

PAPI, DynInst and Hardware Performance Analysis Tools

Lawrence Berkeley National Laboratory
Performance Workshop, August 8th, 2003

Philip Mucci, Research Consultant
Innovative Computing Laboratory/UTK
Performance Evaluation Research Center/NERSC/LBNL

pjmucci@lbl.gov





PAPI provides two standardized APIs to access the underlying performance counter hardware

- A low level interface designed for tool developers and expert users.
- The high level interface is for application engineers.

- Overview of PAPI
 - Features, Functionality and Usage
 - 2.3.4 Release Status
 - Changes coming in 3.0
- Dynamic Instrumentation
- Performance Analysis Tools
- Trends in the field

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.

- To increase application performance and system throughput
- To characterize application and system workload on the CPU
- To stimulate performance tool development and research
- To stimulate research on more sophisticated feedback driven compilation techniques

- Provide a solid foundation for cross platform performance analysis *tools*.
- Loosely *standardize* interface among users, tool developers, vendors and academics.
- Provide implementations for the more popular HPC machines.
- Easy to use, well documented and freely available. (Truly Open Source)

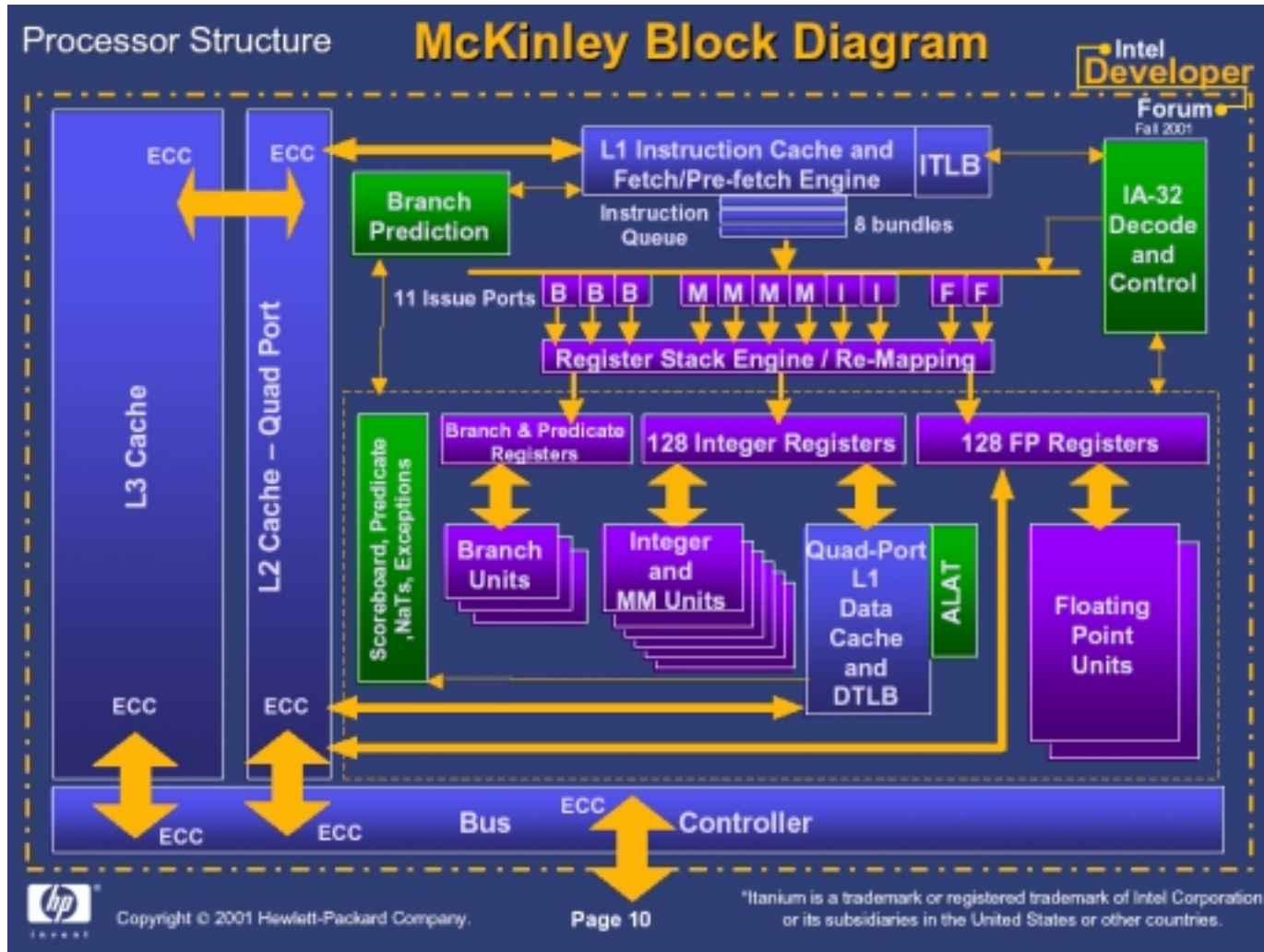
- Engineers often report near linear scalability for large parallel application codes to hundreds of processors.
 - Closer inspection reveals poor per processor performance for their problem class.
 - How does one judge performance?
 - Do we have *good* algorithm? (maybe)
 - Do we have a *good* compiler? (?)
 - We need real performance data from the processor!

- How fast is fast?
 - Single CPU performance in the scientific marketplace often results from good **data locality**, predictable branch behavior, and *instruction pipelines filled with independent instructions*. (Out-of-order scheduling and execution helps the latter 2)
 - Determining the “reuse potential” for a numerically intensive problem often dictates peak performance.
 - Certain ratios of metrics can help diagnose this problem. FP stalls vs. Mem stalls, loads vs. stores, misses vs. references, etc...

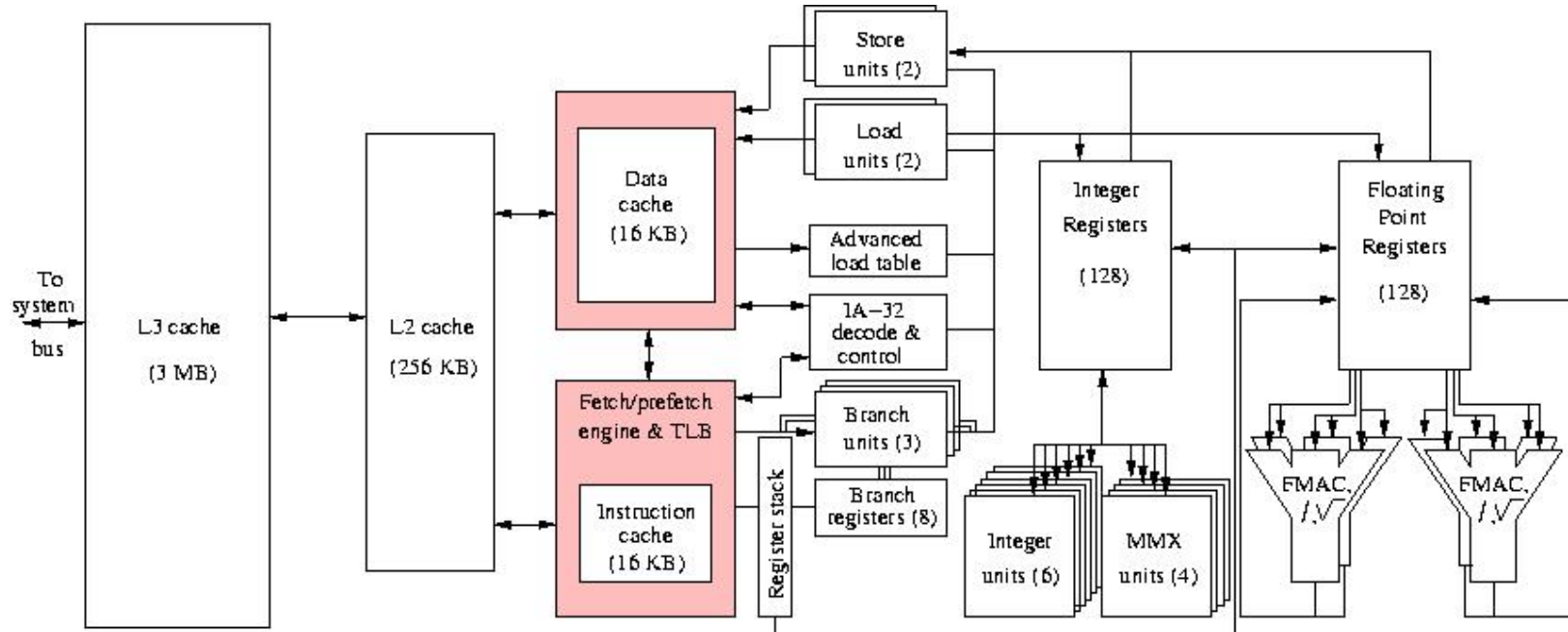
- Portability of an application is of the utmost concern to most HPC environments.
- No common performance tools except **prof / gprof**.
- *Most* commercial tools are based on time.
- HPC have memory and floating point intensive workloads which require detailed analysis.
- Research tools are the most common “solution” among the labs.

- Small number of registers dedicated for performance monitoring functions.
 1. AMD Athlon, 4 counters
 - Power 3, 8 counters
 2. Pentium \leq III, 2 counters
 - Power 4, 8 counters
 3. Pentium IV, 18 counters
 - UltraSparc II, 2 counters
 4. IA64, 4 counters
 - MIPS R14K, 2 counters
 5. Alpha 21x64, 2 counters

Itanium 2 Block Diagram

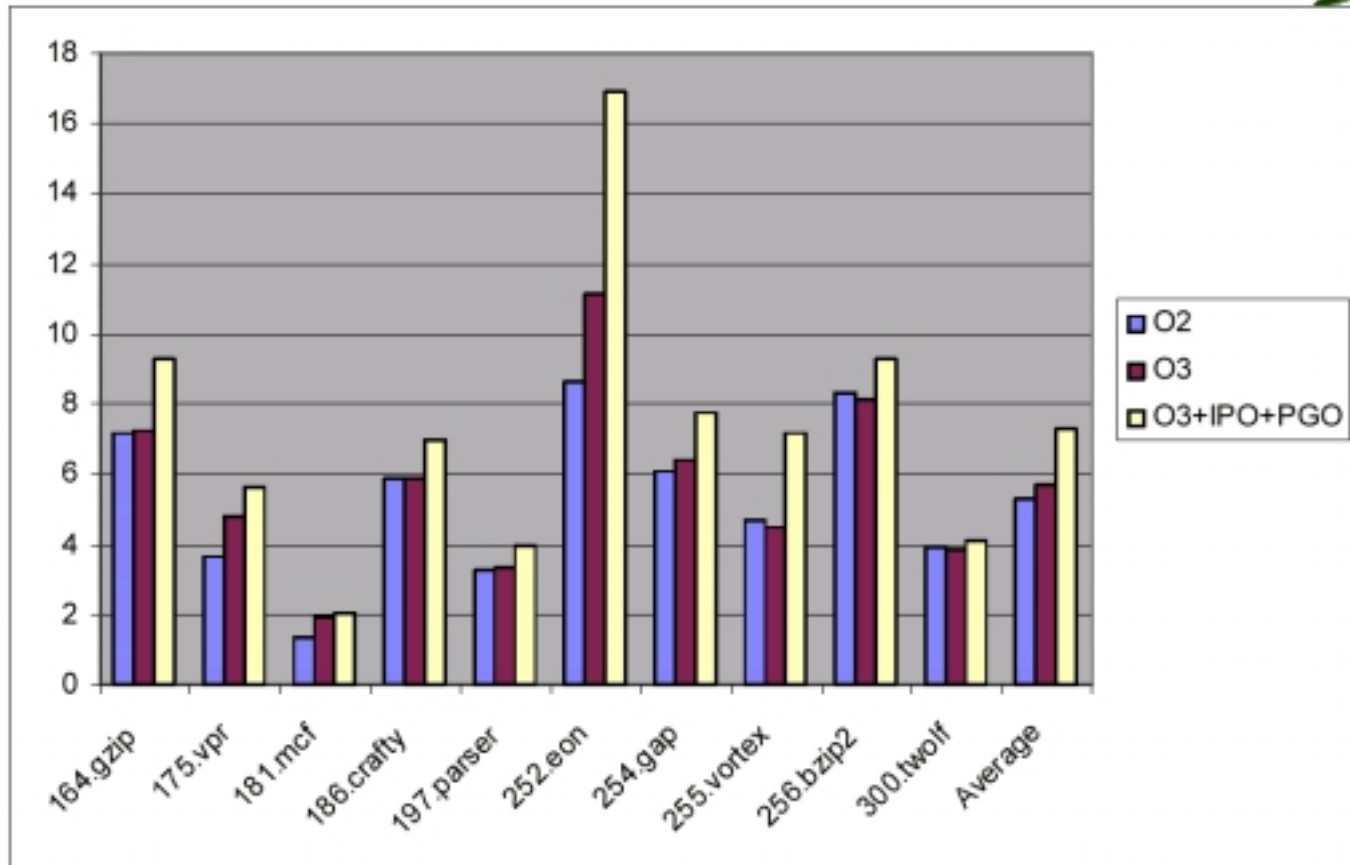


Itanium 2 Block Diagram

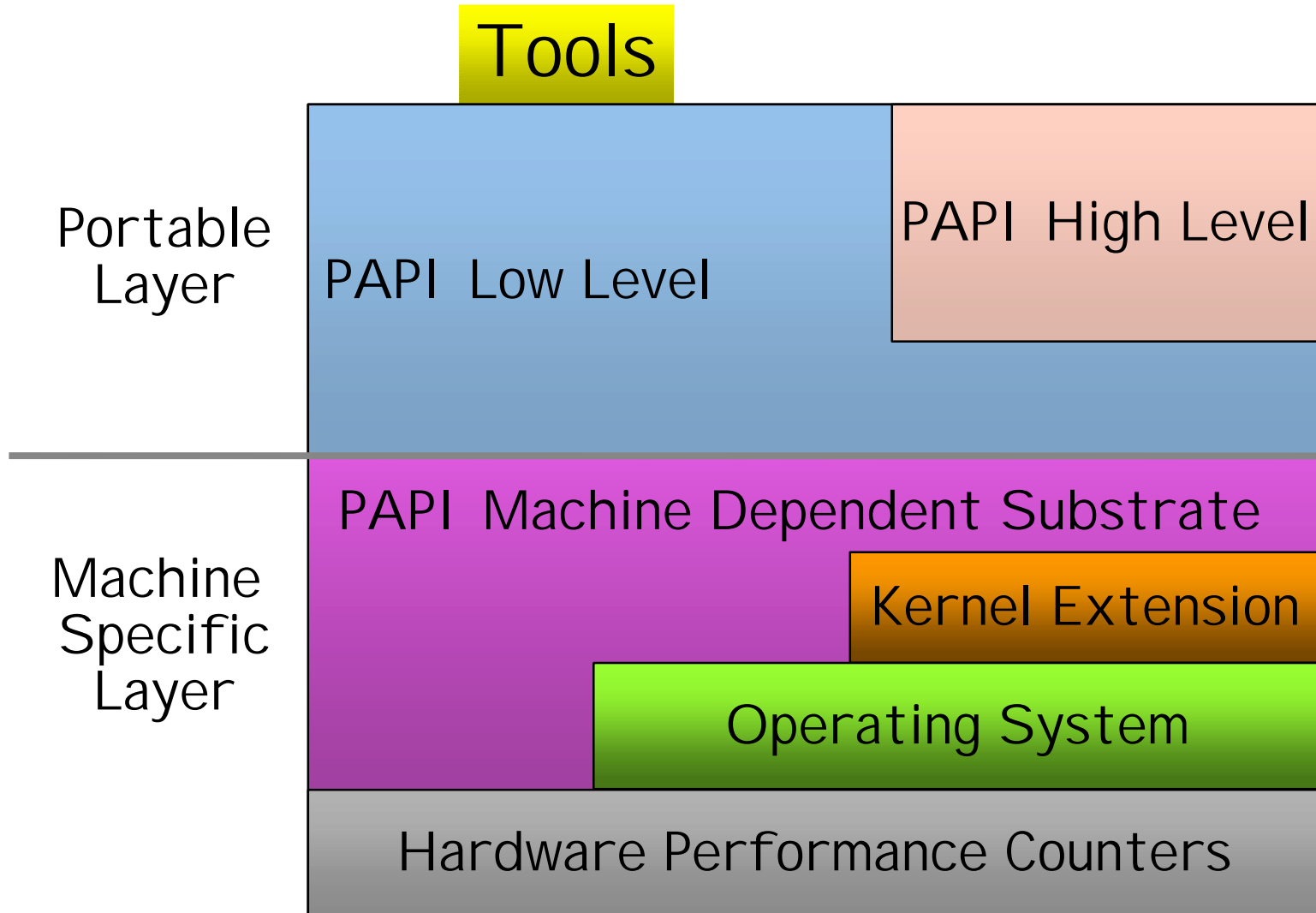


Importance of Optimization

Example: Speed up from Static Compiler Optimization on Itanium-1 in 2002 (SpecInt)



PAPI Implementation



- Proposed standard set of event names deemed most relevant for application performance tuning
- No standardization of the actual definition
- Mapped to native events on a given platform

- PAPI supports approximately 100 preset events.
 - Preset events are mappings from symbolic names to machine specific definitions for a particular hardware event.
 - Example: **PAPI_TOT_CYC**
 - PAPI also supports presets that may be derived from the underlying hardware metrics
 - Example: **PAPI_L1_DCM**

Sample Preset Listing

Test case 8: Available events and hardware information.

```
-----
Vendor string and code      : GenuineIntel (-1)
Model string and code      : Celeron (Mendocino) (6)
CPU revision               : 10.000000
CPU Megahertz              : 366.504944
-----
```

Name	Code	Derived	Description (Note)
PAPI_L1_DCM	0x80000000	No	Level 1 data cache misses
PAPI_L1_ICM	0x80000001	No	Level 1 instruction cache misses
PAPI_L2_DCM	0x80000002	No	Level 2 data cache misses
PAPI_L2_ICM	0x80000003	No	Level 2 instruction cache misses
PAPI_L3_DCM	0x80000004	No	Level 3 data cache misses
PAPI_L3_ICM	0x80000005	No	Level 3 instruction cache misses
PAPI_L1_TCM	0x80000006	Yes	Level 1 cache misses
PAPI_L2_TCM	0x80000007	No	Level 2 cache misses
PAPI_L3_TCM	0x80000008	No	Level 3 cache misses
PAPI_CA_SNP	0x80000009	No	Requests for a snoop
PAPI_CA_SHR	0x8000000a	No	Requests for shared cache line
PAPI_CA_CLN	0x8000000b	No	Requests for clean cache line
PAPI_CA_INV	0x8000000c	No	Requests for cache line inv.



- PAPI supports native events:
 - An event countable by the CPU can be counted even if there is no matching preset PAPI event.
 - The developer uses the same API as when setting up a preset event, but a CPU-specific bit pattern is used instead of the PAPI event definition.

- Meant for application programmers wanting coarse-grained measurements
- As easy to use as the calls present in IRIX.
- Requires no setup code
- Restrictions:
 - Only PAPI preset events may be used
 - Not thread safe (in PAPI 2.3.4)

- `PAPI_num_counters()`
- `PAPI_start_counters(int *cntrs, int alen)`
- `PAPI_stop_counters(long_long *vals, int alen)`
- `PAPI_accum_counters(long_long *vals, int alen)`
- `PAPI_read_counters(long_long *vals, int alen)`
- `PAPI_flops(float *rtime, float *ptime,
 long_long *flpins, float *mflops)`

- Increased efficiency and functionality over the high level PAPI interface
- **Approximately 60 functions**
(http://icl.cs.utk.edu/projects/papi/files/html_man/papi.html#4)
- Thread-safe for all 1:1 thread libraries.
(Native, OpenMP, Pthreads, etc...)
- Supports both presets and native events

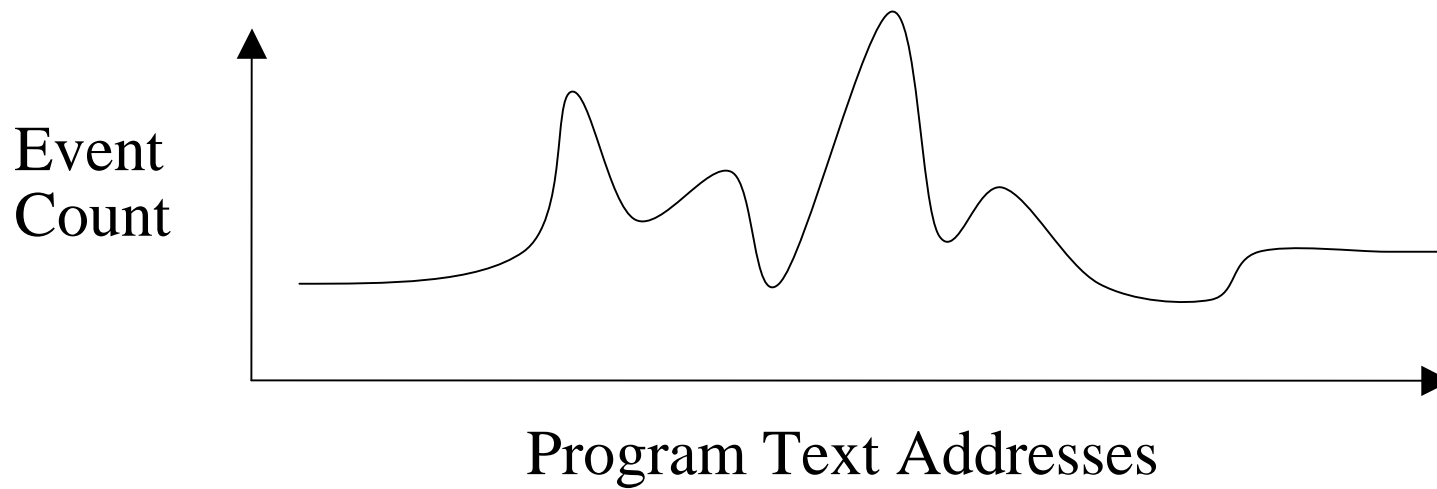
- API Calls for:
 - Counter multiplexing
 - Callbacks on user defined overflow value
 - SVR4 compatible profiling
 - Processor information
 - Address space information
 - Static and dynamic memory information
 - Accurate and low latency timing functions
 - Hardware event inquiry functions
 - Eventset management functions
 - Simple locking operations

- Multiplexing allows simultaneous use of more counters than are supported by the hardware.
 - This is accomplished through timesharing the counter hardware and extrapolating the results.
- Users can enable multiplexing with one API call and then use PAPI normally.

Interrupts on Counter Overflow

- PAPI provides the ability to call user-defined handlers when a specified event exceeds a specified threshold.
- For systems that do not support counter overflow at the hardware level, PAPI emulates this in software at the user level.
 - Code must run a reasonable length of time.

- On overflow of hardware counter, dispatch a signal/interrupt.
- Get the address at which the code was interrupted.
- Store counts of interrupts for *each* address.
- Vendor/GNU **prof** and **gprof** (**-pg** and **-p** compiler options) use interval timers.



- The result: A probabilistic distribution of where the code spent its time and why.

- <http://icl.cs.utk.edu/projects/papi/>
 - Software and documentation
 - Reference materials
 - Papers and presentations
 - Third-party tools
 - Mailing lists

- Supported Platforms
 - IBM 604, 604e, Power 3, 4
 - Intel x86, Pentium IV
 - Intel Itanium I, II
 - Sun UltraSparc I/II/III
 - SGI R10K/R12K/R14K
 - Compaq Alpha 21164/21264 with DADD/DCPI
 - Cray T3E
 - AMD Opteron
 - Windows/x86 (not PIV)
- Enhancements
 - Static/dynamic memory info
 - Multiplexing improvements
 - Lots of bug fixes

- Using lessons learned from years earlier
 - Substrate code: 90% used only 10% of the time
- Complete internal redesign for:
 - Efficiency
 - Robustness
 - Feature Set
 - Elegance
 - Portability

- Multiway multiplexing
 - Use all available counter registers instead of one per time slice.
- Superb performance
 - Example: On Pentium 4, a PAPI_read() costs 230 cycles. (Register read costs 100 cycles)
- Full native and programmable event support
 - Thresholding
 - Instruction matching
 - Per event counting domains

- Third-party interface
 - Allows PAPI to control counters in other processes
- Internal timer/signal/thread abstractions
 - Support signal forwarding
- Additional internal layered API to support robust extensions

- Advanced profiling interface
 - Support profiling on multiple counters
 - Support hardware or operating system assisted profiling
- New sampling interface
 - P4, IA64 provide Event Address Registers of BTB misses, Cache misses, TLB misses, etc...
- Revised memory usage API
 - Process footprint

- System-wide and process wide counting implementation
- Expanded high level API
 - Thread safe
- New language bindings:
 - Java
 - Lisp
 - Matlab

- Initial release expected around SC 2003 (limited scope)
- Additional platforms will be added as they come available:
 - Cray X-1 (partially complete)
 - Nec SX-6
 - Blue Gene (BG/L)

- API for runtime code patching
- permits instrumentation of a program while it executes
- Provides process of index applied to multiple systems
- Includes metadata from inserted code

- Parallel framework based on early DynInst
- Async/Sync operation
- Functions for getting data back to tool
- Integrated with POE
- Available on all HPC platforms (and Windows)
- Breakpoints
- Arbitrary ins. points
- Full Loop, CFG and Basic Block decoding

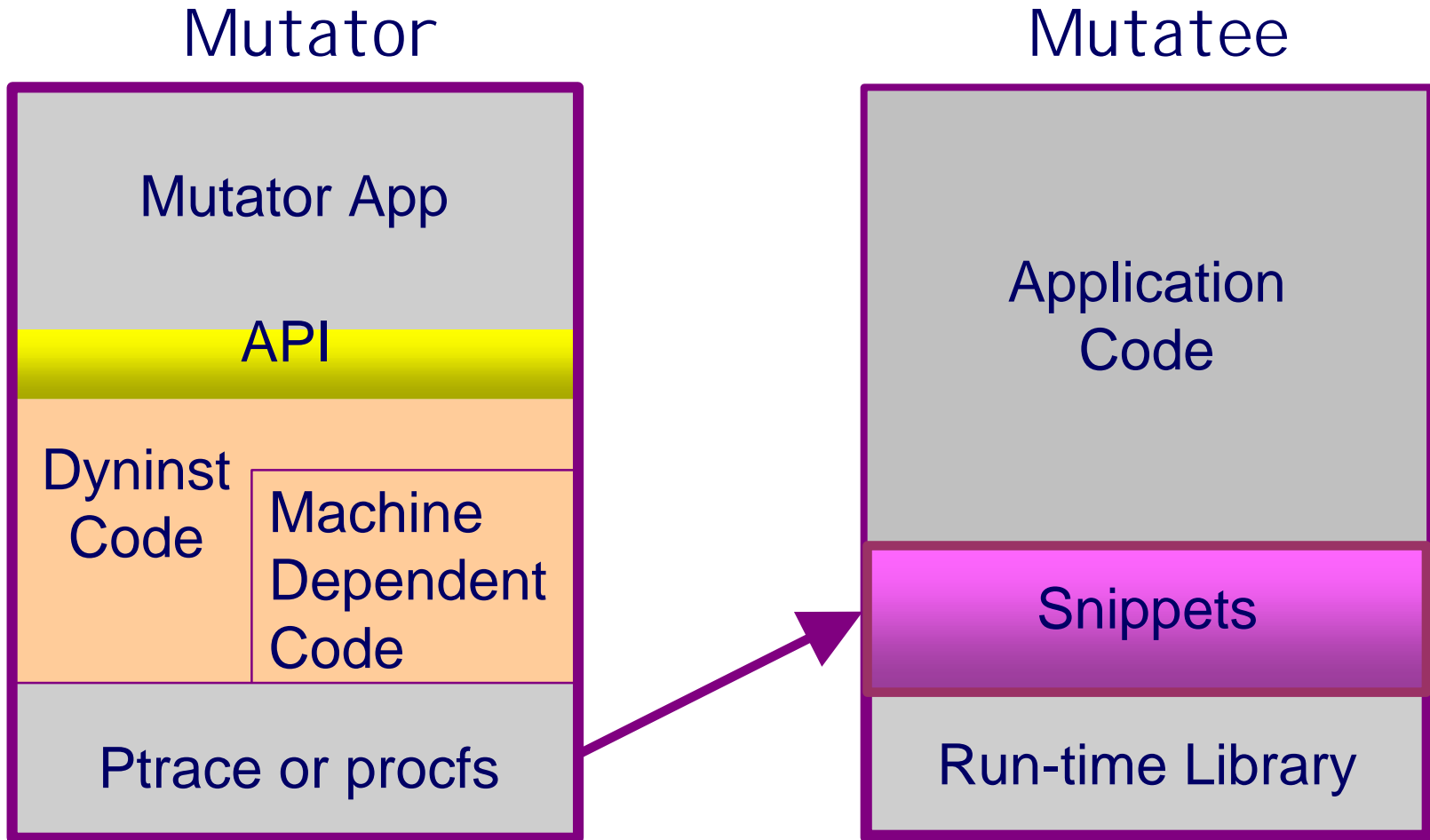
- Popularized by James Larus with EEL: An Executable Editor Library at U. Wisc.
 - <http://www.cs.wisc.edu/~larus/eel.html>
- Technology matured by Dr. Bart Miller and (now Dr.) Jeff Hollingsworth at U. Wisc.
 - DynInst Project at U. Maryland
 - <http://www.dyninst.org/>
 - IBM's DPCL: A Distributed DynInst
 - <http://oss.software.ibm.com/dpcl/>

Why Runtime Code Patching?

- Performance measurement
 - Recording application behavior
- Correctness debugging
 - Fast conditional breakpoints
 - Data breakpoints
- Execution driven simulation
 - Architecture studies
- Testing
 - Code coverage testing
 - Forcing hard to execute paths to be taken

- No forethought needed
 - No user inserted probes
 - No special compiling or linking
 - Start anytime during execution
- Only insert code when needed
 - No wasted checks for “disabled” code
 - Can add new code during execution

Structure of the Dyninst Library



- Provides:
 - Functions for control of mutatee
 - Runtime code generation
 - Information about mutatee
- A set of C++ classes
 - Bpatch_Thread
 - BPatch_image
 - BPatch_snippet
 - BPatch_variableExpr

- Platform Independent Representation
 - Same code can be inserted into apps on any system
- Simple Abstract Syntax Tree
 - Can refer to application state (variables & params)
 - Includes simple looping construct
 - Permits calls to application subroutines
- Type Checking
 - Ensures that snippets are type compatible
 - Based on structural equivalence
 - allows flexibility when adding new code

- Access to local (stack) variables
- Complex types
 - non-integer scalars
 - structures
 - arrays
 - Fortran common blocks
- Correctness debugging
 - print contents of data structures

- New classes added to dyninstAPI
 - BPatch_basicBlock
 - BPatch_sourceBlock
 - BPatch_flowGraph
- Arbitrary Instrumentation points
 - Conservative base trampoline
- Base trampoline deletion

- Code Coverage needs basic block level instrumentation
 - dyninstAPI used to support function level instrumentation for sparc-solaris
 - added arbitrary instrumentation points for SPARC
- More state must be maintained in base trampolines
 - save/restore condition codes before/after arbitrary instrumentation points
 - Sparc arch supports user mode condition-code write/read for version v8plus and later

- Dynamic memory access instrumentation
 - collect low level memory accesses
 - with the flexibility of dynamic instrumentation
- Possible applications
 - offline performance analysis (Sigma etc.)
 - online optimization
 - tools to catch memory errors

- Finding memory access instructions
 - loads, stores, prefetches
- Builds on Arbitrary Instrumentation
- Decoded instruction information
 - type of instruction
 - constants and registers involved in computing
 - the effective address
 - the number of bytes moved
 - available in the mutator before execution
- Memory access snippets
 - effective address in process space
 - byte count
 - available in mutatee at execution time

- Re-running with same modified code requires
 - Parse debug symbols
 - Two processes
 - DyninstAPI shared library
 - Time to re-insert instrumentation
- *Not* a checkpointing mechanism
 - Mutated binary begins at the top of main()
 - Data initialized as in original binary*

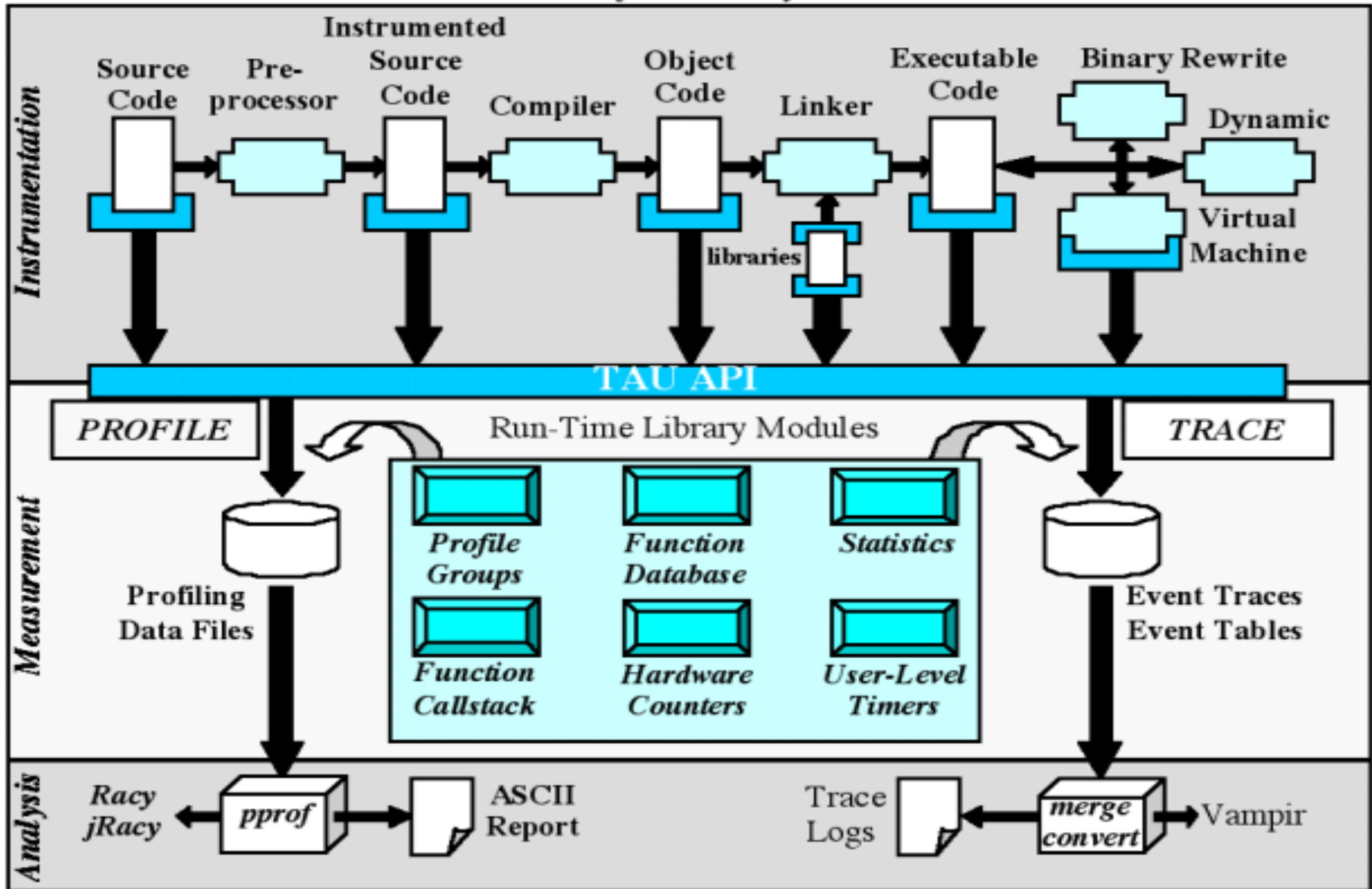
- Supported platforms
 - SPARC (Solaris)
 - x86 (Solaris, Linux, NT)
 - Alpha (Tru64 UNIX)
 - MIPS (IRIX)
 - Power/PowerPC (AIX)
- Software available on the web
 - **`http://www.dyninst.org`**
 - Includes TCL command tool
 - Over 250 sites have downloaded

Paradyn from U. Wisconsin

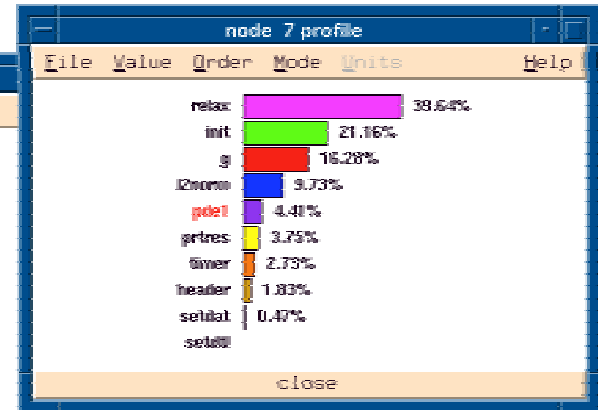
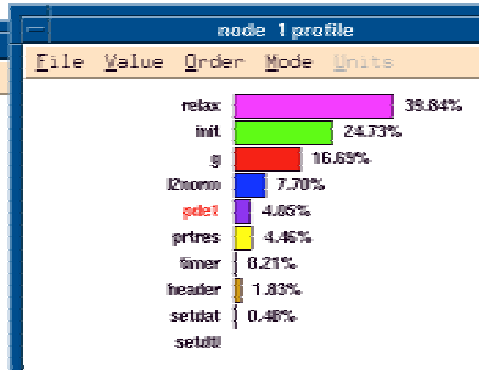
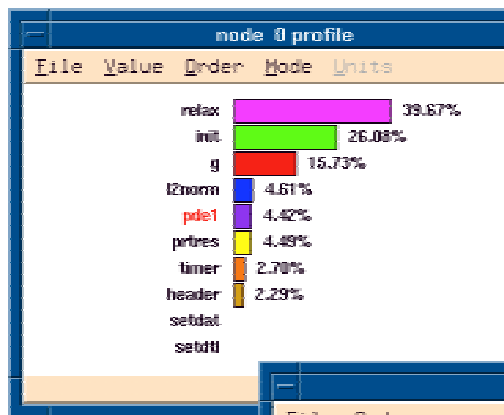
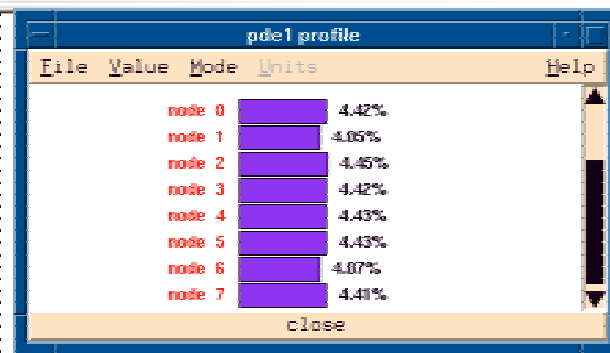
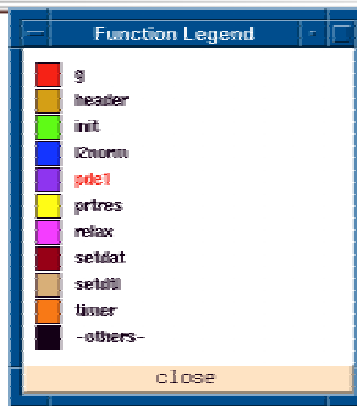
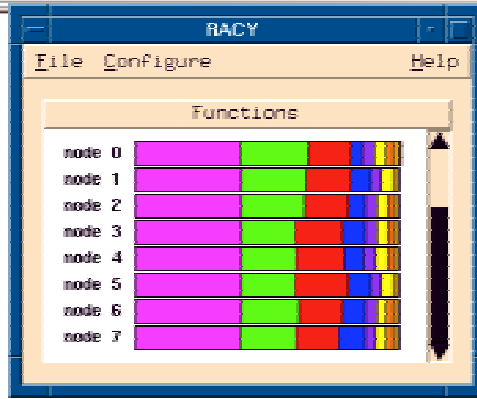
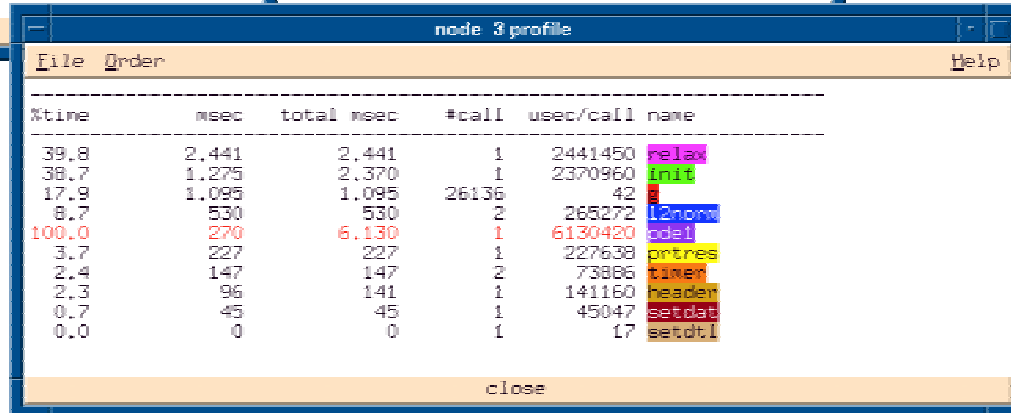
- From Bart Miller Group
- Dynamically generated instructions based on testing
- Different features of visualization plugins
- Work on MPI, PAPI, means to discover bottlenecks
- <http://www.paradyn.org>

- From Allen Malony's Group
- Source or binary based instrumentation
- Supports all forms of parallelism
- Integration with Vampir for trace visualization of MPI, OpenMP and both
- <http://www.cs.uoregon.edu/research/paracomp/tau>

TAU System Architecture



TAU/ParaProf Screenshot

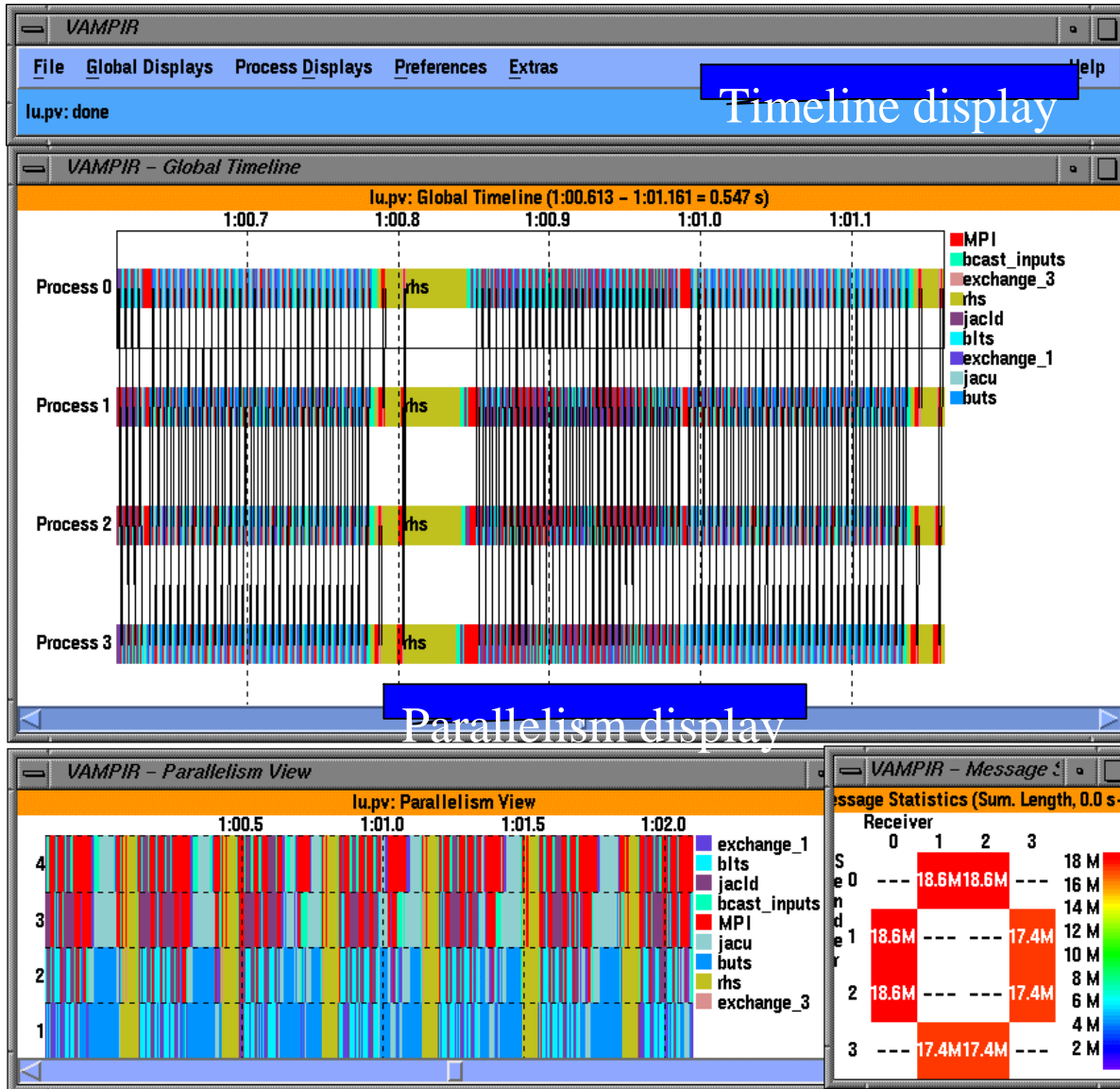



node 3 profile

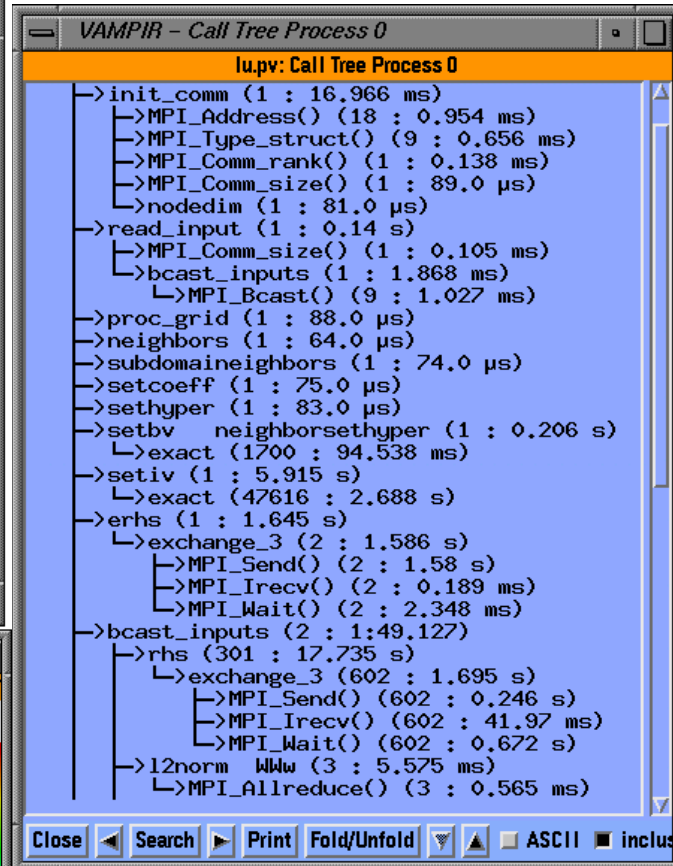
File	Order	%time	msec	total msec	#call	usec/call	name
relax		39.8	2,441	2,441	1	2441450	relax
init		38.7	1,275	2,370	1	2370960	init
g		17.9	1,095	1,095	26136	42	g
l2norm		8.7	530	530	2	265272	l2norm
pde1		100.0	270	6,130	1	6130420	pde1
prtres		3.7	227	227	1	227638	prtres
timer		2.4	147	147	2	73886	timer
header		2.3	96	141	1	141160	header
setdat		0.7	45	45	1	45047	setdat
setdll		0.0	0	0	1	17	setdll

close

Vampir (NAS Parallel Benchmark – LU)

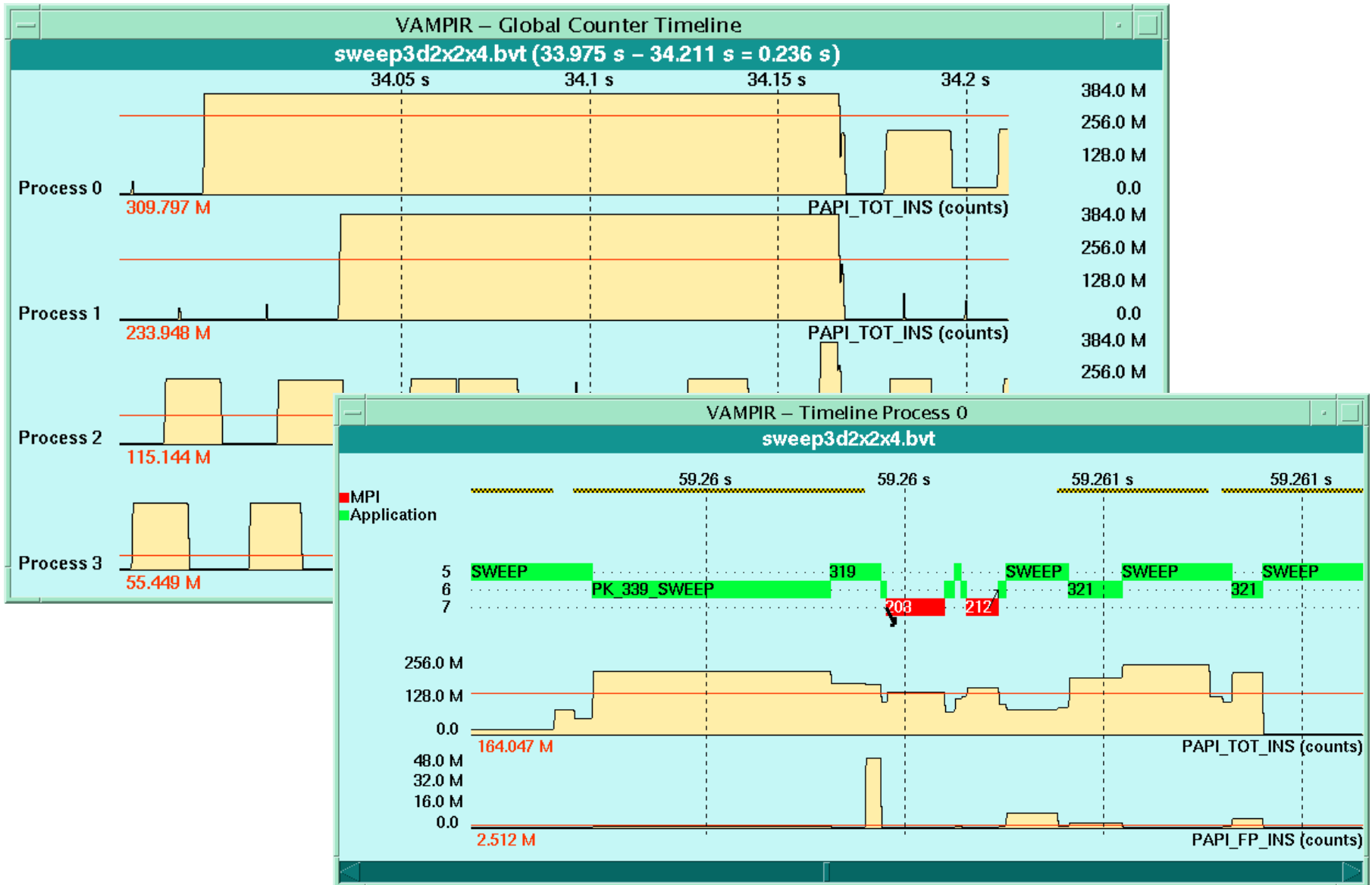


Call tree display

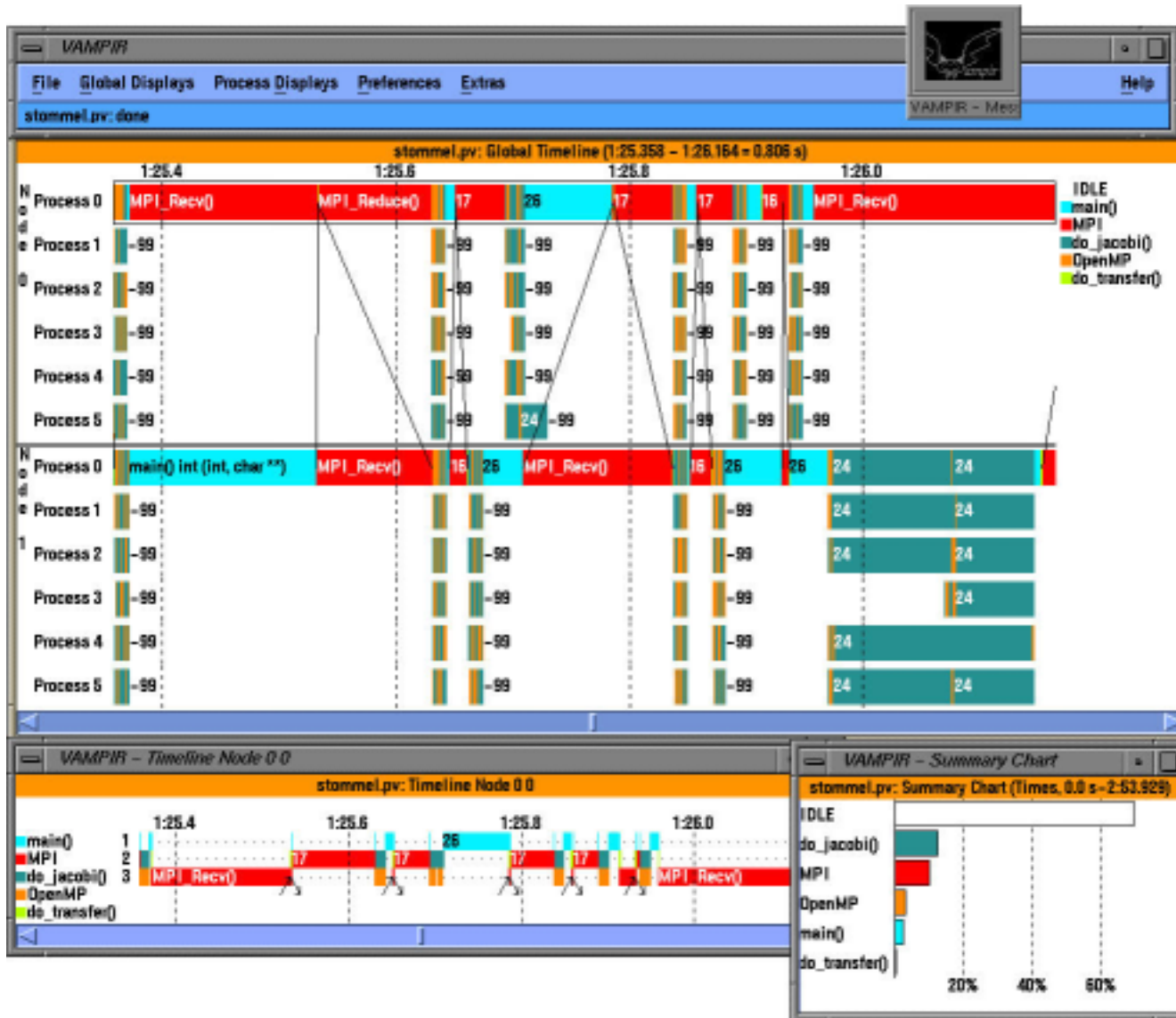


Communications display

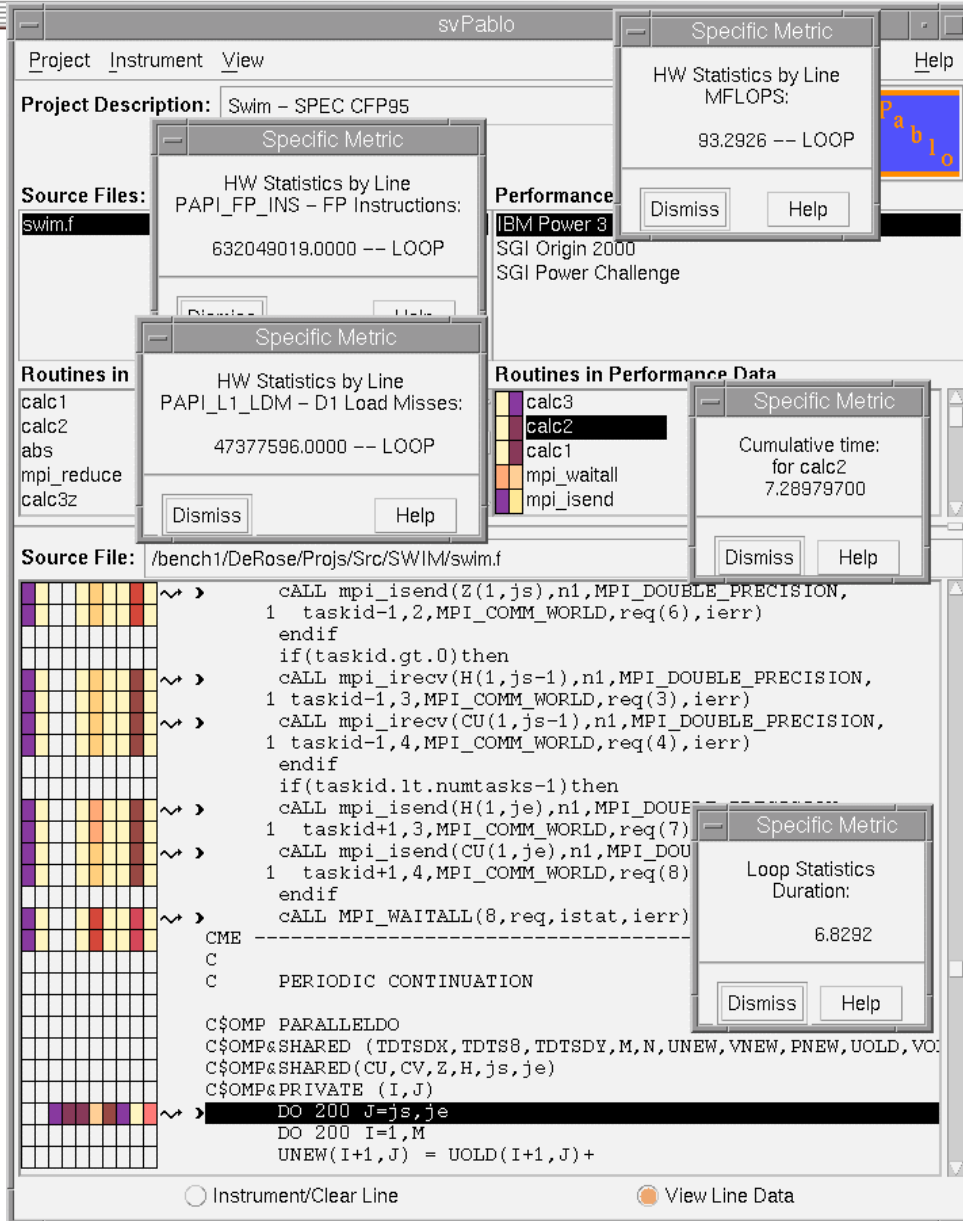
Vampir: PAPI Counter Traces



TAU OpenMP+MPI Vampir Visualization



SvPablo from UIUC



The screenshot shows the SvPablo application window with the following components:

- Project Description:** Swim - SPEC CFP95
- Source Files:** swim.f
- Routines in:** calc1, calc2, abs, mpi_reduce, calc3z
- Source File:** /bench1/DeRose/Projs/Src/SWIM/swim.f
- Performance Data:** IBM Power 3, SGI Origin 2000, SGI Power Challenge
- Routines in Performance Data:** calc3, calc2, calc1, mpi_waitall, mpi_isend
- Source Code:**

```

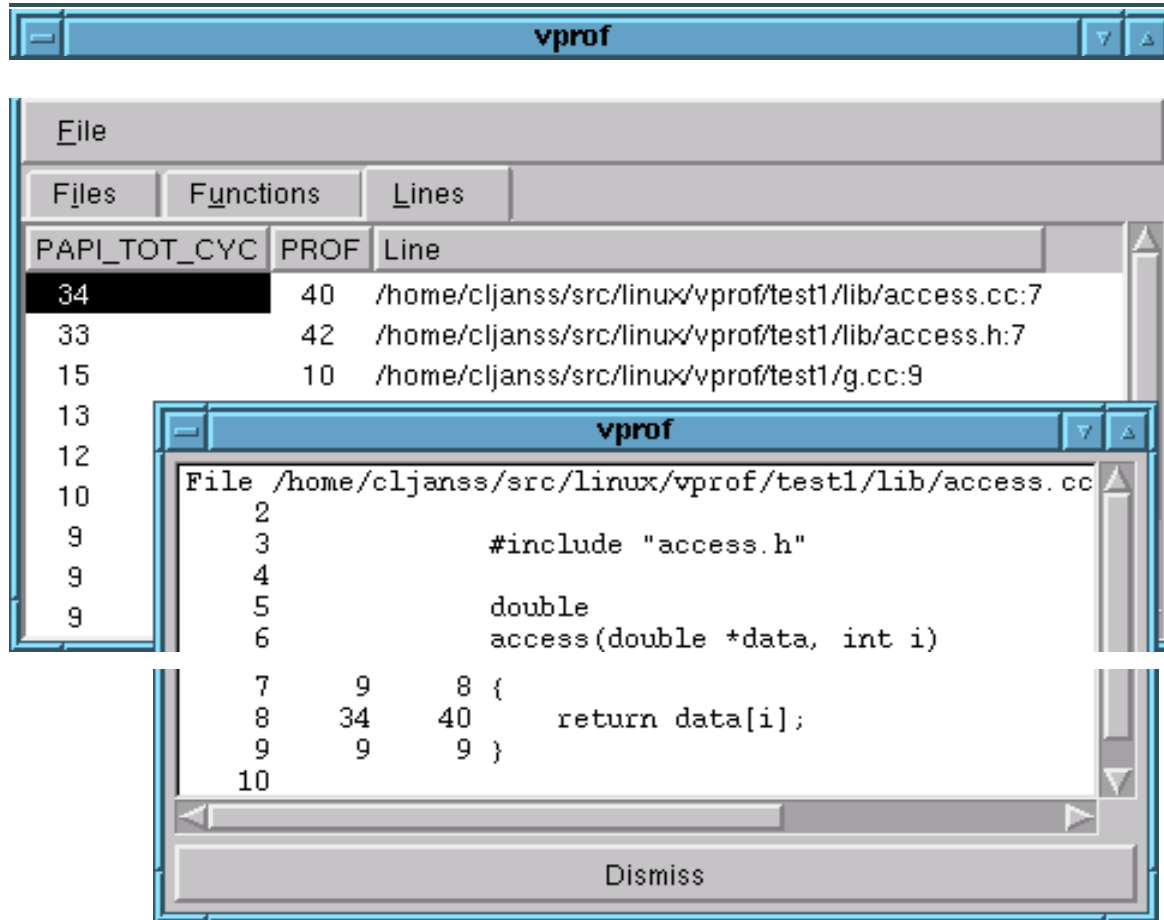
cALL mpi_isend(Z(1,js),n1,MPI_DOUBLE_PRECISION,
1 taskid-1,2,MPI_COMM_WORLD,req(6),ierr)
endif
if(taskid.gt.0)then
cALL mpi_irecv(H(1,js-1),n1,MPI_DOUBLE_PRECISION,
1 taskid-1,3,MPI_COMM_WORLD,req(3),ierr)
cALL mpi_irecv(CU(1,js-1),n1,MPI_DOUBLE_PRECISION,
1 taskid-1,4,MPI_COMM_WORLD,req(4),ierr)
endif
if(taskid.lt.numtasks-1)then
cALL mpi_isend(H(1,je),n1,MPI_DOUE
1 taskid+1,3,MPI_COMM_WORLD,req(7)
cALL mpi_isend(CU(1,je),n1,MPI_DOU
1 taskid+1,4,MPI_COMM_WORLD,req(8)
endif
cALL MPI_WAITALL(8,req,istat,ierr)
CME
C
C PERIODIC CONTINUATION

C$OMP PARALLELDO
C$OMP&SHARED (TDTSDX,TDTS8,TDTS DY,M,N,UNEW,VNEW,PNEW,UOLD,VO
C$OMP&SHARED(CU,CV,Z,H,js,je)
C$OMP&PRIVATE (I,J)
DO 200 J=js,je
DO 200 I=1,M
UNEW(I+1,J) = UOLD(I+1,J)+

```
- Pop-up Windows:**
 - Specific Metric (HW Statistics by Line):** MFLOPS: 93.2926 -- LOOP
 - Specific Metric (HW Statistics by Line):** PAPI_FP_INS - FP Instructions: 632049019.0000 -- LOOP
 - Specific Metric (HW Statistics by Line):** PAPI_L1_LDM - D1 Load Misses: 47377596.0000 -- LOOP
 - Specific Metric (Cumulative time):** for calc2: 7.28979700
 - Specific Metric (Loop Statistics):** Duration: 6.8292

- Source based instrumentation of loops and function calls
- Supports serial and MPI jobs
- Freely available
- Rough F90 parser

Vprof from Sandia National Laboratory



The screenshot shows the vprof application interface. The main window displays a table with columns for PAPI_TOT_CYC, PROF, and Line. The table contains the following data:

PAPI_TOT_CYC	PROF	Line
34	40	/home/cljanss/src/linux/vprof/test1/lib/access.cc:7
33	42	/home/cljanss/src/linux/vprof/test1/lib/access.h:7
15	10	/home/cljanss/src/linux/vprof/test1/g.cc:9

An inset window shows the source code for the file /home/cljanss/src/linux/vprof/test1/lib/access.cc. The code is as follows:

```

2
3     #include "access.h"
4
5     double
6     access(double *data, int i)
7     {
8         9         8 {
9         34     40     return data[i];
9         9         9 }
10

```

The inset window also has a "Dismiss" button at the bottom.

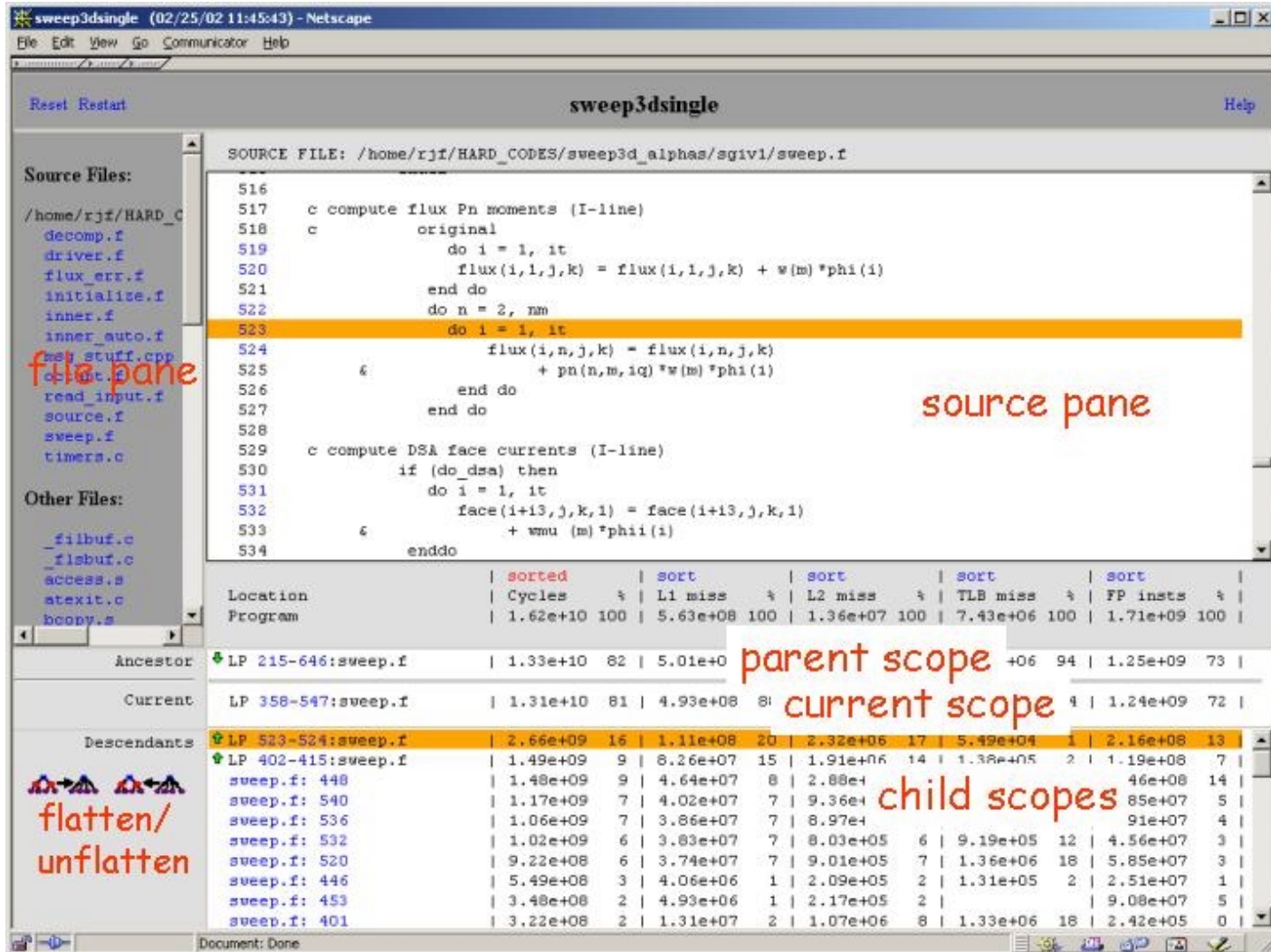
- Based on statistical sampling of the hardware counters
- Must instrument the source
- Ported to other architectures for generalized use
- Parallel codes with some modification
- Not actively supported

<http://aros.ca.sandia.gov/~cljanss/perf/vprof>

- Tools for:
 - Collecting raw statistical profiles
 - Conversion of profiles into platform independent XML
 - Synthesizing browsable representations that correlate metrics with source code
- <http://www.hipersoft.rice.edu/hpctoolkit>

- Collection: papirun/hvprof, equivalent to SGI's “ssrun”
- Loop/CFG recovery from binary: bloop
- Data formatting: papiprof
- Data display and exploration: hpcview
- Call stack profiles: csprof
- Data is aggregated into an XML database
- HPCView is a Java applet that generates dynamic HTML

HPCView Screenshot



The screenshot shows the HPCView interface for a program named 'sweep3dsingle'. The source code pane displays Fortran code with several nested loops. The scope table below the code provides performance metrics for different execution scopes.

Source File: /home/rjf/HARD_CODES/sweep3d_alpha/sgiv1/sweep.f

```

516
517 c compute flux Pn moments (I-line)
518 c     original
519     do i = 1, it
520         flux(i,l,j,k) = flux(i,l,j,k) + w(m)*phi(i)
521     end do
522     do n = 2, nm
523         do i = 1, it
524             flux(i,n,j,k) = flux(i,n,j,k)
525             + pn(n,m,iq)*w(m)*phi(i)
526         end do
527     end do
528
529 c compute DSA face currents (I-line)
530 if (do_dsa) then
531     do i = 1, it
532         face(i+13,j,k,1) = face(i+13,j,k,1)
533         + wmu (m)*phii(i)
534     enddo

```

Scope Table:

Location	sorted	sort	sort	sort	sort	sort	sort
Program	Cycles %	L1 miss %	L2 miss %	TLB miss %	FP insts %		
LP 215-646:sweep.f	1.33e+10	82	5.01e+0	+06	94	1.25e+09	73
LP 358-547:sweep.f	1.31e+10	81	4.93e+08	8	4	1.24e+09	72
LP 523-524:sweep.f	2.66e+09	16	1.11e+08	20	2.32e+06	17	5.49e+04
LP 402-415:sweep.f	1.49e+09	9	8.26e+07	15	1.91e+06	14	1.38e+05
sweep.f: 448	1.48e+09	9	4.64e+07	8	2.88e+		46e+08
sweep.f: 540	1.17e+09	7	4.02e+07	7	9.36e+		85e+07
sweep.f: 536	1.06e+09	7	3.86e+07	7	8.97e+		91e+07
sweep.f: 532	1.02e+09	6	3.83e+07	7	8.03e+05	6	9.19e+05
sweep.f: 520	9.22e+08	6	3.74e+07	7	9.01e+05	7	1.36e+06
sweep.f: 446	5.49e+08	3	4.06e+06	1	2.09e+05	2	1.31e+05
sweep.f: 453	3.48e+08	2	4.93e+06	1	2.17e+05	2	9.08e+07
sweep.f: 401	3.22e+08	2	1.31e+07	2	1.07e+06	8	1.33e+06

Annotations:

- file pane:** Points to the source code pane.
- source pane:** Points to the source code pane.
- parent scope:** Points to the 'LP 215-646:sweep.f' row.
- current scope:** Points to the 'LP 358-547:sweep.f' row.
- child scopes:** Points to the 'LP 523-524:sweep.f' and 'LP 402-415:sweep.f' rows.
- flatten/unflatten:** Points to the tree view of descendant scopes.

- Libraries and tools for machine information, memory information, aggregate counts, derived metrics and statistical profiles
- Targeted for x86 and IA64 systems
- <http://perfsuite.ncsa.uiuc.edu>

- **psinv**: Gather information on a processor and the PAPI events it supports
- **psrun**: Collection of aggregate/derived counts or statistical profiles of unmodified binaries
- **psprocess**: Formatting and output of psrun data into text or HTML

Address <http://perfuite.ncss.ucl.ac.uk/psprocess/psprocess-example-hwpc-vsl.html>

PerfSuite Hardware Performance Report

Derived Metrics	
Graduated instructions per cycle	1.765
Graduated floating point instructions per cycle	0.145
Floating-point percentage of all graduated instructions	8.219%
Graduated loads & stores per cycle	0.219
Graduated loads & stores per floating point instruction	1.514
Data references per instruction	0.029
Ratio of floating point instructions to L1 dcache accesses	2.848
L1 instruction cache miss ratio	0.003
L1 data cache read miss ratio	0.966
L3 data cache miss ratio	0.957
L3 cache data read ratio	0.992
L3 cache instruction miss ratio	0.206
Ratio of mispredicted to correctly predicted branches	0.093
L1 cache data hit rate	0.776
L2 cache data hit rate	0.467
L3 cache data hit rate	0.714
L1 cache line reuse (data)	3.462
L2 cache line reuse (data)	0.877
L3 cache line reuse (data)	2.498
Bandwidth used to L1 cache (MB/s)	1262.361
Bandwidth used to L2 cache (MB/s)	1326.512
Bandwidth used to L3 cache (MB/s)	385.087
Percentage of cycles with no instruction issue	10.41%
Percentage of cycles waiting for memory access	43.14%
MIPS (CPU cycles)	1412.19
MIPS (effective)	1394.349
MFLOPS (CPU cycles)	115.905
MFLOPS (effective)	114.441
Processor utilization	98.74%
PAPI Event	Counter Value



Psrun Statistical Profile Example Output

ICL

Function Summary

Samples	Self %	Total %	Function
1839543	35.01%	35.01%	inl3130
541829	10.31%	45.32%	ns5_core
389741	7.42%	52.74%	inl0100
355349	6.76%	59.51%	spread_q_bsplines
213172	4.06%	63.56%	gather_f_bsplines
200546	3.82%	67.38%	do_longrange
182691	3.48%	70.86%	make_bsplines
149924	2.85%	73.71%	ewald_LRcorrection
112883	2.15%	75.86%	inl3100
105317	2.00%	77.86%	solve_pme
92257	1.76%	79.62%	flincs

- **libperfsuite**: Provides simple wrappers for machine information, process memory usage and high-precision timing
- **libpshwpc**: Provides simple wrappers that are used to collect hardware performance data

```
program mxm  
include 'fperfsuite.h'
```

```
c Initialize libpshwpc
```

```
    call PSF_hwpc_init(ierr)
```

```
c Start performance counting using libpshwpc
```

```
    call PSF_hwpc_start(ierr)
```

```
c Stop hardware performance counting and write the  
c results to a file named 'perf.XXXXXX' (XXXXXX will be  
c replaced by the process ID of the program)
```

```
    call PSF_hwpc_stop('perf', ierr)
```

```
c Shutdown use of libpshwpc and the underlying libraries
```

```
    call PSF_hwpc_shutdown(ierr)
```

- Environment variables and XML input file dictate what gets measured

- Command line utility to gather aggregate counts.
 - PAPI version has been tested on IA32 & IA64
 - User can manually instrument code for more specific information
 - Reports derived metrics like SGI's perfex
- Libhpm for manual instrumentation
- Hpmviz is a GUI to view resulting data

<http://www.ncsa.uiuc.edu/UserInfo/Resources/Software/Tools/HPMToolkit>

hpmviz Screenshot

HPMviz

Label	Excl. Dur. (s)	Incl. Dur. (s)	Count
swim			
-Loop 200	83.837	83.837	1200
-Loop 100	57.339	57.339	1200
-Loop 300	57.237	57.237	1199
-Calc3	2.019	60.261	1198
-loop 110	1.155	1.155	1200
-MPI in Calc3	1.004	1.004	1198
-Calc2	0.891	85.453	1200
-MPI Calc1 end	0.782	0.782	1200
-MPI Calc2 end	0.688	0.688	1200
-Initial	0.205	0.205	1
-Calc1	0.144	59.457	1200
-Calc3z	0.05	0.05	1
-MPI Calc2 start	0.037	0.037	1200
-MPI Calc1 start	0.036	0.036	1200

MPI Calc2 end Metrics	
Task	Value
Thread	29
Count	1200
ExcSec	0.688
IncSec	0.688
Ld/TLB miss	158.069
Total LS	45.39
Instr/LS	1.762
MIPS	116.291
Instr/Cycles	0.342
HW FP/Cycle	0
FPI+FMA	0.022
Mflop/s	0.032
FMA %	0
Comp. Int.	0

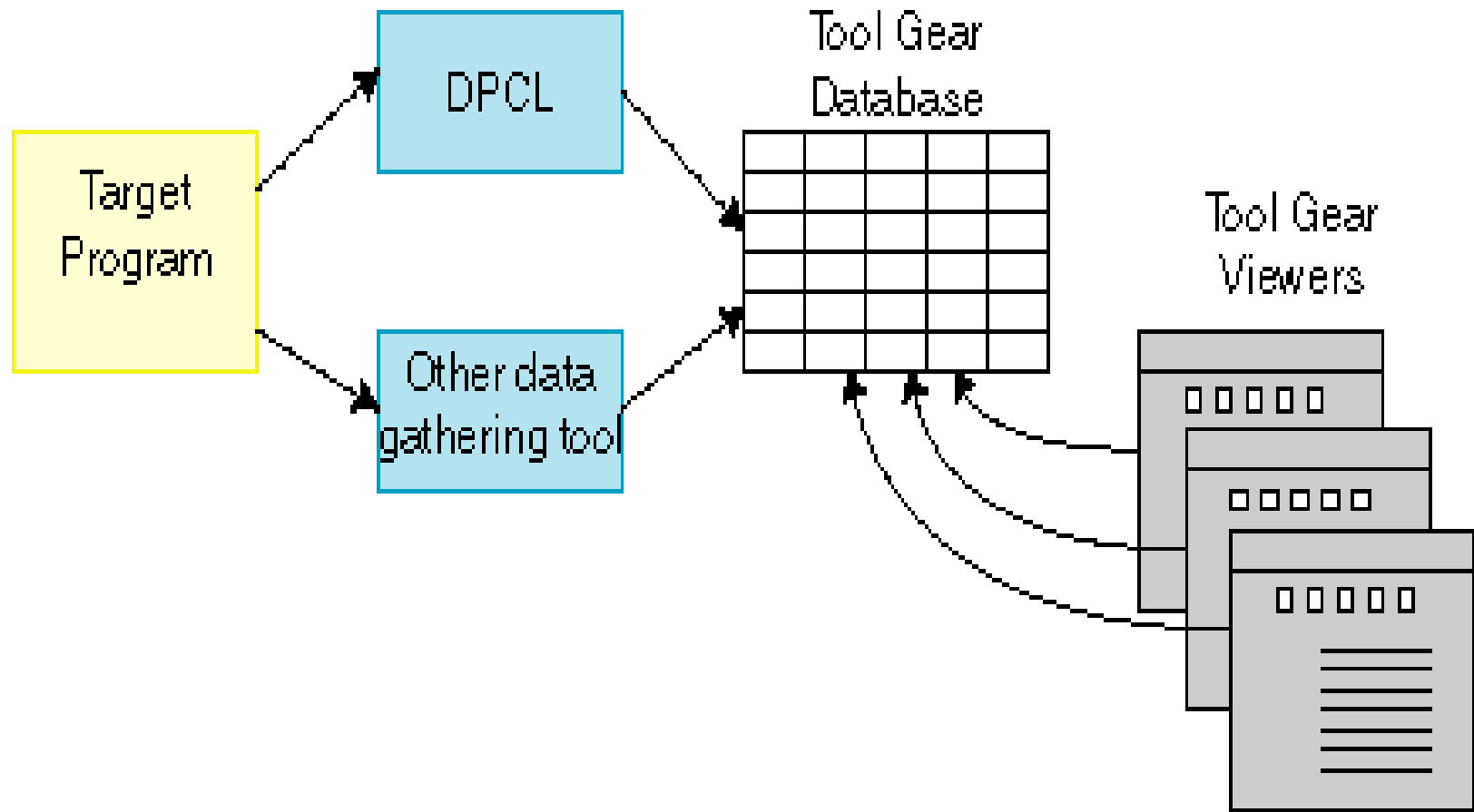
```

call f_hpmstart( 29, "MPI Calc2 end" )
if(taskid.eq.0)then
  call mpi_irecv(VNEW(1:1),ni,MPI_DOUBLE_PRECISION,
  1 numtasks-1,1,MPI_COMM_WORLD,req(1),ierr)
endif
if(taskid.eq.numtasks-1)then
  call mpi_irecv(UNEW(1,n+1),ni,MPI_DOUBLE_PRECISION,
  1 0,2,MPI_COMM_WORLD,req(2),ierr)
  call mpi_irecv(PNEW(1,n+1),ni,MPI_DOUBLE_PRECISION,
  1 0,3,MPI_COMM_WORLD,req(3),ierr)
endif
if(taskid.eq.numtasks-1)then
  call mpi_isend(VNEW(1,n+1),ni,MPI_DOUBLE_PRECISION,
  1 0,1,MPI_COMM_WORLD,req(4),ierr)
endif
if(taskid.eq.0)then
  call mpi_isend(UNEW(1,1),ni,MPI_DOUBLE_PRECISION,
  1 numtasks-1,2,MPI_COMM_WORLD,req(5),ierr)
  call mpi_isend(PNEW(1,1),ni,MPI_DOUBLE_PRECISION,
  1 numtasks-1,3,MPI_COMM_WORLD,req(6),ierr)
endif
if(taskid.eq.0.or.taskid.eq.numtasks-1)then
  call MPI_WAITALL(6,req,istat,ierr)
endif
do 215 i=1,M
  UNEW(i+1,N+1) = UNEW(i+1,1)
  VNEW(i,i) = VNEW(i,N+1)
  PNEW(i,N+1) = PNEW(i,1)
c 215 CONTINUE
  UNEW(i,N+1) = UNEW(M+1,1)
  VNEW(M+1,1) = VNEW(1,N+1)
  PNEW(M+1,N+1) = PNEW(1,1)
  if(taskid.eq.numtasks-1)then
    call mpi_irecv(UNEW(1,N+1),1,MPI_DOUBLE_PRECISION,
    1 0,2,MPI_COMM_WORLD,req(1),ierr)
    call mpi_irecv(PNEW(M+1,N+1),1,MPI_DOUBLE_PRECISION,
    1 0,3,MPI_COMM_WORLD,req(2),ierr)
  endif
  if(taskid.eq.0)then
    call mpi_irecv(VNEW(M+1,1),1,MPI_DOUBLE_PRECISION,
    1 numtasks-1,1,MPI_COMM_WORLD,req(3),ierr)
  endif
  if(taskid.eq.0)then
    call mpi_isend(UNEW(M+1,1),1,MPI_DOUBLE_PRECISION,
    1 numtasks-1,2,MPI_COMM_WORLD,req(4),ierr)
    call mpi_isend(PNEW(1,1),1,MPI_DOUBLE_PRECISION,
    1 numtasks-1,3,MPI_COMM_WORLD,req(5),ierr)
  
```

```
#include "libhpm.h"  
hpmInit( tasked, "my program" );  
hpmStart( 1, "outer call" );  
do_work();  
hpmStart( 2, "computing meaning of life" );  
do_more_work();  
hpmStop( 2 );  
hpmStop( 1 );  
hpmTerminate( taskID );
```

- Dynamic instrumentation and analysis suite from LLNL
- Based on DPCL from IBM
 - Tested only on AIX
- Qt Front end can theoretically accept data from any source
- GUI displays instrumentable points
- Instrumented points update display with data in real time
- http://www.llnl.gov/CASC/tool_gear

ToolGear Architecture



ToolGear Screenshot: Instrumentation

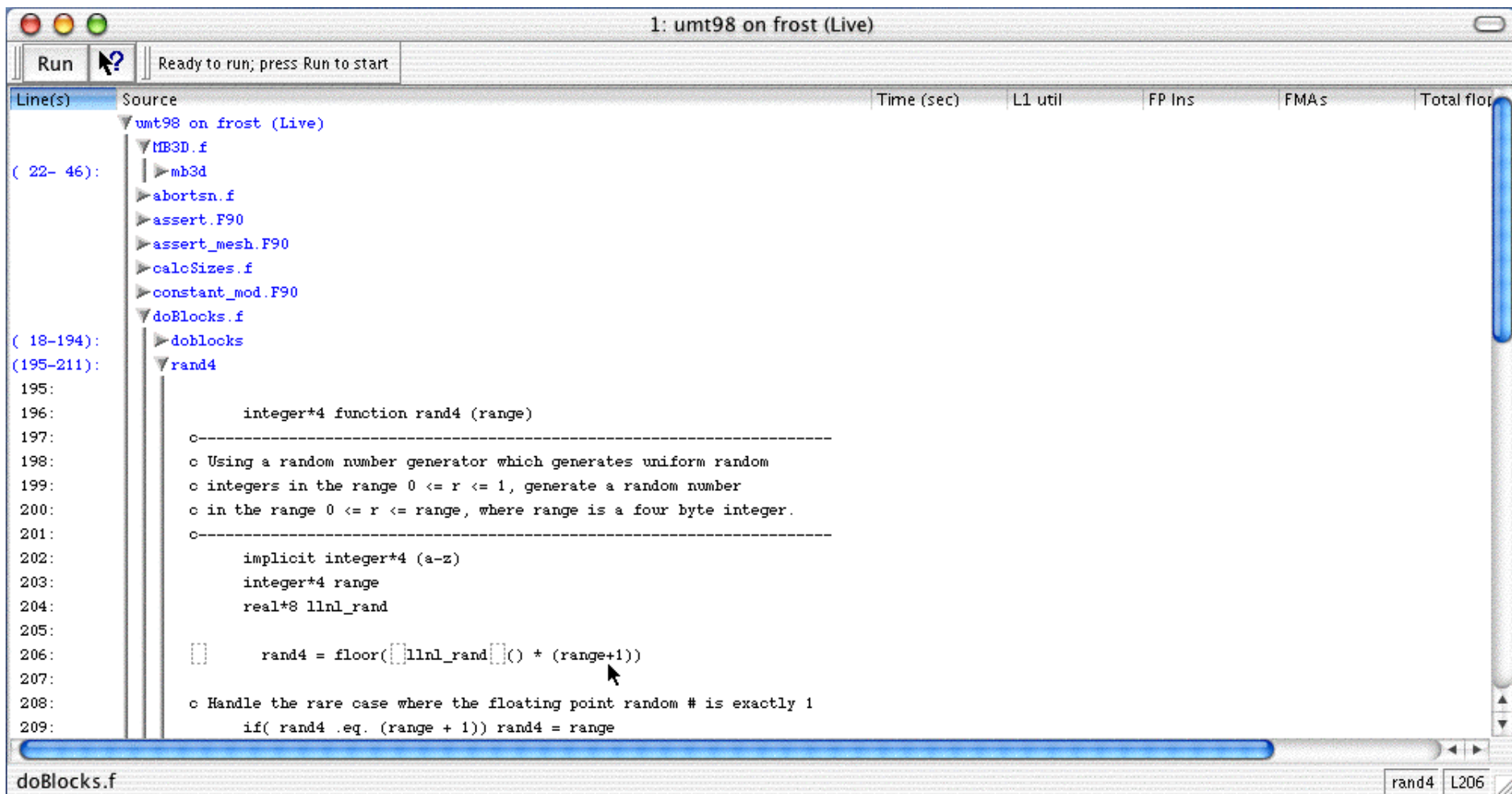
1: testcpmod_mpi on snow (Live)

Run Program has terminated

Line(s)	Source	L1 util	FP Ins	FMA's	Total flop	FLOP/sec
38:	for(i = 0; i < 10; i++) {					
39:	/* Tiled */					
40:	init_array();					
41:						
42:	printf("Doing %d flops of tiled test\n", FLOPS);					
43:	do_tiled_cache_test(FLOPS);	0.991517	1.04858e+07	1.04858e+07	2.09716e+07	1.46422e+08
44:						
45:	/* Untiled */					
46:	init_array();					
47:						
48:	printf("Doing %d flops of untiled test\n", FLOPS);					
49:	do_untiled_cache_test(FLOPS);	0.937468	1.04858e+07	1.04858e+07	2.09716e+07	1.43855e+08
50:						
51:	/* Indexed */					
52:	init_array();					
53:						
54:	printf("Doing %d flops of indirect address test\n", FLOPS);					
55:	do_indirect_address_test(FLOPS);	0.512543	1.04861e+07	1.04861e+07	2.09722e+07	3.26941e+07
56:	printf("Done with series %i\n", i);					
57:	fprintf(stderr, "Done with series %i\n", i);					
58:						

8 data pts: Max 0.93747 (Rank 0/Thread 1) Min 0.937464 (7/1) Mean 0.937468 StdDev 1.90709e-06 Sum 7.49974

Mean L49



The screenshot shows the ToolGear IDE interface. At the top, there is a window title "1: umt98 on frost (Live)" and a "Run" button with a question mark icon. Below the toolbar, a status bar indicates "Ready to run; press Run to start".

The main area is divided into two panes. The left pane shows a tree view of the source files:

- umt98 on frost (Live)
 - MB3D.f
 - mb3d
 - aborts.n.f
 - assert.F90
 - assert_mesh.F90
 - calcSizes.f
 - constant_mod.F90
 - doBlocks.f
 - doblocks
 - rand4

The right pane shows the source code for the selected file, `rand4`. The code is as follows:

```

integer*4 function rand4 (range)
c-----
c Using a random number generator which generates uniform random
c integers in the range 0 <= r <= 1, generate a random number
c in the range 0 <= r <= range, where range is a four byte integer.
c-----
implicit integer*4 (a-z)
integer*4 range
real*8 llnl_rand

rand4 = floor(llnl_rand() * (range+1))

c Handle the rare case where the floating point random # is exactly 1
if( rand4 .eq. (range + 1)) rand4 = range
  
```

The status bar at the bottom of the code editor shows the current file is `doBlocks.f`, the current function is `rand4`, and the current line is `L206`.

1: pi3 on snow (Live)

Run ? Program has terminated

Line(s)	Source	Count	Max time	Min time	Mean time	App %	MPI %
	pi3 on snow (Live)	16	6.57	0.518	2.541	23.41	100
	pi3.f	16	6.57	0.518	2.541	23.41	100
(1- 71):	main	16	6.57	0.518	2.541	23.41	100
43:	call MPI_BCAST(n,1,MPI_INTEGER,0,MPI_COMM_WORLD,ierr)	8	2.19	0.228	0.611	5.63	24.06
55:	call MPI_REDUCE(myp_i,pi,1,MPI_DOUBLE_PRECISION,MPI_SUM,0,	8	4.38	0.29	1.93	17.78	75.94
1:	c pi3.f						
2:	c slightly modified from the MPICH pi3 example code						
3:	c*****						
4:	c pi.f - compute pi by integrating $f(x) = 4/(1 + x^2)$						
5:	c						

pi3.f main L2

- A portable tool to dynamically instrument serial and parallel programs for the purpose of performance analysis.
- Simple and intuitive command line interface like GDB.
- Java/Swing GUI.
- Instrumentation is done through the run-time insertion of function calls to specially developed performance probes.



Why the “Dyna” in DynaProf?

- Instrumentation:
 - Functions are contained in shared libraries.
 - Calls to those functions are generated at run-time.
 - Those calls are dynamically inserted into the program’s address space.
- Built on DynInst and DPCL
- Can choose the mode of instrumentation, currently:
 - Function Entry/Exit
 - Call site Entry/Exit
 - One-shot

- Make collection of run-time performance data easy by:
 - Avoiding instrumentation and recompilation
 - Avoiding perturbation of compiler optimizations
 - Providing complete language independence
 - Allowing multiple insert/remove instrumentation cycles

No source code required!

DynaProf Goals 2

- Using the same tool with different probes
- Providing useful and meaningful probe data
- Providing different kinds of probes
- Allowing custom probe development Make collection of run-time performance data easy by:

No source code required!

- perfometerprobe
 - Visualize hardware event rates in “real-time”
- papiprobe
 - Measure any combination of PAPI presets and native events
- wallclockprobe
 - Highly accurate elapsed wallclock time in microseconds.
- The latter 2 probes report:
 - Inclusive
 - Exclusive
 - 1 Level Call Tree

- Probes export a few functions with loosely standardized interfaces.
- Easy to roll your own.
 - If you can code a timer, you can write a probe.
- DynaProf detects thread model.
- Probes dictate how the data is recorded and visualized.

- For threaded code, use the same probe!
- Dynaprof detects threads and loads a special version of the probe library.
- Each probe specifies what to do when a new thread is discovered.
- Each thread gets the same instrumentation.

- Can count any PAPI preset or Native event accessible through PAPI
- Can count multiple events
- Supports PAPI multiplexing
- Supports multithreading
 - AIX: SMP, OpenMP, Pthreads
 - Linux: SMP, OpenMP, Pthreads

- Counts microseconds using RTC
- Supports multithreading
 - AIX: SMP, OpenMP, Pthreads
 - Linux: SMP, OpenMP, Pthreads

- The wallclock and PAPI probes produce very similar data.
- Both use a parsing script written in Perl.
 - wallclockrpt <file>
 - papiproberpt <file>
- Produce 3 profiles
 - Inclusive: $T_{function} = T_{self} + T_{children}$
 - Exclusive: $T_{function} = T_{self}$
 - 1-Level Call Tree: $T_{child} = \text{Inclusive } T_{function}$

```
./dynaprof
(dynaprof) load tests/swim
(dynaprof) list
DEFAULT_MODULE
swim.F
libm.so.6
libc.so.6
(dynaprof) list swim.F
MAIN__
inital_
calc1_
calc2_
calc3z_
calc3_
(dynaprof) list swim.F MAIN__
Entry
  Call s_wsle
  Call do_lio
  Call e_wsle
  Call s_wsle
  Call do_lio
  Call e_wsle
  Call calc3_

(dynaprof) use probes/papiprobe
Module papiprobe.so was loaded.
Module libpapi.so was loaded.
Module libperfctr.so was loaded.
(dynaprof) instr module swim.F calc*
swim.F, inserted 4 instrumentation points
(dynaprof) run
papiprobe: output goes to
/home/mucci/dynaprof/tests/swim.1671
```

```
(dynaprof) use probes/papiprobe PAPI_TOT_CYC, PAPI_TOT_INS
Module papiprobe.so was loaded.
Module libpapi.so was loaded.
Module libperfctr.so was loaded.
(dynaprof) instr function swim.F calc*
Swim.F, inserted 3 instrumentation points
(dynaprof) instr
calc1_
calc2_
calc3_
calc3z_
```


SWIM Benchmark: Cycles &

Instructions

Exclusive Profile of Metric PAPI_TOT_INS.

Name	Percent	Total	Calls
TOTAL	100	1.723e+09	1
calc2	38.28	6.598e+08	120
calc1	32.31	5.567e+08	120
calc3	22.33	3.847e+08	118
unknown	7.084	1.221e+08	1

Inclusive Profile of Metric PAPI_TOT_INS.

Name	Percent	Total	SubCalls
TOTAL	100	1.723e+09	0
calc2	39.42	6.793e+08	1680
calc1	35.28	6.08e+08	1800
calc3	22.87	3.942e+08	1652

1-Level Inclusive Call Tree of Metric PAPI_TOT_INS.

Parent/-Child	Percent	Total	Calls
TOTAL	100	1.723e+09	1
calc1	100	6.08e+08	120
- fsav	0.02065	1.255e+05	120
- mpi_irecv	0.03132	1.904e+05	120
- mpi_isend	0.05911	3.593e+05	120
- mpi_isend	0.06434	3.912e+05	120
-mpi_waitall	0.9013	5.479e+06	120
- mpi_irecv	0.03132	1.904e+05	120
- mpi_irecv	0.03132	1.904e+05	120
- mpi_isend	0.05356	3.256e+05	120
- mpi_isend	0.05079	3.088e+05	120
-mpi_waitall	6.813	4.142e+07	120
- mpi_irecv	0.03132	1.904e+05	120
- mpi_irecv	0.03132	1.904e+05	120
- mpi_isend	0.07504	4.562e+05	120
- mpi_isend	0.06757	4.108e+05	120
-mpi_waitall	0.161	9.791e+05	120

Exclusive Profile of Metric PAPI_TOT_CYC.

Name	Percent	Total	Calls
TOTAL	100	3.181e+09	1
calc2	34.85	1.108e+09	120
calc1	33.48	1.065e+09	120
calc3	26.1	8.301e+08	118
unknown	5.568	1.771e+08	1

Inclusive Profile of Metric PAPI_TOT_CYC.

Name	Percent	Total	SubCalls
TOTAL	100	3.181e+09	0
calc2	35.98	1.144e+09	1680
calc1	35.61	1.133e+09	1800
calc3	26.88	8.55e+08	1652

1-Level Inclusive Call Tree of Metric PAPI_TOT_CYC.

Parent/-Child	Percent	Total	Calls
TOTAL	100	3.181e+09	1
calc1	100	1.133e+09	120
- fsav	0.03432	3.887e+05	120
- mpi_irecv	0.07356	8.332e+05	120
- mpi_isend	0.0663	7.51e+05	120
- mpi_isend	0.0739	8.371e+05	120
-mpi_waitall	0.7189	8.143e+06	120
- mpi_irecv	0.1646	1.864e+06	120
- mpi_irecv	0.03407	3.859e+05	120
- mpi_isend	0.1867	2.115e+06	120
- mpi_isend	0.06067	6.872e+05	120
-mpi_waitall	4.22	4.78e+07	120
- mpi_irecv	0.03979	4.506e+05	120
- mpi_irecv	0.03008	3.407e+05	120
- mpi_isend	0.1014	1.148e+06	120
- mpi_isend	0.07568	8.573e+05	120
-mpi_waitall	0.1076	1.219e+06	120



SWIM Benchmark: Instructions

ICL

per Cycle

Exclusive Profile of Metric PAPI_TOT_INS.

Name	Percent	Total	Calls
TOTAL	100	1.723e+09	1
calc2	38.28	6.598e+08	120
calc1	32.31	5.567e+08	120
calc3	22.33	3.847e+08	118
unknown	7.084	1.215e+09	1

Inclusive Profile of Metric PAPI_TOT_INS.

Name	Percent	Total	SubCalls
TOTAL	100	1.723e+09	0
calc2	39.42	6.793e+08	1680
calc1	35.28	6.08e+08	1800
calc3	22.87	3.942e+08	1652

1-Level Inclusive Call Tree of Metric PAPI_TOT_INS.

Parent/-Child	Percent	Total	Calls
TOTAL	100	1.723e+09	1
calc1	100	6.08e+08	1800
- fsav	0.02065	1.215e+09	1
- mpi_irecv	0.03132	1.215e+09	1
- mpi_isend	0.05911	3.847e+08	3
- mpi_isend	0.06434	3.847e+08	3
-mpi_waitall	0.9013	5.567e+08	5
- mpi_irecv	0.03132	1.215e+09	1
- mpi_irecv	0.03132	1.215e+09	1
- mpi_isend	0.05356	3.231e+08	3
- mpi_isend	0.05079	3.08e+08	3
-mpi_waitall	6.813	4.142e+08	12
- mpi_irecv	0.03132	1.904e+08	1
- mpi_irecv	0.03132	1.904e+08	1
- mpi_isend	0.07504	4.562e+05	1
- mpi_isend	0.06757	4.108e+05	1
-mpi_waitall	0.161	9.791e+05	1

Exclusive Profile of Metric PAPI_TOT_CYC.

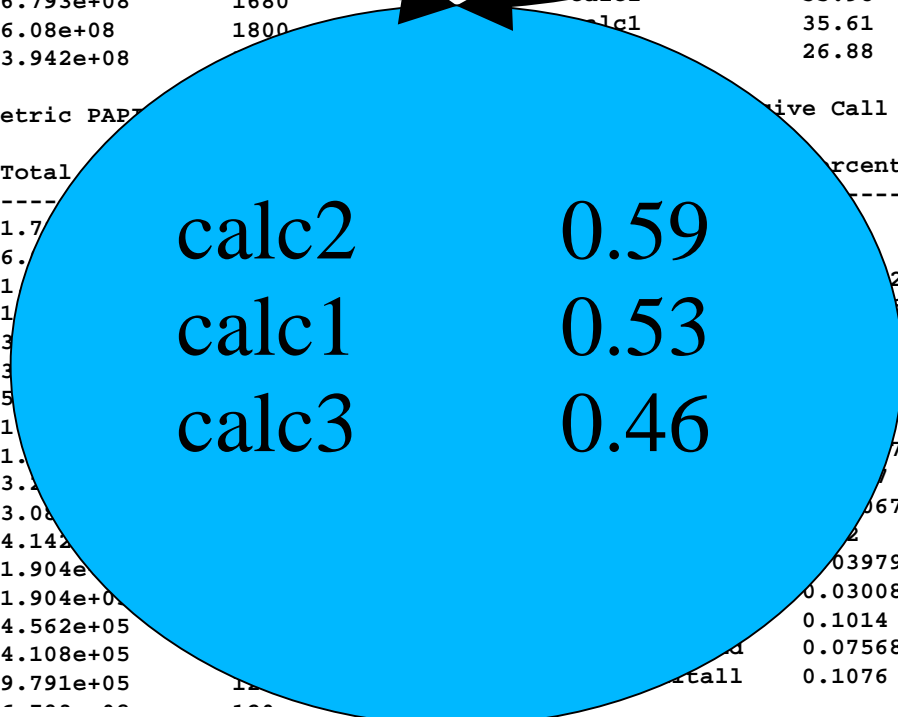
Name	Percent	Total	Calls
TOTAL	100	3.181e+09	1
calc2	34.85	1.108e+09	120
calc1	33.48	1.067e+09	120
calc3	26.1	8.301e+08	118
unknown	5.568	1.771e+09	1

Inclusive Profile of Metric PAPI_TOT_CYC.

Name	Percent	Total	SubCalls
TOTAL	100	3.181e+09	0
calc2	35.98	1.144e+09	1680
calc1	35.61	1.133e+09	1800
calc3	26.88	8.55e+08	1652

1-Level Inclusive Call Tree of Metric PAPI_TOT_CYC.

Parent/-Child	Percent	Total	Calls
TOTAL	100	3.181e+09	1
calc1	100	1.133e+09	120
- fsav	0.02065	3.887e+05	120
- mpi_irecv	0.03132	8.332e+05	120
- mpi_isend	0.05911	7.51e+05	120
- mpi_isend	0.06434	8.371e+05	120
-mpi_waitall	0.9013	8.143e+06	120
- mpi_irecv	0.03132	1.864e+06	120
- mpi_irecv	0.03132	3.859e+05	120
- mpi_isend	0.05356	2.115e+06	120
- mpi_isend	0.05079	6.872e+05	120
-mpi_waitall	6.813	4.78e+07	120
- mpi_irecv	0.03132	4.506e+05	120
- mpi_irecv	0.03132	3.407e+05	120
- mpi_isend	0.07504	1.148e+06	120
- mpi_isend	0.06757	8.573e+05	120
-mpi_waitall	0.161	1.219e+06	120



- Displays module tree for instrumentation
- Simple selection of probes and instrumentation points
- Single-click execution of common DynaProf commands
- Coupling of probes and visualizers (e.g. Perfometer)

- It's a bit rough
- Supported Platforms
 - Using DynInst 3.0
 - Linux 2.x
 - AIX 4.3/5
 - Using DPCL (formal MPI support)
 - AIX 4.3
 - AIX 5
- Available as a development snapshot from:
- Includes:
 - Java/Swing GUI
 - User's Guide
 - Probe libraries

<http://www.cs.utk.edu/~mucci/dynaprof>

Thanks!

Philip J. Mucci

LBNL 50B-3207

510 486-8616

pjmucci@lbl.gov

mucci@cs.utk.edu