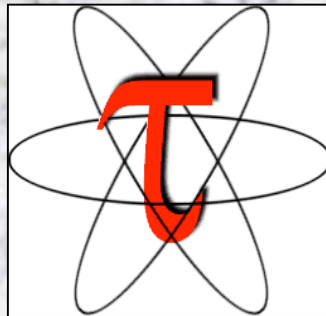


TAU Parallel Performance System

DOD UGC 2004 Tutorial



Part 2: TAU Components and Usage



Tutorial Outline – Part 2

TAU Components and Usage

- ❑ Configuration
- ❑ Instrumentation
 - Source, library, dynamic, multi-level
- ❑ Measurement
 - Profiling and tracing
- ❑ Analysis
 - ParaProf
 - Vampir
- ❑ Examples of use



How To Use TAU?

- ❑ Instrumentation
 - Application code and libraries
 - Selective instrumentation
- ❑ Install, compile, and link with TAU measurement library
 - Configure TAU system
 - Multiple configurations for different measurements options
 - Does not require change in instrumentation – just relink
 - Selective measurement control
- ❑ Execute “experiments” to produce performance data
 - Performance data generated at end or during execution
- ❑ Use analysis tools to look at performance results



Using TAU in Practice

- ❑ **Install TAU**
 - % configure ; make clean install
- ❑ **Instrument application**
 - TAU Profiling API
- ❑ **Typically modify application makefile**
 - Include TAU's stub makefile, modify variables
- ❑ **Set environment variables**
 - Directory where profiles/traces are to be stored
- ❑ **Execute application**
 - % mpirun -np <procs> a.out
- ❑ **Analyze performance data**
 - ParaProf, vampir, pprof, paraver ...

TAU System Configuration



□ configure [OPTIONS]

- `{-c++=<CC>, -cc=<cc>}` Specify C++ and C compilers
- `{-pthread, -sproc}` Use pthread or SGI sproc threads
- `-openmp` Use OpenMP threads
- `-jdk=<dir>` Specify Java instrumentation (JDK)
- `-opari=<dir>` Location of Opari OpenMP tool
- `-papi=<dir>` Location of PAPI
- `-pdt=<dir>` Location of PDT
- `-dyninst=<dir>` Location of DynInst Package
- `-mpi[inc/lib]=<dir>` Specify MPI library instrumentation
- `-python[inc/lib]=<dir>` Specify Python instrumentation
- `-epilog=<dir>` Specify location of EPILOG



Configuring TAU

```
% configure [options]  
% make clean install
```

- ❑ Creates `<arch>/lib/Makefile.tau<options>` stub Makefile
- ❑ Creates `<arch>/lib/libTau<options.a [.so]` libraries
- ❑ Defines a single configuration of TAU
- ❑ Attempts to automatically detect architecture



Examples: TAU Configuration

- ❑ Use TAU with xlc_r and pthread library under AIX
Enable TAU profiling (default)
 - ./configure -c++=xlc_r -pthread
- ❑ Enable both TAU profiling and tracing
 - ./configure -TRACE -PROFILE
- ❑ Use IBM's xlc_r and xlc_r compilers with PAPI, PDT, MPI packages and multiple counters for measurements
 - ./configure -c++=xlc_r -cc=xlc_r
-papi=/usr/local/packages/papi
-pdt=/usr/local/pdtoolkit-3.0 -arch=ibm64
-mpiinc=/usr/lpp/ppe.poe/include
-mpilib=/usr/lpp/ppe.poe/lib -MULTIPLECOUNTERS
- ❑ Typically configure multiple measurement libraries



Instrumentation Alternatives

- ❑ Manual instrumentation at the source
 - Use TAU API appropriate for source language
- ❑ Automatic source-level instrumentation
 - Source rewriting
 - Directive rewriting (e.g., for OpenMP)
- ❑ Library instrumentation
 - Typically done at source level
 - Wrapper interposition library (e.g., PMPI)
- ❑ Binary Instrumentation
 - Pre-execution or runtime binary rewriting
- ❑ Dynamic runtime instrumentation

TAU Measurement API for C/C++



□ Initialization and runtime configuration

- TAU_PROFILE_INIT(*argc*, *argv*);
TAU_PROFILE_SET_NODE(*myNode*);
TAU_PROFILE_SET_CONTEXT(*myContext*);
TAU_PROFILE_EXIT(*message*);
TAU_REGISTER_THREAD();

□ Function and class methods

- TAU_PROFILE(*name*, *type*, *group*);

□ Template

- TAU_TYPE_STRING(*variable*, *type*);
TAU_PROFILE(*name*, *type*, *group*);
CT(*variable*);

□ User-defined timing

- TAU_PROFILE_TIMER(*timer*, *name*, *type*, *group*);
TAU_PROFILE_START(*timer*);
TAU_PROFILE_STOP(*timer*);

TAU Measurement API for C/C++ (continued)



□ User-defined events

- TAU_REGISTER_EVENT(**variable**, **event_name**);
TAU_EVENT(**variable**, **value**);
TAU_PROFILE_STMT(**statement**);

□ Mapping

- TAU_MAPPING(**statement**, **key**);
TAU_MAPPING_OBJECT(**funcIdVar**);
TAU_MAPPING_LINK(**funcIdVar**, **key**);
- TAU_MAPPING_PROFILE (**funcIdVar**);
TAU_MAPPING_PROFILE_TIMER(**timer**, **funcIdVar**);
TAU_MAPPING_PROFILE_START(**timer**);
TAU_MAPPING_PROFILE_STOP(**timer**);

□ Reporting

- TAU_REPORT_STATISTICS();
TAU_REPORT_THREAD_STATISTICS();



Example: Manual Instrumentation (C++)

```
#include <TAU.h>
int main(int argc, char **argv)
{
    TAU_PROFILE("int main(int, char **)", " ", TAU_DEFAULT); /* name,type,group */
    TAU_PROFILE_INIT(argc, argv);
    TAU_PROFILE_SET_NODE(0); /* for sequential programs */
    foo();
    return 0;
}
int foo(void)
{
    TAU_PROFILE("int foo(void)", " ", TAU_DEFAULT); // measures entire foo()
    TAU_PROFILE_TIMER(t, "foo(): for loop", "[23:45 file.cpp]", TAU_USER);
    TAU_PROFILE_START(t);
    for(int i = 0; i < N ; i++){
        work(i);
    }
    TAU_PROFILE_STOP(t);
    // other statements in foo ...
}
```



Example: Manual Instrumentation (C)

```
#include <TAU.h>
int main(int argc, char **argv)
{
    TAU_PROFILE_TIMER(tmain, "int main(int, char **)", " ", TAU_DEFAULT);
    TAU_PROFILE_INIT(argc, argv);
    TAU_PROFILE_SET_NODE(0); /* for sequential programs */
    TAU_PROFILE_START(tmain);
    foo();
    ...
    TAU_PROFILE_STOP(tmain);
    return 0;
}
int foo(void)
{
    TAU_PROFILE_TIMER(t, "foo()", " ", TAU_USER);
    TAU_PROFILE_START(t);
    for(int i = 0; i < N ; i++){
        work(i);
    }
    TAU_PROFILE_STOP(t);
}
```



Example: Manual Instrumentation (F90)

```
PROGRAM SUM_OF_CUBES
  integer profiler(2)
  save profiler
  INTEGER :: H, T, U
  call TAU_PROFILE_INIT()
  call TAU_PROFILE_TIMER(profiler, 'PROGRAM SUM_OF_CUBES')
  call TAU_PROFILE_START(profiler)
  call TAU_PROFILE_SET_NODE(0)
  ! This program prints all 3-digit numbers that
  ! equal the sum of the cubes of their digits.
  DO H = 1, 9
    DO T = 0, 9
      DO U = 0, 9
        IF (100*H + 10*T + U == H**3 + T**3 + U**3) THEN
          PRINT "(3I1)", H, T, U
        ENDIF
      END DO
    END DO
  END DO
  call TAU_PROFILE_STOP(profiler)
END PROGRAM SUM_OF_CUBES
```



Instrumenting Multithreaded Applications

```
#include <TAU.h>
void * threaded_function(void *data)
{
    TAU_REGISTER_THREAD(); // Before any other TAU calls
    TAU_PROFILE("void * threaded_function", " ", TAU_DEFAULT);
    work();
}
int main(int argc, char **argv)
{
    TAU_PROFILE("int main(int, char **)", " ", TAU_DEFAULT);
    TAU_PROFILE_INIT(argc, argv);
    TAU_PROFILE_SET_NODE(0); /* for sequential programs */
    pthread_attr_t attr;
    pthread_t tid;

    pthread_attr_init(&attr);
    pthread_create(&tid, NULL, threaded_function, NULL);
    return 0;
}
```



Compiling: TAU Makefiles

- ❑ Include TAU Stub Makefile (<arch>/lib) in the user's Makefile
- ❑ 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_DISABLE** TAU's dummy F90 stub library
- ❑ **Note:** Not including TAU_DEFS in CFLAGS disables instrumentation in C/C++ programs (**TAU_DISABLE** for f90)



Example: Including TAU Makefile

```
include /usr/tau/sgi64/lib/Makefile.tau-pthread-kcc
CXX = $(TAU_CXX)
CC  = $(TAU_CC)
CFLAGS = $(TAU_DEFS) $(TAU_INCLUDE)
LIBS = $(TAU_LIBS)
OBJS = ...
TARGET= a.out
TARGET: $(OBJS)
    $(CXX) $(LDFLAGS) $(OBJS) -o $@ $(LIBS)
.cpp.o:
    $(CC) $(CFLAGS) -c $< -o $@
```



Example: Including TAU Makefile (F90)

```
include $PET_HOME/PTOOLS/tau-2.13.5/rs6000/lib/Makefile.tau-pdt
F90 = $(TAU_F90)
FFLAGS = -I<dir>
LIBS = $(TAU_LIBS) $(TAU_CXXLIBS)
OBJS = ...
TARGET= a.out
TARGET: $(OBJS)
    $(F90) $(LDFLAGS) $(OBJS) -o $@ $(LIBS)
.f.o:
    $(F90) $(FFLAGS) -c $< -o $@
```



Example: Using TAU's malloc Wrapper Library

```
include $PET_HOME/PTOOLS/tau-2.13.5/rs6000/lib/Makefile.tau-pdt
CC=$(TAU_CC)
CFLAGS=$(TAU_DEFS) $(TAU_INCLUDE) $(TAU_MEMORY_INCLUDE)
LIBS = $(TAU_LIBS)
OBJS = f1.o f2.o ...
TARGET= a.out
TARGET: $(OBJS)
    $(F90) $(LDFLAGS) $(OBJS) -o $@ $(LIBS)
.c.o:
    $(CC) $(CFLAGS) -c $< -o $@
```



TAU's malloc() / free() wrapper

- ❑ Used to capture measurements of memory usage

```
#include <TAU.h>
#include <malloc.h>
int main(int argc, char **argv)
{
    TAU_PROFILE("int main(int, char **)", " ", TAU_DEFAULT);

    int *ary = (int *) malloc(sizeof(int) * 4096);

    // TAU's malloc wrapper library replaces this call automatically
    // when $(TAU_MEMORY_INCLUDE) is used in the Makefile.

    ...
    free(ary);
    // other statements in foo ...
}
```

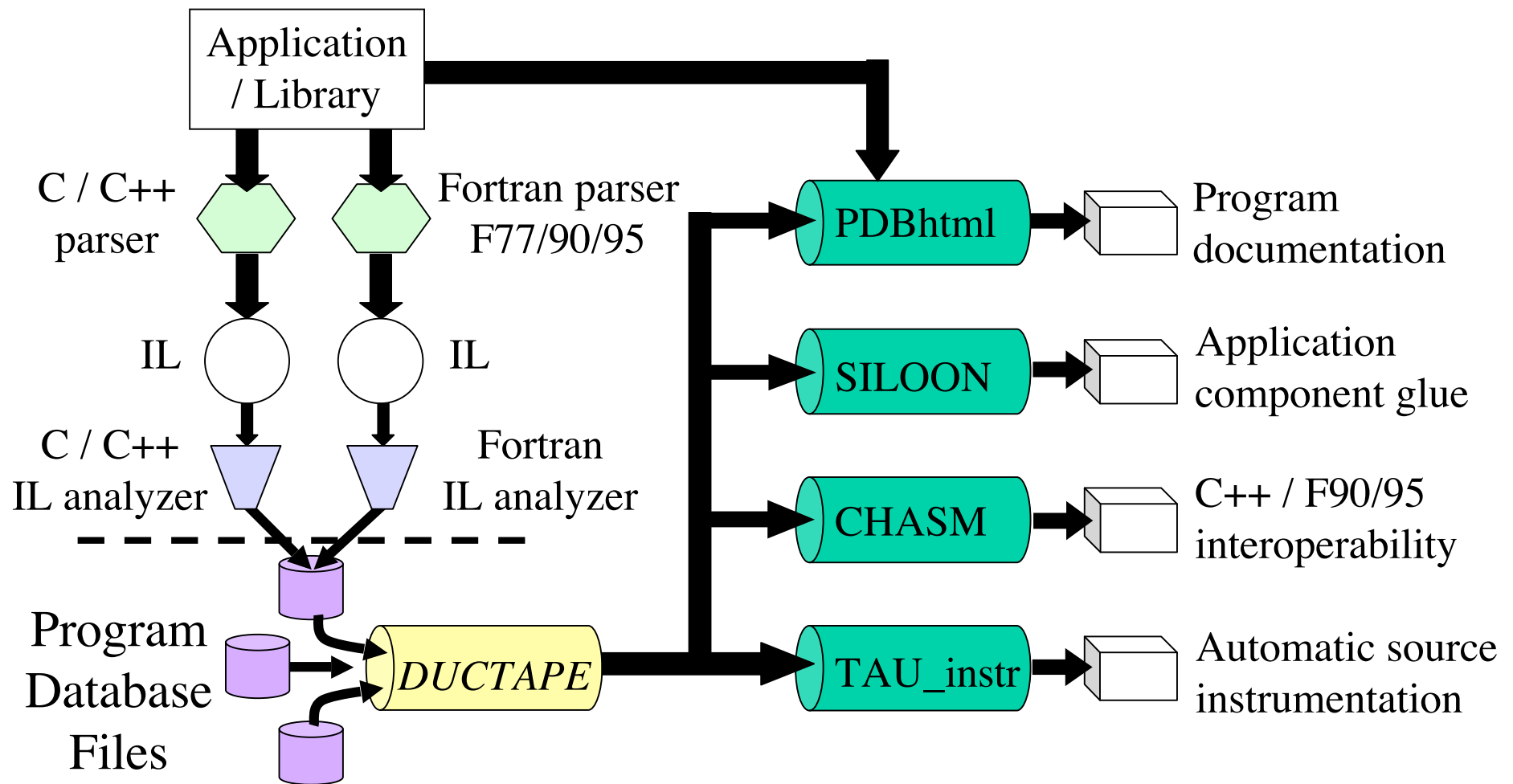


Program Database Toolkit (PDT)

- ❑ Program code analysis framework
 - develop source-based analysis tools
- ❑ *High-level interface* to source code information
- ❑ *Integrated toolkit* for source code parsing, database creation, and database query
 - Commercial grade front-end parsers
 - Portable IL analyzer, database format, and access API
 - Open software approach for tool development
- ❑ Multiple source languages
- ❑ Implement automatic performance instrumentation tools
 - *tau_instrumentor* for automatic source instrumentation



Program Database Toolkit (PDT)





PDT Components

□ Language front end

- Edison Design Group (EDG): C, C++, Java
- Mutek Solutions Ltd.: F77, F90
- Cleanscape FortranLint F95 parser/analyzer
- Creates an intermediate-language (IL) tree

□ IL Analyzer

- Processes the intermediate language (IL) tree
- Creates “program database” (PDB) formatted file

□ DUCTAPE (Bernd Mohr, ZAM, Germany)

- C++ program Database Utilities and Conversion Tools
Application Environment
- C++ library to process the PDB for PDT applications

□ Intel/KAI C++ headers for std. C++ library



Contents of PDB files

- ❑ Source file names
- ❑ Routines, Classes, Methods, Templates, Macros, Modules
- ❑ Parameters, signature
- ❑ Entry and exit point information (return)
- ❑ Location information for all of the above
- ❑ Static callgraph
- ❑ Header file inclusion tree
- ❑ Statement-level information
 - Loops, if-then-else, switch ...



PDT 3.2 Functionality

- ❑ C++ statement-level information implementation
 - for, while loops, declarations, initialization, assignment...
 - PDB records defined for most constructs
- ❑ DUCTAPE
 - Processes PDB 1.x, 2.x, 3.x uniformly
- ❑ PDT applications
 - XMLgen
 - PDB to XML converter
 - Used for CHASM and CCA tools
 - PDBstmt
 - Statement callgraph display tool



PDT 3.2 Functionality (continued)

- ❑ **Cleanscape Flint parser fully integrated for F90/95**
 - Flint parser (f95parse) is very robust
 - Produces PDB records for TAU instrumentation (stage 1)
 - Linux (x86, IA-64, Opteron, Power4), HP Tru64, IBM AIX, Cray X1,T3E, Solaris, SGI, Apple, Windows, Power4 Linux (IBM Blue Gene/L compatible)
 - Full PDB 2.0 specification (stage 2) [SC'04]
 - Statement level support (stage 3) [SC'04]
- ❑ **PDT 3.2 released in June 3, 2004**
 - Important bug fixes
- ❑ **<http://www.cs.uoregon.edu/research/paracomp/pdtoolkit>**



Configuring PDT

Step I: Configure PDT:

```
% configure -arch=IRIX64 -CC  
% make clean; make install
```

Builds <pdt_dir>/<arch>/bin/cxxparse, cparse, f90parse and f95parse

Builds <pdt_dir>/<arch>/lib/libpdb.a. See <pdt_dir>/README file

Step II: Configure TAU with PDT for auto-instrumentation of source code:

```
% configure -arch=IRIX64 -c++=CC -cc=cc  
  -pdt=/usr/contrib/TAU/pdtoolkit-3.0  
% make clean; make install
```

Builds <taudir>/<arch>/bin/tau_instrumentor,

<taudir>/<arch>/lib/Makefile.tau<options> and libTau<options>.a

See <taudir>/INSTALL file



TAU Makefile for PDT

```
include /usr/tau/include/Makefile
CXX = $(TAU_CXX)
CC  = $(TAU_CC)
PDTPARSE = $(PDTDIR)/$(CONFIG_ARCH)/bin/cxxparse
TAUINSTR = $(TAUROOT)/$(CONFIG_ARCH)/bin/tau_instrumentor
CFLAGS = $(TAU_DEFS) $(TAU_INCLUDE)
LIBS = $(TAU_LIBS)
OBJS = ...
TARGET= a.out
TARGET: $(OBJS)
    $(CXX) $(LDFLAGS) $(OBJS) -o $@ $(LIBS)
.cpp.o:
    $(PDTPARSE) $<
    $(TAUINSTR) $*.pdb $< -o $*.inst.cpp -f select.dat
    $(CC) $(CFLAGS) -c $*.inst.cpp -o $@
```



tau_instrumentor: A PDT Instrumentation Tool

```
% 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
```



tau_reduce: Rule-Based Overhead Analysis

- Analyze the performance data to determine events with high (relative) overhead performance measurements
- Create a select list for excluding those events
- Rule grammar (used in *tau_reduce* tool)

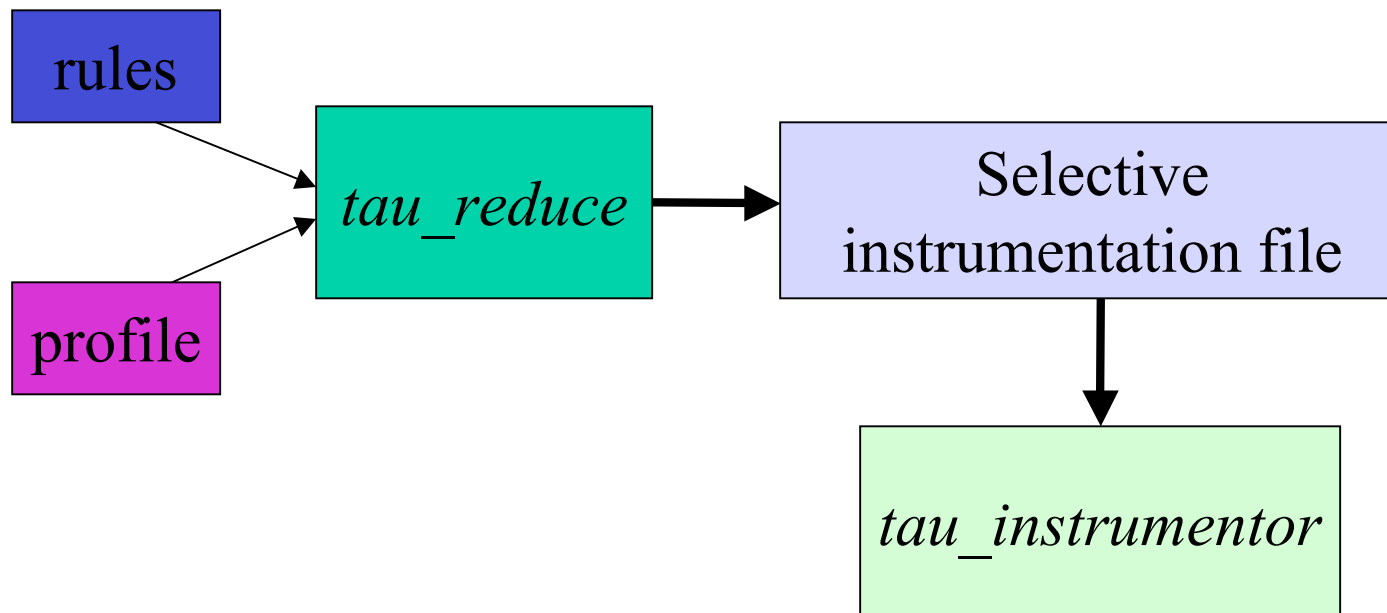
[GroupName:] Field Operator Number

- *GroupName* indicates rule applies to events in group
- *Field* is a event metric attribute (from profile statistics)
 - numcalls, numsubs, percent, usec, cumusec, count [PAPI], totalcount, stdev, usecs/call, counts/call
- *Operator* is one of >, <, or =
- *Number* is any number
- Compound rules possible using & between simple rules



Iterative Instrumentation Process

- ❑ Reads profile files and rules
- ❑ Creates selective instrumentation file
 - Specifies which routines should be excluded
 - Input to *tau_instrumentor*





Examples: Instrumentation Rules

- ❑ #Exclude all events that are members of TAU_USER
#and use less than 1000 microseconds
 - TAU_USER:usec < 1000
- ❑ #Exclude all events that have less than 100
#microseconds and are called only once
 - usec < 1000 & numcalls = 1
- ❑ #Exclude all events that have less than 1000 usecs per
#call OR have a (total inclusive) percent less than 5
 - usecs/call < 1000
percent < 5
- ❑ Scientific notation can be used
 - usec>1000 & numcalls>400000 & usecs/call<30 & percent>25



Instrumentation Control

- ❑ Selection of which performance events to observe
 - Could depend on scope, type, level of interest
 - Could depend on instrumentation overhead
- ❑ How is selection supported in instrumentation system?
 - No choice
 - Include / exclude routine and file lists (TAU)
 - Environment variables
 - Static vs. dynamic
- ❑ Problem: Controlling instrumentation of small routines
 - High relative measurement overhead
 - Significant intrusion and possible perturbation



Example: tau_reduce

- ❑ *tau_reduce* implements overhead reduction in TAU
- ❑ Consider *klargest* example
 - Find *k*th largest element in a *N* elements
 - Compare two methods: *quicksort*, *select_kth_largest*
- ❑ *i = 2324, N = 1000000 (uninstrumented)*
 - *quicksort*: (wall clock) = 0.188511 secs
 - *select_kth_largest*: (wall clock) = 0.149594 secs
 - Total: (P3/1.2GHz, *time*) = 0.340u 0.020s 0:00.37
- ❑ Execution with all routines instrumented
- ❑ Execution with rule-based selective instrumentation
 - `usec>1000 & numcalls>400000 & usecs/call<30 & percent>25`



Reducing Instrumentation on One Processor

Before selective instrumentation reduction

NODE 0;CONTEXT 0;THREAD 0:

%Time	Exclusive msec	Inclusive msec	#Call	#Subrs	Inclusive usec/call	Name
100.0	13	4,982	1	4	4982030	int main
93.5	3,223	4,659	4.20241E+06	1.40268E+07	1	void quicksort
62.9	0.00481	3,134	5	5	626839	int kth_largest_qs
36.4	137	1,813	28	450057	64769	int select_kth_largest
33.6	150	1,675	449978	449978	4	void sort_5elements
28.8	1,435	1,435	1.02744E+07	0	0	void interchange
0.4	20	20	1	0	20668	void setup
0.0	0.0118	0.0118	49	0	0	int ceil

After selective instrumentation reduction

NODE 0;CONTEXT 0;THREAD 0:

%Time	Exclusive msec	Inclusive total msec	#Call	#Subrs	Inclusive usec/call	Name
100.0	14	383	1	4	383333	int main
50.9	195	195	5	0	39017	int kth_largest_qs
40.0	153	153	28	79	5478	int select_kth_largest
5.4	20	20	1	0	20611	void setup
0.0	0.02	0.02	49	0	0	int ceil



TAU's MPI Wrapper Interposition Library

- ❑ Uses standard MPI Profiling Interface
 - Provides name shifted interface
 - *MPI_Send* \Leftrightarrow *PMPI_Send*
 - weak bindings
- ❑ Instrument MPI wrapper library
 - Use TAU measurement API
- ❑ Interpose TAU's MPI wrapper library
 - Replace **-lmpi** by **-lTauMpi -lpmpi -lmpi**
- ❑ No change to the source code!
 - Just **re-link** the application to generate performance data



Using MPI Wrapper Interposition Library

Step I: Configure TAU with MPI:

```
% configure -mpiinc=/usr/include -mpilib=/usr/lib64  
-arch=sgi64 -c++=CC -cc=cc  
-pdt=/usr/contrib/TAU/pdtoolkit-3.0  
% make clean; make install
```

**Builds <taudir>/<arch>/lib/libTauMpi<options>,
<taudir>/<arch>/lib/Makefile.tau<options> and libTau<options>.a**



MPI Library Instrumentation (MPI_Send)

```
int  MPI_Send(...) /* TAU redefines MPI_Send */
...
{
  int  returnVal, typesize;
  TAU_PROFILE_TIMER(tautimer, "MPI_Send()", " ", TAU_MESSAGE);
  TAU_PROFILE_START(tautimer);
  if (dest != MPI_PROC_NULL) {
    PMPI_Type_size(datatype, &typesize);
    TAU_TRACE_SENDMSG(tag, dest, typesize*count);
  }
  /* Wrapper calls PMPI_Send */
  returnVal = PMPI_Send(buf, count, datatype, dest, tag, comm);
  TAU_PROFILE_STOP(tautimer);
  return returnVal;
}
```



Including TAU's Stub Makefile (C, C++)

```
include /usr/tau/sgi64/lib/Makefile.tau-mpi
CXX = $(TAU_CXX)
CC  = $(TAU_CC)
CFLAGS = $(TAU_DEFS) $(TAU_INCLUDE) $(TAU_MPI_INCLUDE)
LIBS = $(TAU_MPI_LIBS) $(TAU_LIBS)
LD_FLAGS = $(TAU_LDFLAGS)
OBJS = ...
TARGET= a.out
TARGET: $(OBJS)
    $(CXX) $(LDFLAGS) $(OBJS) -o $@ $(LIBS)
.cpp.o:
    $(CC) $(CFLAGS) -c $< -o $@
```



Including TAU's Stub Makefile (Fortran)

```
include $PET_HOME/PTOOLS/tau-2.13.5/rs6000/lib/Makefile.tau-mpi-pdt
F90 = $(TAU_F90)
CC = $(TAU_CC)
LIBS = $(TAU_MPI_LIBS) $(TAU_LIBS) $(TAU_CXXLIBS)
LD_FLAGS = $(TAU_LDFLAGS)
OBJS = ...
TARGET= a.out
TARGET: $(OBJS)
    $(CXX) $(LDFLAGS) $(OBJS) -o $@ $(LIBS)
.f.o:
    $(F90) $(FFLAGS) -c $< -o $@
```




Including TAU's Stub Makefile with PAPI

```
include $PET_HOME/PTOOLS/tau-2.13.5/rs6000/lib/Makefile.tau-  
papiwallclock-multiplecounters-papivirtual-mpi-papi-pdt  
CC = $(TAU_CC)  
LIBS = $(TAU_MPI_LIBS) $(TAU_LIBS) $(TAU_CXXLIBS)  
LD_FLAGS = $(TAU_LDFLAGS)  
OBJS = ...  
TARGET= a.out  
TARGET: $(OBJS)  
        $(CXX) $(LDFLAGS) $(OBJS) -o $@ $(LIBS)  
.f.o:  
        $(F90) $(FFLAGS) -c $< -o $@
```



TAU Makefile for PDT with MPI and F90

```
include $PET/PTOOLS/tau-2.13.5/rs6000/lib/Makefile.tau-mpi-pdt
FCOMPILE = $(TAU_F90) $(TAU_MPI_INCLUDE)
PDTF95PARSE = $(PDTDIR)/$(PDTARCHDIR)/bin/f95parse
TAUINSTR = $(TAUROOT)/$(CONFIG_ARCH)/bin/tau_instrumentor
PDB=merged.pdb
COMPILE_RULE= $(TAU_INSTR) $(PDB) $< -o $*.inst.f -f sel.dat;\
    $(FCOMPILE) $*.inst.f -o $@;
LIBS = $(TAU_MPI_FLIBS) $(TAU_LIBS) $(TAU_CXXLIBS)
OBJS = f1.o f2.o f3.o ...
TARGET= a.out
TARGET: $(PDB) $(OBJS)
    $(TAU_F90) $(LDFLAGS) $(OBJS) -o $@ $(LIBS)
$(PDB): $(OBJS:.o=.f)
    $(PDTF95PARSE) $(OBJS:.o=.f) $(TAU_MPI_INCLUDE) -o$(PDB)
# This expands to f95parse *.f -I.../mpi/include -omerged.pdb
.f.o:
    $(COMPILE_RULE)
```



Instrumentation of OpenMP Constructs



- ❑ **OpenMP Pragma And Region Instrumentor**
- ❑ Source-to-Source translator to insert *POMP* calls around OpenMP constructs and API functions
- ❑ **Supports**
 - Fortran77 and Fortran90, OpenMP 2.0
 - C and C++, OpenMP 1.0
 - *POMP* Extensions
 - Preserves source information (**#line line file**)
 - Measurement library implementations
 - EPILOG, TAU POMP, DPOMP (IBM)
- ❑ **Work in Progress**
 - Investigating standardization through OpenMP Forum



POMP OpenMP Performance Tool Interface

□ OpenMP Instrumentation

- OpenMP Directive Instrumentation
- OpenMP Runtime Library Routine Instrumentation

□ POMP Extensions

- Runtime Library Control (**init**, **finalize**, **on**, **off**)
- (Manual) User Code Instrumentation (**begin**, **end**)
- Conditional Compilation (**#ifdef _POMP, !\$P**)
- Conditional / Selective Transformations
(**[no]instrument**)



Example: !\$OMP PARALLEL DO

```
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)
```



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: Opari Directive Instrumentation

```
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);  
#line 261 "stommel.c"
```



Example: TAU POMP Implementation

```
TAU_GLOBAL_TIMER(tfor, "for enter/exit",  
                 "[OpenMP]", OpenMP);  
  
void pomp_for_enter(OMPRegDescr* r) {  
    #ifdef TAU_AGGREGATE_OPENMP_TIMINGS  
        TAU_GLOBAL_TIMER_START(tfor)  
    #endif  
    #ifdef TAU_OPENMP_REGION_VIEW  
        TauStartOpenMPRegionTimer(r);  
    #endif  
}  
  
void pomp_for_exit(OMPRegDescr* r) {  
    #ifdef TAU_AGGREGATE_OPENMP_TIMINGS  
        TAU_GLOBAL_TIMER_STOP(tfor)  
    #endif  
    #ifdef TAU_OPENMP_REGION_VIEW  
        TauStopOpenMPRegionTimer(r);  
    #endif  
}
```




OPARI: Makefile Template (C, C++)

```
OMPCC = ...          # insert C OpenMP compiler here
OMPCXX = ...         # insert C++ OpenMP compiler here

.c.o:
    opari $<
    $(OMPCC) $(CFLAGS) -c $*.mod.c

.cc.o:
    opari $<
    $(OMPCXX) $(CXXFLAGS) -c $*.mod.cc

opari.init:
    rm -rf opari.rc

opari.tab.o:
    opari -table opari.tab.c
    $(CC) -c opari.tab.c

myprog: opari.init myfile*.o ... opari.tab.o
    $(OMPCC) -o myprog myfile*.o opari.tab.o -lpomp

myfile1.o: myfile1.c myheader.h
myfile2.o: ...
```



OPARI: Makefile Template (Fortran)

```
OMPF77 = ...           # insert f77 OpenMP compiler here
OMPF90 = ...           # insert f90 OpenMP compiler here

.f.o:
    opari $<
    $(OMPF77) $(CFLAGS) -c $*.mod.F

.f90.o:
    opari $<
    $(OMPF90) $(CXXFLAGS) -c $*.mod.F90

opari.init:
    rm -rf opari.rc

opari.tab.o:
    opari -table opari.tab.c
    $(CC) -c opari.tab.c

myprog: opari.init myfile*.o ... opari.tab.o
    $(OMPF90) -o myprog myfile*.o opari.tab.o -lpomp

myfile1.o: myfile1.f90
myfile2.o: ...
```



OPARI: Basic Usage (F90)

- ❑ Reset OPARI state information
 - `rm -f opari.rc`
- ❑ Call OPARI for each input source file
 - `opari file1.f90`
 - ...
 - `opari fileN.f90`
- ❑ Generate OPARI runtime table, compile it with ANSI C
 - `opari -table opari.tab.c`
 - `cc -c opari.tab.c`
- ❑ Compile modified files `*.mod.f90` using OpenMP
- ❑ Link the resulting object files, the OPARI runtime table `opari.tab.o` and the TAU POMP RTL



Using Opari with TAU

Step I: Configure KOJAK/opari

[Download from <http://www.fz-juelich.de/zam/kojak/>]

```
% cd kojak-0.99; cp mf/Makefile.defs.sgi Makefile.defs;  
  edit Makefile
```

```
% make
```

Builds opari

Step II: Configure TAU with Opari (used here with MPI and PDT)

```
% configure -opari=/usr/contrib/TAU/kojak-0.99/opari  
  -mpiinc=/usr/include -mpilib=/usr/lib64  
  -arch=sgi64 -c++=CC -cc=cc  
  -pdt=/usr/contrib/TAU/pdtoolkit-3.0  
% make clean; make install
```



Dynamic Instrumentation

- ❑ TAU uses DyninstAPI for runtime code patching
- ❑ *tau_run* (mutator) loads measurement library
- ❑ Instruments mutatee
 - Application binary
 - Uses TAU-developed instrumentation specification
- ❑ 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-4.0.2/core; make
```

Set DyninstAPI environment variables (including LD_LIBRARY_PATH)

Step II: Configure TAU with Dyninst

```
% configure -dyninst=/usr/local/dyninstAPI-4.0.2
```

```
% 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

NOTE: All compilers and platforms are not yet supported (work in progress)



Using TAU with Python Applications

Step I: Configure TAU with Python

```
% configure -pythoninc=/usr/include/python2.2/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>]
```



Example: Python Manual Instrumentation

□ Python measurement API and dynamic library

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

import pytau
from time import sleep

x = pytau.profileTimer('`Timer A`')
pytau.start(x)

print " Sleeping for 5 seconds "
sleep(5)

pytau.stop(x)

Running:
% setenv PYTHONPATH <tau>/<arch>/lib
% ./application.py
```




Example: Python Automatic Instrumentation

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

import tau
from time import sleep

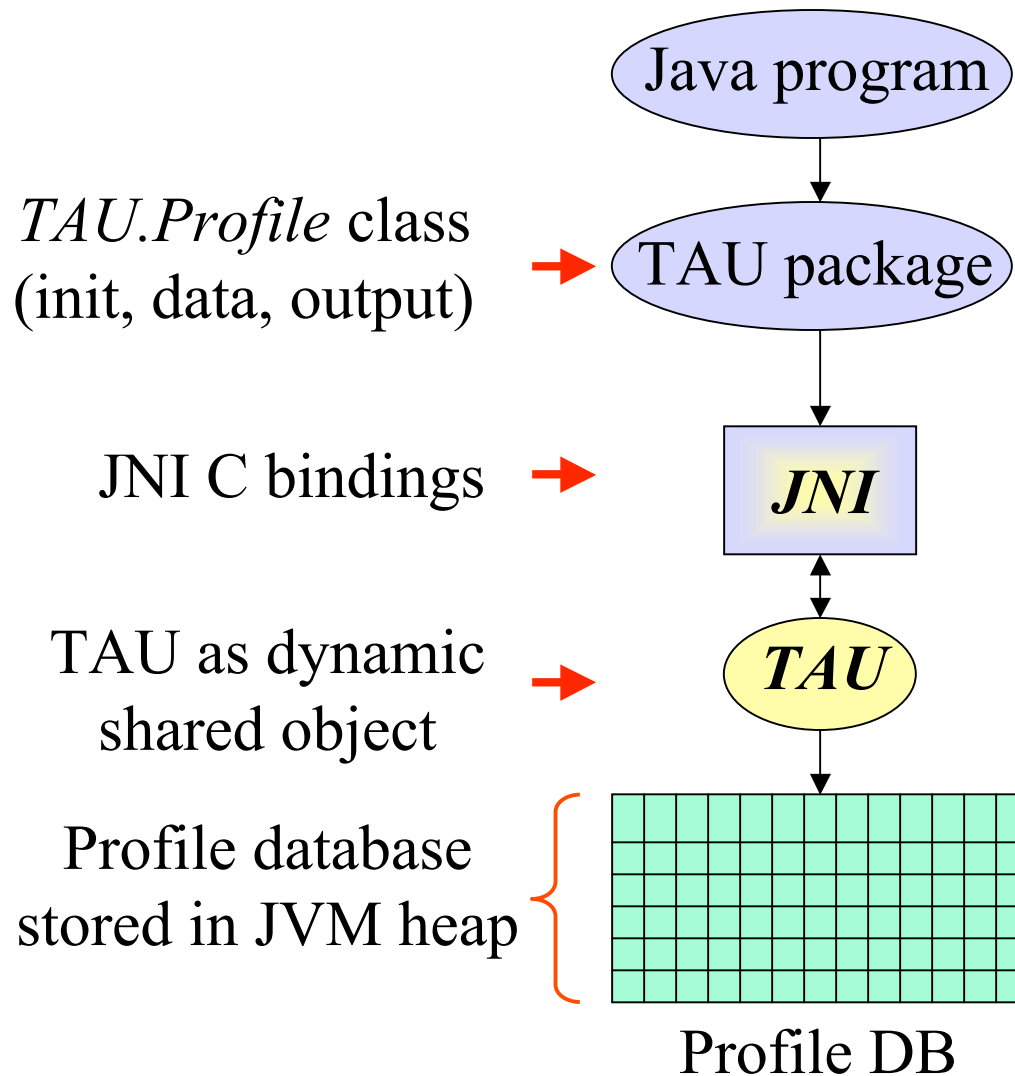
def f2():
    print " In f2: Sleeping for 2 seconds "
    sleep(2)
def f1():
    print " In f1: Sleeping for 3 seconds "
    sleep(3)

def OurMain():
    f1()
tau.run('OurMain()')
```

Running:

```
% setenv PYTHONPATH <tau>/<arch>/lib
% ./auto.py
Instruments OurMain, f1, f2, print...
```

TAU Java Source Instrumentation Architecture



- ❑ Any code section can be measured
- ❑ Portability
- ❑ Measurement options
 - Profiling, tracing
- ❑ Limitations
 - Source access only
 - Lack of thread information
 - Lack of node information



Java Source-Level Instrumentation

- ❑ TAU Java package
- ❑ User-defined events
- ❑ *TAU.Profile* class for new “timers”
 - Start/Stop
- ❑ Performance data output at end

```
emacs@neutron.cs.uoregon.edu
Buffers Files Tools Edit Search Mule Java Help
import TAU.*;
import mpi.*;

public class Life {

    static TAU.Profile blocktimer= new TAU.Profile("Life compute local block info",\
    "", "TAU_DEFAULT", TAU.Profile.TAU_DEFAULT);

    static TAU.Profile updatetimer = new TAU.Profile("Life main update loop", "", "\
TAU_DEFAULT", TAU.Profile.TAU_DEFAULT);

    // .. other static data
    static public void main(String [] args) throws MPIException {
        MPI.Init(args) ;

        Cartcomm p = MPI.COMM_WORLD.Create_cart(dims, periods, false) ;

        /* Compute local `blockSizeX', `blockBaseX', `blockSizeY', `blockBaseY'. */
        blocktimer.Start();
        {
            // Code to compute blockSizeX, blockBaseX, blockSizeY, blockBaseY
        }
        blocktimer.Stop();

        updatetimer.Start();
        for(int iter = 0 ; iter < NITER ; iter++) {
            // Shift this block's upper x edge into next neighbour's lower ghost edge
            p.Sendrecv(block, blockSizeX * sY, 1, edgeXType, dstX[0], 0,
                block, 0, 1, edgeXType, srcX[0], 0) ;

            // other synchronization operations and loops
            dumpBoard() ;
        }
        updatetimer.Stop();

        MPI.Finalize();
    }
}

--:** LifeBenchmark.java (Java)--L8--Top-----
```

Virtual Machine Performance Instrumentation



□ Integrate performance system with VM

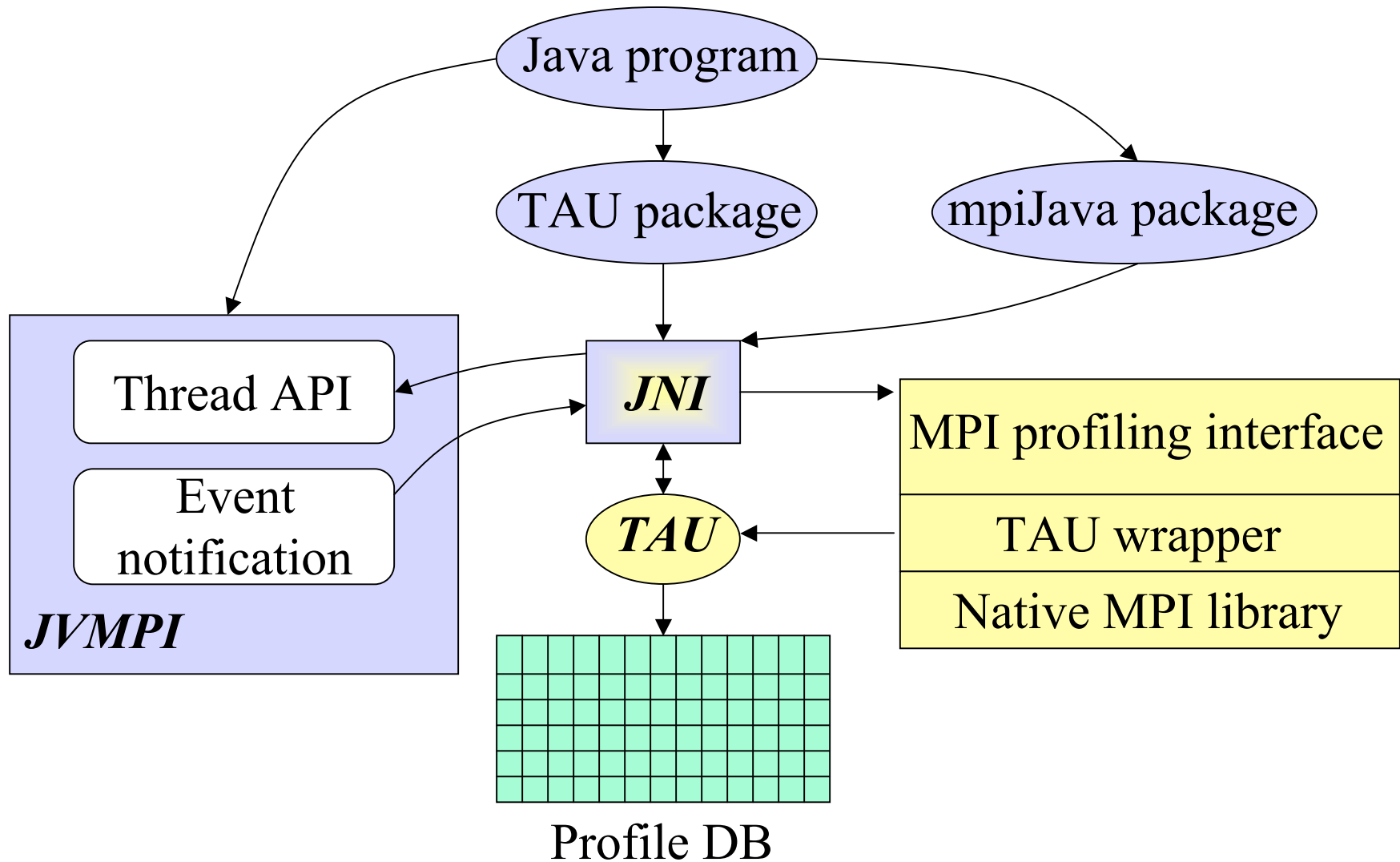
- Captures robust performance data (e.g., thread events)
- Maintain features of environment
 - portability, concurrency, extensibility, interoperation
- Allow use in optimization methods

□ 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



TAU Java Instrumentation Architecture



TAU Measurement



□ Performance information

- High-resolution *timer library* (real-time / virtual clocks)
- General *software counter library* (user-defined events)
- Hardware performance counters
 - *PAPI* (Performance API) (UTK, Ptools Consortium)
 - consistent, portable API

□ Measurement types

- Parallel profiling
 - includes multiple counters, callpaths, performance mapping
- Parallel tracing

□ Support for online performance data access



Multi-Threading Performance Measurement

□ General issues

- Thread identity and per-thread data storage
- Performance measurement support and synchronization
- Fine-grained parallelism
 - different forms and levels of threading
 - greater need for efficient instrumentation

□ TAU general threading and measurement model

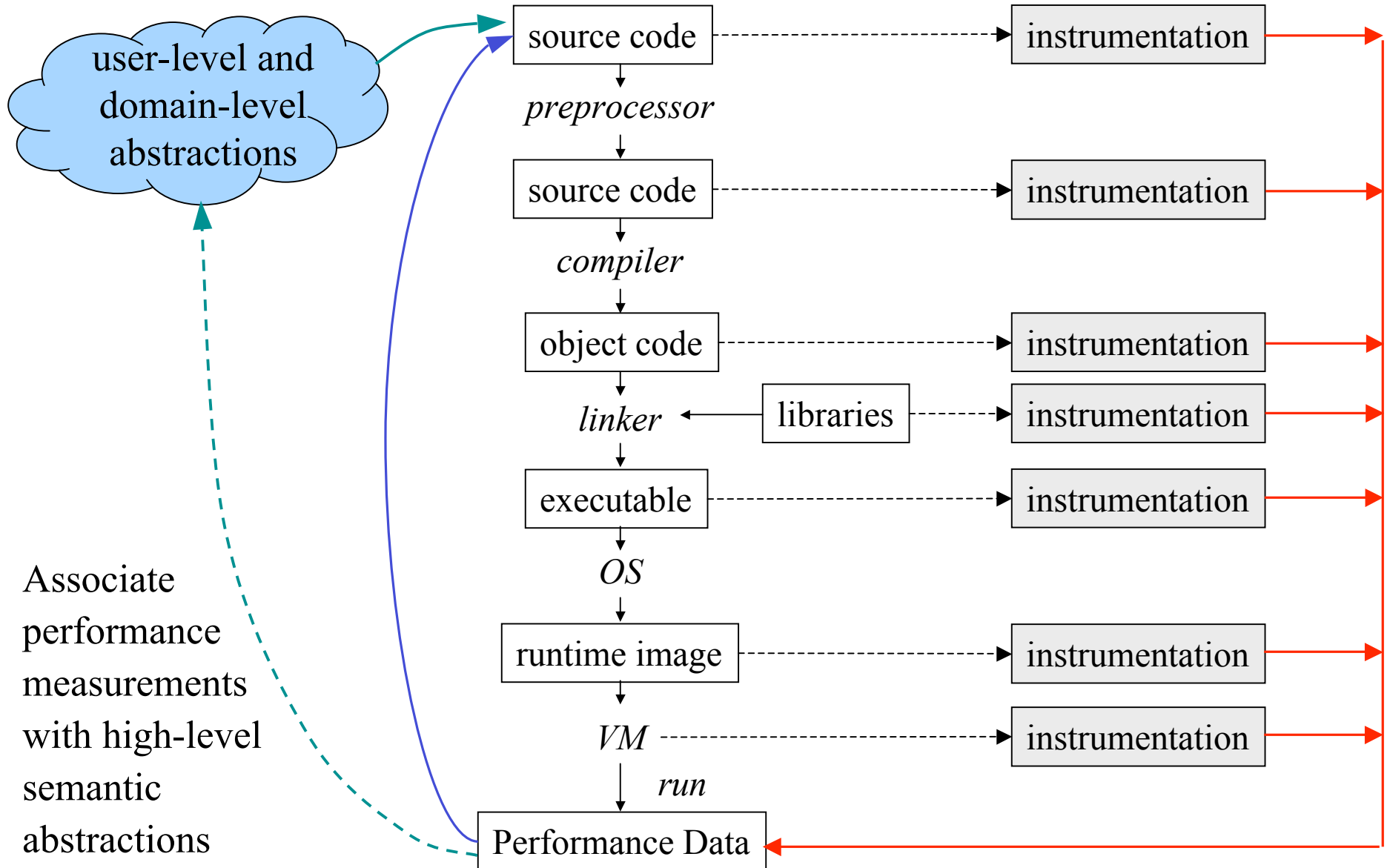
- Common thread layer and measurement support
- Interface to system specific libraries (reg, id, sync)

□ Target different thread systems with core functionality

- Pthreads, Windows, Java, OpenMP



Semantic Performance Mapping



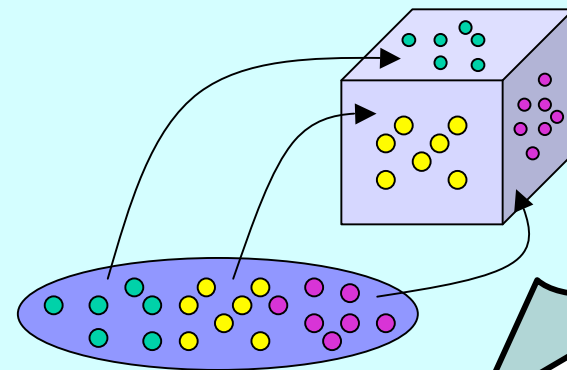
Associate performance measurements with high-level semantic abstractions



Hypothetical Mapping Example

- Particles distributed on surfaces of a 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++){
        /* particles on this face */
        int particles_on_this_face = num(face);
        for (int i=last; i < particles_on_this_face; i++) {
            /* particle properties are a function of face */
            P[i] = ... f(face);
            ...
        }
        last+= particles_on_this_face;
    }
}
```





Hypothetical Mapping Example (continued)

```
int ProcessParticle(Particle *p) {
    /* perform some computation on p */
}
int main() {
    GenerateParticles();
    /* create a list of particles */
    for (int i = 0; i < N; i++)
        /* iterates over the list */
        ProcessParticle(P[i]);
}
```

- ❑ How much time is spent processing *face i* particles?
- ❑ What is the distribution of performance among faces?
- ❑ How is this determined if execution is parallel?

Semantic Entities/Attributes/Associations (SEAA)



- New dynamic mapping scheme
 - Entities defined at any level of abstraction
 - Attribute entity with semantic information
 - Entity-to-entity associations
- Two association types (implemented in TAU API)
 - Embedded
 - External
- “The Role of Performance Mapping”
 - Dr. Sameer Shende
 - Ph.D. thesis

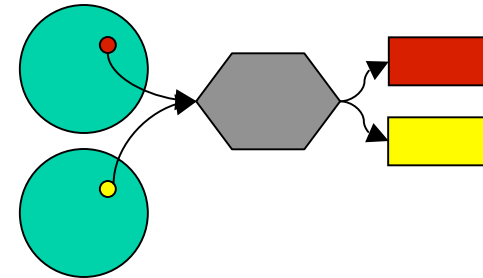


Mapping Associations

□ Embedded association

○ Embedded

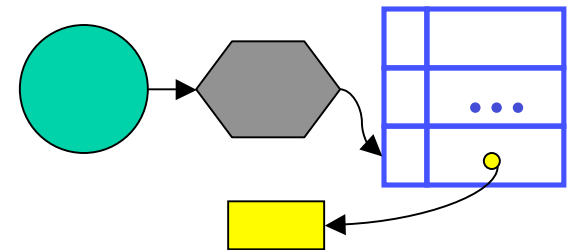
extends associated object to store performance measurement entity



□ External association

○ External

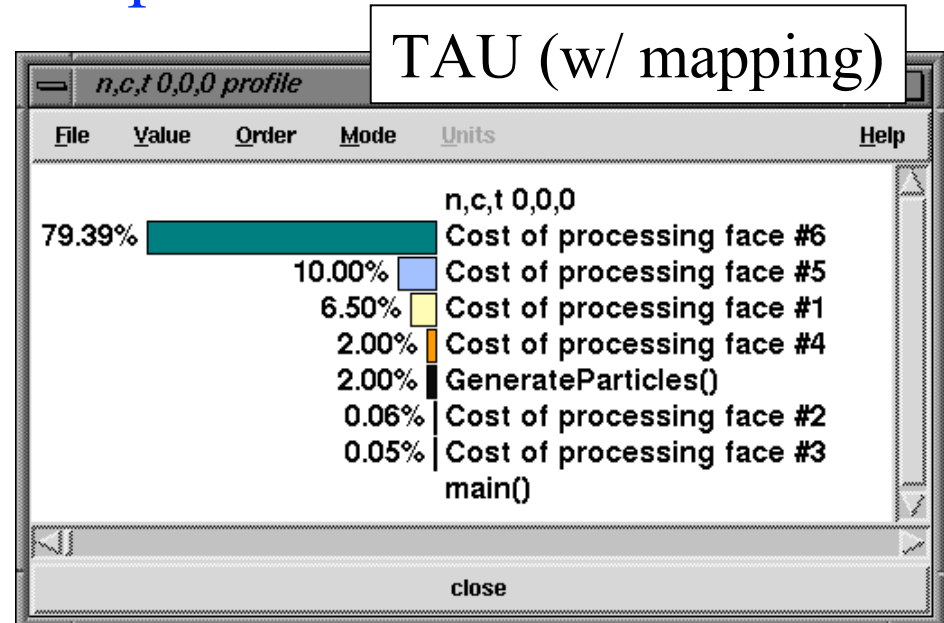
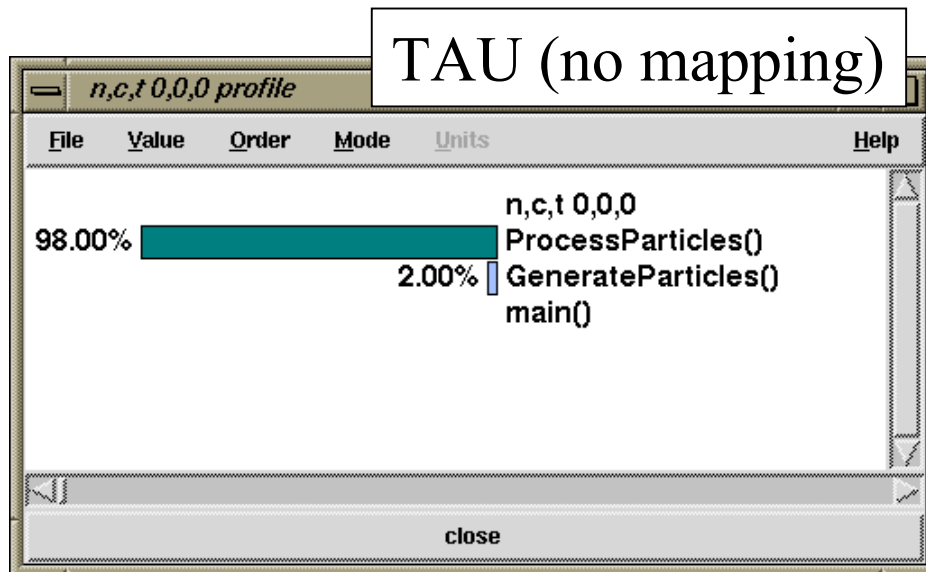
creates an external look-up table using address of object as key to locate performance measurement entity





No Performance Mapping versus Mapping

- Typical performance tools report performance with respect to routines
- Does not provide support for mapping
- Performance tools with SEAA mapping can observe performance with respect to scientist's programming and problem abstractions





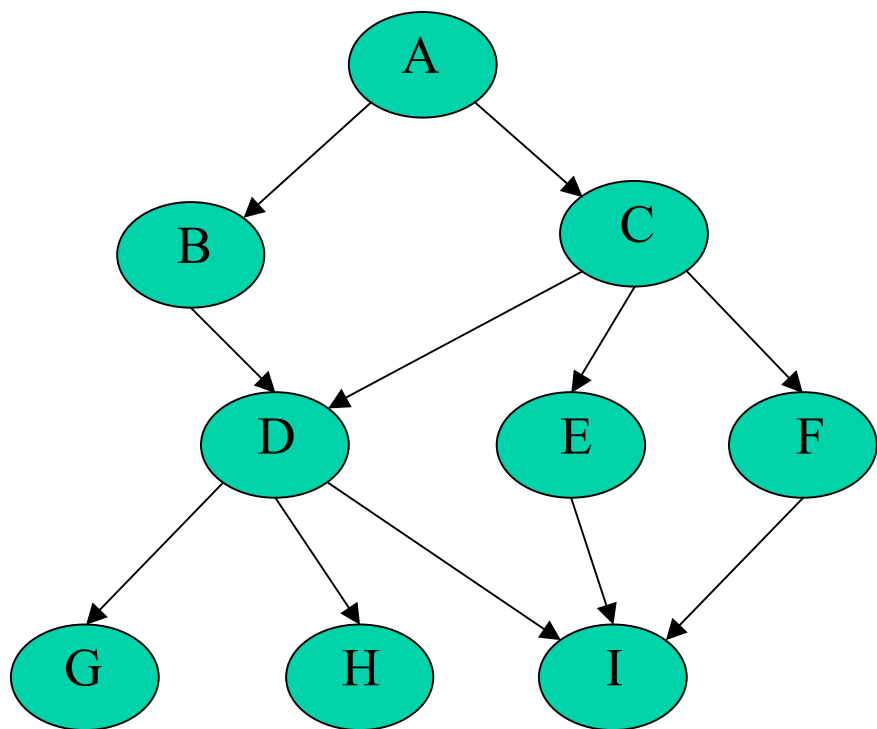
Performance Mapping and Callpath Profiling

- ❑ Associate performance with “significant” entities (events)
- ❑ Source code points are important
 - Functions, regions, control flow events, user events
- ❑ Execution process and thread entities are important
- ❑ Some entities are more abstract, harder to measure
- ❑ Consider callgraph (callpath) profiling
 - Measure time (metric) along an edge (path) of callgraph
 - Incident edge gives parent / child view
 - Edge sequence (path) gives parent / descendant view
- ❑ Problem: Callpath profiling when callgraph is unknown
 - Determine callgraph dynamically at runtime
 - Map performance measurement to dynamic call path state



Callgraph (Callpath) Profiling

- Measure time (metric) along an edge (path) of callgraph
 - Incident edge gives parent / child view
 - Edge sequence (path) gives parent / descendant view



- 1-level callpath
 - Immediate descendant
 - $A \rightarrow B$, $E \rightarrow I$, $D \rightarrow H$
 - $C \rightarrow H$?
- k -level callpath
 - k call descendant
 - 2-level: $A \rightarrow D$, $C \rightarrow I$
 - 2-level: $A \rightarrow I$?
 - 3-level: $A \rightarrow H$



k-Level Callpath Implementation in TAU

- ❑ TAU maintains a performance event (routine) callstack
- ❑ Profiled routine (child) looks in callstack for parent
 - Previous profiled performance event is the parent
 - A *callpath profile structure* created first time parent calls
 - TAU records parent in a *callgraph map* for child
 - String representing k-level callpath used as its key
 - “a()=>b()=>c()” : name for time spent in “c” when called by “b” when “b” is called by “a”
- ❑ Map returns pointer to callpath profile structure
 - k-level callpath is profiled using this profiling data
 - Set environment variable **TAU_CALLPATH_DEPTH** to depth
- ❑ Build upon TAU’s performance mapping technology
- ❑ Measurement is independent of instrumentation
- ❑ Use `-PROFILECALLPATH` to configure TAU



Running Applications

```
% set path=($path <taudir>/<arch>/bin)
% set path=($path $PET_HOME/PTOOLS/tau-2.13.5/src/rs6000/bin)
% setenv LD_LIBRARY_PATH $LD_LIBRARY_PATH\:<taudir>/<arch>/lib
```

For PAPI (1 counter, if multiplecounters is not used):

```
% setenv PAPI_EVENT PAPI_L1_DCM (Level 1 Data cache misses)
```

For PAPI (multiplecounters):

```
% setenv COUNTER1 PAPI_FP_INS (Floating point instructions)
```

```
% setenv COUNTER2 PAPI_TOT_CYC (Total cycles)
```

```
% setenv COUNTER3 P_VIRTUAL_TIME (Virtual time)
```

```
% setenv COUNTER4 LINUX_TIMERS (Wallclock time)
```

(NOTE: PAPI_FP_INS and PAPI_L1_DCM cannot be used together on Power4. Other restrictions may apply to no. of counters used.)

```
% mpirun -np <n> <application>
```

```
% llsubmit job.sh
```

```
% paraprof (for performance analysis)
```



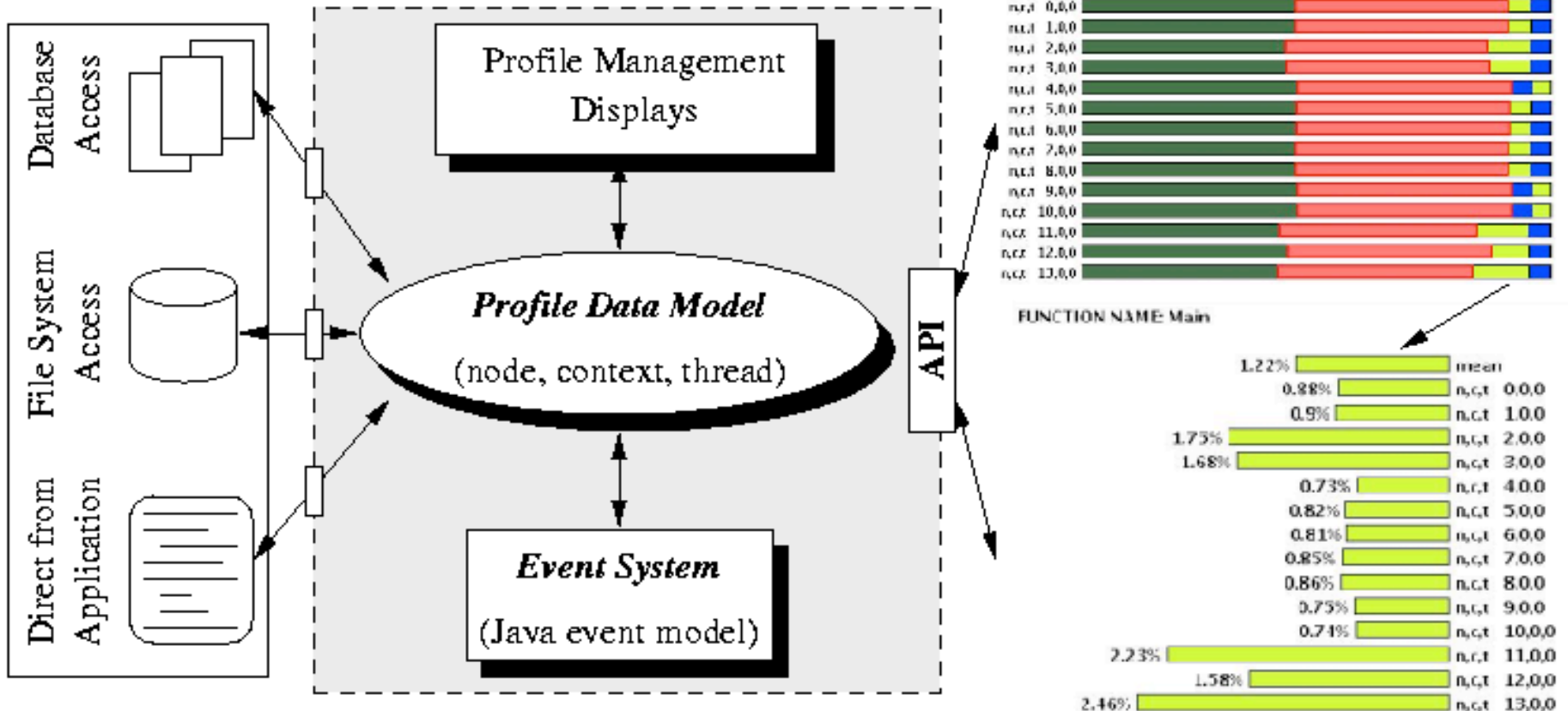
TAU Performance Analysis

- ❑ Analysis of parallel profile and trace measurement
- ❑ Parallel profile analysis
 - *Pprof*: parallel profiler with text-based display
 - *ParaProf*: graphical, scalable parallel profile analysis
- ❑ Parallel trace analysis
 - Format conversion (ALOG, VTF 3.0, Paraver, EPILOG)
 - Trace visualization using *Vampir* (Pallas/Intel)
 - Parallel profile generation from trace data



ParaProf Framework Architecture

- Portable, extensible, and scalable tool for profile analysis
- Try to offer “best of breed” capabilities to analysts
- Build as profile analysis framework for extensibility





ParaProf Manager

The screenshot shows the ParaProf Manager application window. The title bar reads "ParaProf Manager". The menu bar contains "File", "Options", and "Help".

The left pane shows a tree view under "Applications":

- Applications
 - Stan
 - Update Meta Data in DB
 - Add Trial
 - Default Exp
 - Default Trial
 - Time
 - Runtime Applications
 - DB Applications

Field	Value
Name	Default Exp
Application ID	0
Experiment ID	0
User Data	
System Name	
System Machine Type	
System Arch.	
System OS	
System Memory Size	
System Processor Amount	
System L1 Cache Size	
System L2 Cache Size	
System User Data	
Configuration Prefix	
Configuration Architecture	
Configuration CPP	
Configuration CC	
Configuration JDK	
Configuration Profile	
Configuration User Data	
Compiler CPP Name	
Compiler CPP Version	
Compiler CC Name	
Compiler CC Version	
Compiler Java Dir. Path	
Compiler Java Version	
Compiler User Data	

- ❑ Powerful manager for control of data sources
 - Directly from files
 - Profile database
 - Runtime (online)
- ❑ Conveniences to facilitate working with data



ParaProf Manager (continued)

- Data management windows
 - Loading flat files from disk
 - Generating new derived metrics
 - Database interface

The screenshot displays the ParaProf Manager application window. The main interface is divided into several sections:

- Applications Tree:** A hierarchical tree view on the left showing 'Applications' with sub-items like 'Star', 'Default Exp', 'Default Trial', 'Runtime Applications', and 'DB Applications'. A context menu is open over the 'Default Trial' folder, listing options: 'Update Meta Data in DB', 'Add Trial', 'CPU_TIME', 'GET_TIME_OF_DAY', 'packed_DP_uop_all', and 'packed_DP_uop_all / GET_TIME_OF_DAY'. The last option is selected.
- Table:** A table on the right showing trial information:

Field	Value
Name	packed_DP_uop_all / GET_T...
Application ID	0
Experiment ID	0
Trial ID	0
Metric ID	3
- Arguments:** Two text input fields for 'Argument 1' (0:0:0:3) and 'Argument 2' (0:0:0:2), with a 'Divide' dropdown menu and an 'Apply operation' button.
- Database Configuration Window:** A separate window with a 'Password' field (masked with asterisks), a 'Config File' field (./bertie/Robert/Code/Data/multiplecount), and 'Cancel' and 'Ok' buttons.
- Load Trial Window:** A separate window with a 'Trial Type' dropdown (pprof.dat), a 'Dir. Location' field (bertie/Robert/Code/Data/multiplecount), and 'Cancel' and 'Ok' buttons.

Red arrows indicate the flow of information: from the selected application in the tree to the 'Load Trial' window, from the 'DB Applications' folder to the 'Database Configuration' window, and from the table to a box labeled 'Trial information'.



ParaProf Derived Metrics

The screenshot shows the ParaProf Manager application window. The title bar reads "ParaProf Manager". The menu bar contains "File" and "Help".

The left pane is titled "Standard Applications" and contains a tree view:

- Standard Applications
 - Default App
 - Experiments
 - Default Exp
 - Trials
 - Default Trial : 512proc/samrai/taudata/neutron
 - 0000 - P_WALL_CLOCK_TIME
 - 0001 - PAPI_FP_INS
 - 0002 - PAPI_FP_INS / P_WALL_CLOCK_TIME

Below the tree view are sections for "Runtime Applications" and "DB Applications".

The right pane is titled "ParaProf Manager" and contains the following text:

Clicking on different values causes ParaProf to display the clicked on metric.

The sub-window below allow you to generate new metrics based on those that were gathered during the run. The operand number options for Operand A and B correspond the numbers prefixing the values.

At the bottom of the window is a section titled "Apply operations here!" with the following controls:

- Op A:
- Op B:
- Operation:
-



ParaProf Profile Analysis Displays

textual profile

%Time	Time	total Time	#calls	#subrs	total Time
100.0	8.757522992004E8	1.1717499847999E9	1.0	253.0	1.17174998
21.8	2.553087943998E8	2.553087943998E8			
0.9	1.0481863198E7	1.0481863198E7			
2.7	5075165.6002	3.16434583997E7			
0.4	5025481.6	5025481.6			
0.3	4027989.5978	4027989.5978			
0.2	2935568.0	2935568.0			
0.2	2130867.201	2130867.201			
0.1	1642712.0005	1642712.0005			
0.1	1632580.8013	1632580.8013			

legend

- main() void (int, char **)
- MPI_Reduce()
- MPI_Waitsome()
- MPI_Scheduler::execute()
- MPI_Init_thread()
- MPI_Allreduce()
- MPI_Barrier()
- MPI_Type_indexed()

full profile with display adjustment

Bar Multiple: 0 5 10 15 20 25 30 35 40

thread display

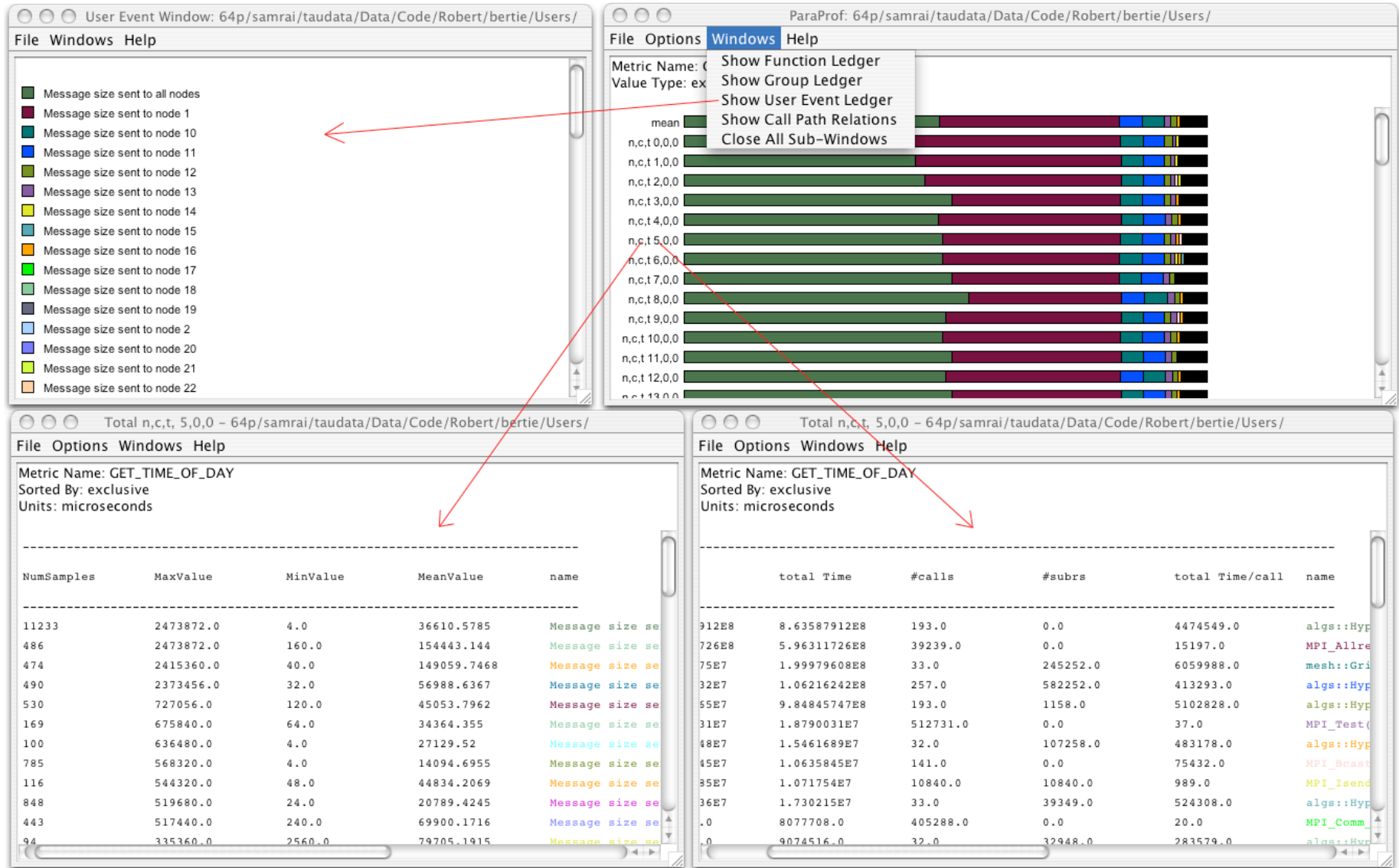
Metric Name	Value
main	74.74%
MPI_Reduce	21.79%
MPI_Waitsome	0.89%
MPI_Scheduler::execute	0.43%
MPI_Init_thread	0.43%
MPI_Allreduce	0.34%
MPI_Barrier	0.25%
MPI_Type_indexed	0.18%
SerialMPM::comp	0.14%
MPI_Probe	0.14%

event display

Metric Name	Value
mean	72.66%
n,c,t 223,0,0	78.19%
n,c,t 240,0,0	77.62%
n,c,t 264,0,0	77.46%

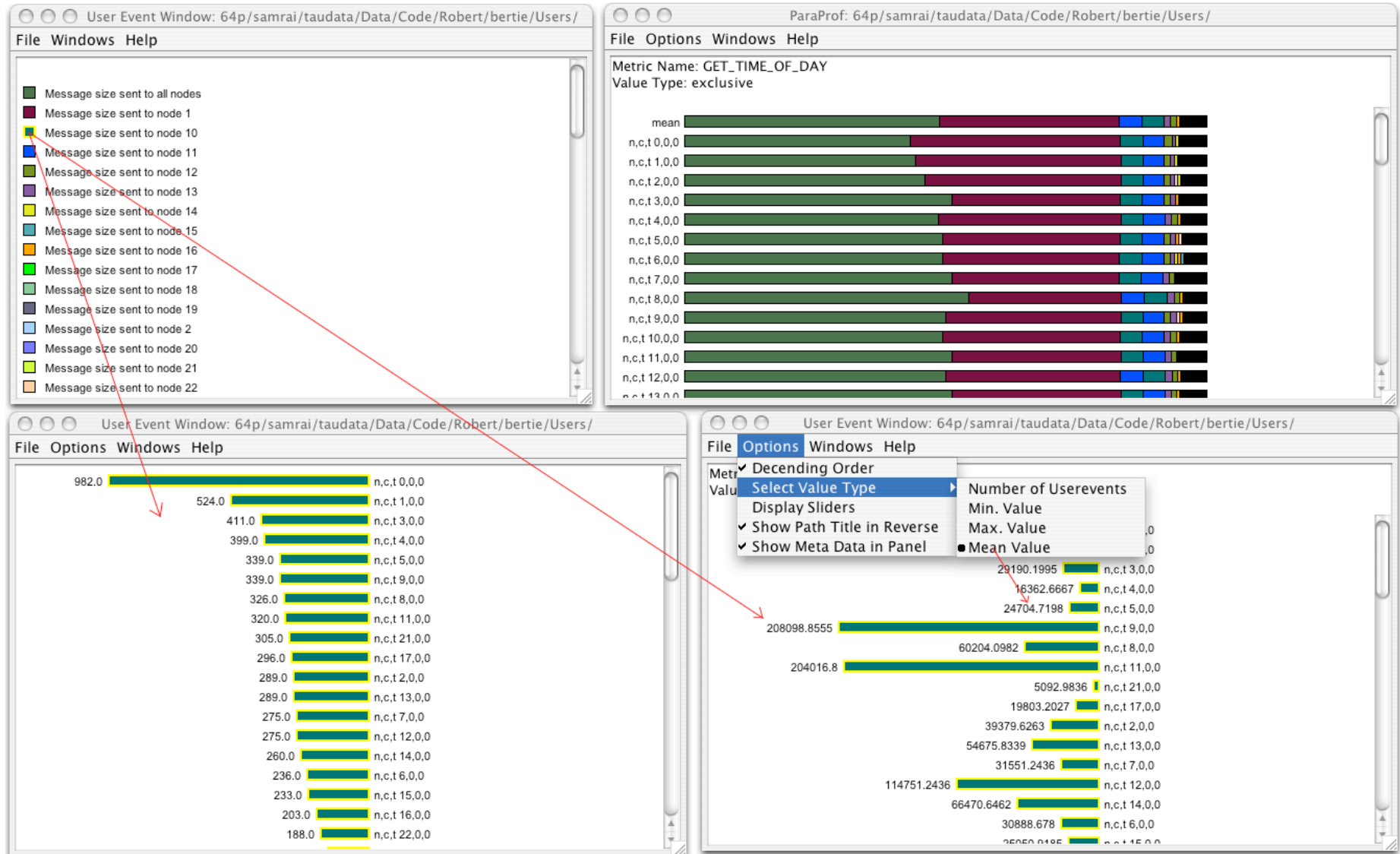


ParaProf User Event Display (MPI message size)





ParaProf User Event Details (MPI message size)





Using TAU's Malloc Wrapper Library for C/C++

NumSamples	MaxValue	MinValue	MeanValue	name
1	40016.0	40016.0	40016.0	malloc size <file=main.cpp, line=252>
1	40016.0	40016.0	40016.0	free size <file=main.cpp, line=298>
12	30000.0	240.0	5590.0	malloc size <file=select.cpp, line=80>
12	30000.0	240.0	5590.0	malloc size <file=select.cpp, line=81>
3	30000.0	6000.0	17000.0	free size <file=select.cpp, line=107>
3	30000.0	6000.0	17000.0	free size <file=select.cpp, line=109>
1	8000.0	8000.0	8000.0	malloc size <file=main.cpp, line=258>
1	8000.0	8000.0	8000.0	free size <file=main.cpp, line=299>
7	6000.0	600.0	2228.5714	free size <file=select.cpp, line=118>
7	6000.0	600.0	2228.5714	free size <file=select.cpp, line=119>
2	240.0	240.0	240.0	free size <file=select.cpp, line=126>
2	240.0	240.0	240.0	free size <file=select.cpp, line=128>



ParaProf Profile Analysis Features

- ❑ Inter-window event management
 - Full event propagation
 - Hyperlinked displays
- ❑ Window configuration and help management
 - Popup menus
 - Full preference control
 - Data view control
- ❑ Maturation of profile performance data views
- ❑ Java-based implementation
 - Extensible
- ❑ Performance database connectivity



ParaProf Enhancements

- ❑ Readers completely separated from the GUI
- ❑ Access to performance profile database
- ❑ Profile translators
 - *mpiP, papiprof, dynaprof*
- ❑ Callgraph display
 - *prof / gprof* style with hyperlinks
- ❑ Integration of 3D performance plotting library
- ❑ Scalable profile analysis
 - Statistical histograms, cluster analysis, ...
- ❑ Generalized programmable analysis engine
- ❑ Cross-experiment analysis



Callpath Profiling Example (NAS LU v2.3)

```
% configure -PROFILECALLPATH -SGITIMERS -arch=sgi64
-mpiinc=/usr/include -mpilib=/usr/lib64 -useropt=-O2
```

Mean Total Stat Window: /tmp_mnt/inf/research/parallel/sameer/gar/rs/demo/tau2/examples/NPB2.3/bin/lpprof.dat

File Options Windows Help

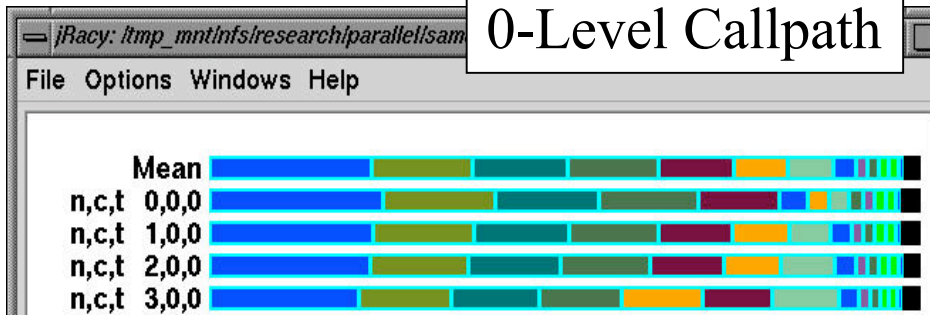
%time	msec	total msec	#call	#subrs	usec/call	name
30.7	13,037	17,495	301	602	58126	bcast_inputs => rhs
30.7	13,037	17,495	301	602	58126	rhs
17.3	8,195	9,838	9300	18600	1058	bcast_inputs => buts
17.3	8,195	9,838	9300	18600	1058	buts
21.0	7,669	11,998	9300	18600	1290	bcast_inputs => blts
21.0	7,669	11,998	9300	18600	1290	blts
12.8	7,320	7,320	9300	0	787	bcast_inputs => jacld
12.8	7,320	7,320	9300	0	787	jacld
10.6	6,049	6,049	9300	0	651	bcast_inputs => jacu
10.6	6,049	6,049	9300	0	651	jacu
7.7	4,385	4,385	18600	0	236	MPI_Recv()
7.7	4,385	4,385	18600	0	236	exchange_1 => MPI_Recv()
6.5	3,700	3,700	606	0	6106	MPI_Wait()
6.5	3,700	3,700	604	0	6126	exchange_3 => MPI_Wait()
95.8	1,882	54,609	2.25	37517	24270831	bcast_inputs
95.7	1,882	54,604	1.25	37508	43683863	applu => bcast_inputs
1.8	1,012	1,012	1	0	1012219	MPI_Finalize()
1.8	1,012	1,012	1	0	1012219	applu => MPI_Finalize()
10.5	996	5,971	37200	37200	161	exchange_1
1.5	882	882	19206	0	46	MPI_Send()
1.7	833	950	1	44686.5	950499	applu => setiv



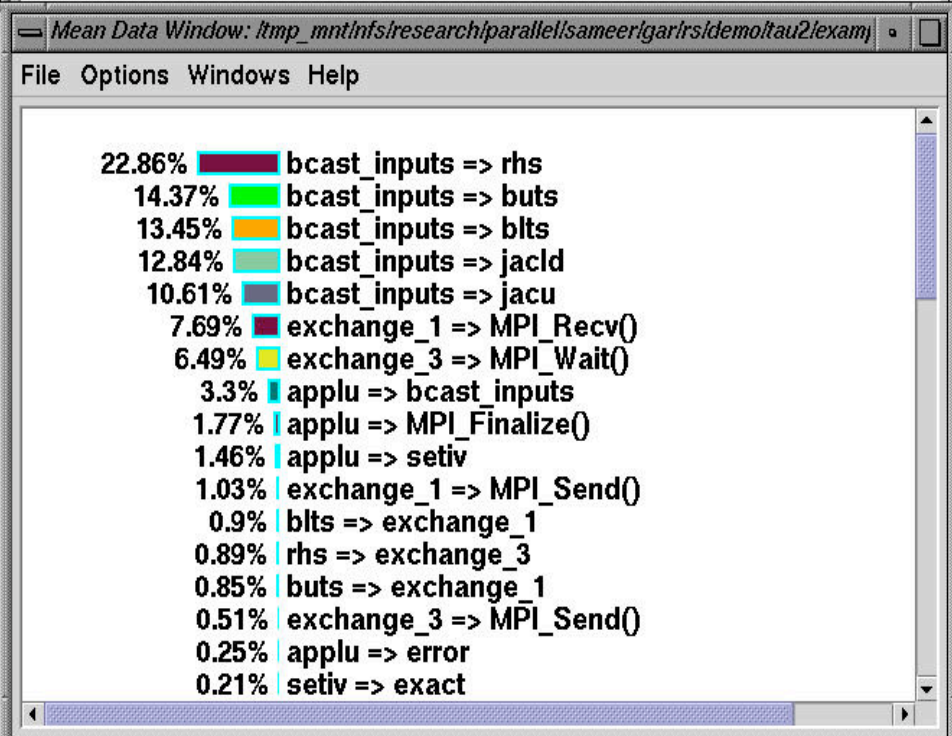
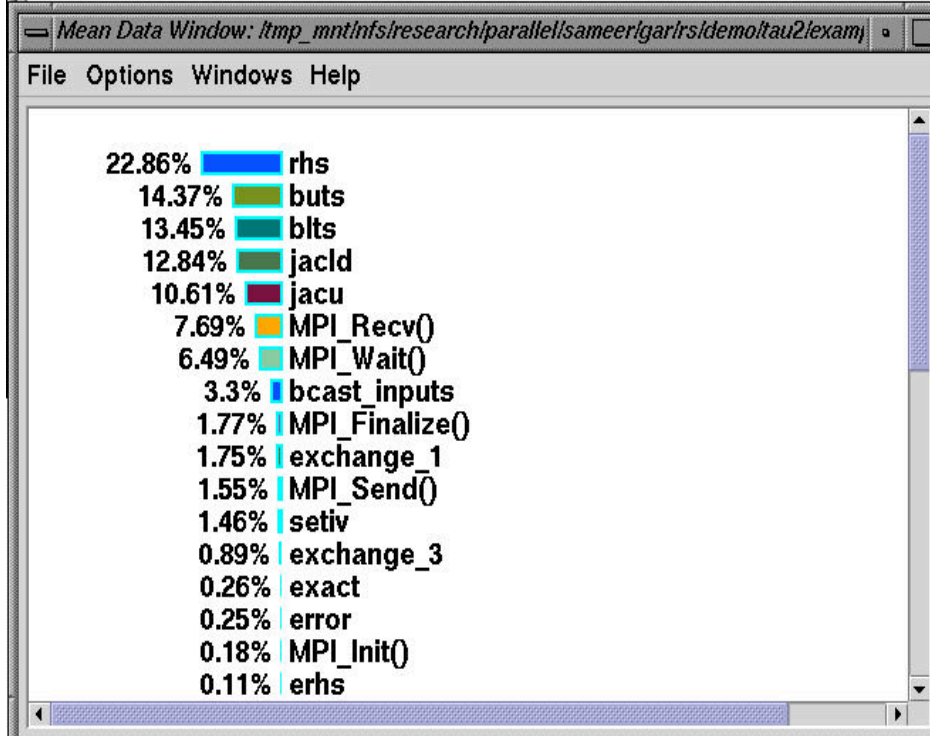
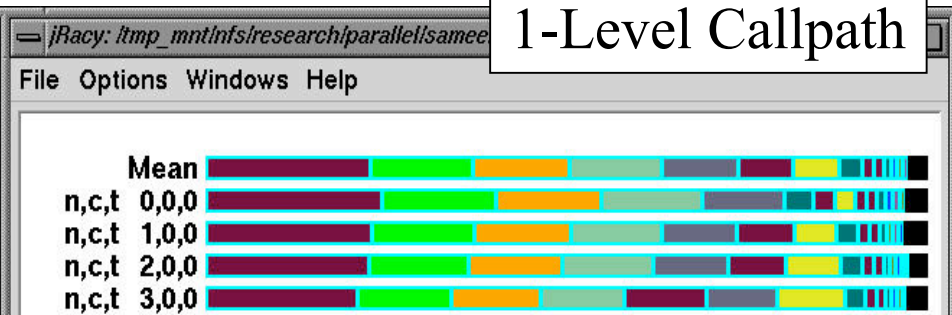
Callpath Parallel Profile Display

0-level and 1-level callpath grouping

0-Level Callpath

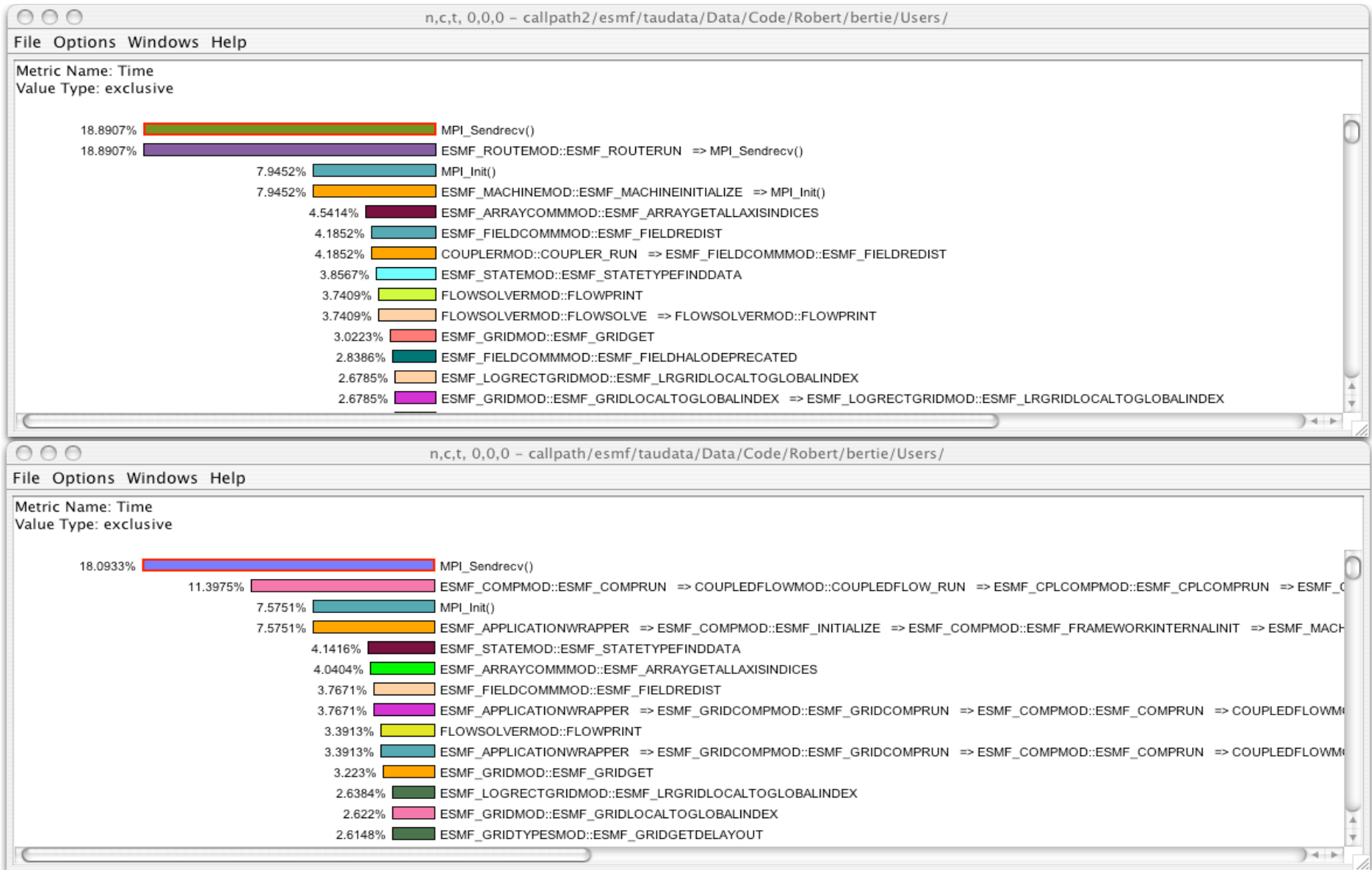


1-Level Callpath





Callpath Profiling Screenshot





Callpath Profiling Parent/Child Relations

```
Call Path Data Relations - callpath/esmf/Data/Code/Robert/bertie/Users/
File Options Windows Help

ESMF_COMPMOD::ESMF_INITIALIZE [1]
--> ESMF_COMPMOD::ESMF_FRAMEWORKINTERNALINIT [3]
    ESMF_MACHINEMOD::ESMF_MACHINEINITIALIZE [5]
        ESMF_DELAYOUTMOD::ESMF_DELAYOUTCREATEDDEFAULTID [13]

ESMF_COMPMOD::ESMF_FRAMEWORKINTERNALINIT [3]
--> ESMF_MACHINEMOD::ESMF_MACHINEINITIALIZE [5]
    MPI_Init() [7]
    MPI_Comm_size() [9]
    MPI_Comm_rank() [11]

ESMF_MACHINEMOD::ESMF_MACHINEINITIALIZE [5]
--> MPI_Init() [7]

ESMF_MACHINEMOD::ESMF_MACHINEINITIALIZE [5]
    ESMF_DELAYOUTMOD::ESMF_DELAYOUTCREATEDDEFAULTID [13]
    ESMF_DELAYOUTMOD::ESMF_DELAYOUTCREATEFROMDELIST [181]
--> MPI_Comm_size() [9]

ESMF_MACHINEMOD::ESMF_MACHINEINITIALIZE [5]
ESMF_DELAYOUTMOD::ESMF_DELAYOUTCREATEDDEFAULTID [13]
```




Callpath Profiling Screenshot (cont.)

Call Path Data n,c,t, 0,0,0 - callpath/esmf/Data/Code/Robert/bertie/Users/

File Options Windows Help

Metric Name: Time
Sorted By: exclusive
Units: microseconds

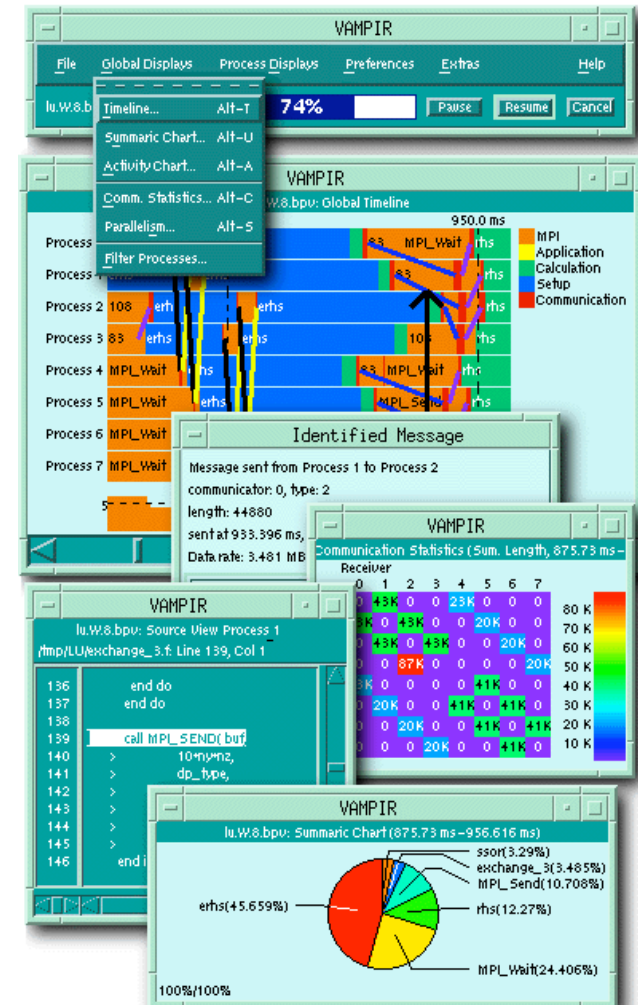
--> 483085.0	1538110.0	55811	ESMF_STATEMOD::ESMF_STATEGETFIELD [444]
176.0	634	77435	INJECTORMOD::INJECTOR_INIT2 [436]
177.0	636	77435	FLWSOLVERMOD::FLOW_INIT2 [461]
108242.0	351	77435	FLWSOLVERMOD::FLWSOLVE [497]
219049.0	704	77435	COUPLERMOD::COUPLER_RUN [674]
112038.0	357078.0	12600/50414	INJECTORMOD::INJECTOR_RUN [719]
--> 439682.0	1414389.0	50414	ESMF_STATEMOD::ESMF_STATEISNEEDED [438]
446373.0	724341.0	50414/12600	ESMF_STATEMOD::ESMF_STATETYPEFINDDATA [337]
250366.0	250366.0	50414/12600	ESMF_STATEMOD::ESMF_NEEDEQ [442]
--> 424991.0	2660567.0	1800	FLWSOLVERMOD::FLOWVEL [597]
--> 417305.0	417305.0	77435	ESMF_DELAYOUTMOD::ESMF_DELAYOUTGETDEID [32]
> 416249.0	416249.0	77435	ESMF_DELAYOUTMOD::ESMF_DELAYOUTGETDEID [32]

Show Function Details
 Find Function
 Change Function Color
 Reset to Generic Color



Vampir Trace Visualization

- Visualization and Analysis of MPI Programs
- Originally developed by Forschungszentrum Jülich
- Current development by Technical University Dresden
- Distributed by PALLAS, Germany



<http://www.pallas.de/pages/vampir.htm>



Using TAU with Vampir

```
include $PET_HOME/PTOOLS/tau-  
2.13.5/rs6000/lib/Makefile.tau-mpi-pdt-trace  
F90 = $(TAU_F90)  
LIBS = $(TAU_MPI_LIBS) $(TAU_LIBS) $(TAU_CXXLIBS)  
OBJS = ...  
TARGET= a.out  
TARGET: $(OBJS)  
    $(CXX) $(LDFLAGS) $(OBJS) -o $@ $(LIBS)  
.f.o:  
    $(F90) $(FFLAGS) -c $< -o $@
```



Using TAU with Vampir

- ❑ **Configure TAU with `-TRACE` option**

```
% configure -TRACE -SGITIMERS ...
```

- ❑ **Execute application**

```
% mpirun -np 4 a.out
```

- ❑ **This generates TAU traces and event descriptors**

- ❑ **Merge all traces using `tau_merge`**

```
% tau_merge *.trc app.trc
```

- ❑ **Convert traces to Vampir Trace format using `tau_convert`**

```
% tau_convert -pv app.trc tau.edf app.pv
```

- ❑ **Load generated trace file in Vampir**

```
% vampir app.pv
```



Case Study: SIMPLE Performance Analysis

- SIMPLE hydrodynamics benchmark
 - C code with MPI message communication
 - Multiple instrumentation methods
 - source-to-source translation (PDT)
 - MPI wrapper library level instrumentation (PMPI)
 - pre-execution binary instrumentation (DyninstAPI)
 - Alternative measurement strategies
 - statistical profiles of software actions
 - statistical profiles of hardware actions (PCL, PAPI)
 - program event tracing
 - choice of time source
 - gettimeofday, physical clock, CPU, process virtual



SIMPLE Source Instrumentation (Preprocessed)

- PDT automatically generates instrumentation code
 - Names events with full function signatures

```
int compute_heat_conduction(  
    double theta_hat[X][Y], double deltat, double new_r[X][Y],  
    double new_z[X][Y], double new_alpha[X][Y],  
    double new_rho[X][Y], double theta_l[X][Y],  
    double Gamma_k[X][Y], double Gamma_l[X][Y])  
{  
    TAU_PROFILE("int compute_heat_conduction(  
        double (*)[259], double, double (*)[259],  
        double (*)[259], double (*)[259], double (*)[259],  
        double (*)[259], double (*)[259], double (*)[259])",  
        " ", TAU_USER);  
    ...  
}
```

- Similarly for all other routines in SIMPLE program



MPI Library Instrumentation (MPI_Send)

- Uses MPI profiling interposition library (PMPI)

```
int MPI_Send(...)
...
{
    int returnVal, typesize;
    TAU_PROFILE_TIMER(tautimer, "MPI_Send()", " ", TAU_MESSAGE);
    TAU_PROFILE_START(tautimer);
    if (dest != MPI_PROC_NULL) {
        PMPI_Type_size(datatype, &typesize);
        TAU_TRACE_SENDSMSG(tag, dest, typesize*count);
    }
    returnVal = PMPI_Send(buf, count, datatype, dest, tag, comm);
    TAU_PROFILE_STOP(tautimer);
    return returnVal;
}
```



MPI Library Instrumentation (MPI_Recv)

```
int MPI_Recv(...)
...
{
    int returnVal, size;
    TAU_PROFILE_TIMER(tautimer, "MPI_Recv()", " ", TAU_MESSAGE);
    TAU_PROFILE_START(tautimer);
    returnVal = PMPI_Recv(buf, count, datatype, src, tag, comm,
        status);
    if (src != MPI_PROC_NULL && returnVal == MPI_SUCCESS) {
        PMPI_Get_count(status, MPI_BYTE, &size);
        TAU_TRACE_RECVMSG(status->MPI_TAG, status->MPI_SOURCE,
            size);
    }
    TAU_PROFILE_STOP(tautimer);
    return returnVal;
}
```




Multi-Level Instrumentation (Profiling)

four processes

event legend

Profile per process

%time	msec	total msec	#call	#subrs	usec/call	name
49.9	19,295	31,066	206388	4.12776E+06	151	polynomial
17.6	10,946	10,946	3	9	3648787	net_accept
17.6	10,930	10,930	3.71498E+06	0	3	power
13.3	6,684	8,277	273	546	30320	socket_recv
3.1	1,952	1,954	561	1122	3484	net_recv
3.9	1,554	2,399	10	317580	239940	compute_viscos
2.1	1,316	1,318	628	1256	2100	net_send
35.8	1,117	22,279	10	119070	2227940	compute_temperature

global profile

Function	Percentage
polynomial	31.0%
net_accept	17.0%
power	17.0%
socket_recv	10.0%
net_recv	3.0%
compute_viscosity_interior	2.0%

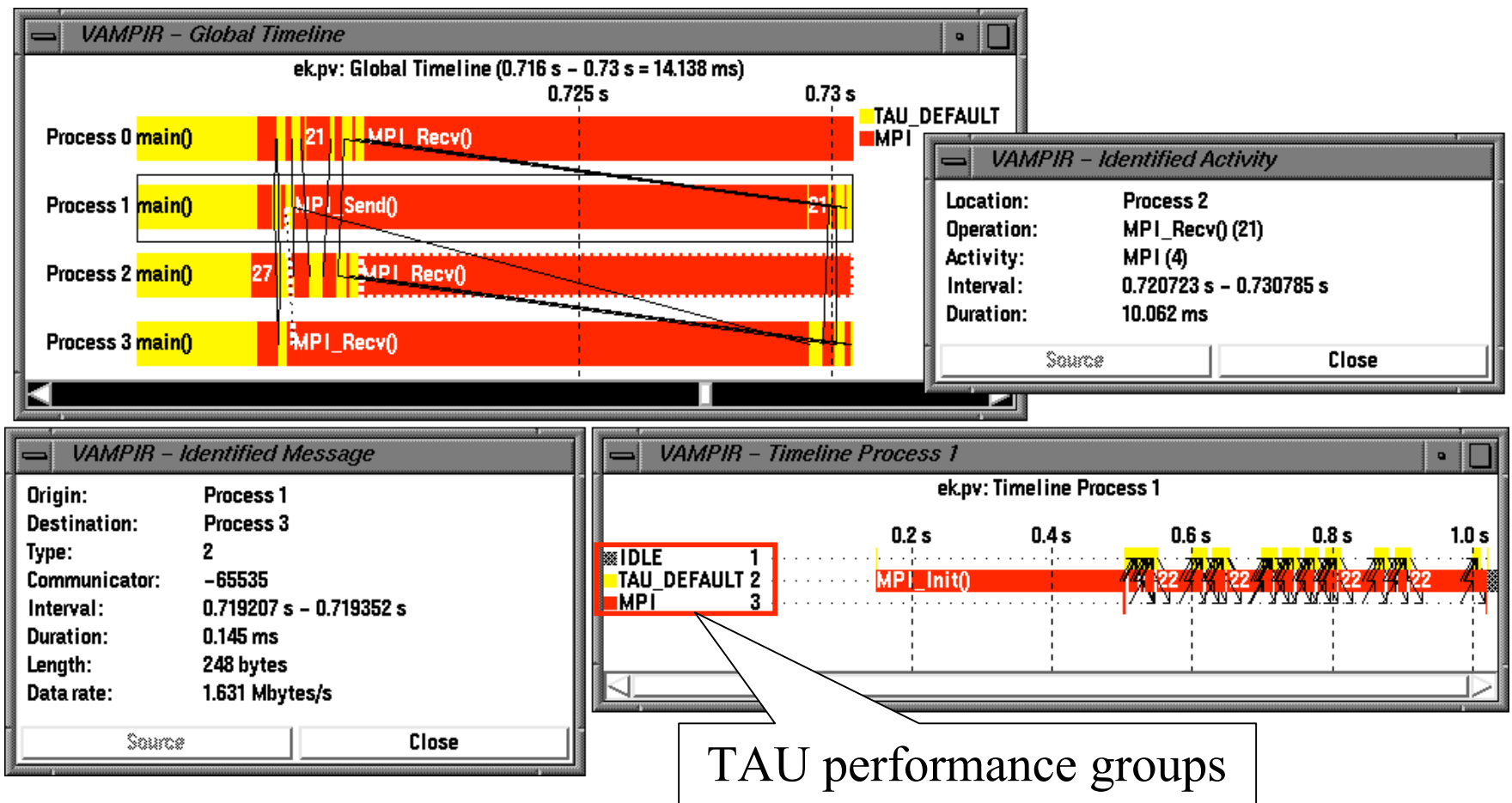
Racy Help Window

n,c,t stands for: Node, Context and Thread.
Using the right mouse button, double click here to display more detailed data about this thread.



Multi-Level Instrumentation (Tracing)

- Relink with TAU library configured for tracing
 - No modification of source instrumentation required!





Dynamic Instrumentation of SIMPLE

- ❑ Uses DynInstAPI for runtime code patching
- ❑ Mutator loads measurement library, instruments mutatee
 - One mutator (*tau_run*) per executable image
 - `mpirun -np <n> tau.shell`

