots





Profile Snapshot Views

Only show main loop



POINT

Percentage
 breakdown



SC '09: Productive Performance Engineering of Petascale Applications with POINT and VI-HPS $^{\rm 48}_{\rm 48}$



Snapshot Replay in ParaProf



PerfExplorer – Runtime Breakdown

000

X TAU/PerfExplorer: Total Runtime Breakdown

<u>File</u><u>H</u>elp

Total Runtime Breakdown for S3D (Jaguar, ORNL):Harness Scaling Study: GET_TIME_OF_DAY



 $\blacksquare \text{ DERIVATIVE_X_COMM [[derivative_x.pp.f90] {53,14}]} \blacksquare \text{ Loop: CHEMRIN_M::REACTION_RATE_BOUNDS [[cnemRin_m.pp.f90] {574,3}-[386,7]]} \\ \blacksquare \text{ Loop: DERIVATIVE_X_CALC [{derivative_x.pp.f90} {432,10}-[441,15]] \blacksquare \text{ Loop: DERIVATIVE_X_CALC [{derivative_x.pp.f90} {566,19}-[589,24]]} \\ \blacksquare \text{ Loop: DERIVATIVE_X_CALC [{derivative_x.pp.f90} {566,1$

Loop: DERIVATIVE_Y_CALC [{derivative_y.pp.f90} {431,10}-{440,15}]

Loop: INTEGRATE [(integrate_erk.pp.f90) {73,3}-{93,13}]

Loop: RHSF [{rhsf.pp.f90} {537,3}-{543,16}] Loop: RHSF [{rhsf.pp.f90} {545,3}-{551,16}]

Loop: THERMCHEM_M::CALC_INV_AVG_MOL_WT [{thermchem_m.pp.f90} { 127,5}-{ 129,9}]

Loop: THERMCHEM_M::CALC_SPECENTH_ALLPTS [{thermchem_m.pp.f90} {506,3}-{512,8}]

Loop: THERMCHEM_M::CALC_TEMP [{thermchem_m.pp.f90} { 175,5 }-{216,9 }]

Loop: TRANSPORT_M::COMPUTECOEFFICIENTS [{mixavg_transport_m.pp.f90} {492,5}-{520,9}]

Loop: TRANSPORT_M::COMPUTEHEATELUX [(mixavg_transport_m.pp.f90) {782,5}-{790,19}]

Loop: TRANSPORT_M::COMPUTESPECIESDIFFFLUX [(mixavg_transport_m.pp.f90] {630,5]-{656,19}]

Loop: VARIABLES_M::GET_MASS_FRAC [{variables_m.pp.f90} {96,3}-{99,7}]5[] MPI_Comm_compare() - MPI_Wait()

📕 READWRITE_SAVEFILE_DATA [(io.pp.f90) {544,14}] 📕 RHSF [(rhsf.pp.f90) {1,12}] 👘 WRITE_SAVEFILE [(io.pp.f90) {240,14}] 📕 other

PerfExplorer – Relative Comparisons

- Total execution time
- Timesteps per second
- Relative efficiency
- Relative efficiency per event
- Relative speedup
- Relative speedup per event
- Group fraction of total
- Runtime breakdown

- Correlate events with total runtim
- Relative efficiency per phase
- Relative speedup per phase
- Distribution visualizations





PerfExplorer – Correlation Analysis

600

PerfExplorer Client



PerfExplorer – Correlation Analysis

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







PerfExplorer – Cluster Analysis



PerfExplorer – Cluster Analysis

 Four significant events automatically selected



55

PerfExplorer – Performance Regression





Other Projects in TAU

- TAU Portal
 - Support collaborative performance study
- Kernel-level system measurements (KTAU)
 - Application to OS noise analysis and I/O system analysis
- TAU performance monitoring
 - TAUoverSupermon and TAUoverMRNet
- PerfExplorer integration and expert-based analysis
 - OpenUH compiler optimizations
 - Computational quality of service in CCA
- Eclipse CDT and PTP integration
- Performance tools integration (NSF POINT project) SC '09: Productive Performance Engineering of Petascale Applications with POINT and VI-HPS



Using TAU

- Install TAU
 - % configure [options]; make clean install
- Modify application makefile and choose TAU configuration
 - Select TAU's stub makefile
 - Change name of compiler in makefile
- Set environment variables
 - Directory where profiles/traces are to be stored/counter selection
 - TAU options

- Execute application
 - % mpirun –np <procs> a.out;
- Analyze performance data
 - paraprof, vampir, pprof, paraver ...



Application Build Environment

- Minimize impact on user's application build procedures
- Handle 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 AU COMPILER



Configuring TAU

- TAU can measure several metrics with profiling and tracing approaches
- Different tools can also be invoked to instrument programs for TAU measurement
- Each configuration of TAU produces a measurement library for an architecture
- Each measurement configuration of TAU also creates a corresponding stub makefile that can be used to compile programs
- Typically configure multiple measurement libraries



SC '09: Productive Performance Engineering of Petascale Applications with POINT and VI-HPS



TAU Measurement System Configuration

- configure [OPTIONS]
 - {-c++=<CC>, -cc=<cc>}
 - -pdt=<dir>
 - opari=<dir>
 - -papi=<dir>
 - vampirtrace=<dir>
 - -mpi[inc/lib]=<dir>
 - -dyninst=<dir>
 - shmem[inc/lib]=<dir>
 - -python[inc/lib]=<dir>
 - -tag=<name>
 - -epilog=<dir>
 - -slog2
 - -otf=<dir>
 - arch=<architecture>
 (bgl_xt3)
 - {-pthread, -sproc}
 - openmp

POINT

- -jdk=<dir>
- fortran=[vendor]

Specify C++ and C compilers Specify location of PDT Specify location of Opari OpenMP tool Specify location of PAPI Specify location of VampirTrace Specify MPI library instrumentation Specify location of DynInst Package Specify PSHMEM library instrumentation Specify Python instrumentation Specify a unique configuration name Specify location of EPILOG Build SLOG2/Jumpshot tracing package Specify location of OTF trace package itecture> Specify architecture explicitly (bgl, xt3,ibm64,ibm64linux...) Use pthread or SGI sproc threads Use OpenMP threads Specify Java instrumentation (JDK) Specify Fortran compiler

VI-HRS

TAU Measurement System Configuration

- configure [OPTIONS]
 - TRACE Generate binary TAU traces
 - -PROFILE (default)
 - - PROFILECALLPATH
 - - PROFILEPHASE

POINT

- - PROFILEMEMORY

- Generate profiles (summary)
- Generate call path profiles
 - Generate phase based profiles
- 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 Configuration – Examples

- Configure using PDT and MPI for x86_64 Linux ./configure –pdt=/usr/pkgs/pkgs/pdtoolkit-3.14 -mpiinc=/usr/pkgs/mpich/include -mpilib= /usr/pkgs/mpich/lib -mpilibrary='-Impich -L/usr/gm/lib64 -Igm -Ipthread -Idl'
- Use PAPI counters (one or more) with C/C++/F90 automatic instrumentation for Cray CNL. Also instrument the MPI library. Use PGI compilers.

./configure -arch=craycnl -cc=cc -c++=CC -fortran=pgi -papi= /opt/xt-tools/papi/3.2.1 -mpi -MULTIPLECOUNTERS; make clean install

Stub makefiles

/usr/pkgs/tau/x86_64/lib/Makefile.tau-mpi-pdt-pgi /usr/pkgs/tau/x86_64/lib/Makefile.tau-multiplecounterspoi-papi-pdt-pgile Performance Engineering of Petascale Applications with POINT and VI-HPS

Stub Makefiles Configuration Parameters

- TAU scripts use stub makefiles to select performance measurements
- Variables:
 - TAU_CXX Specify the C++ compiler used by TAU
 - TAU_CC, TAU_F90
 - TAU_DEFS
 - TAU_LDFLAGS
 - TAU_INCLUDE
 - TAU_LIBS
 - TAU_SHLIBS
 - TAU_MPI_LIBS
 - TAU_MPI_FLIBS
 - TAU_FORTRANLIBS
 - TAU_CXXLIBS
 - TAU_INCLUDE_MEMORY
 - TAU_DISABLE
 - TAU_COMPILER



Defines used by TAU (add to CFLAGS) Linker options (add to LDFLAGS) Header files include path (add to CFLAGS) Statically linked TAU library (add to LIBS) Dynamically linked TAU library TAU's MPI wrapper library for C/C++ TAU's MPI wrapper library for F90 Must be linked in with C++ linker for F90 Must be linked in with F90 linker EMORY Use TAU's malloc/free wrapper lib TAU's dummy F90 stub library

Specify the C, F90 compilers

Instrument using tau_compiler.sh script



TAU Measurement Configuration

- % cd /opt/tau-2.18/x86_64/lib; Is Makefile.*
 - 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-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 /opt/tau-2.187x86_64/lib/Makefile.tau-mpi-pdt



-PROFILE Option

- Generates flat profiles
 - One for each MPI process
 - It is the default option.
- Uses wallclock time
 gettimeofday() sys call

POINT

 Calculates exclusive, inclusive time spent in each timer and number of calls



Generating a Flat Profile with MPI

```
% paraprof app.ppk
```



Generating a Loop-level Profile

```
% setenv TAU MAKEFILE /opt/tau-2.18/x86 64
                      /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=(/opt/tau-2.18/x86 64/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
```



Compiler-based Instrumentation





-MULTIPLECOUNTERS 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 ...
- 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

POINT

papi_native_avail and papi_avail are two useful tools



Generate a PAPI profile

```
% setenv TAU MAKEFILE /opt/tau-2.18/x86 64
                          /lib/Makefile.tau-multiplecounters-papi-mpi-pdt
% setenv TAU OPTIONS `-optTauSelectFile=select.tau -optVerbose'
% cat select.tau
 BEGIN INSTRUMENT SECTION
  loops routine="#"
 END INSTRUMENT SECTION
% set path=(/opt/tau-2.18/x86 64/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
% 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.
```



-PROFILECALLPATH Option

- Generates profiles that show the calling order (edges and 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 environment variable
 - -callpath in the name of the stub Makefile name or setting TAU_CALLPATH= 1 at runtime

POINT AU V2: 09 POuctive Performance Engineering of Petascale Applications with POINT and VI-HPS 72


-DEPTHLIMIT 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.tauicpc-mpi-depthlimit-pdt





Generate a Callpath Profile

```
% setenv TAU MAKEFILE /opt/tau-2.18/x86 64
                       /lib/Makefile.tau-callpath-mpi-pdt
% set path=(/opt/tau-2.18/x86 64/bin $path)
% make F90=tau f90.sh
(Or edit Makefile and change F90=tau f90.sh)
% setenv TAU CALLPATH DEPTH 100
NOTE: In TAU v2.18.1+ you may simply use:
% setenv TAU CALLPATH 1
to generate the callpath profiles without any recompilation.
% qsub run.job
% paraprof --pack app.ppk
  Move the app.ppk file to your desktop.
% paraprof app.ppk
(Windows -> Thread -> Call Graph)
```



-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-icpcmpi-pdt-trace



Generate a Trace File

```
% setenv TAU MAKEFILE /opt/tau-2.18/x86 64
               /lib/Makefile.tau-mpi-pdt-trace
% set path=(/opt/tau-2.18/x86 64/bin $path)
% make F90=tau f90.sh
(Or edit Makefile and change F90=tau f90.sh)
% asub 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
VAMPTR:
% tau2otf tau.trc tau.edf app.otf -n 4 -z
(4 streams, compressed output trace)
% vampir app.otf
(or vng client with vngd server)
```

Instrumentation Specification

```
% tau instrumentor
Usage : tau instrumentor <pdbfile> <sourcefile> [-o <outputfile>] [-noinline]
[-g groupname] [-i headerfile] [-c|-c++|-fortran] [-f <instr reg 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
```



Outer Loop Level Instrumentation

BEGIN_INSTRUMENT_SECTION						
<pre>loops file="loop test.cpp" routine="multiply"</pre>						
# it also understands # as the wildcard in routine name						
# and * and ? wildcards in file name.						
# You can also specify the full						
# name of the routine as is found in profile files						
# name of the foutine as is found in profile files.						
#100ps file="100p_test.cpp" routine="double multiply#"						
END_INSTRUMENT_SECTION						
	_					
% pprof						
NODE 0	;CONTEXT 0;THR	EAD 0:				
 %Time	Exclusive	Inclusive	 #Call	#Subrs	Inclusive	Name
	msec	total msec			usec/call	
100.0	0.12	25,162	1	1	25162827	int main(int, char **)
100.0	0.175	25,162	1	4	25162707	double multiply()
90.5	22,778	22,778	1	0	22778959	Loop: double multiply()[
file =	<loop_test.cp< td=""><td><pre>p> line,col =</pre></td><td><23,3> to <30</td><td>,3>]</td><td></td><td></td></loop_test.cp<>	<pre>p> line,col =</pre>	<23,3> to <30	,3>]		
9.3	2,345	2,345	1	0	2345823	Loop: double multiply()[
file =	<loop_test.cp< td=""><td><pre>p> line,col =</pre></td><td><38,3> to <46</td><td>,7>]</td><td></td><td></td></loop_test.cp<>	<pre>p> line,col =</pre>	<38,3> to <46	,7>]		
0.1	33	33	1	0	33964	Loop: double
multip	ly()[file = <	loop_test.cpp>	line, col = <	16,10> to	<21,12>]	



Support Acknowledgements

Office of

- Department of Energy (DOE)
 - Office of Science
 - MICS, Argonne National Lab
 - ASC/NNSA
 - University of Utah ASC/NNSA Level 1
 - ASC/NNSA, Lawrence Livermore National La
- Department of Defense (DoD)
 - HPC Modernization Office (HPCMO)
- NSF Software Development for Cyberinfrastructure (SD)
- Research Centre Juelich
- Los Alamos National Laboratory
- **TU Dresden**
- ParaTools, Inc. POIN1 SC '09: Productive Performance Engineering of Petascale Applications with POINT and VI-HPS

ParaTools

For more information

- TAU Website: http://tau.uoregon.edu
 - Software
 - Release notes
 - Documentation



