# *Parallel Performance Evaluation using the TAU Performance System Project*

## Workshop on Performance Tools for Petascale Computing

9:30 – 10:30am, Tuesday, July 17, 2007, Snowbird, UT

### Sameer S. Shende

sameer@cs.uoregon.edu

http://www.cs.uoregon.edu/research/tau

Performance Research Laboratory

University of Oregon

UNIVERSITY

OF OREGON

# *Acknowledgements*

- Dr. Allen D. Malony, Professor
- Alan Morris, Senior software engineer
- Wyatt Spear, Software engineer
- Scott Biersdorff, Software engineer
- Kevin Huck, Ph.D. student
- Aroon Nataraj, Ph.D. student
- Brad Davidson, Systems administrator

# *Outline*

□ Overview of features

□ Instrumentation

□ Measurement

□ Analysis tools

  ○ Parallel profile analysis (ParaProf)

  ○ Performance data management (PerfDMF)

  ○ Performance data mining (PerfExplorer)

□ Application examples

□ Kernel monitoring and KTAU

# TAU Performance System

□ *T*uning and *A*nalysis *U*tilities (15+ year project effort)
□ **Performance system framework for HPC systems**
  ○ Integrated, scalable, flexible, and parallel
□ Targets a general complex system computation model
  ○ *Entities*: nodes / contexts / threads
  ○ *Multi-level*: system / software / parallelism
  ○ Measurement and analysis abstraction
□ **Integrated toolkit for performance problem solving**
  ○ Instrumentation, measurement, analysis, and visualization
  ○ Portable performance profiling and tracing facility
  ○ Performance data management and data mining
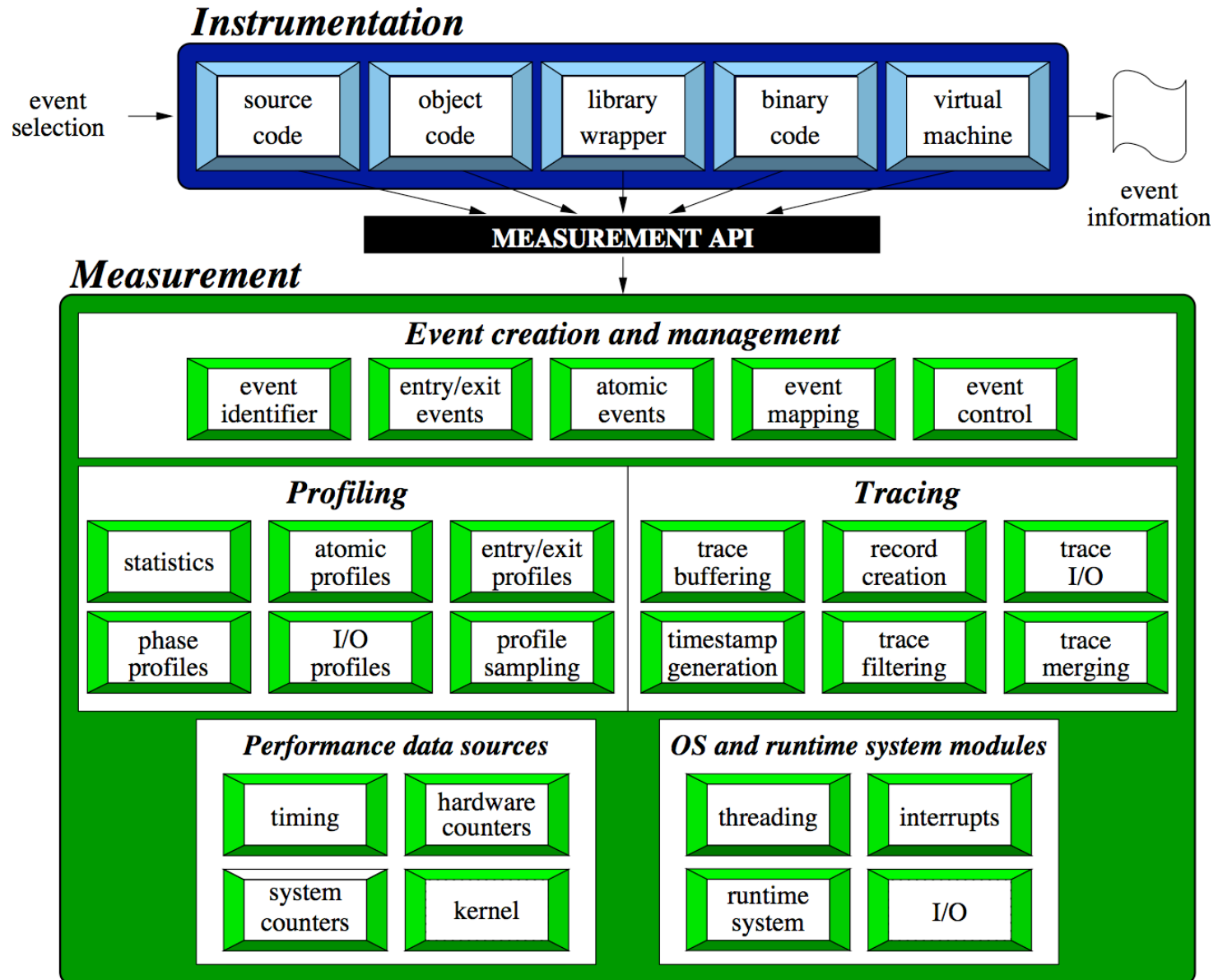□ *Partners*: LLNL, ANL, LANL, Research Center Jülich

# *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
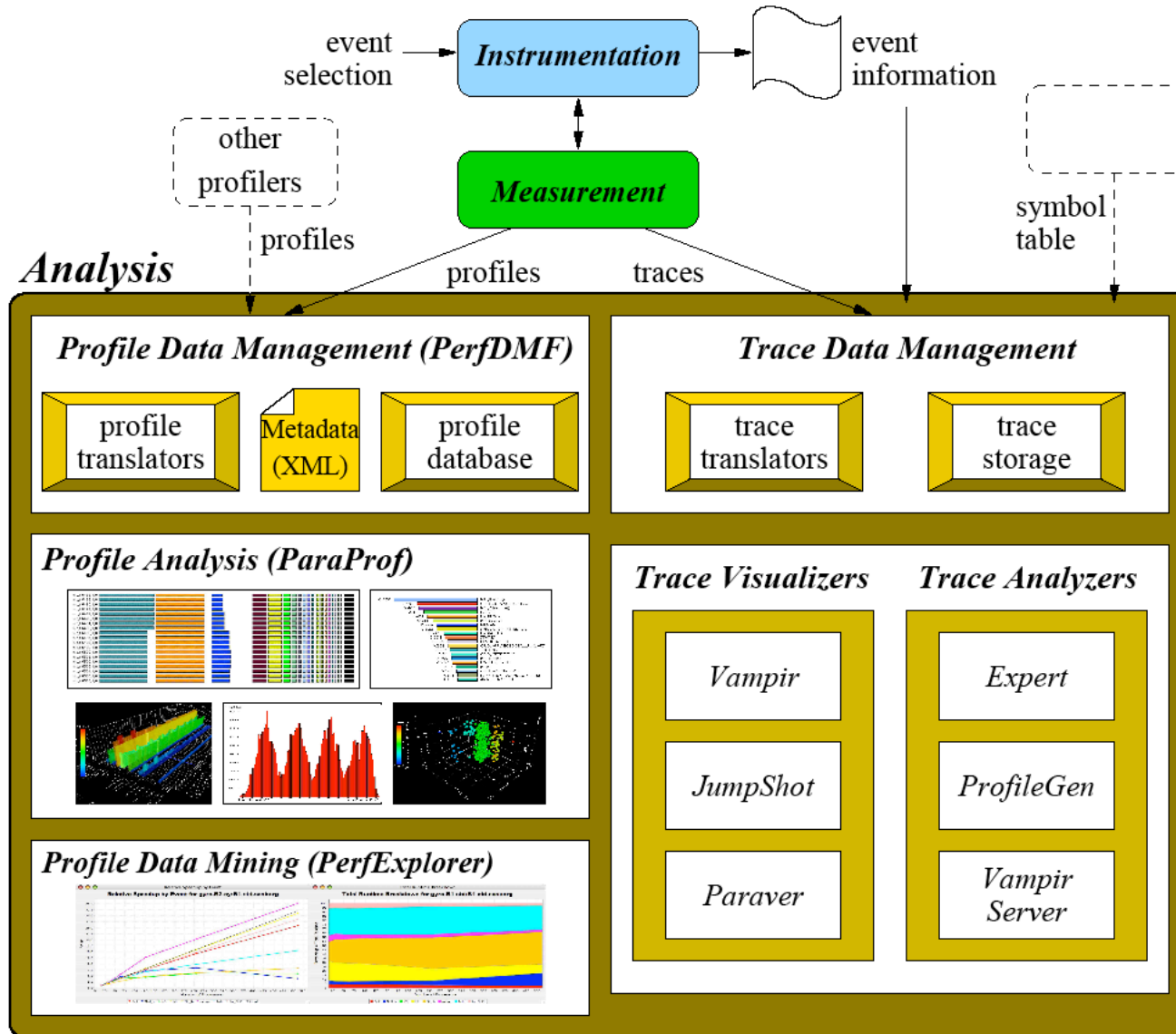- **Scalable (very large) parallel performance analysis**

# TAU Performance System Architecture

## Instrumentation

event selection →

| source code | object code | library wrapper | binary code | virtual machine |
|---|---|---|---|---|

→ event information

**MEASUREMENT API**

## Measurement

### Event creation and management

| event identifier | entry/exit events | atomic events | event mapping | event control |
|---|---|---|---|---|

### Profiling

| statistics | atomic profiles | entry/exit profiles |
|---|---|---|
| phase profiles | I/O profiles | profile sampling |

### Tracing

| trace buffering | record creation | trace I/O |
|---|---|---|
| timestamp generation | trace filtering | trace merging |

### Performance data sources

| timing | hardware counters |
|---|---|
| system counters | kernel |

### OS and runtime system modules

| threading | interrupts |
|---|---|
| runtime system | I/O |

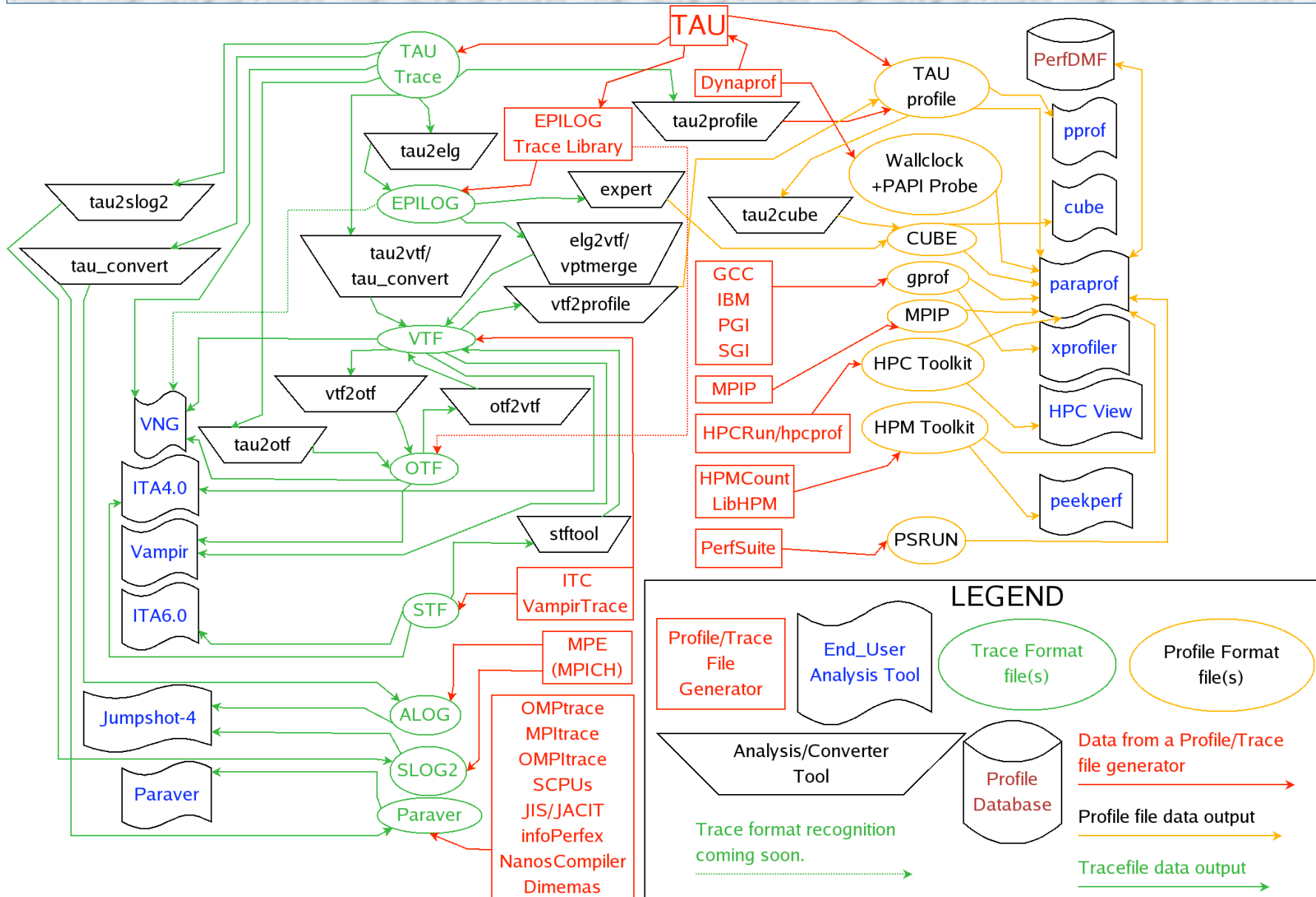# TAU Performance System Architecture

# Building Bridges to Other Tools: TAU

# *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 for hardware performance counters (PAPI)

❑ Support definition of "semantic" entities for mapping

❑ Support for event groups (aggregation, selection)

❑ Instrumentation optimization
- ○ Eliminate instrumentation in lightweight routines

# *PAPI*

- **Performance Application Programming Interface**
  - The purpose of the PAPI project is to design, standardize and implement a portable and efficient API to access the hardware performance monitor counters found on most modern microprocessors.
- Parallel Tools Consortium project started in 1998
- Developed by University of Tennessee, Knoxville
- http://icl.cs.utk.edu/papi/

# 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)

□ **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*)

□ `TAU_COMPILER` to automate instrumentation process
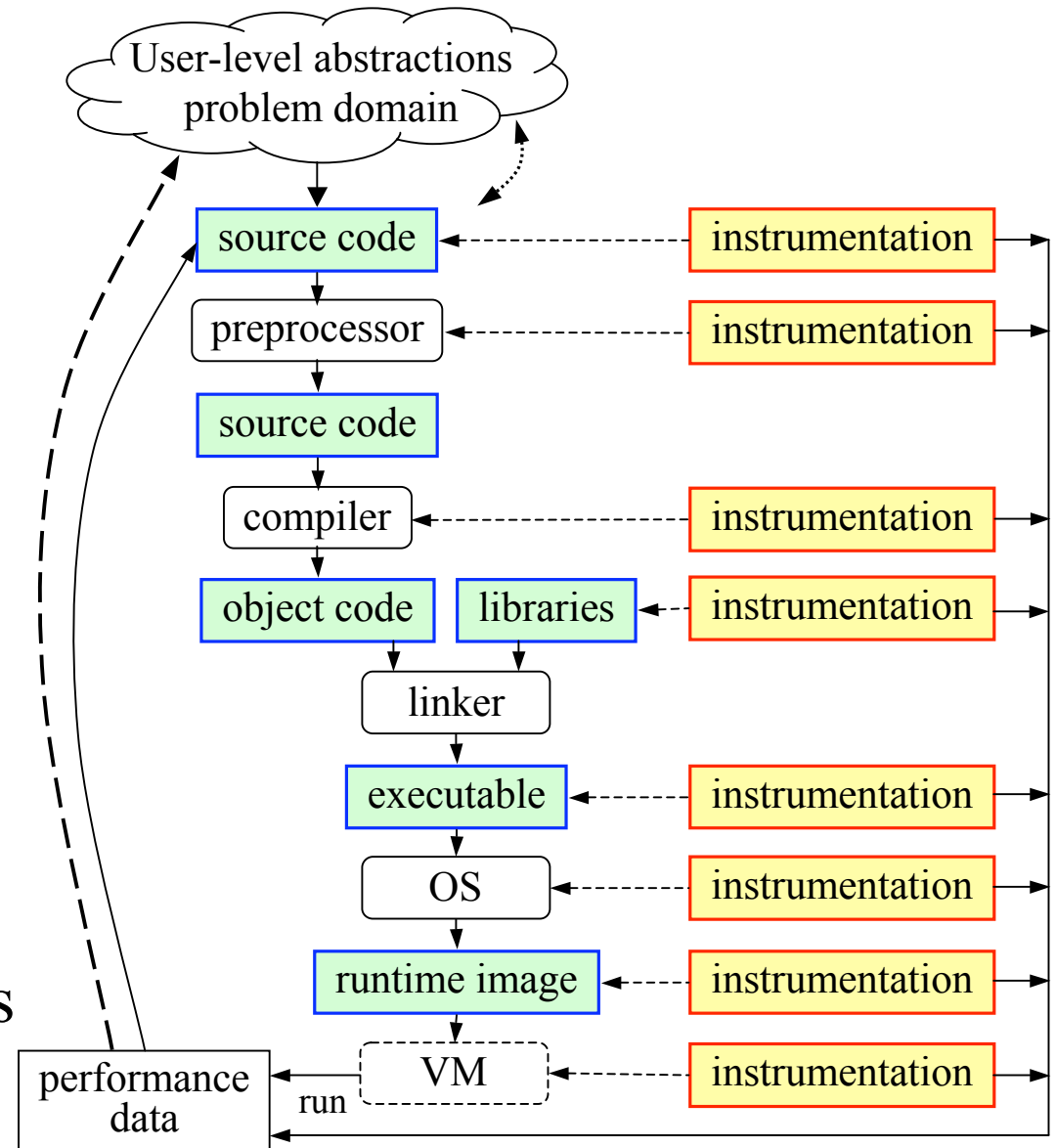
# *Using TAU: A brief Introduction*

❒ To instrument source code using PDT

   ○ Choose an appropriate TAU stub makefile in <arch>/lib:

   **% setenv TAU_MAKEFILE
   /usr/tau-2.x/xt3/lib/Makefile.tau-mpi-pdt-pgi**

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

   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)**

# *Multi-Level Instrumentation and Mapping*

- □ **Multiple interfaces**
- □ **Information sharing**
  - ○ Between interfaces
- □ **Event selection**
  - ○ Within/between levels
- □ **Mapping**
  - ○ Associate performance data with high-level semantic abstractions

# *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 for CCA component software

# *TAU Measurement Mechanisms*

□ **Parallel profiling**

  ○ Function-level, block-level, statement-level

  ○ Supports user-defined events and mapping events

  ○ TAU parallel profile stored (dumped) during execution

  ○ Support for flat, callgraph/callpath, phase profiling

  ○ Support for memory profiling (headroom, malloc/leaks)

  ○ Support for tracking I/O (wrappers, Fortran instrumentation of read/write/print calls)

□ **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) phases

# *Performance Analysis and Visualization*

- ❏ Analysis of parallel profile and trace measurement
- ❏ **Parallel profile analysis**
  - ○ *ParaProf*: parallel profile analysis and presentation
  - ○ *ParaVis*: parallel performance visualization package
  - ○ Profile generation from trace data (*tau2profile*)
- ❏ Performance data management framework (*PerfDMF*)
- ❏ **Parallel trace analysis**
  - ○ Translation to *VTF* (V3.0), *EPILOG, OTF* formats
  - ○ Integration with *VNG* (Technical University of Dresden)
- ❏ Online parallel analysis and visualization
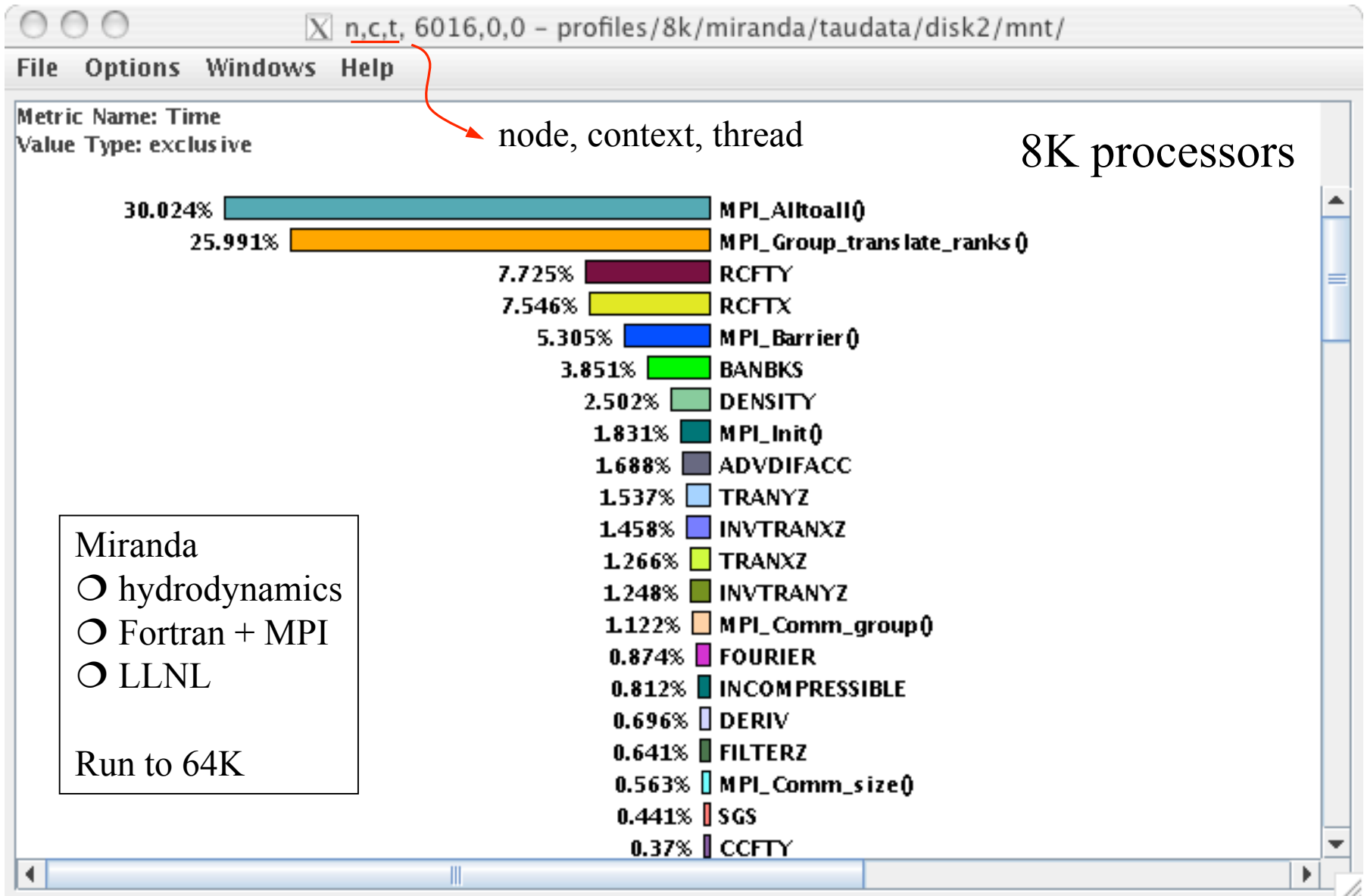- ❏ Integration with *CUBE* browser (KOJAK, UTK, FZJ)

# *ParaProf Parallel Performance Profile Analysis*

# *ParaProf – Flat Profile (Miranda, BG/L)*

X n,c,t, 6016,0,0 – profiles/8k/miranda/taudata/disk2/mnt/

File   Options   Windows   Help

Metric Name: Time
Value Type: exclusive

node, context, thread

8K processors

| | |
|---|---|
| 30.024% | MPI_Alltoall() |
| 25.991% | MPI_Group_translate_ranks() |
| 7.725% | RCFTY |
| 7.546% | RCFTX |
| 5.305% | MPI_Barrier() |
| 3.851% | BANBKS |
| 2.502% | DENSITY |
| 1.831% | MPI_Init() |
| 1.688% | ADVDIFACC |
| 1.537% | TRANYZ |
| 1.458% | INVTRANXZ |
| 1.266% | TRANXZ |
| 1.248% | INVTRANYZ |
| 1.122% | MPI_Comm_group() |
| 0.874% | FOURIER |
| 0.812% | INCOMPRESSIBLE |
| 0.696% | DERIV |
| 0.641% | FILTERZ |
| 0.563% | MPI_Comm_size() |
| 0.441% | SGS |
| 0.37% | CCFTY |

Miranda
O hydrodynamics
O Fortran + MPI
O LLNL

Run to 64K

# *ParaProf – Stacked View (Miranda)*

# *ParaProf – Callpath Profile (Flash)*

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

File Options Windows Help

Metric Name: Time
Value Type: exclusive

26.474% MODULEHYDROSWEEP::HYDRO_SWEEP
26.474% FLASH => EVOLVE => HYDRO::HYDRO_3D => MODULEHYDROSWEEP::HYDRO_SWEEP
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% MPI_Ssend()
3.678% FLASH => EVOLVE => HYDRO::HYDRO_3D => MODULEHYDROSWEEP::HYDRO_SWEEP => MODULEEOS3D::EOS3D
3.536% MPI_Allreduce()
2.727% MPI_Waitall()
2.242% MODULEUPDATE_SOLN::UPDATE_SOLN
2.242% FLASH => EVOLVE => HYDRO::HYDRO_3D => MODULEHYDROSWEEP::HYDRO_SWEEP => MODULEUPDATE_SOLN::UPDATE_SOLN
2.059% AMR_GUARDCELL_CC_SRL
1.703% FLASH => EVOLVE => HYDRO::HYDRO_3D => MODULEHYDROSWEEP::HYDRO_SWEEP => MESH_GUARDCELL => AMR_GUARDCELL_SRL => AMR_
1.56% FLASH => EVOLVE => HYDRO::HYDRO_3D => MODULEHYDROSWEEP::HYDRO_SWEEP => MESH_GUARDCELL => AMR_GUARDCELL_SRL => AMR_
1.406% FLASH => EVOLVE => MESH_UPDATE_GRID_REFINEMENT => MESH_REFINE_DEREFINE => AMR_REFINE_DEREFINE => AMR_MORTON_ORDER => /
1.361% FLASH => TIMESTEP => MPI_Allreduce()
1.319% AMR_RESTRICT_UNK_FUN
1.272% AMR_PROLONG_GEN_UNK_FUN
1.093% FLASH => EVOLVE => HYDRO::HYDRO_3D => MODULEHYDROSWEEP::HYDRO_SWEEP => MESH_GUARDCELL => AMR_GUARDCELL_C_TO_F => A
1.077% ABUNDANCE_RESTRICT
1.077% FLASH => EVOLVE => HYDRO::HYDRO_3D => MODULEHYDROSWEEP::HYDRO_SWEEP => ABUNDANCE_RESTRICT
1.064% DBASETREE::DBASENEIGHBORBLOCKLIST
1% FLASH => EVOLVE => HYDRO::HYDRO_3D => MODULEHYDROSWEEP::HYDRO_SWEEP => MESH_GUARDCELL => AMR_RESTRICT => AMR_RESTRI
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
0.916% MPI_Barrier()
0.807% FLASH => EVOLVE => HYDRO::HYDRO_3D => MODULEHYDROSWEEP::HYDRO_SWEEP => MESH_GUARDCELL => TOT_BND => DBASETREE::DBAS
0.806% AMR_PROLONG_UNK_FUN
0.735% AMR_DIAGONAL_PATCH
0.699% DIFFUSE
0.699% FLASH => EVOLVE => HYDRO::HYDRO_3D => MODULEHYDROSWEEP::HYDRO_SWEEP => DIFFUSE
0.671% AMR_RESTRICT_RED
0.671% FLASH => EVOLVE => HYDRO::HYDRO_3D => MODULEHYDROSWEEP::HYDRO_SWEEP => MESH_FLUX_CONSERVE => AMR_FLUX_CONSERVE_UDT
0.657% FLASH => 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% FLASH => 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
0.508% TOT_BND
0.454% FLASH => EVOLVE => MESH_UPDATE_GRID_REFINEMENT => MARK_GRID_REFINEMENT => MODULEEOS3D::EOS3D

Flash
○ thermonuclear
  flashes
○ Fortran + MPI
○ Argonne
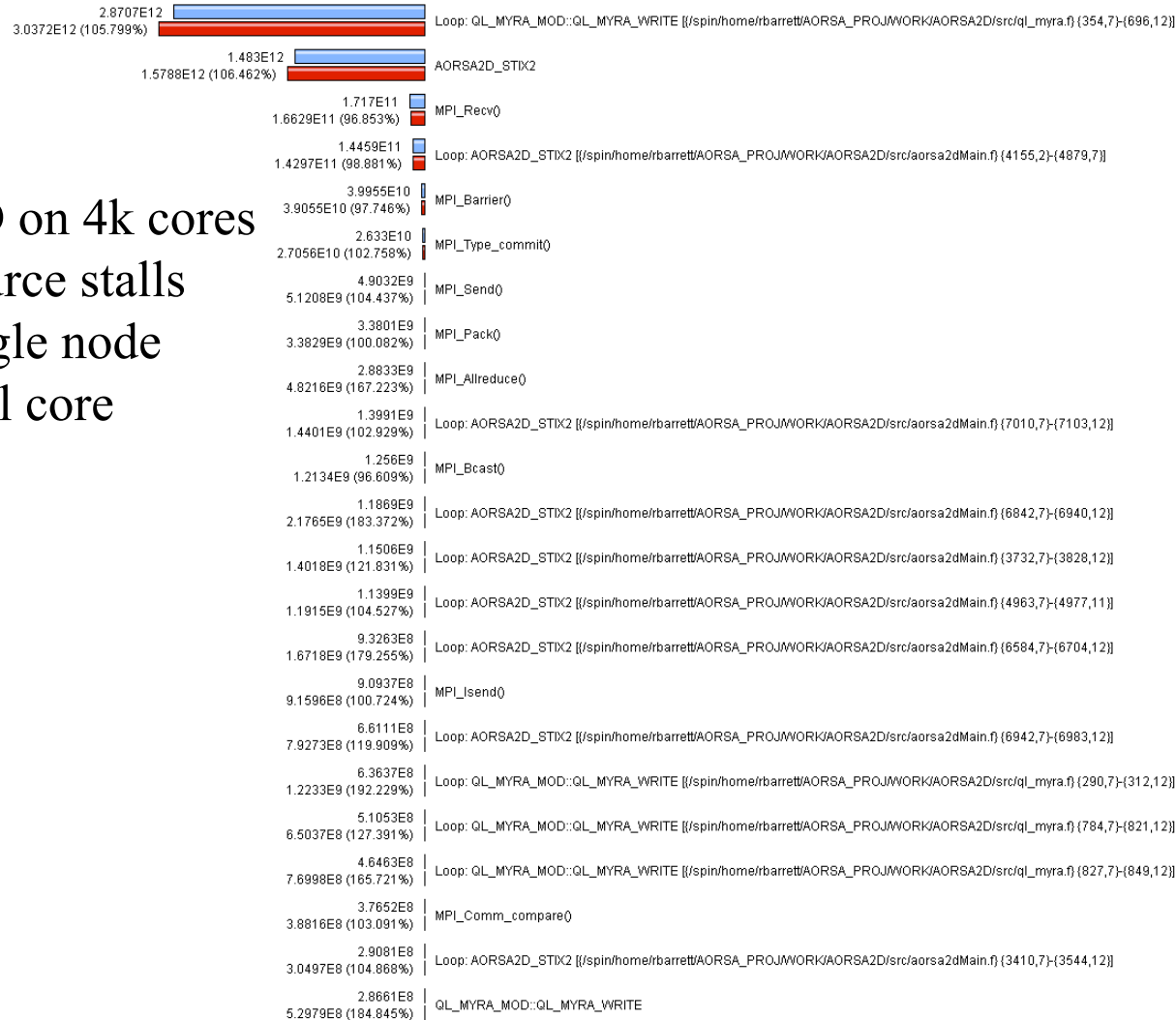
# *Comparing Effects of MultiCore Processors*

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

```
2.8707E12            Loop: QL_MYRA_MOD::QL_MYRA_WRITE [{/spin/home/rbarrett/AORSA_PROJ/WORK/AORSA2D/src/ql_myra.f} {354,7}-{696,12}]
3.0372E12 (105.799%)

1.483E12             AORSA2D_STIX2
1.5788E12 (106.462%)

1.717E11             MPI_Recv()
1.6629E11 (96.853%)

1.4459E11            Loop: AORSA2D_STIX2 [{/spin/home/rbarrett/AORSA_PROJ/WORK/AORSA2D/src/aorsa2dMain.f} {4155,2}-{4879,7}]
1.4297E11 (98.881%)

3.9955E10            MPI_Barrier()
3.9055E10 (97.746%)

2.633E10             MPI_Type_commit()
2.7056E10 (102.758%)

4.9032E9             MPI_Send()
5.1208E9 (104.437%)

3.3801E9             MPI_Pack()
3.3829E9 (100.082%)

2.8833E9             MPI_Allreduce()
4.8216E9 (167.223%)

1.3991E9             Loop: AORSA2D_STIX2 [{/spin/home/rbarrett/AORSA_PROJ/WORK/AORSA2D/src/aorsa2dMain.f} {7010,7}-{7103,12}]
1.4401E9 (102.929%)

1.256E9              MPI_Bcast()
1.2134E9 (96.609%)

1.1869E9             Loop: AORSA2D_STIX2 [{/spin/home/rbarrett/AORSA_PROJ/WORK/AORSA2D/src/aorsa2dMain.f} {6842,7}-{6940,12}]
2.1765E9 (183.372%)

1.1506E9             Loop: AORSA2D_STIX2 [{/spin/home/rbarrett/AORSA_PROJ/WORK/AORSA2D/src/aorsa2dMain.f} {3732,7}-{3828,12}]
1.4018E9 (121.831%)

1.1399E9             Loop: AORSA2D_STIX2 [{/spin/home/rbarrett/AORSA_PROJ/WORK/AORSA2D/src/aorsa2dMain.f} {4963,7}-{4977,11}]
1.1915E9 (104.527%)

9.3263E8             Loop: AORSA2D_STIX2 [{/spin/home/rbarrett/AORSA_PROJ/WORK/AORSA2D/src/aorsa2dMain.f} {6584,7}-{6704,12}]
1.6718E9 (179.255%)

9.0937E8             MPI_Isend()
9.1596E8 (100.724%)

6.6111E8             Loop: AORSA2D_STIX2 [{/spin/home/rbarrett/AORSA_PROJ/WORK/AORSA2D/src/aorsa2dMain.f} {6942,7}-{6983,12}]
7.9273E8 (119.909%)

6.3637E8             Loop: QL_MYRA_MOD::QL_MYRA_WRITE [{/spin/home/rbarrett/AORSA_PROJ/WORK/AORSA2D/src/ql_myra.f} {290,7}-{312,12}]
1.2233E9 (192.229%)

5.1053E8             Loop: QL_MYRA_MOD::QL_MYRA_WRITE [{/spin/home/rbarrett/AORSA_PROJ/WORK/AORSA2D/src/ql_myra.f} {784,7}-{821,12}]
6.5037E8 (127.391%)

4.6463E8             Loop: QL_MYRA_MOD::QL_MYRA_WRITE [{/spin/home/rbarrett/AORSA_PROJ/WORK/AORSA2D/src/ql_myra.f} {827,7}-{849,12}]
7.6998E8 (165.721%)

3.7652E8             MPI_Comm_compare()
3.8816E8 (103.091%)

2.9081E8             Loop: AORSA2D_STIX2 [{/spin/home/rbarrett/AORSA_PROJ/WORK/AORSA2D/src/aorsa2dMain.f} {3410,7}-{3544,12}]
3.0497E8 (104.868%)

2.8661E8             QL_MYRA_MOD::QL_MYRA_WRITE
5.2979E8 (184.845%)
```
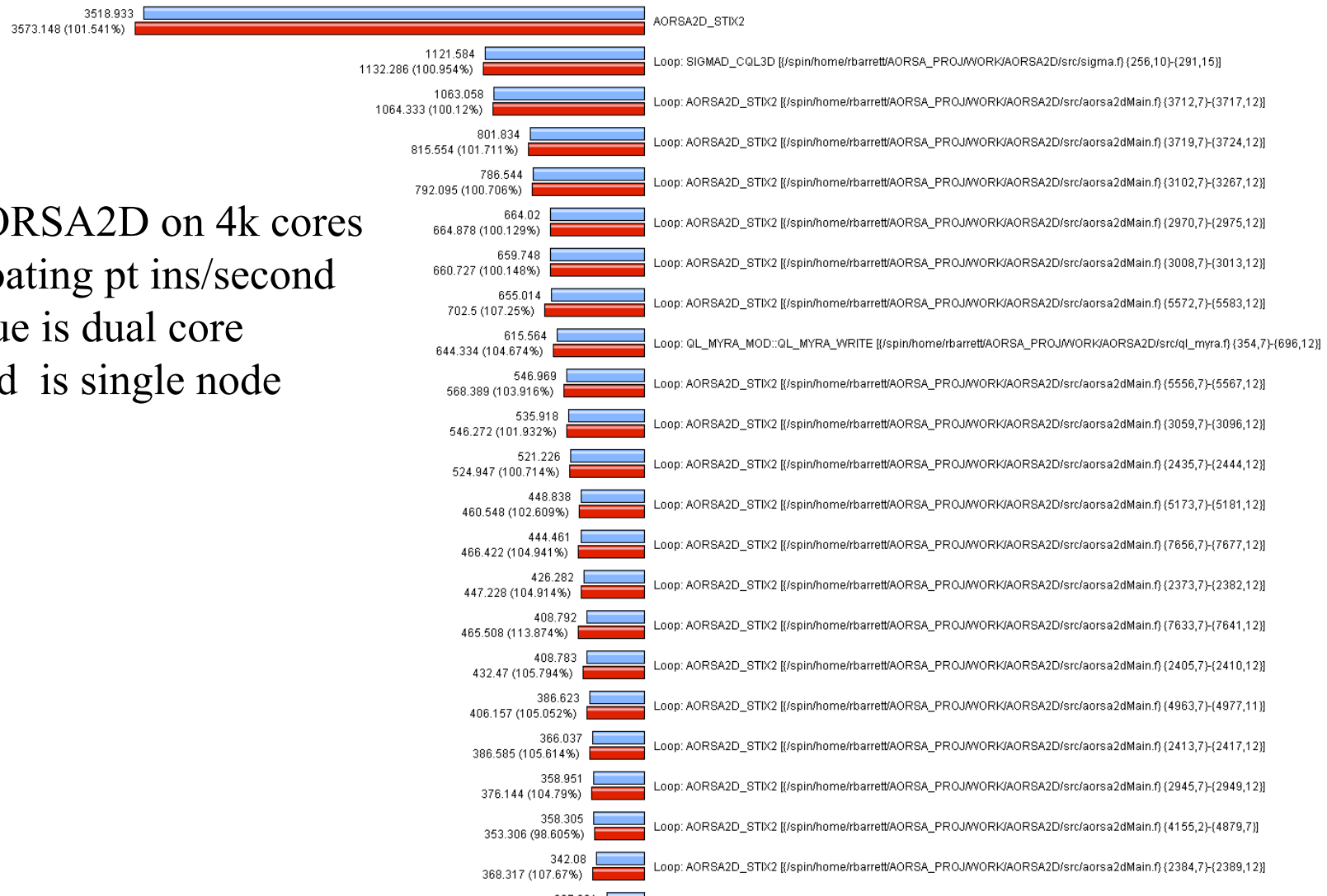
❑ AORSA2D on 4k cores
❑ PAPI resource stalls
❑ Blue is single node
❑ Red  is dual core

# *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

3518.933
3573.148 (101.541%) — AORSA2D_STIX2

1121.584
1132.286 (100.954%) — Loop: SIGMAD_CQL3D [{/spin/home/rbarrett/AORSA_PROJ/WORK/AORSA2D/src/sigma.f} {256,10}-{291,15}]

1063.058
1064.333 (100.12%) — Loop: AORSA2D_STIX2 [{/spin/home/rbarrett/AORSA_PROJ/WORK/AORSA2D/src/aorsa2dMain.f} {3712,7}-{3717,12}]

801.834
815.554 (101.711%) — Loop: AORSA2D_STIX2 [{/spin/home/rbarrett/AORSA_PROJ/WORK/AORSA2D/src/aorsa2dMain.f} {3719,7}-{3724,12}]

786.544
792.095 (100.706%) — Loop: AORSA2D_STIX2 [{/spin/home/rbarrett/AORSA_PROJ/WORK/AORSA2D/src/aorsa2dMain.f} {3102,7}-{3267,12}]

664.02
664.878 (100.129%) — Loop: AORSA2D_STIX2 [{/spin/home/rbarrett/AORSA_PROJ/WORK/AORSA2D/src/aorsa2dMain.f} {2970,7}-{2975,12}]

659.748
660.727 (100.148%) — Loop: AORSA2D_STIX2 [{/spin/home/rbarrett/AORSA_PROJ/WORK/AORSA2D/src/aorsa2dMain.f} {3008,7}-{3013,12}]

655.014
702.5 (107.25%) — Loop: AORSA2D_STIX2 [{/spin/home/rbarrett/AORSA_PROJ/WORK/AORSA2D/src/aorsa2dMain.f} {5572,7}-{5583,12}]

615.564
644.334 (104.674%) — Loop: QL_MYRA_MOD::QL_MYRA_WRITE [{/spin/home/rbarrett/AORSA_PROJ/WORK/AORSA2D/src/ql_myra.f} {354,7}-{696,12}]

546.969
568.389 (103.916%) — Loop: AORSA2D_STIX2 [{/spin/home/rbarrett/AORSA_PROJ/WORK/AORSA2D/src/aorsa2dMain.f} {5556,7}-{5567,12}]

535.918
546.272 (101.932%) — Loop: AORSA2D_STIX2 [{/spin/home/rbarrett/AORSA_PROJ/WORK/AORSA2D/src/aorsa2dMain.f} {3059,7}-{3096,12}]

521.226
524.947 (100.714%) — Loop: AORSA2D_STIX2 [{/spin/home/rbarrett/AORSA_PROJ/WORK/AORSA2D/src/aorsa2dMain.f} {2435,7}-{2444,12}]

448.838
460.548 (102.609%) — Loop: AORSA2D_STIX2 [{/spin/home/rbarrett/AORSA_PROJ/WORK/AORSA2D/src/aorsa2dMain.f} {5173,7}-{5181,12}]

444.461
466.422 (104.941%) — Loop: AORSA2D_STIX2 [{/spin/home/rbarrett/AORSA_PROJ/WORK/AORSA2D/src/aorsa2dMain.f} {7656,7}-{7677,12}]

426.282
447.228 (104.914%) — Loop: AORSA2D_STIX2 [{/spin/home/rbarrett/AORSA_PROJ/WORK/AORSA2D/src/aorsa2dMain.f} {2373,7}-{2382,12}]

408.792
465.508 (113.874%) — Loop: AORSA2D_STIX2 [{/spin/home/rbarrett/AORSA_PROJ/WORK/AORSA2D/src/aorsa2dMain.f} {7633,7}-{7641,12}]

408.783
432.47 (105.794%) — Loop: AORSA2D_STIX2 [{/spin/home/rbarrett/AORSA_PROJ/WORK/AORSA2D/src/aorsa2dMain.f} {2405,7}-{2410,12}]

386.623
406.157 (105.052%) — Loop: AORSA2D_STIX2 [{/spin/home/rbarrett/AORSA_PROJ/WORK/AORSA2D/src/aorsa2dMain.f} {4963,7}-{4977,11}]

366.037
386.585 (105.614%) — Loop: AORSA2D_STIX2 [{/spin/home/rbarrett/AORSA_PROJ/WORK/AORSA2D/src/aorsa2dMain.f} {2413,7}-{2417,12}]

358.951
376.144 (104.79%) — Loop: AORSA2D_STIX2 [{/spin/home/rbarrett/AORSA_PROJ/WORK/AORSA2D/src/aorsa2dMain.f} {2945,7}-{2949,12}]

358.305
353.306 (98.605%) — Loop: AORSA2D_STIX2 [{/spin/home/rbarrett/AORSA_PROJ/WORK/AORSA2D/src/aorsa2dMain.f} {4155,2}-{4879,7}]

342.08
368.317 (107.67%) — Loop: AORSA2D_STIX2 [{/spin/home/rbarrett/AORSA_PROJ/WORK/AORSA2D/src/aorsa2dMain.f} {2384,7}-{2389,12}]

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

# *ParaProf – Scalable Histogram View (Miranda)*



8k processors

16k processors

# *ParaProf – 3D Full Profile (Miranda)*

ParaProf Visualizer

File   Options   Windows   Help

○ Triangle Mesh

○ Bar Plot

○ Scatter Plot

**Height Metric**

Exclusive ▾   Time ▾

**Color Metric**

Exclusive ▾   Time ▾

Function    MPI_Barrier()

◀ ▌▌▌ ▶

Thread    16:0:0

◀ ▌▌ ▶

Height value   1.2229E8 microseconds

Color value    1.2229E8 microseconds

Mesh Plot | Axes | ColorScale | Render

Orientation

☑ Show Axes      NW ○      ○ NE

SE ◉      ○ SW

1.6734E8

1.6734E8

1.2551E8

8.3672E7

4.1836E7

16k processors

# *ParaProf – 3D Scatterplot (S3D – XT4 only)*

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



6400 cores

I/O takes less time on one node (rank 0)

- Events (exclusive time metric)
  - MPI_Barrier(), two loops
  - write operation

# *S3D Scatter Plot: Visualizing Hybrid XT3+XT4*



6400 cores

☐ Red nodes are XT4, blue are XT3

❏ Gap represents XT3 nodes

# *Visualizing S3D Profiles in ParaProf*



❏ Gap represents XT3 nodes

○ MPI_Wait takes less time, other routines take more time

# *Profile Snapshots in ParaProf*

❑ Profile snapshots are parallel profiles recorded at runtime

❑ Used to highlight profile changes during execution

Initialization

Checkpointing

Finalization

# *Profile Snapshots in ParaProf*

□ Filter snapshots (only show main loop iterations)

# *Profile Snapshots in ParaProf*

❏ Breakdown as a percentage

# Snapshot replay in ParaProf



All windows dynamically update

# *Profile Snapshots in ParaProf*

❑ Follow progression of various displays through time

❑ 3D scatter plot shown below



T = 0s → T = 11s

# *New automated metadata collection*

# *Performance Data Management: Motivation*

❑ Need for robust processing and storage of multiple profile performance data sets

❑ Avoid developing independent data management solutions
  - Waste of resources
  - Incompatibility among analysis tools

❑ Goals:
  - Foster multi-experiment performance evaluation
  - Develop a common, reusable foundation of performance data storage, access and sharing
  - A core module in an analysis system, and/or as a central repository of performance data

# *PerfDMF Approach*

- ❑ Performance <u>D</u>ata <u>M</u>anagement <u>F</u>ramework
- ❑ Originally designed to address critical TAU requirements
- ❑ Broader goal is to provide an open, flexible framework to support common data management tasks
- ❑ Extensible toolkit to promote integration and reuse across available performance tools
  - ○ Supported profile formats:
    TAU, CUBE, Dynaprof, HPC Toolkit, HPM Toolkit, gprof, mpiP, psrun (PerfSuite), others in development
  - ○ Supported DBMS:
    PostgreSQL, MySQL, Oracle, DB2, Derby/Cloudscape

# *PerfDMF Architecture*



K. Huck, A. Malony, R. Bell, A. Morris, "**Design and Implementation of a Parallel Performance Data Management Framework**," ICPP 2005.

# *Recent PerfDMF Development*

❑ Integration of XML metadata for each profile

  ❍ Common Profile Attributes

  ❍ Thread/process specific Profile Attributes

  ❍ Automatic collection of runtime information

  ❍ Any other data the user wants to collect can be added

   ➢ Build information
   ➢ Job submission information

  ❍ Two methods for acquiring metadata:

   ➢ TAU_METADATA() call from application
   ➢ Optional XML file added when saving profile to PerfDMF

  ❍ TAU Metadata XML schema is simple, easy to generate from scripting tools (no XML libraries required)

# *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)



K. Huck and A. Malony, "**PerfExplorer: A Performance Data Mining Framework For Large-Scale Parallel Computing**," SC 2005.

# PerfExplorer Analysis Methods

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

# *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

# PerfExplorer: Cross Experiment Analysis for S3D

# *PerfExplorer: S3D Total Runtime Breakdown*



WRITE_SAVEFILE

MPI_Wait

12,000 cores!

# *TAU Plug-Ins for Eclipse: Motivation*

❑ 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

# *Adding TAU to 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

# *TAU Eclipse Plug-In Features*

❑ **Performance data collection**

- ○ Graphical selection of TAU stub makefiles and compiler options
- ○ Automatic instrumentation, compilation and execution of target C, C++ or Fortran projects
- ○ Selective instrumentation via source editor and source outline views
- ○ Full integration with the Parallel Tools Platform (PTP) parallel launch system for performance data collection from parallel jobs launched within Eclipse

❑ **Performance data management**

- ○ Automatically place profile output in a PerfDMF database or upload to TAU-Portal
- ○ Launch ParaProf on profile data collected in Eclipse, with performance counters linked back to the Eclipse source editor

# *TAU Eclipse Plug-In Features*



PerfDMF

# *Choosing PAPI Counters with TAU's in Eclipse*

# *Future Plug-In Development*

❒ Integration of additional TAU components

    ○ Automatic selective instrumentation based on previous experimental results

    ○ Trace format conversion from within Eclipse

❒ Trace and profile visualization within Eclipse

❒ Scalability testing interface

❒ Additional user interface enhancements

# *KTAU Project*

□ Trend toward Extremely Large Scales

  ○ System-level influences are increasingly dominant performance bottleneck contributors

  ○ Application sensitivity at scale to the system (e.g., OS noise)

  ○ Complex I/O path and subsystems another example

  ○ Isolating system-level factors non-trivial

□ OS Kernel instrumentation and measurement is important to understanding system-level influences

□ But can we closely correlate observed application and OS performance?

□ KTAU / TAU (Part of the ANL/UO ZeptoOS Project)

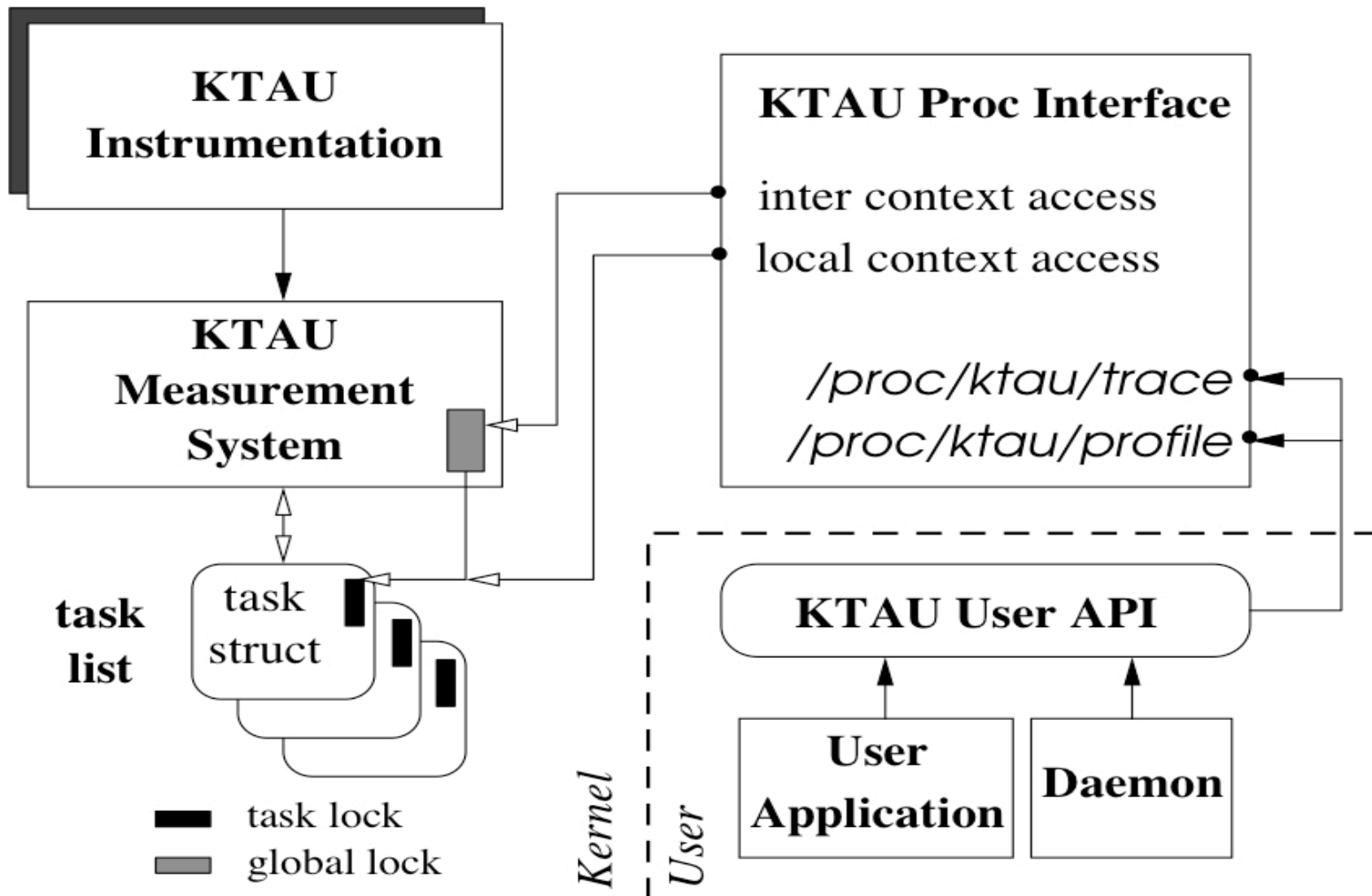  ○ Integrated methodology and framework to measure whole-system performance

# *Applying KTAU+TAU*

❏ How does *real* OS-noise affect *real* applications on target platforms?

   ○ Requires a tightly coupled performance measurement & analysis approach provided by KTAU+TAU

   ○ Provides an estimate of application slowdown due to Noise (and in particular, different noise-components - IRQ, scheduling, etc)

   ○ Can empower both application and the middleware and OS communities.

   ○ A. Nataraj, A. Morris, A. Malony, M. Sottile, P. Beckman, "The Ghost in the Machine : Observing the Effects of Kernel Operation on Parallel Application Performance", SC'07.

❏ Measuring and analyzing complex, multi-component I/O subsystems in systems like BG(L/P) (work in progress).

**KTAU Instrumentation**

**KTAU Measurement System**

task list

task struct

■ task lock
▬ global lock

**KTAU Proc Interface**

inter context access
local context access

/proc/ktau/trace
/proc/ktau/profile

*Kernel* | *User*

**KTAU User API**

**User Application**

**Daemon**

A. Nataraj, A. Malony, S. Shende, and A. Morris, "**Kernel-level Measurement for Integrated Performance Views: the KTAU Project**," *Cluster 2006*, distinguished paper.

# *TAU: Interoperability*

❑ **What we can offer other tools:**
  ○ Automated source-level instrumentation (tau_instrumentor, PDT)
  ○ ParaProf 3D profile browser
  ○ PerfDMF database, PerfExplorer cross-experiment analysis tool
  ○ Eclipse/PTP plugins for performance evaluation tools
  ○ Conversion of trace and profile formats
  ○ Kernel-level performance tracking using KTAU
  ○ Support for most HPC platforms, compilers, MPI-1,2 wrappers

❑ **What help we need from other projects:**
  ○ Common API for compiler instrumentation
    ➢ Scalasca/Kojak and VampirTrace compiler wrappers
    ➢ Intel, Sun, GNU, Hitachi, PGI, …
  ○ Support for sampling for hybrid instrumentation/sampling measurement
    ➢ HPCToolkit, PerfSuite
  ○ Portable, robust binary rewriting system that requires no root previleges
    ➢ DyninstAPI
  ○ Scalable communication framework for runtime data analysis
    ➢ MRNet, Supermon

# *Support Acknowledgements*

- □ US 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 Lab
- □ US Department of Defense (DoD)
- □ NSF Software and Tools for High-End Computing
- □ Research Centre Juelich
- □ TU Dresden
- □ Los Alamos National Laboratory
- □ ParaTools, Inc.

# *TAU Transport Substrate - Motivations*

- Transport Substrate
  - Enables movement of measurement-related data
  - TAU, in the past, has relied on shared file-system
- Some Modes of Performance Observation
  - Offline / Post-mortem observation and analysis
    - least requirements for a specialized transport
  - Online observation
    - long running applications, especially at scale
    - dumping to file-system can be suboptimal
  - Online observation with feedback into application
    - in addition, requires that the transport is bi-directional
- Performance observation problems and requirements are a function of the mode

# Requirements

- Improve performance of transport
  - NFS can be slow and variable
  - Specialization and remoting of FS-operations to front-end
- Data Reduction
  - At scale, cost of moving data too high
  - Sample in different domain (node-wise, event-wise)
- Control
  - Selection of events, measurement technique, target nodes
  - What data to output, how often and in what form?
  - Feedback into the measurement system, feedback into application
- Online, distributed processing of generated performance data
  - Use compute resource of transport nodes
  - Global performance analyses within the topology
  - Distribute statistical analyses
- Scalability, most important - All of above at very large scales

# *Approach and Prototypes*

- ❒ Measurement and measured data transport de-coupled
  - ○ Earlier, no such clear distinction in TAU
- ❒ Created abstraction to separate and hide transport
  - ○ *TauOutput*
- ❒ Did not create a custom transport for TAU(as yet)
  - ○ Use existing monitoring/transport capabilities
- ❒ TAUover: Supermon (Sottile and Minnich, LANL) and MRNET (Arnold and Miller, UWisc)
- ❒ A. Nataraj, M.Sottile, A. Morris, A. Malony, S. Shende "TAUoverSupermon: Low-overhead Online Parallel Performance Monitoring", Europar'07.
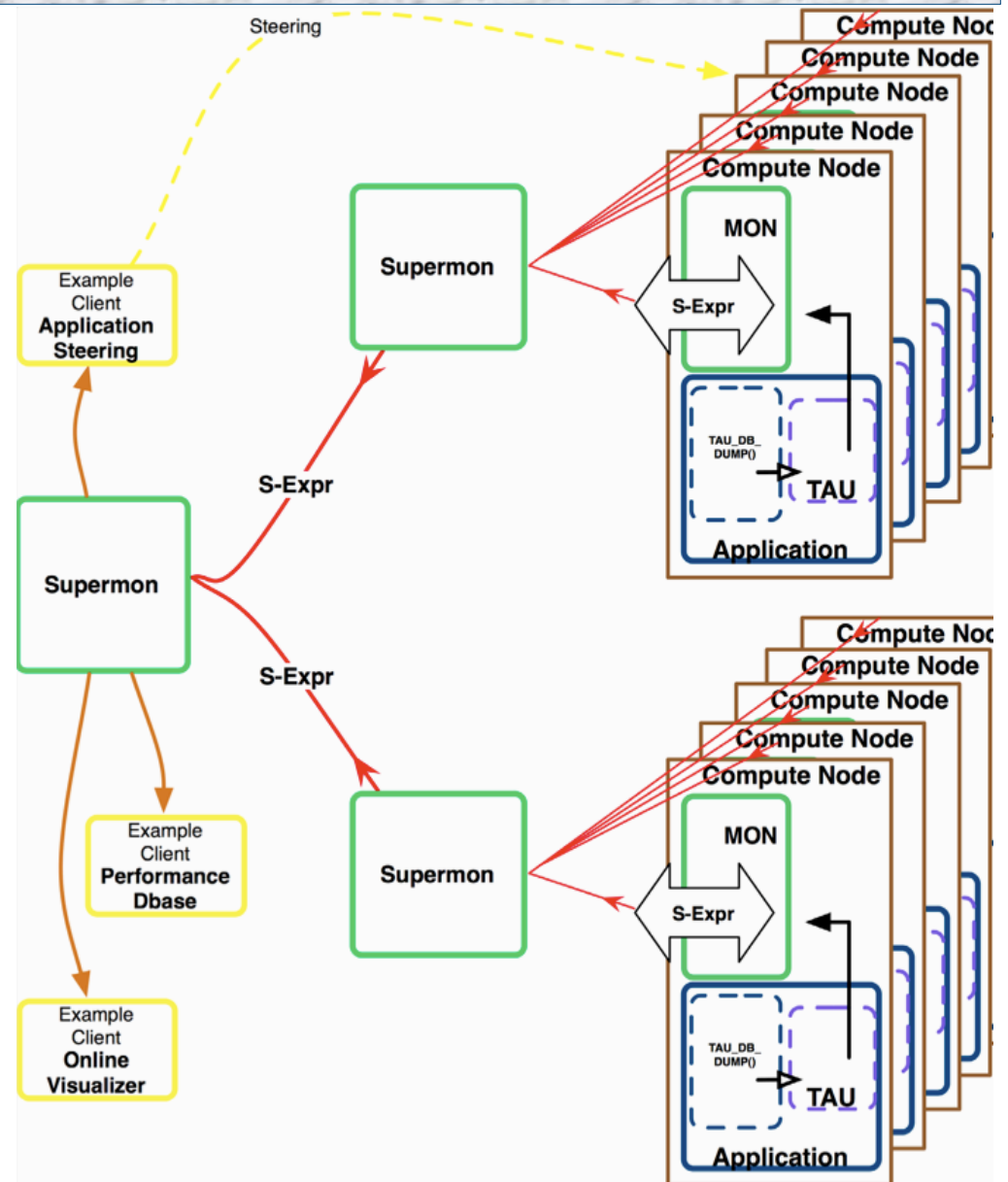
# *Rationale*

❒ Moved away from NFS

❒ Separation of concerns

  ○ Scalability, portability, robustness

  ○ Addressed independent of TAU

❒ Re-use existing technologies where appropriate

❒ Multiple bindings

  ○ Use different solutions best suited to particular platform

❒ Implementation speed

  ○ Easy, fast to create adapter that binds to existing transport

# *Substrate Architecture - High-level*

- Components
  - Front-End (FE)
  - Intermediate Nodes
  - Back-End (BE)

- NFS, Supermon, MRNet API

- Push-Pull model of data retrieval

- Figure shows *ToS* high-level view

# Substrate Architecture - Back-End

- ❒ Application calls into TAU
  - ○ Per-Iteration explicit call to output routine
  - ○ Periodic calls using alarm
- ❒ TauOutput object invoked
  - ○ Configuration specific: compile or runtime
  - ○ One per thread
- ❒ TauOutput mimics subset of FS-style operations
  - ○ Avoids changes to TAU code
  - ○ If required rest of TAU can be made aware of output type
- ❒ Non-blocking *recv* for control
- ❒ Back-end pushes, Sink pulls

**TAU-Instrumented Application**
**The MRNET Back-End**

Control
MRNET
Stream

Non-Blocking
check for
Control
Information

Data
MRNET
Stream

Packetize

TauMrnetOutput

**TauMrnetOutput Implementation**

Instrumented
application code
calls into TAU

TAU

file | smon | mrnet

TAU_DB_DUMP()

Transport