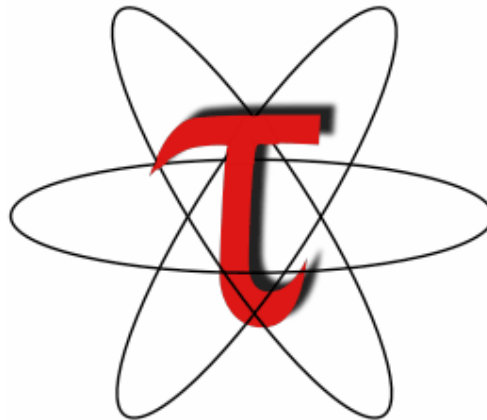# TAU Performance System
## Alan Morris, Sameer Shende, Allen D. Malony
## University of Oregon

*{amorris, sameer, malony}@cs.uoregon.edu*

UNIVERSITY
OF OREGON

# *Acknowledgements*

❑ Pete Beckman, ANL

❑ Holger Brunst and Wolfgang Nagel [TU Dresden]

❑ Bernd Mohr [Research Center Juelich, Germany]

❑ Aroon Nataraj, U. Oregon

❑ Suravee Suthikulpanit, U. Oregon

# *Outline*

- Overview of features
  - Instrumentation
  - Measurement (Profiling, Tracing)
  - Analysis tools
- New features in TAU
  - Runtime MPI shared library instrumentation
  - Workload characterization
- New features for BG/L
  - PAPI now supported
  - Open Trace Format (OTF), tau2otf
  - I/O node Linux kernel profiling with TAU (KTAU)

# *TAU Performance System*

- <u>T</u>uning and <u>A</u>nalysis <u>U</u>tilities (13+ year project effort)
- Performance system framework for HPC systems
    - Integrated, scalable, portable, flexible, and parallel
- Integrated toolkit for performance problem solving
    - Automatic instrumentation
    - Highly configurable measurement system with support for many flavors of profiling and tracing
    - Portable analysis and visualization tools
    - Performance data management and data mining
- http://www.cs.uoregon.edu/research/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)

❑ 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*)
- ○ Runtime Linking (LD_PRELOAD)

# *Automatic Instrumentation*

❑ We now provide compiler wrapper scripts

　○ Simply replace `mpxlf90` 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 = mpxlf90_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 $<
```

# *Profiling Options*

- ❒ Flat profiles
  - ○ Time (or counts) spent in each routine (nodes in callgraph).
  - ○ Exclusive/inclusive time, no. of calls, child calls
  - ○ Support for hardware counters (PAPI, PCL), multiple counters.
- ❒ Callpath Profiles
  - ○ Flat profiles, **plus**
  - ○ 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.
  - ○ Configurable callpath depth limit (`TAU_CALLPATH_DEPTH` environment variable)
- ❒ Phase based profiles
  - ○ 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` during "`IO`" phase

# *ParaProf – Manager Window*



performance database

derived performance metrics

# *ParaProf – Full Profile (Miranda)*

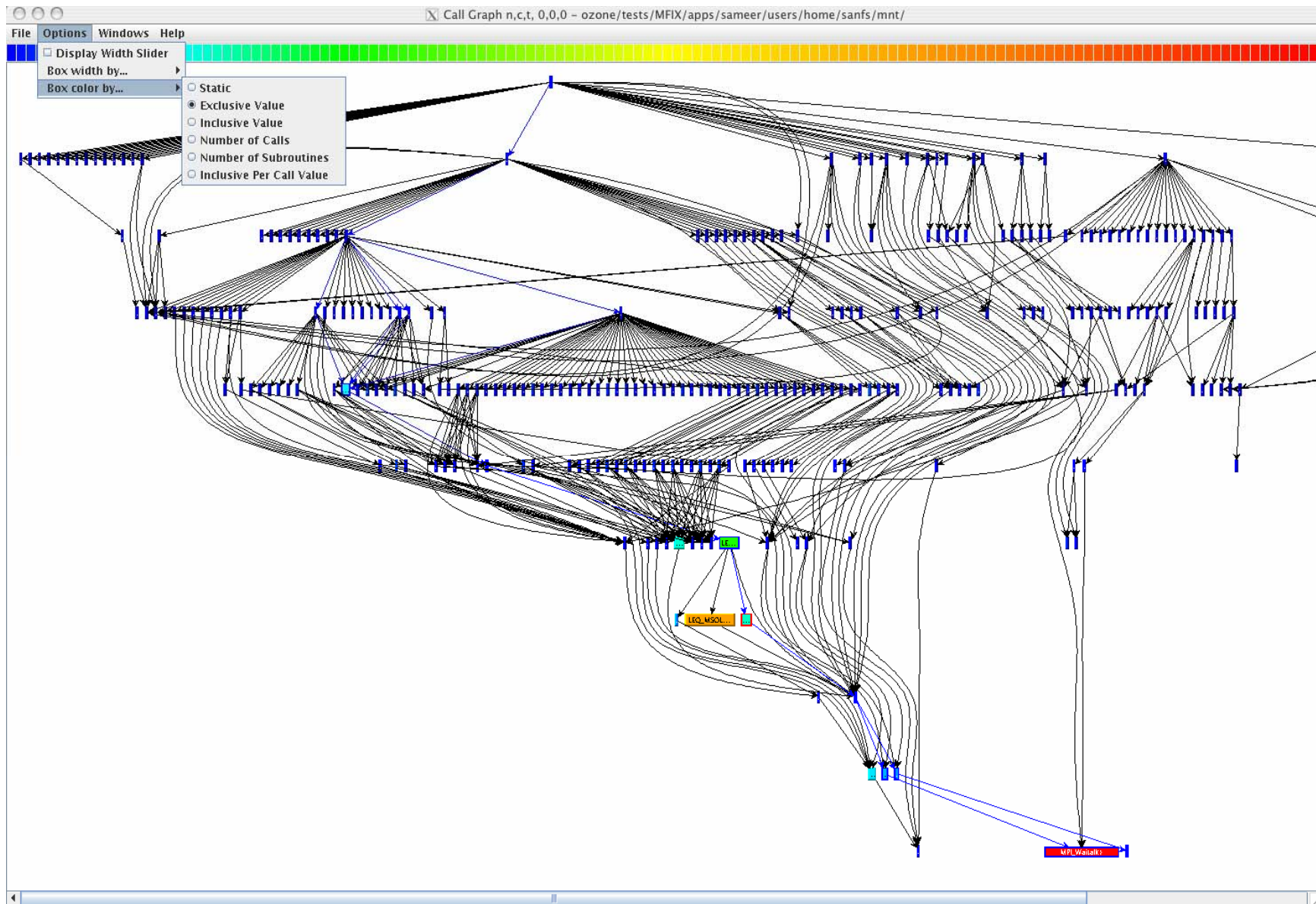# *ParaProf - Statistics Table (Uintah)*

Thread Statistics: n,c,t, 3,0,0 - /home/amorris/data/packed/uintah16.ppk

File   Options   Windows   Help

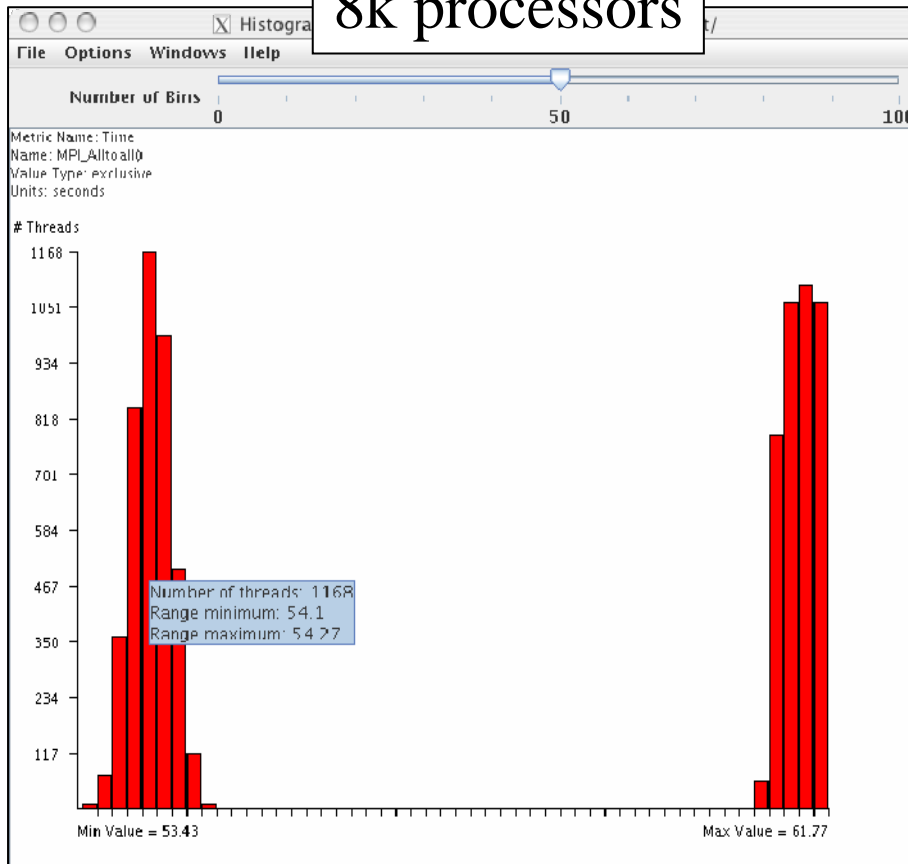| Name △ | P_WALL_CLOCK_TIME | Calls | Child Calls |
|---|---|---|---|
| main() void (int, char **) | 0.015 | 1 | 14 |
| Uintah::ProcessorGroup *Uintah::Parallel::getRootProcessor( | 0 | 1 | 0 |
| Uintah::SimpleSimulationController &Uintah::SimpleSimulatior | 0 | 1 | 0 |
| Uintah::SimulationController &Uintah::SimulationController::Si | 0 | 1 | 0 |
| bool Uintah::Parallel::usingMPI() | 0 | 1 | 0 |
| int Uintah::Parallel::getMPIRank() | 0 | 1 | 0 |
| void Uintah::OnDemandDataWarehouse::~OnDemandDataW | 0 | 2 | 0 |
| void Uintah::Parallel::determineIfRunningUnderMPI(int, char | 0.002 | 1 | 0 |
| void Uintah::Parallel::finalizeManager(Uintah::Parallel::Circur | 0.011 | 1 | 1 |
| void Uintah::Parallel::initializeManager(int &, char **&, const | 0.001 | 1 | 3 |
| MPI_Comm_rank() | 0 | 1 | 0 |
| MPI_Comm_size() | 0 | 1 | 0 |
| MPI_Init_thread() | 6.327 | 1 | 39 |
| void Uintah::Parallel::noThreading() | 0 | 1 | 0 |
| void Uintah::SimpleSimulationController::run() Uintah::Simple | 0.074 | 1 | 154 |
| MPIScheduler::actuallyCompile() | 0.109 | 2 | 44 |
| MPIScheduler::execute() | 27.68 | 11 | 3,460 |
| MPI_Reduce() | 0.001 | 40 | 40 |
| Uintah::DataWarehouse::ScrubMode Uintah::OnDemandD | 0 | 21 | 0 |
| Uintah::OnDemandDataWarehouse &Uintah::OnDemandD | 0 | 11 | 0 |
| bool Uintah::OnDemandDataWarehouse::timestepAborte | 0 | 10 | 0 |
| bool Uintah::OnDemandDataWarehouse::timestepRestar | 0 | 10 | 0 |
| bool Uintah::SimpleSimulationController::needRecompile | 0 | 10 | 0 |
| void Uintah::OnDemandDataWarehouse::get(Uintah::Red | 0.001 | 10 | 30 |
| void Uintah::OnDemandDataWarehouse::override(const L | 0.001 | 20 | 40 |

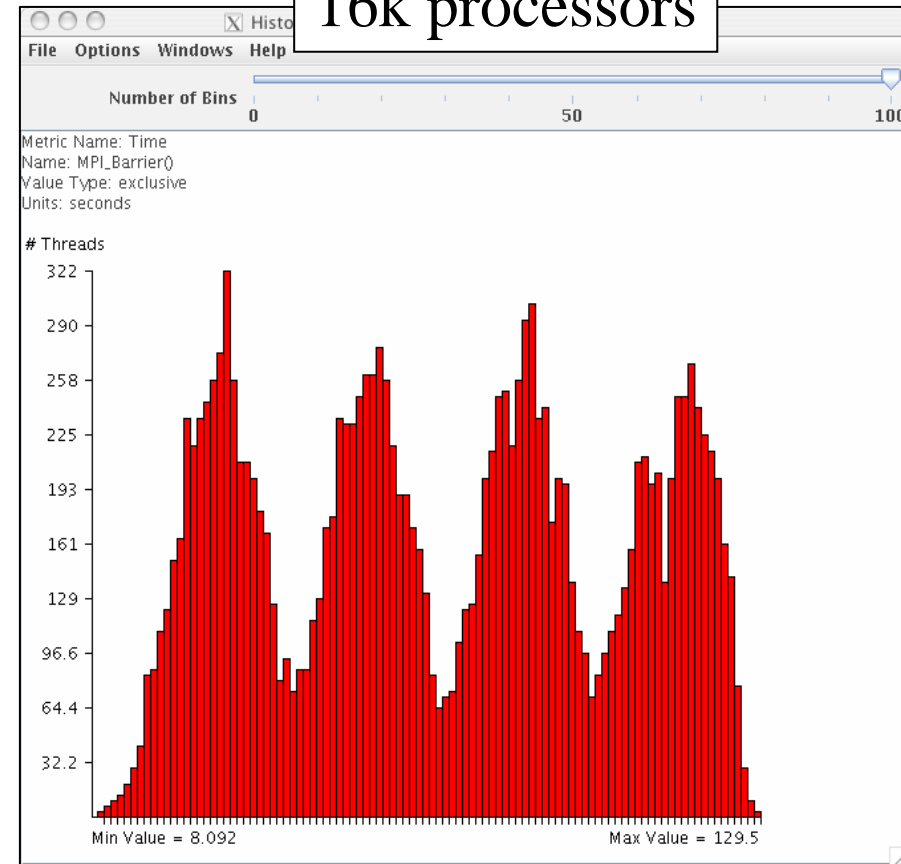# *ParaProf –Callgraph View (MFIX)*

# *ParaProf – Histogram View (Miranda)*
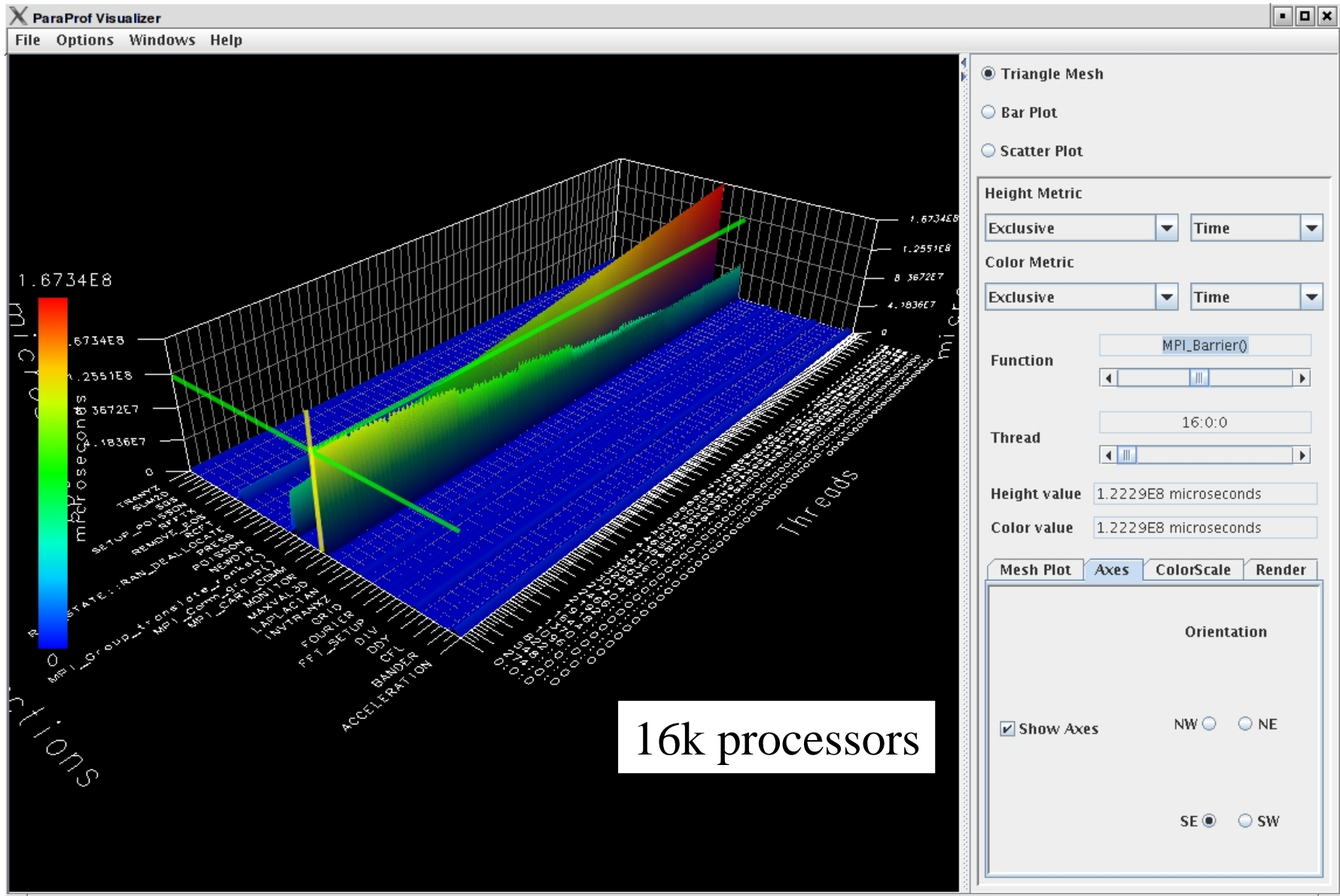
❑ Scalable 2D displays



8k processors

16k processors
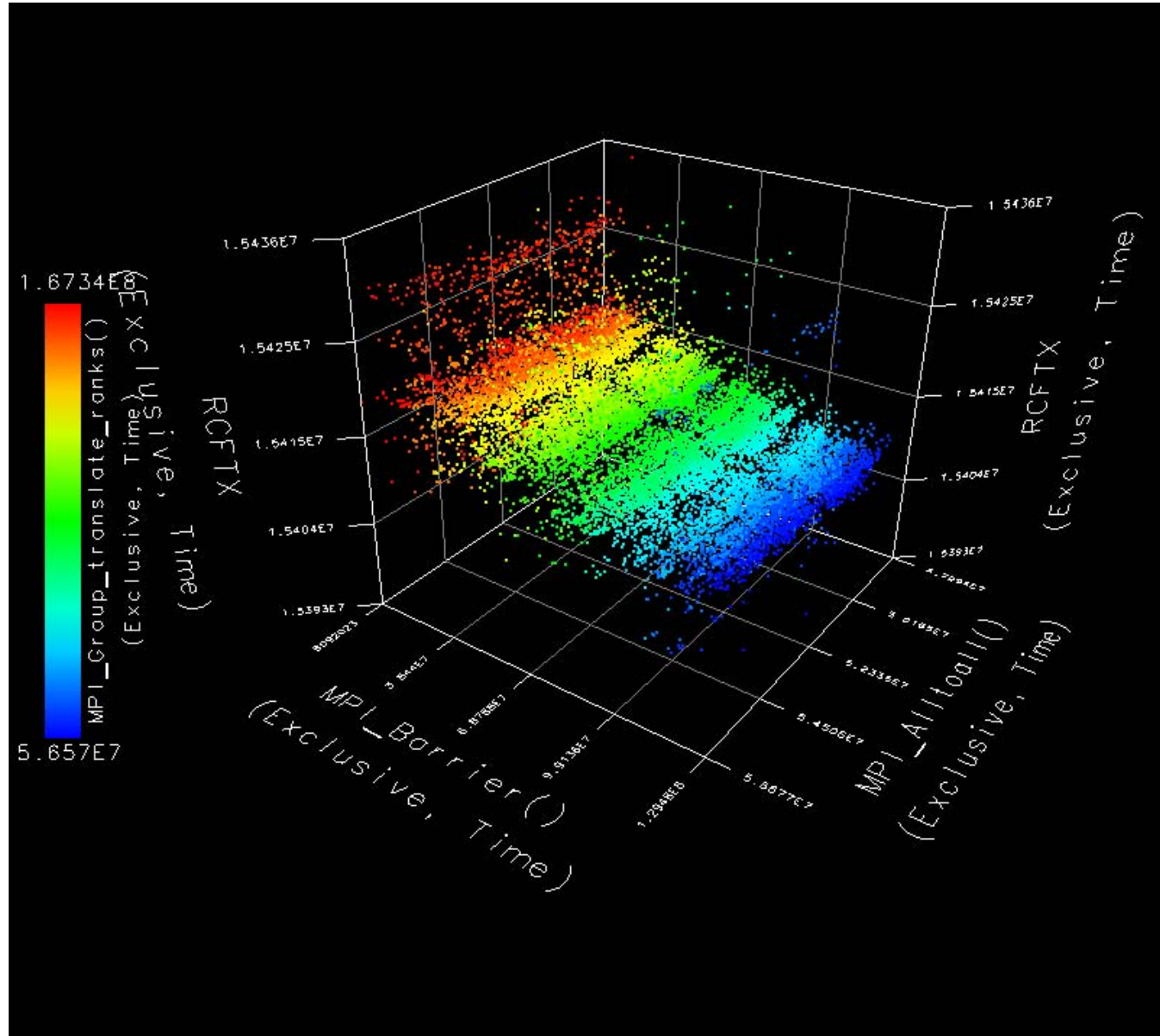
# *ParaProf – 3D Full Profile (Miranda)*



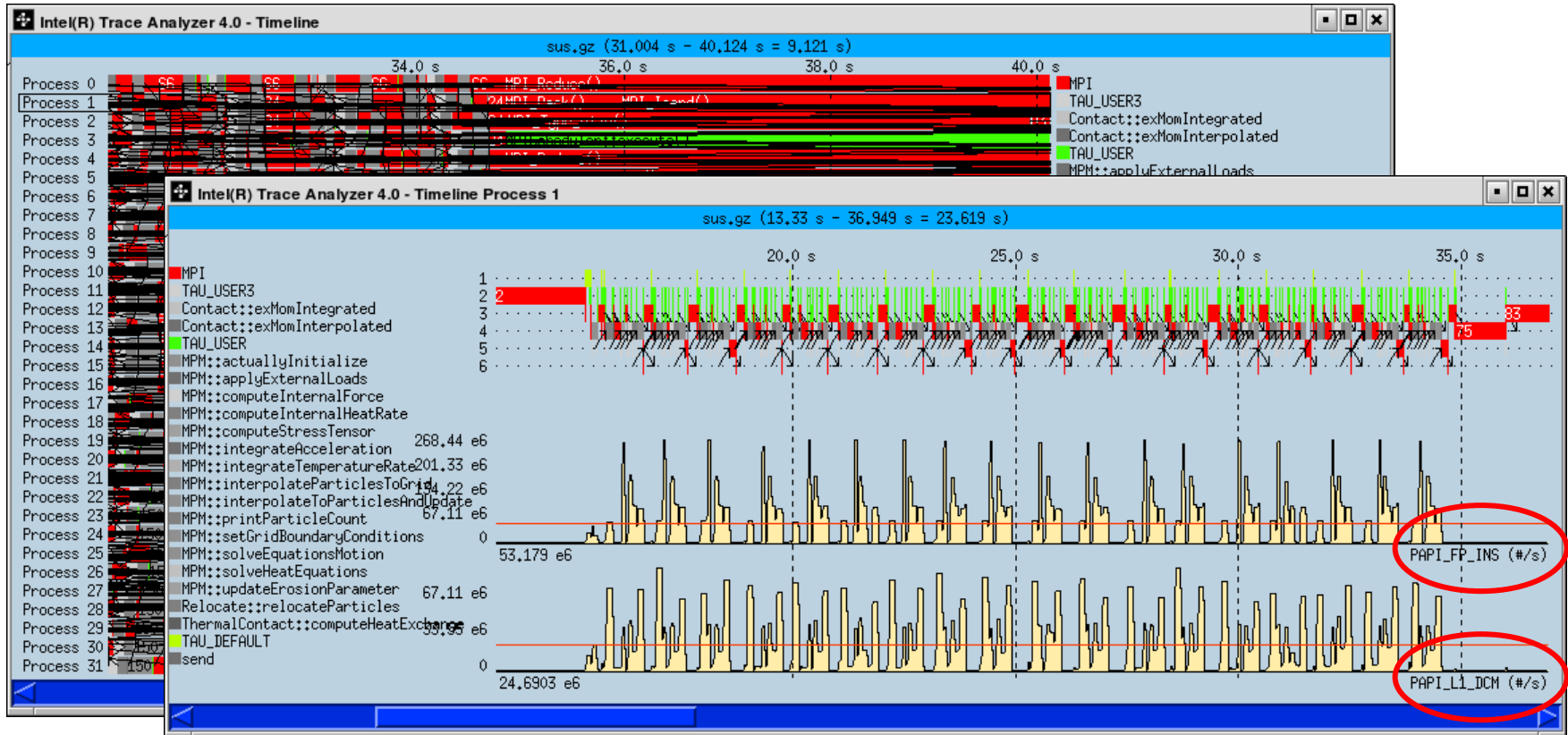16k processors

# *ParaProf – 3D Scatterplot (Miranda)*

- Each point is a "thread" of execution
- Relation between four routines shown at once

# *Tracing (Vampir)*

- ❑ Trace analysis provides in-depth understanding of temporal event and message passing relationships
- ❑ Traces can even store hardware counters

# *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

❒ Soon on AIX using `MPI_EUILIB/MPI_EUILIBPATH`

❒ Simply compile TAU with MPI support and prefix your MPI program with `tau_load.sh`

```
% mpirun –np 4 tau_load.sh a.out
```

❒ Requires shared library MPI

# *Workload Characterization*

- 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(int 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*

❏ Simple example, send/receive squared message sizes (0-32MB)

```c
#include <stdio.h>
#include <mpi.h>

int main(int argc, char **argv) {
  int rank, size, i, j;
  int buffer[16*1024*1024];
  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<16*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*

☐ Use `tau_load.sh` to instrument MPI routines (SGI Altix)

```
% icc mpi.c –lmpi

% mpirun –np 2 tau_load.sh a.out
```



SGI MPI

Intel MPI (SGI Altix)

# *Workload Characterization*

☐ MPI Results (NAS Parallel Benchmark 3.1, LU class D on 16 processors of SGI Altix)

# *Workload Characterization*

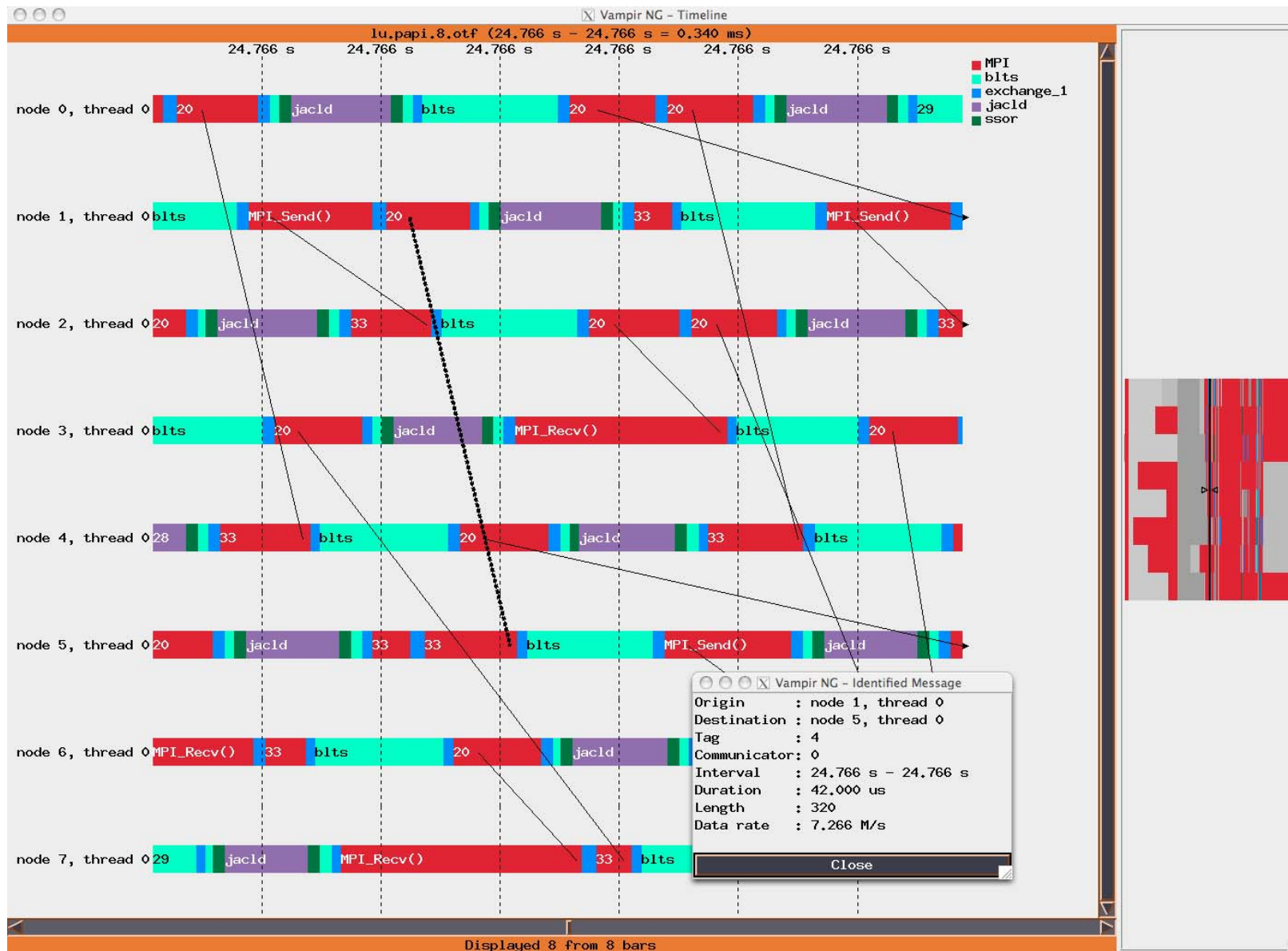❑ Two different message sizes (~3.3MB and ~4K)

# *Vampir, VNG, and OTF*

- ❏ Commercial trace based tools developed at ZiH, T.U. Dresden
  - ○ Wolfgang Nagel, Holger Brunst and others…
- ❏ Vampir Trace Visualizer (aka Intel ® Trace Analyzer v4.0)
  - ○ Sequential program
- ❏ Vampir Next Generation (VNG)
  - ○ Client (vng) runs on a desktop, server (vngd) on a cluster
  - ○ Parallel trace analysis
  - ○ Orders of magnitude bigger traces (more memory)
- ❏ Open Trace Format (OTF)
  - ○ Hierarchical trace format, efficient streams based parallel access with VNGD
  - ○ Replacement for proprietary formats such as STF
  - ○ Tracing library available on IBM BG/L platform
  - ○ Open Source release of OTF by SC06
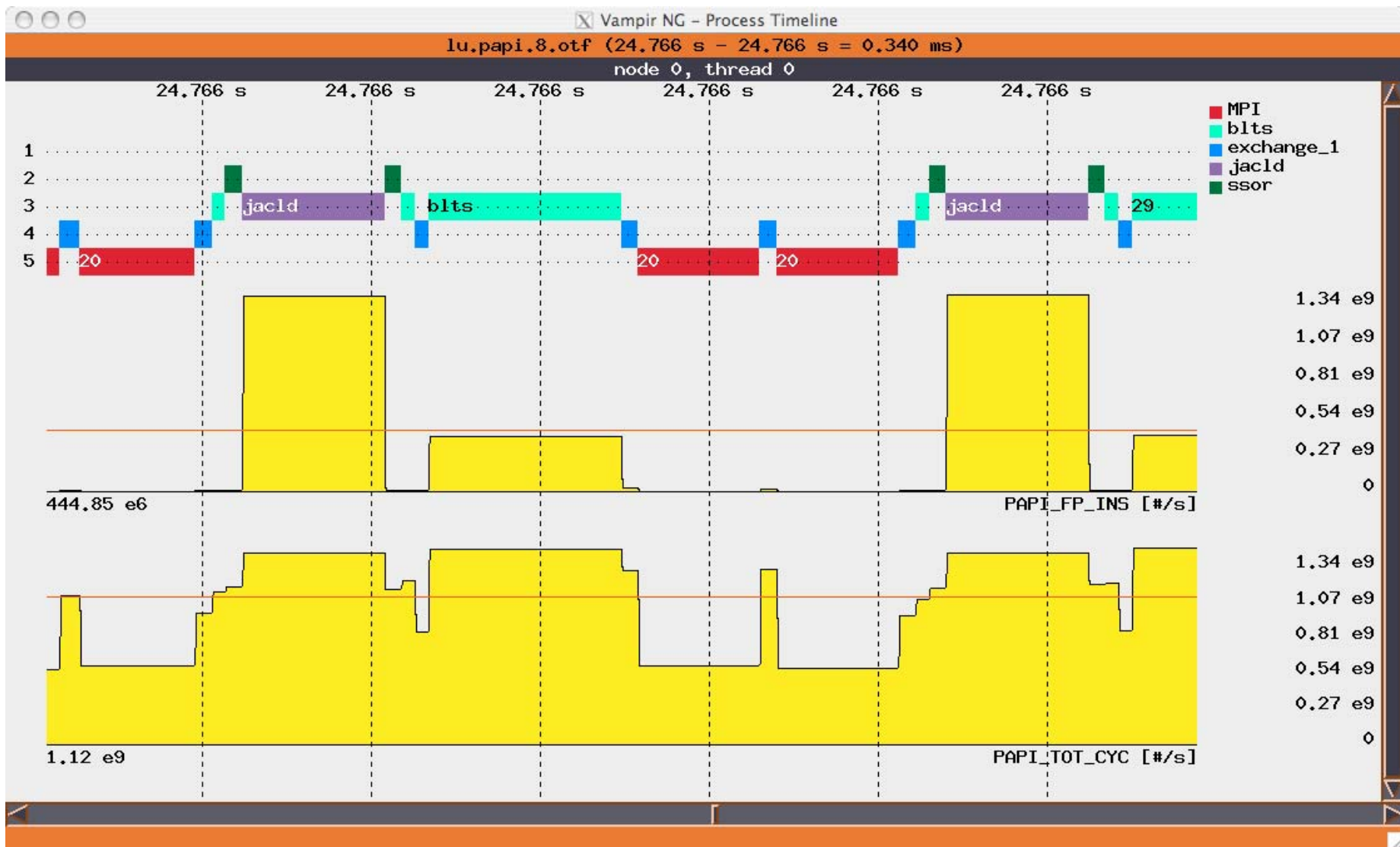- ❏ Development of OTF supported by LLNL contract

**http://www.vampir-ng.de**
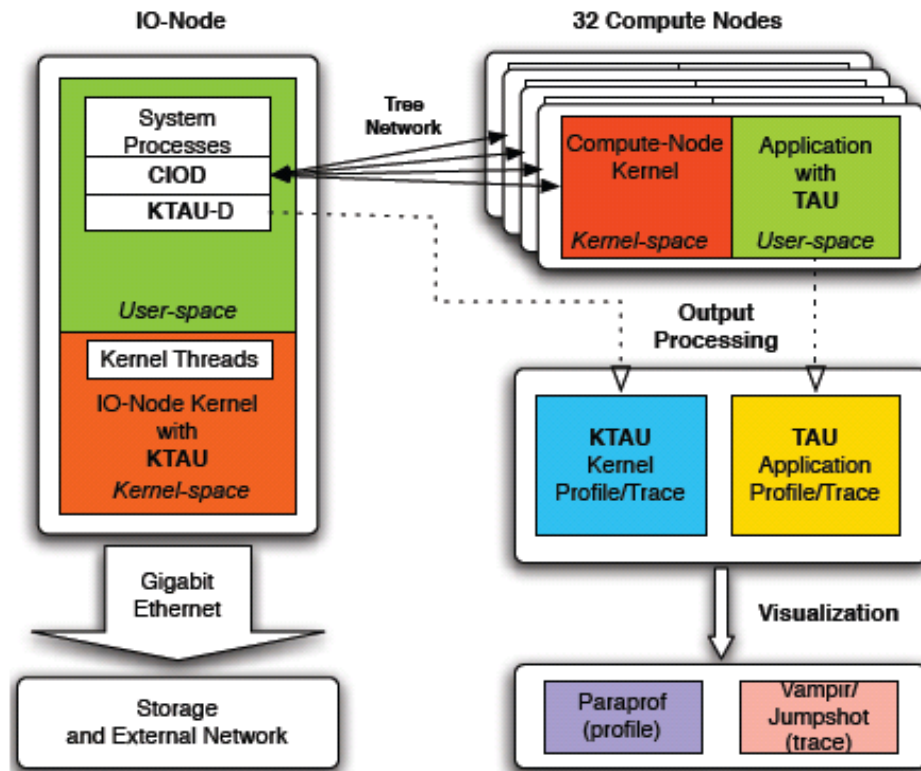
# *VNG Timeline Display (Miranda on BGL)*

# *VNG Process Timeline with PAPI Counters*

# *KTAU on BG/L*

- ❑ KTAU designed for Linux Kernel profiling
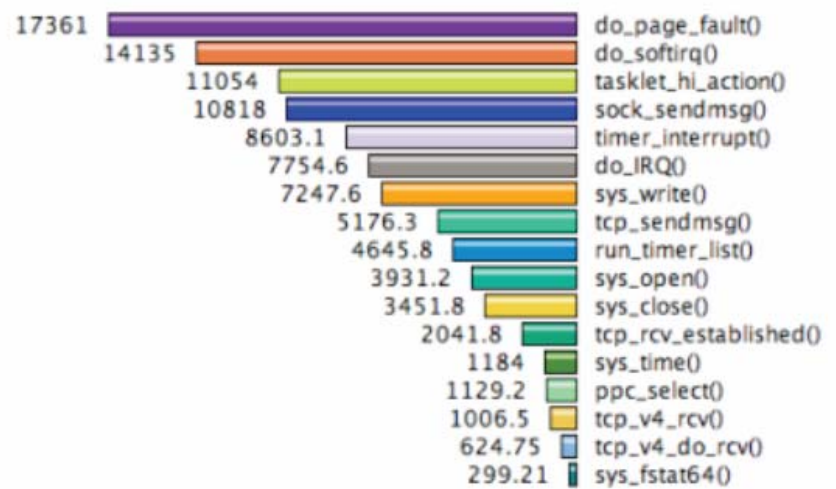- ❑ Provides merged application/system profile
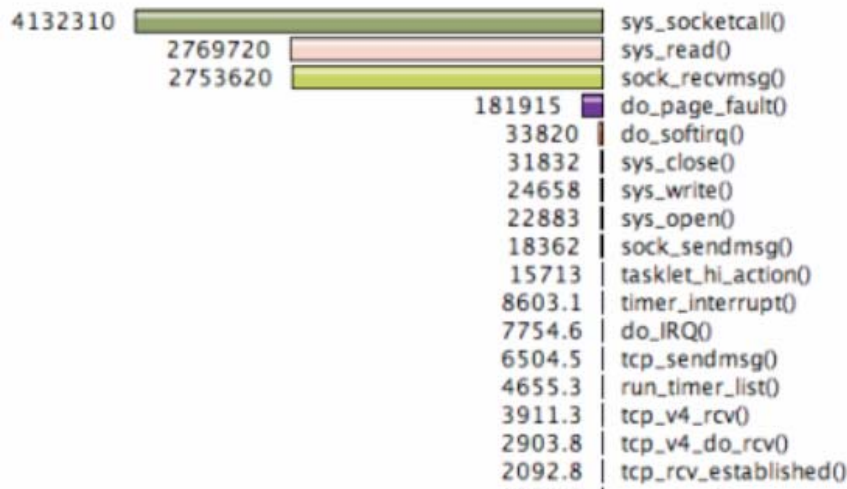- ❑ Runs on I/O-Node of BG/L

□ Current status

- Detailed I/O Node kernel profiling/tracing
- KTAU integrated into ZeptoOS build system
- KTAU-Daemon (KTAU-D) on I/O Node
  - ➢ Monitors system-wide and/or individual processes
- Visualization of trace/profile of ZeptoOS and CIOD
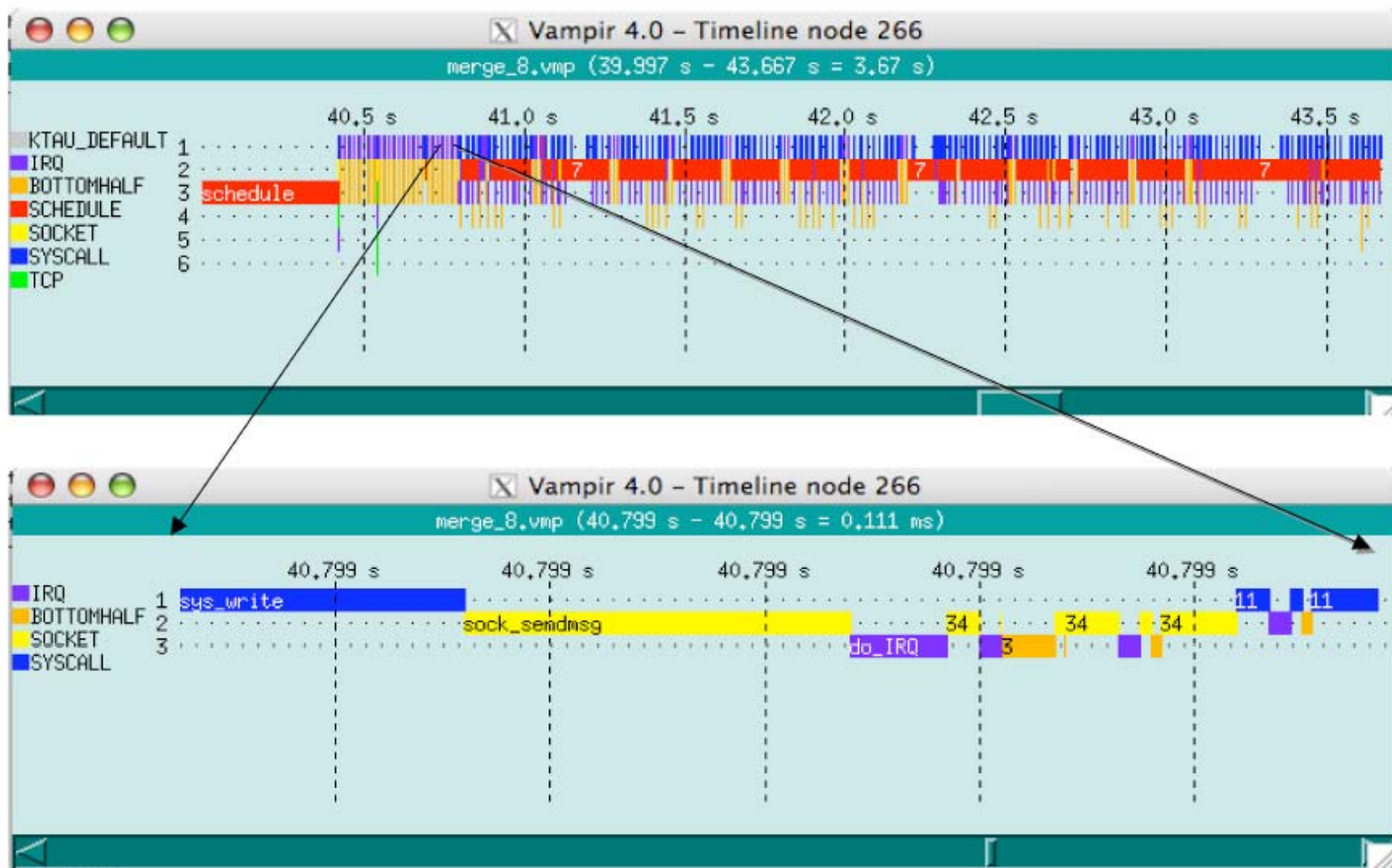  - ➢ Vampir/JumpShot (trace), and Paraprof (profile)

- Example of I/O Node profile data
- Numbers in microseconds, inclusive left, exclusive right

# KTAU on BG/L, Trace Data

# *Support Acknowledgements*

- Department of Energy (DOE)
  - Office of Science contracts
  - University of Utah ASC Level 1 sub-contract
  - LLNL ASC/NNSA Level 3 contract
  - LLNL ParaTools/GWT contract
- NSF
  - High-End Computing Grant
- T.U. Dresden, GWT
  - Dr. Wolfgang Nagel and Holger Brunst
- Research Centre Juelich
  - Dr. Bernd Mohr
- Los Alamos National Laboratory contracts