

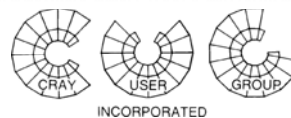


The Supercomputer Company

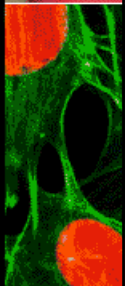
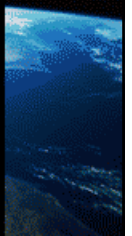
The New Generation of Cray Performance Tools

Luiz DeRose
Sr. Principal Engineer
Tools Manager

Albuquerque, NM



CUG 2005



Outline

- Previous status
- The Cray tools strategy
- Cray Tools
 - HWPC
 - CrayPat
 - Cray Apprentice²
- Future directions
- Conclusions



The Starting Point

- Missing a complete solution
 - Cray X1
 - Very good measurement infrastructure
 - Very poor user interface
 - XT3 & XD1
 - No in-house solution
 - Dependence on external providers
- Needed integration across platforms



The Cray Tools Strategy

- Must be easy to use
- Integrated performance tools solution
- Strategy based on the three main steps normally used for application optimization and tuning:
 1. Debug application
 - Nobody really cares how fast the program can compute the wrong answer
 - TotalView is our debugger of choice on all platforms
 2. Single processor and vector optimization
 - Make the common case fast
 3. Parallel processing and I/O optimization
 - Communication / barriers / etc.



Single Processor Optimization

- Answer the following questions:
 1. Does my program have performance problems?
 - Time / Resource / Hardware Counters measurement
 - Provides overall view of the program execution
 - `% pat_hwpc a.out`



Single Processor Optimization (2)

- Answer the following questions:
 2. Where are the main bottlenecks?
 - Sample based profiler (Flat, Call graph, Call path)
 - Event Profiler with runtime summarization



Single Processor Optimization (3)

- Answer the following questions:
 3. What are the causes of performance problems?
 - Compiler feedback
 - Canal / Loopmarks
 - HW counters based Instrumentation library
 - Run-time summarization
 - Hand or automatic instrumentation



Parallel Processing Optimization

- Answer the following questions:
 1. Does my program have communication or synchronization problems?
 - Communication Profiler



Parallel Processing Optimization (2)

- Answer the following questions:
 2. What are the causes of parallel performance problems?
 - Task/Thread tracing
 - Tracing library
 - » CrayPat
 - Visualization GUI
 - Cray Apprentice²



Approach

- Take advantage of existing tools infrastructure
 - Improve usability
 - Develop components for missing functionality
 - Port to all Cray platforms
- Close interaction and collaboration with tools builders in Universities, Labs, and Research Centers
- Close interaction with application developers for quick feedback, targeting functionality enhancements



Collaborations

- Interactions with lead researchers in the performance analysis and optimization field
 - Jack Dongarra's group at UTK
 - PAPI port to all Cray systems
 - Bernd Mohr at Forschungszentrum Jülich
 - KOJAK expert system ported to X1
 - Investigation of CAF and UPC measurement and analysis
 - Automatic performance analysis
 - Jeff Vetter's group at ORNL
 - MPI Profiler (mpiP) port to all Cray platforms
 - Data management techniques in performance analysis
 - Development of filters
 - Hierarchical analysis



Cray Performance Tools

- PAT_HWPC
 - Hardware Counters based measurement for whole application
- CrayPat
 - Performance Analysis Toolkit
- Cray Apprentice²
 - Performance visualization



PAT_HWPC (X1 Series only)

- Collects hardware performance counters information for an application
- No instrumentation required
 - `pat_hwpc [options] a.out`
 - Accepts various hardware counters groups
- Results in report with raw counts and derived metrics for the whole execution
 - Hardware counters summed across threads per process



LIBHWPC (XT3 and XD1)

- Instrumentation library for Fortran, C, and C++
- For each instrumented region provides:
 - Total count & duration (wall clock time & user time)
 - Hardware performance counters information
 - Derived metrics
- Supports:
 - Multiple instrumentation points
 - Nested instrumentation
 - Multiple calls to an instrumented point
- Uses PAPI Version 3
 - PAPI available on base distribution



CrayPat

- Performance Analysis Toolkit
 - pat_build
 - Utility that automatically instruments the application
 - No source code modification required
 - pat runtime library
 - Transparent to the user
 - Collects performance data during execution
 - Writes performance file
 - pat_report
 - Performance analysis utility
 - pat_run
 - pat_help



pat_build

- `pat_build <options> a.out instrumented.a.out`
- Automatic instrumentation
 - User functions
 - MPI, SHMEM, OpenMP, UPC, CAF, pThreads
 - System calls
 - Dynamic heap
 - Raw I/O, buffered I/O, flexible file I/O
- API available for fine grain instrumentation
- “`instrumented.a.out`” is executed the same way as the original application
- No recompilation needed!



CrayPat Functionality: Instrumentation

- Capture of:
 - Software state
 - Thread
 - Call Stack
 - parameter values passed into a function
 - Hardware state
 - Program counter (PC)
 - Hardware performance counters (HWPC)
 - Time stamps
 - Real time clock
 - HWPC Cycle counter



Performance Collection

- Trigger mechanism:
 - Activated by external agent (sampling)
 - Timer
 - Hardware counters overflow
 - (not supported on the XT3 yet)
 - Activated internally (event tracing or synchronous)
 - Code inserted through instrumentation
- Data recording:
 - Summarized during runtime
 - Stored in the form of a profile
 - Trace records for each event
 - Stored in the form of a trace file



pat_report: Flat Profile

Time%	Cum.Time%	Time	Traces	Function
7.4%	93.8%	3.737219	24000	MPI_Waitall
1.4%	96.7%	0.707859	56383	MPI_Isend
0.9%	98.9%	0.455456	56383	MPI_Irecv
0.1%	99.8%	0.051986	56383	mpi_
0.1%	99.9%	0.040626	4	InitalCalc1Calc2#88
0.1%	100.0%	0.026122	24000	mpi_waitall_
0.0%	100.0%	0.009164	4800	Calc1Calc2#10
0.0%	100.0%	0.007347	4800	Calc2#20
0.0%	100.0%	0.006329	4792	Calc3#30
0.0%	100.0%	0.001980	60	MPI_Reduce
0.0%	100.0%	0.000083	60	mpi_reduce_
0.0%	100.0%	0.000030	60	MPI_Type_get_true_extent
0.0%	100.0%	0.000016	4	main
0.0%	100.0%	0.000011	4	MPI_Finalize
0.0%	100.0%	0.000005	4	mpi_comm_rank_
0.0%	100.0%	0.000005	4	mpi_comm_size_
0.0%	100.0%	0.000003	4	MPI_Comm_size
0.0%	100.0%	0.000003	4	_Exit
0.0%	100.0%	0.000003	4	MPI_Comm_rank
0.0%	100.0%	0.000003	4	MPI_Init



pat_report Call Tree Profile

Time%	Cum.Time%	Time	Traces	Calltree
100.0%	100.0%	50.313549	316932	Total

98.5%	98.5%	49.541495	316864	main

98.5%	98.5%	49.541465	316852	MAIN_

39.0%	39.0%	19.645725	115200	calc2_

36.4%	36.4%	18.336529	4800	pat_region_begin_
				Do 200#21
1.1%	37.6%	0.559553	19200	mpi_waitall_

1.1%	37.5%	0.550431	9600	MPI_Waitall
0.0%	37.6%	0.009122	9600	mpi_waitall_(exclusive)
=====				
0.6%	38.2%	0.300126	43200	mpi_isend_

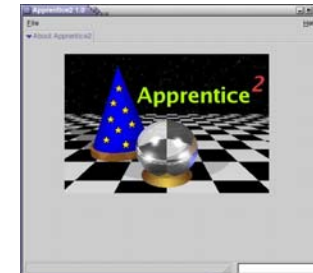
0.6%	38.1%	0.280955	21600	MPI_Isend
0.0%	38.2%	0.019171	21600	mpi_isend_(exclusive)
=====				
0.5%	38.6%	0.246678	4800	calc2_(exclusive)
0.4%	39.0%	0.202839	43200	mpi_irecv_

0.4%	39.0%	0.180679	21600	MPI_Irecv
0.0%	39.0%	0.022161	21600	mpi_irecv_(exclusive)
=====				
29.5%	68.5%	14.844131	110400	calc1_

23.7%	62.8%	11.933617	4800	pat_region_begin_
				Do 100#11
4.5%	67.3%	2.282388	19200	mpi_waitall_



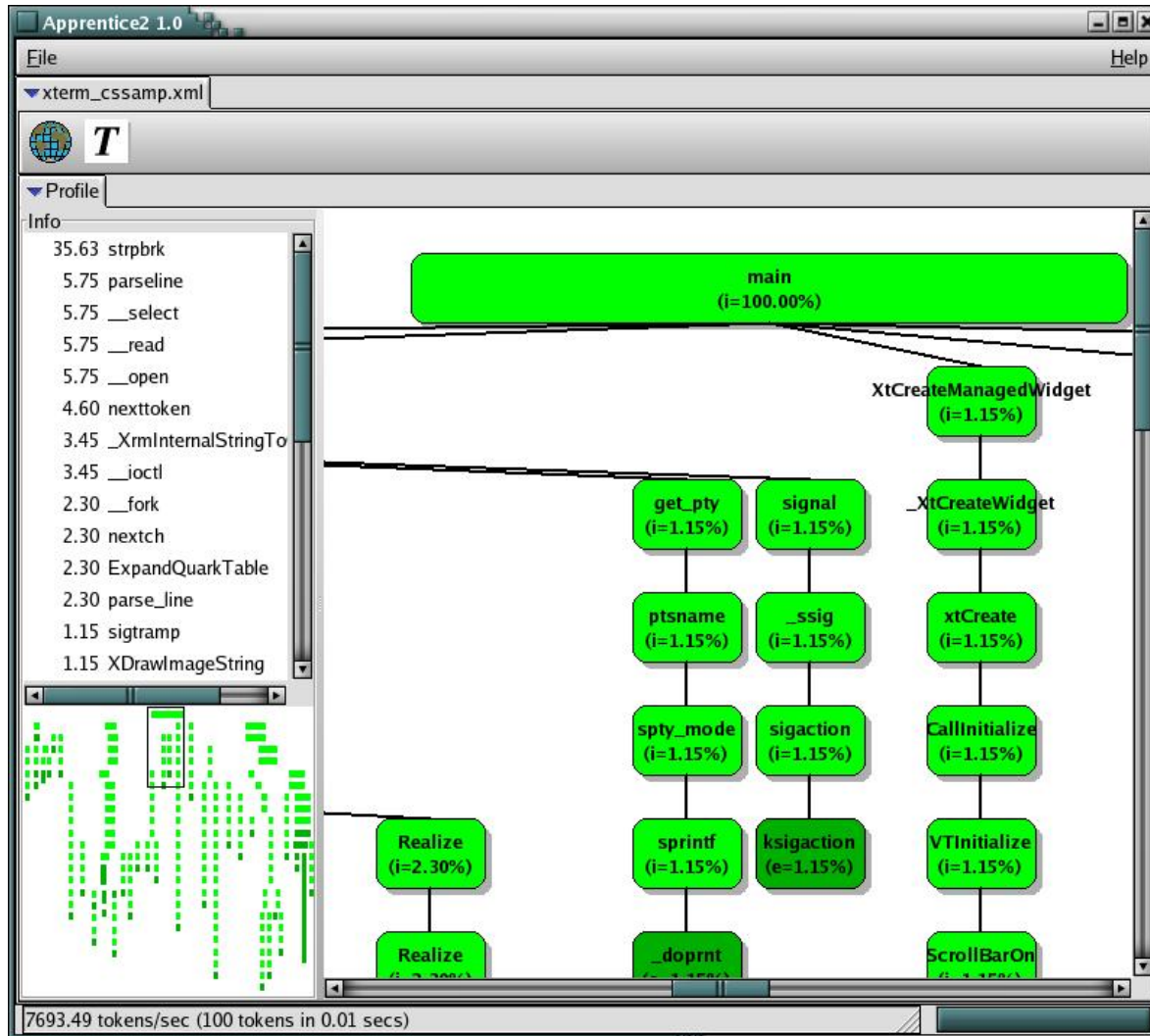
Cray Apprentice²



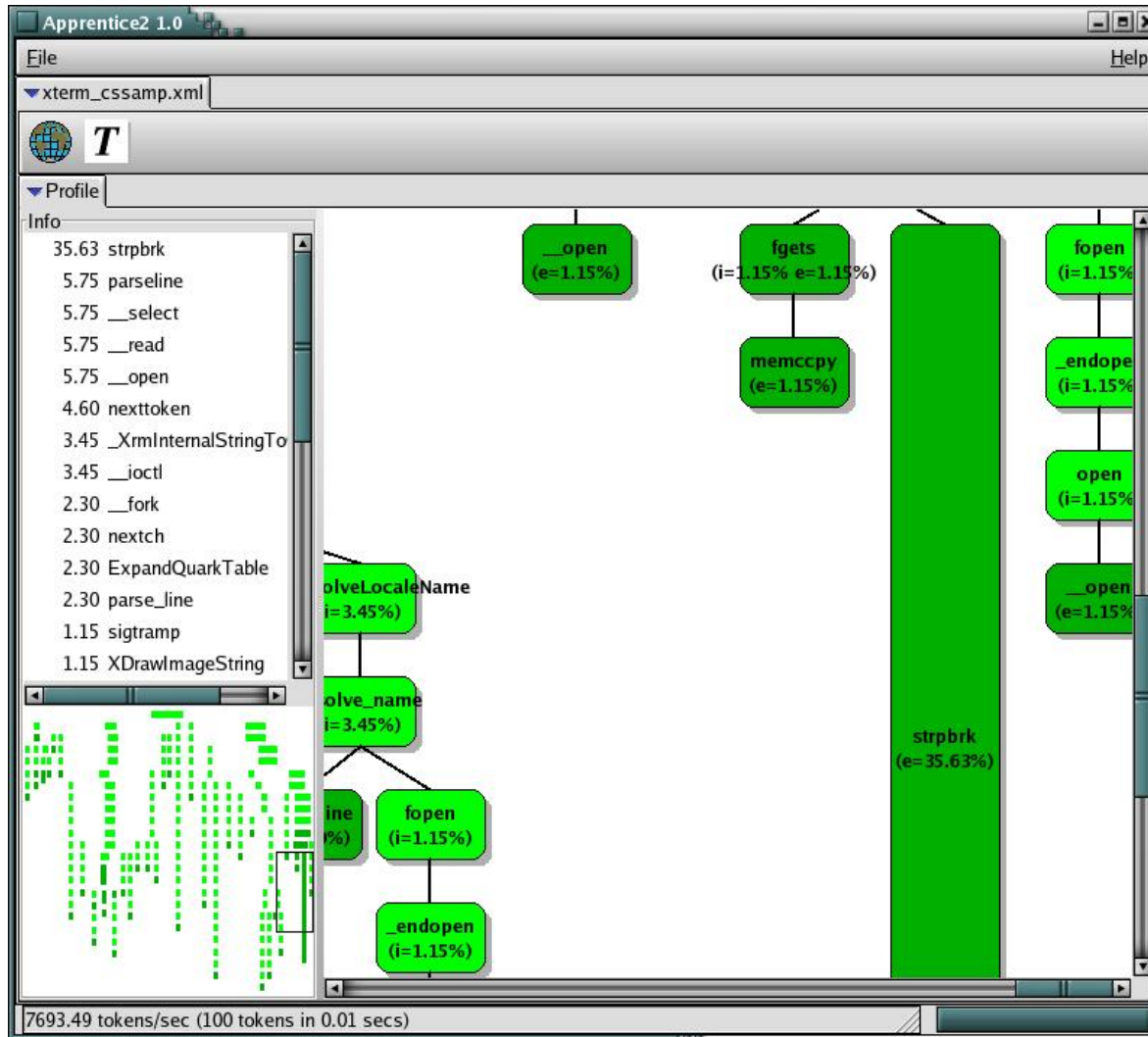
- Call graph profile
- Communication statistics
- Time-line view
 - Communication
 - I/O
- Activity view
- Pair-wise communication statistics
- Text reports
- Source code mapping
- Target to help identify and correct:
 - Load imbalance
 - Excessive serialization
 - Excessive communication
 - Network contention
 - Poor use of the memory hierarchy
 - Poor functional unit use



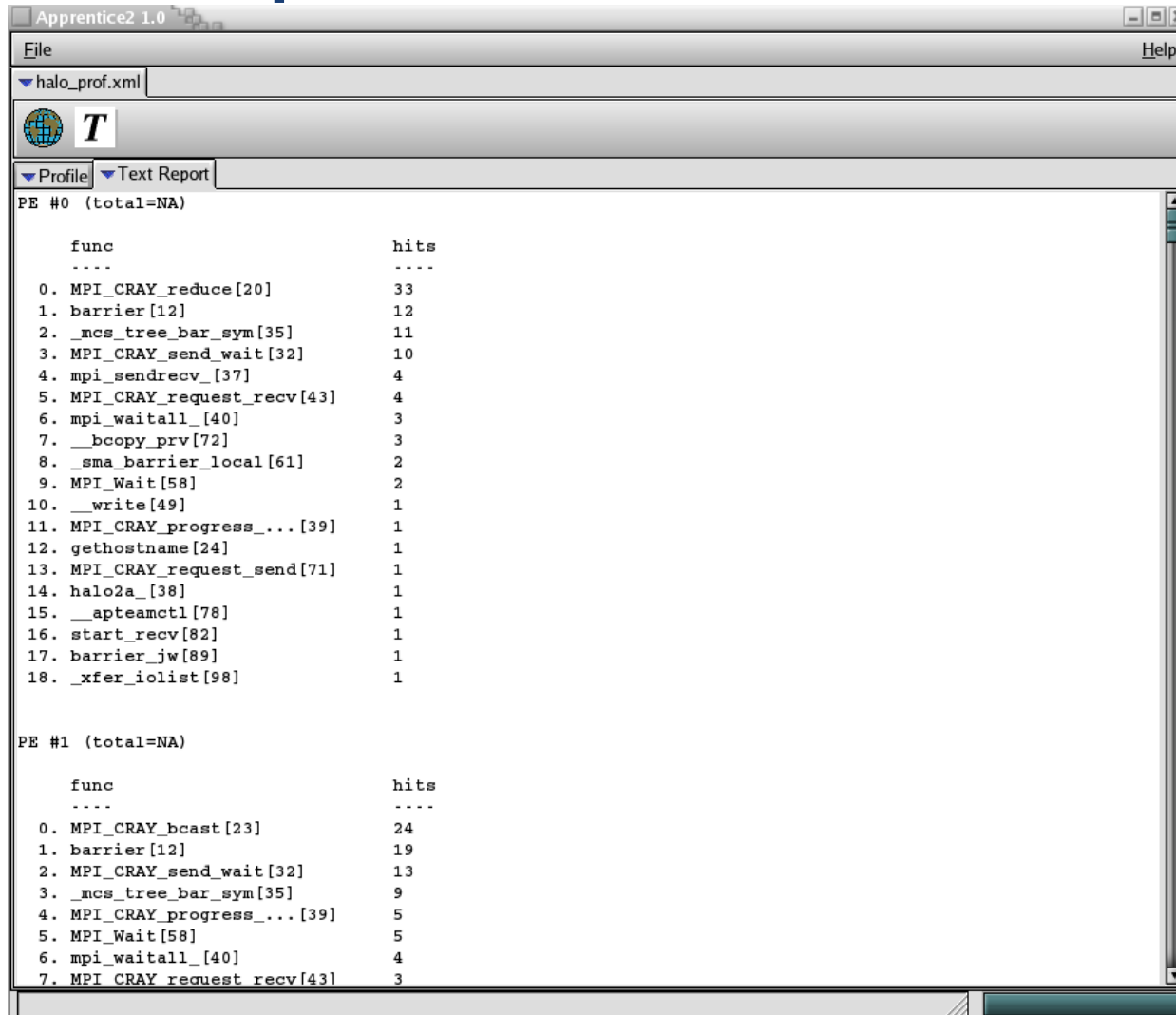
Call Graph Profile View



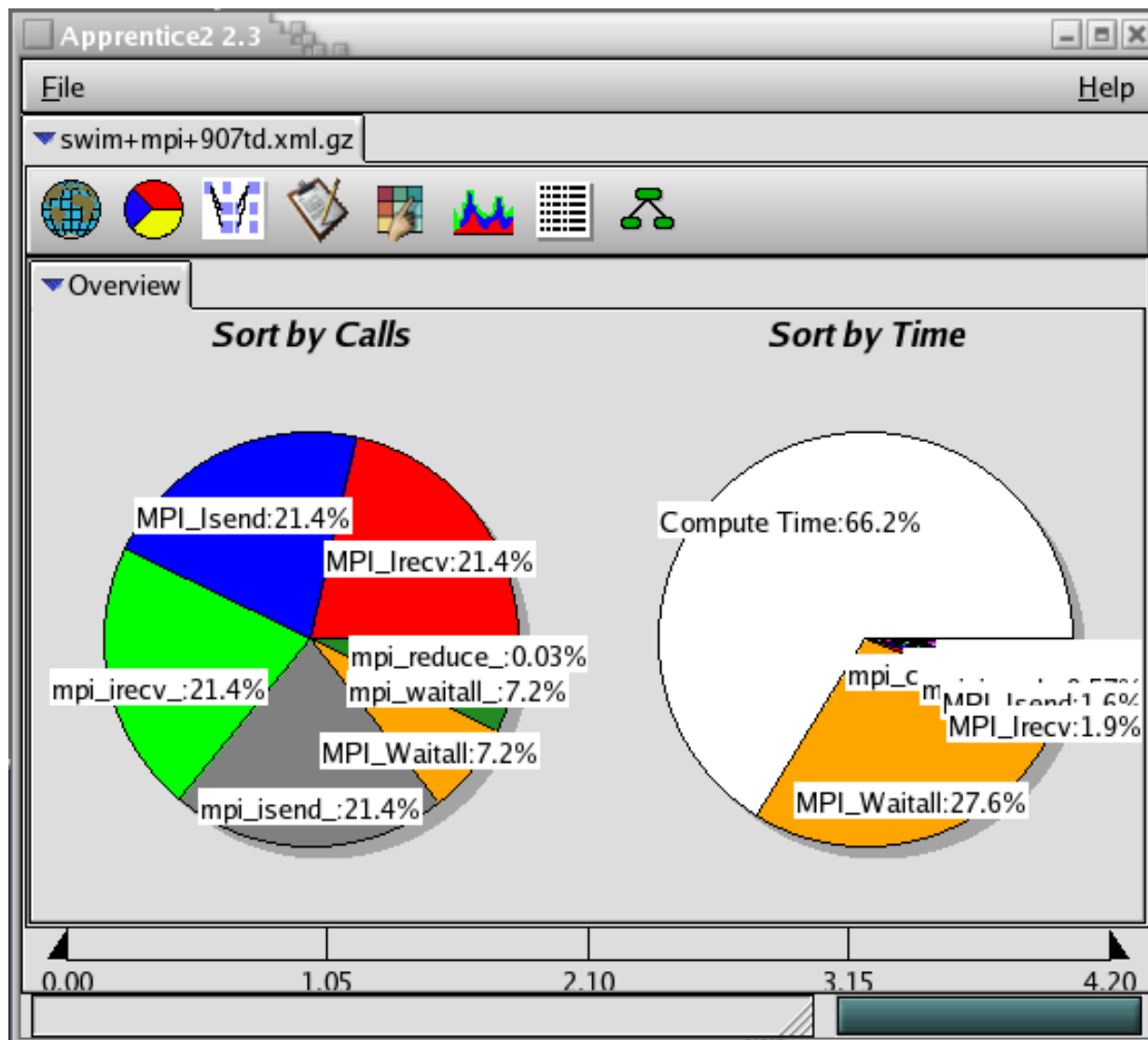
Call Graph Profile View



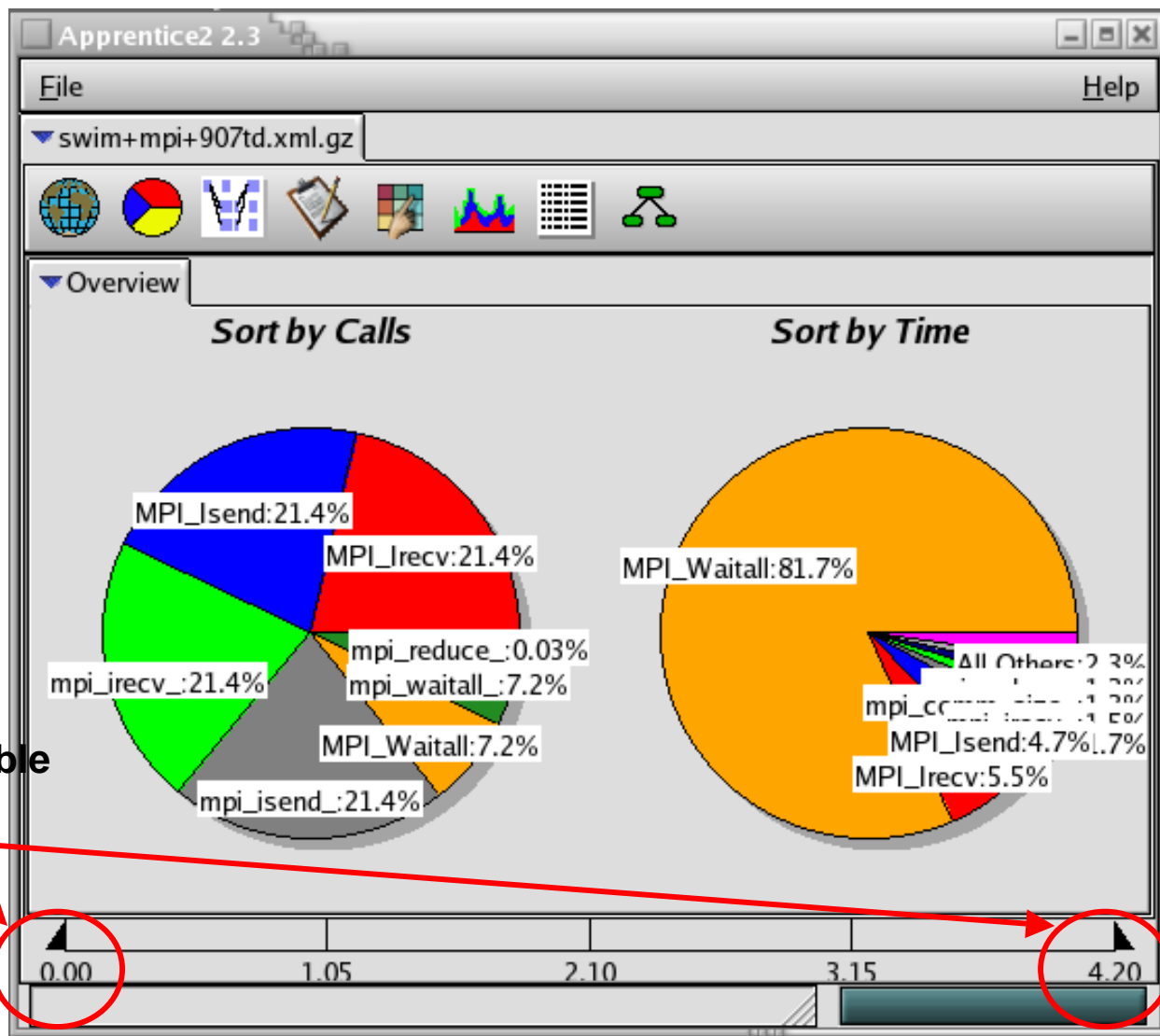
Call Graph Profile View



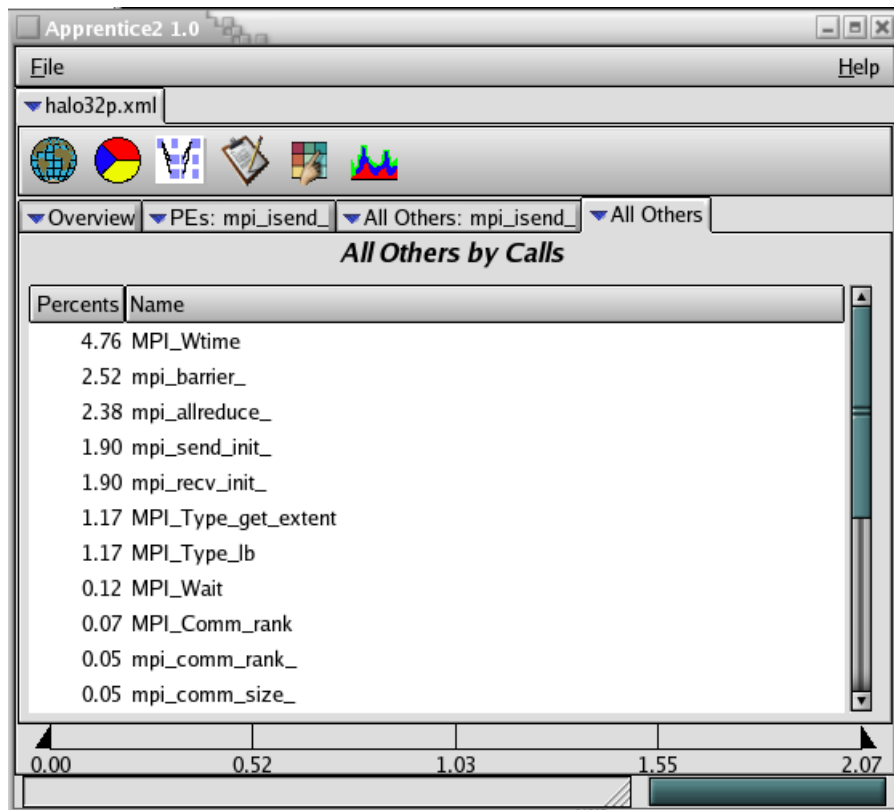
Communication Statistics



Communication Statistics

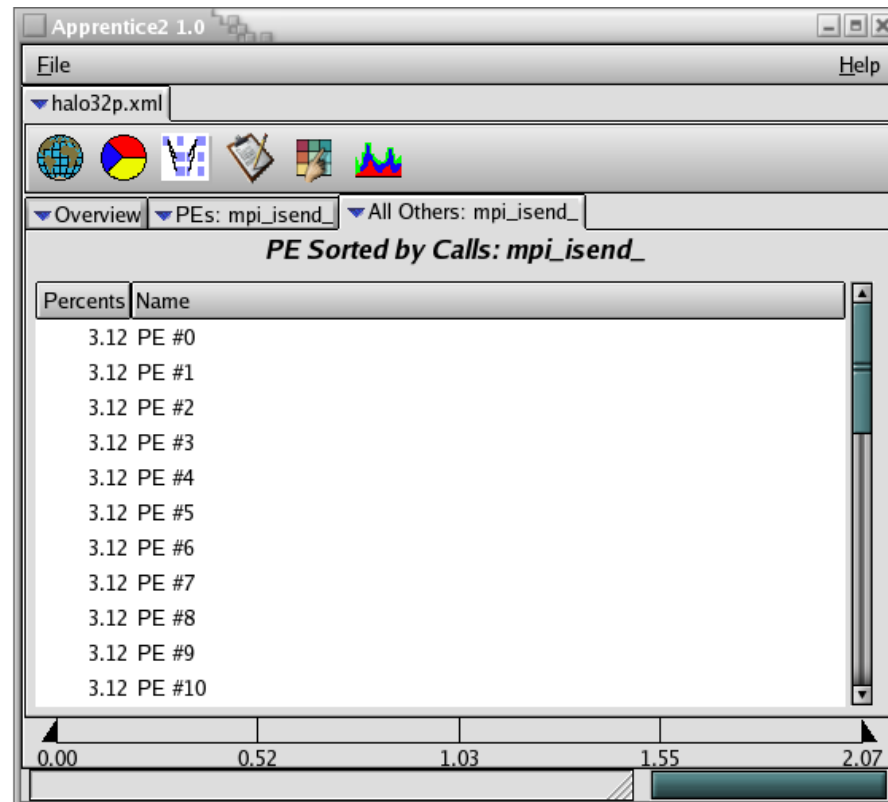


Communication Statistics

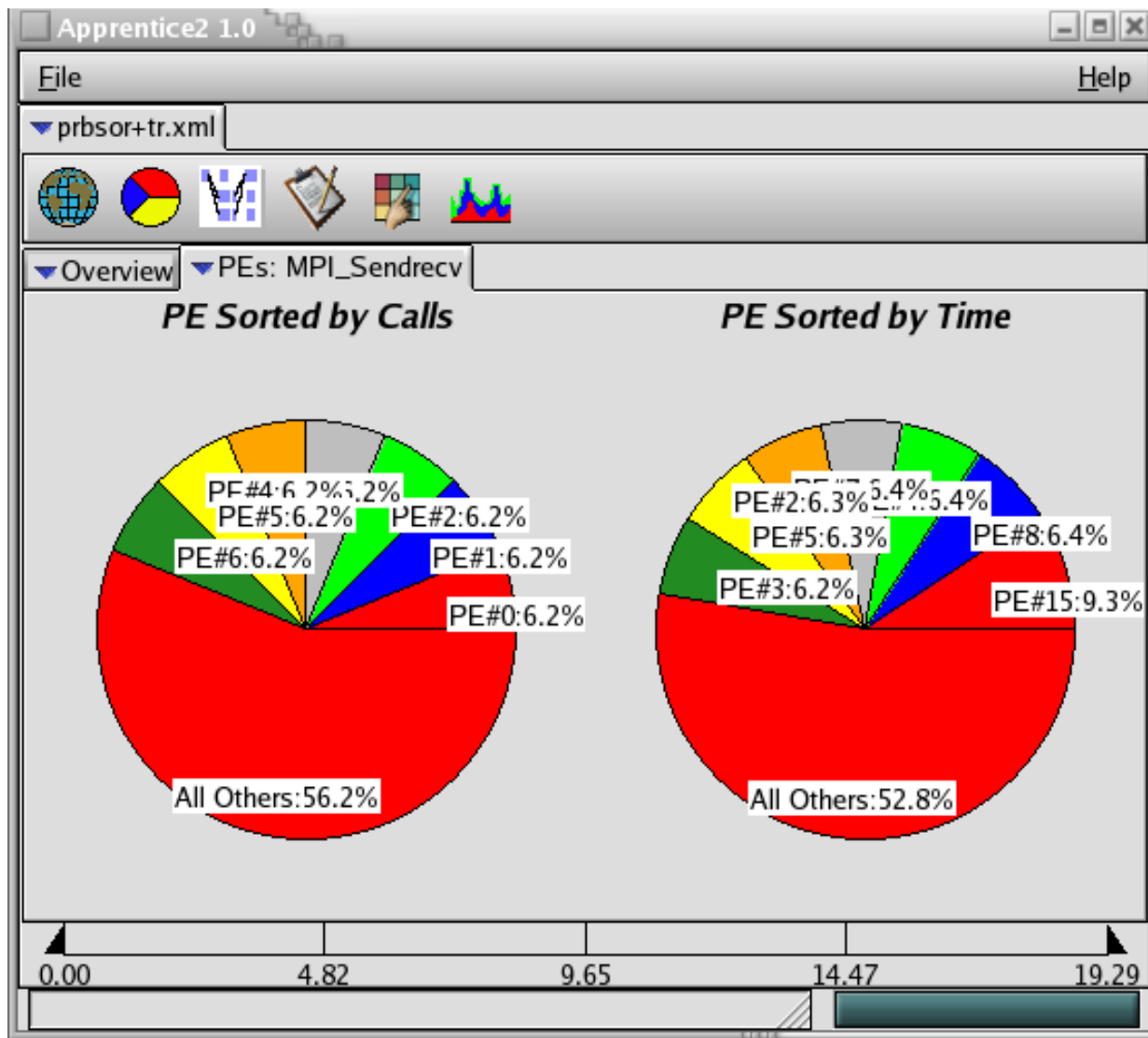


All Others %

PE Distribution (mpi_isend)



Statistics by PE



Statistics by PE

Profile

Apprentice2 1.0

File

prbsor+tr.xml

Overview PEs: MPI_Sendrecv Callers

Hits	Caller	Line	Source
1000	main	290	/scratch/scr4tb/derose/prbsor/prbsor.c
1000	main	286	/scratch/scr4tb/derose/prbsor/prbsor.c
1000	relax	76	/scratch/scr4tb/derose/prbsor/prelax.c
1000	relax	72	/scratch/scr4tb/derose/prbsor/prelax.c

0.00 4.82 9.65 14.47 19.29

Apprentice2 1.0

File Help

prbsor+tr.xml

Overview PEs: MPI_Sendrecv Callers Traffic Report Text Report

```

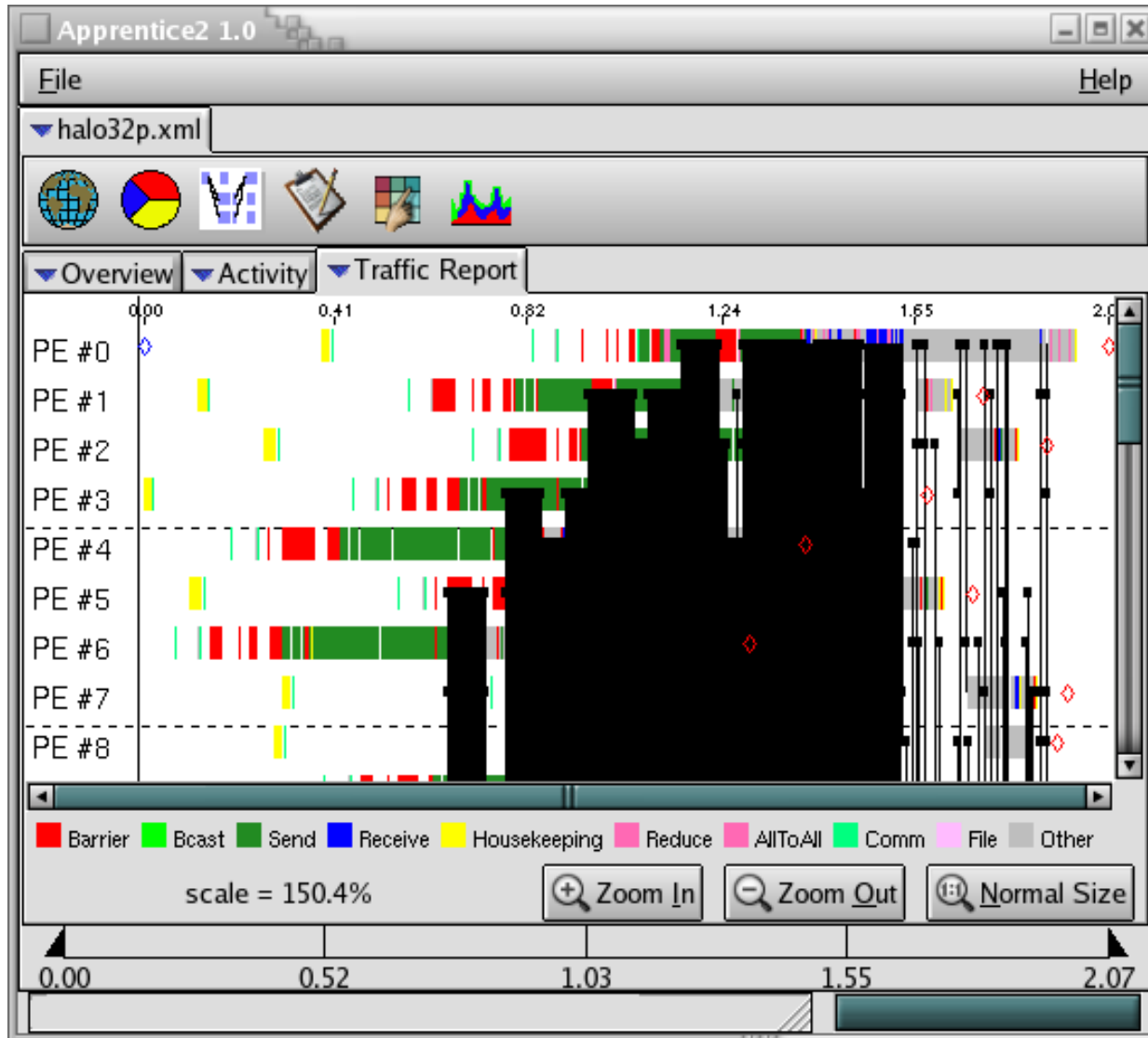
PE #0 (total=NA)
func           time           calls
-----
0. MPI_Sendrecv[28] 2.788269      4000
1. MPI_Reduce[26]  0.721298      1000
2. MPI_Bcast[12]   0.021593        5
3. MPI_Init[23]    0.006056        1
4. ioctl[4]        0.003599        4
5. MPI_Type_get_extent[29] 0.000879        49
6. MPI_Type_lb[30] 0.000383        49
7. MPI_Get_processor_...[21] 0.000134        1
8. MPI_Wait[32]    0.000067        4
9. MPI_Finalize[19] 0.000065        1
10. MPI_Comm_dup[14] 0.000044        1
11. MPI_Comm_rank[16] 0.000040        4
12. MPI_Comm_size[17] 0.000034        3
13. MPI_Comm_free[15] 0.000015        1
    
```

0.00 4.82 9.65 14.47 19.29

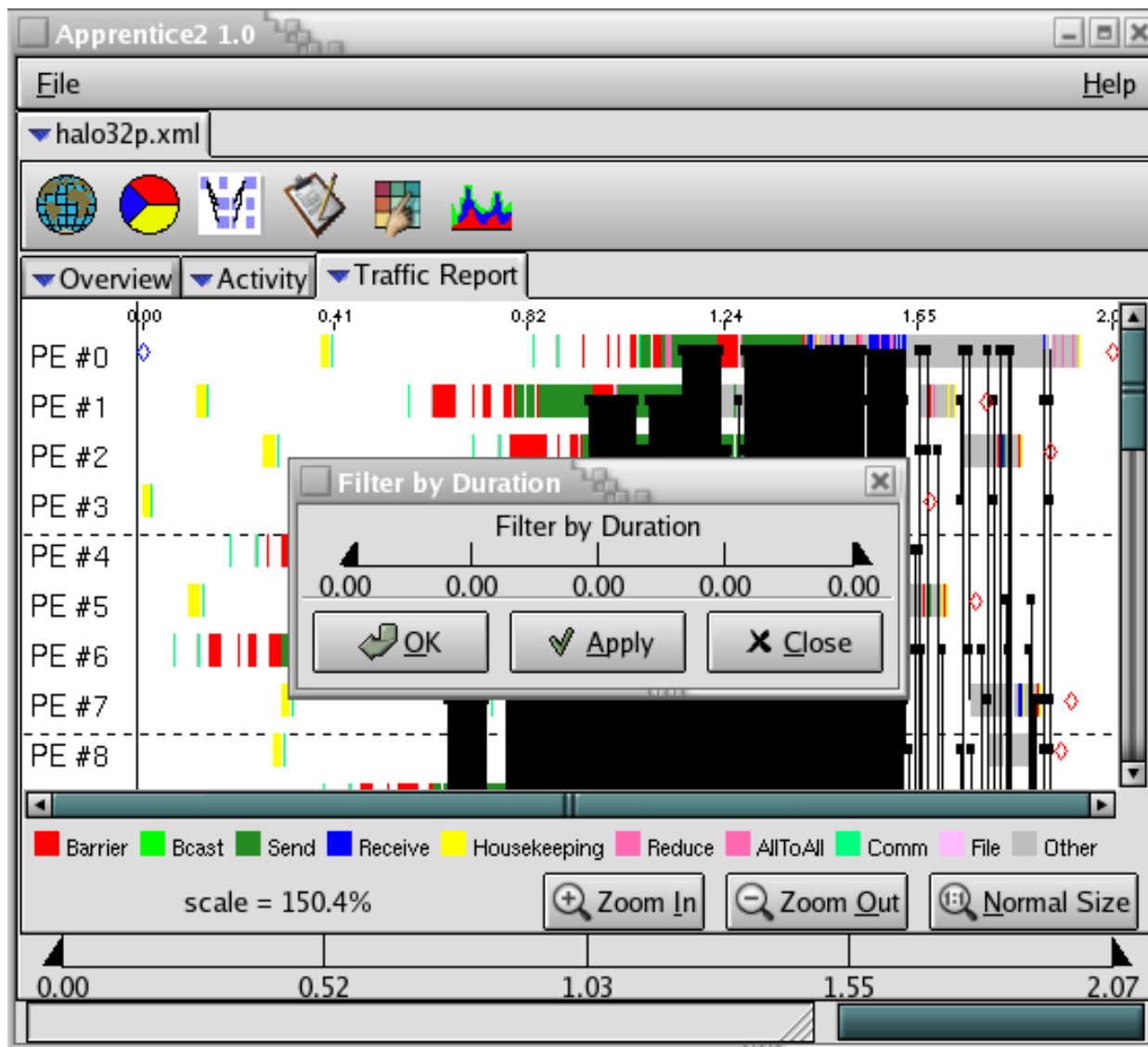
Source code mapping



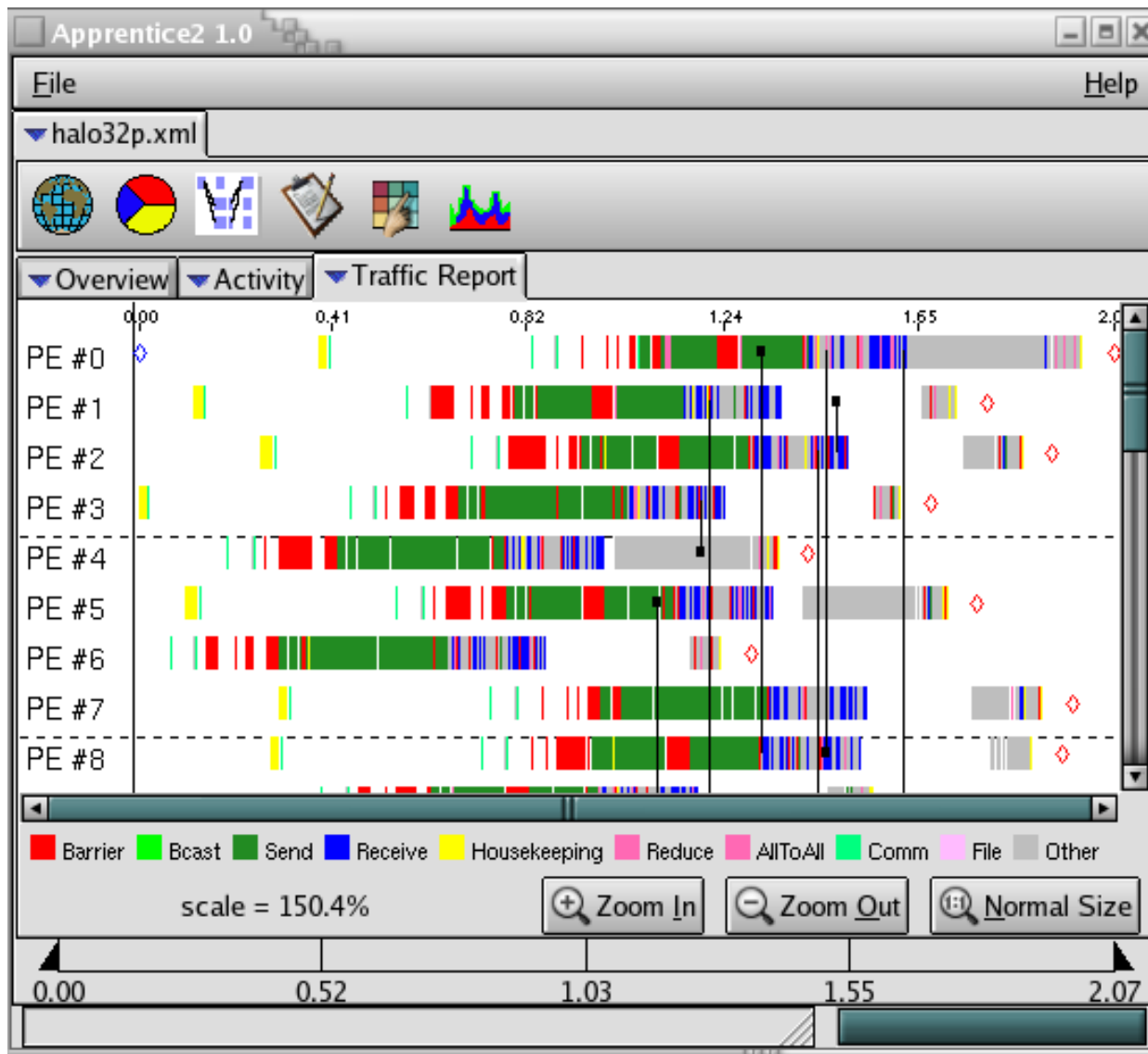
Timeline View



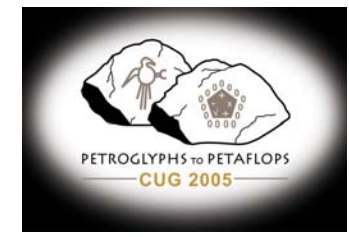
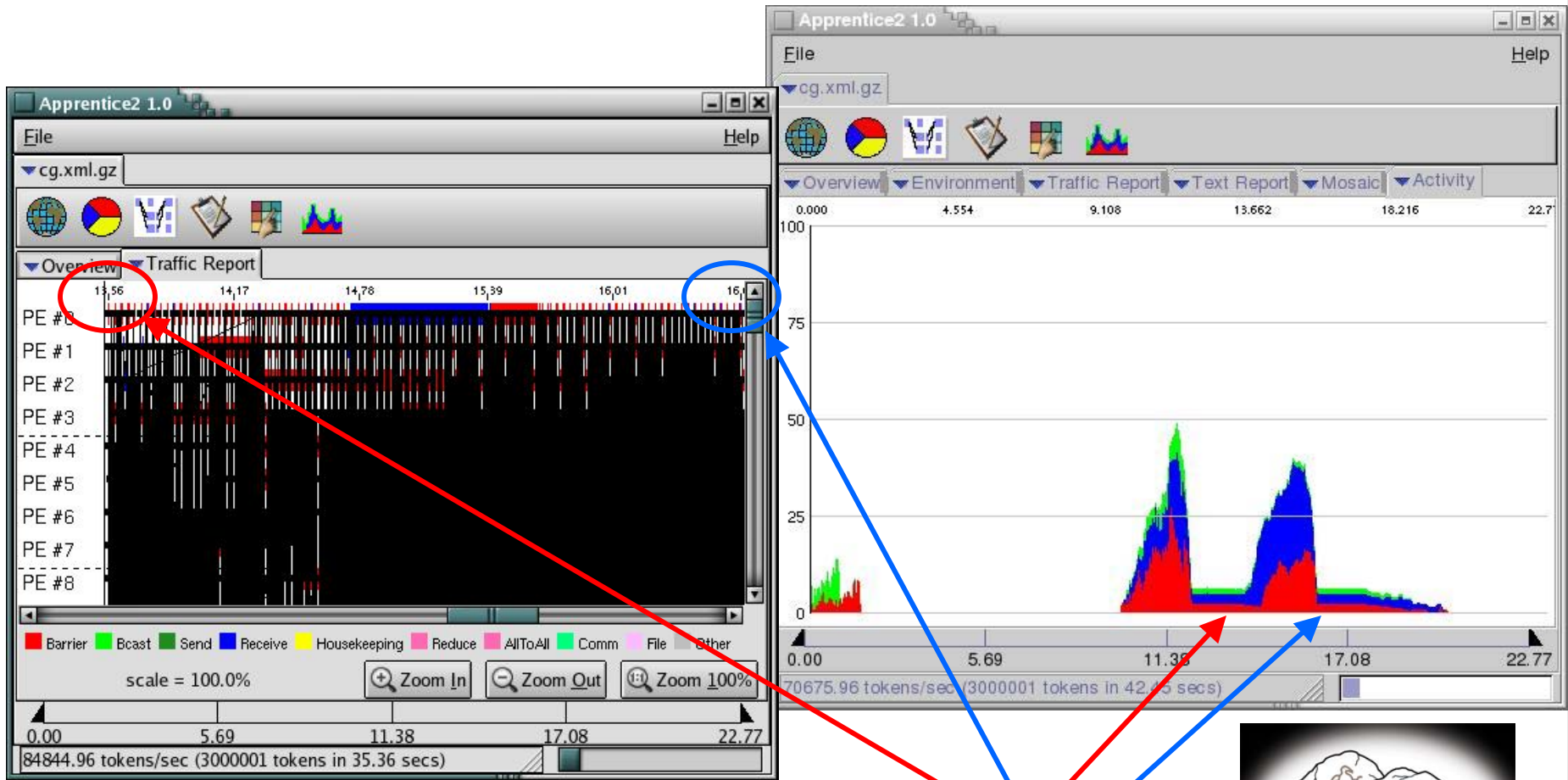
Timeline View



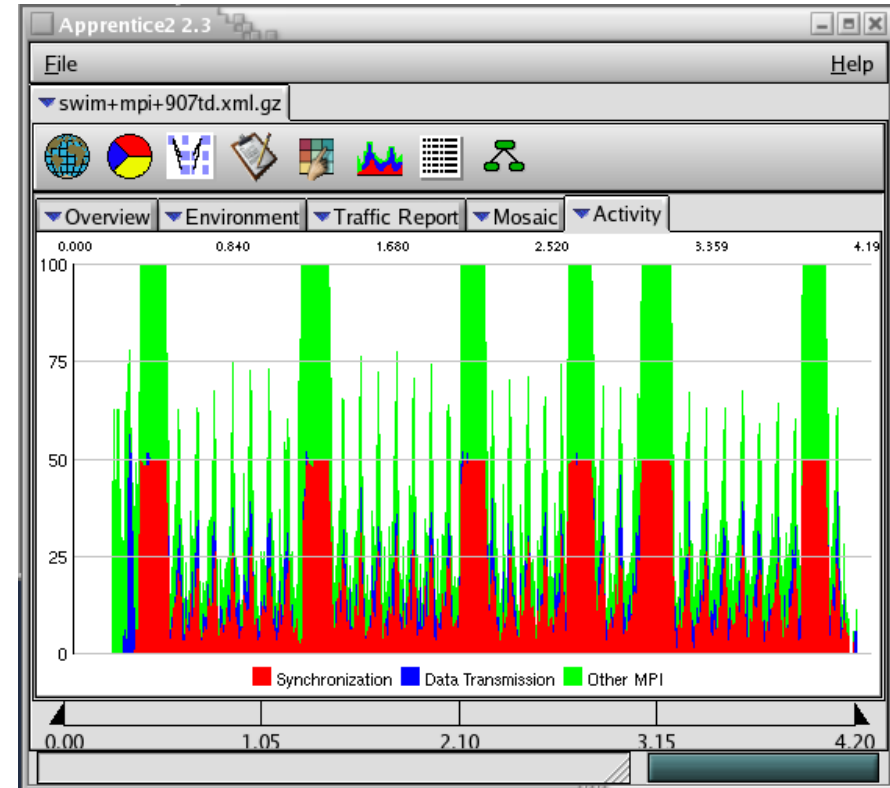
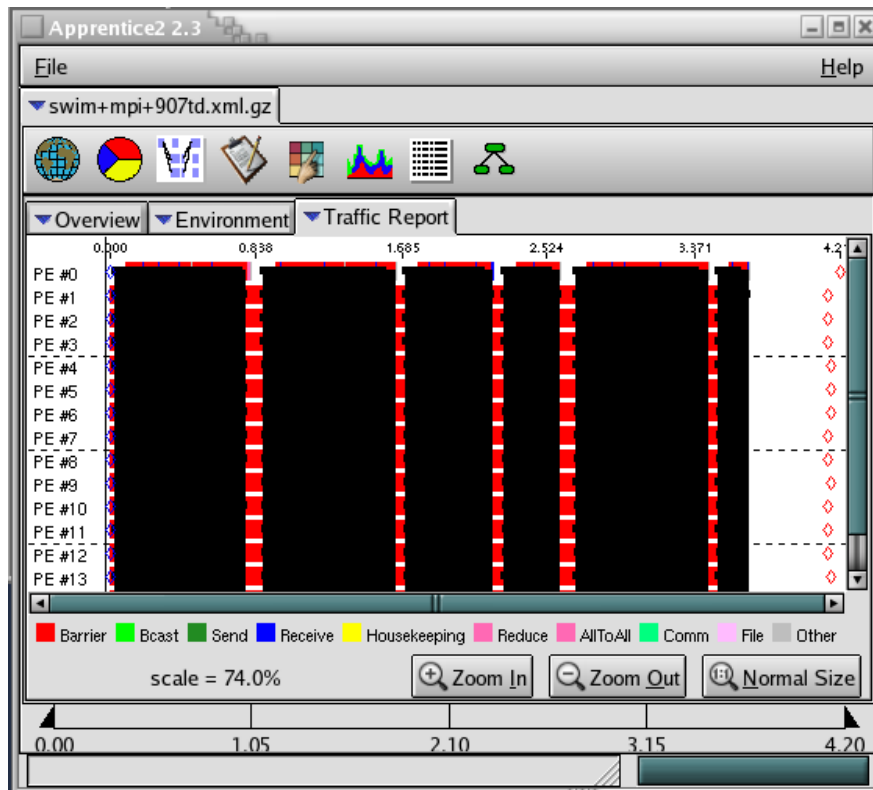
Timeline View



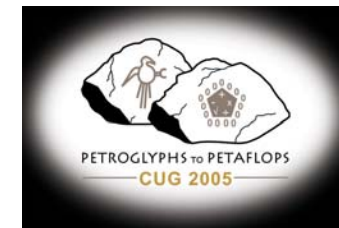
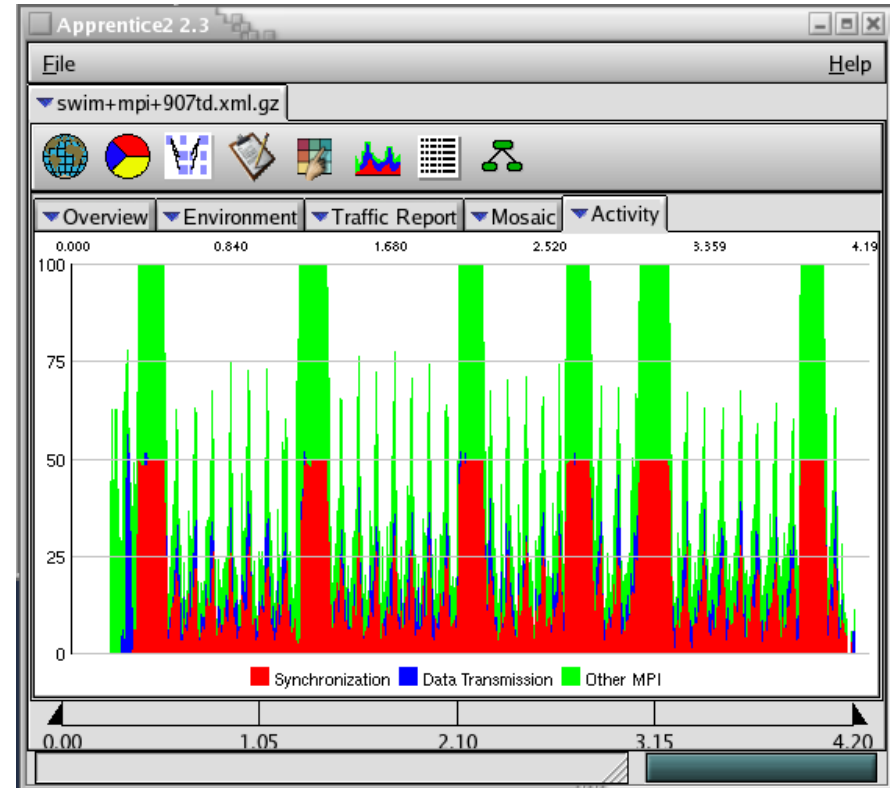
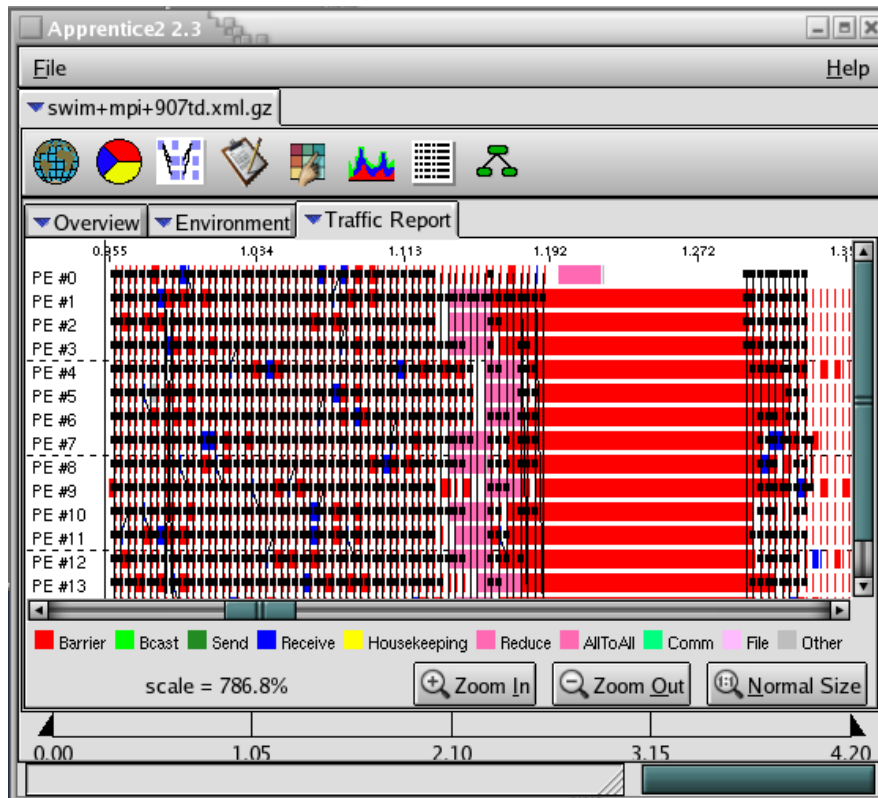
Activity Panel



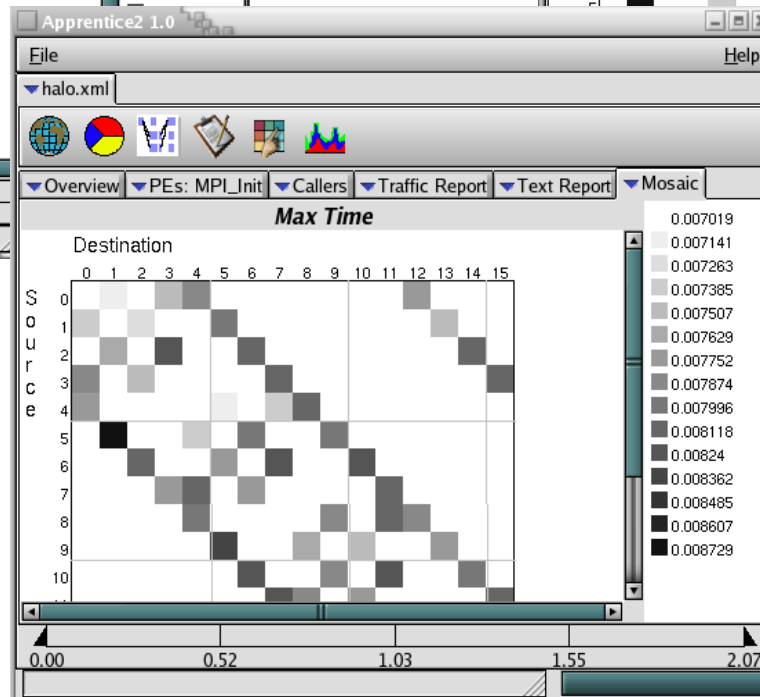
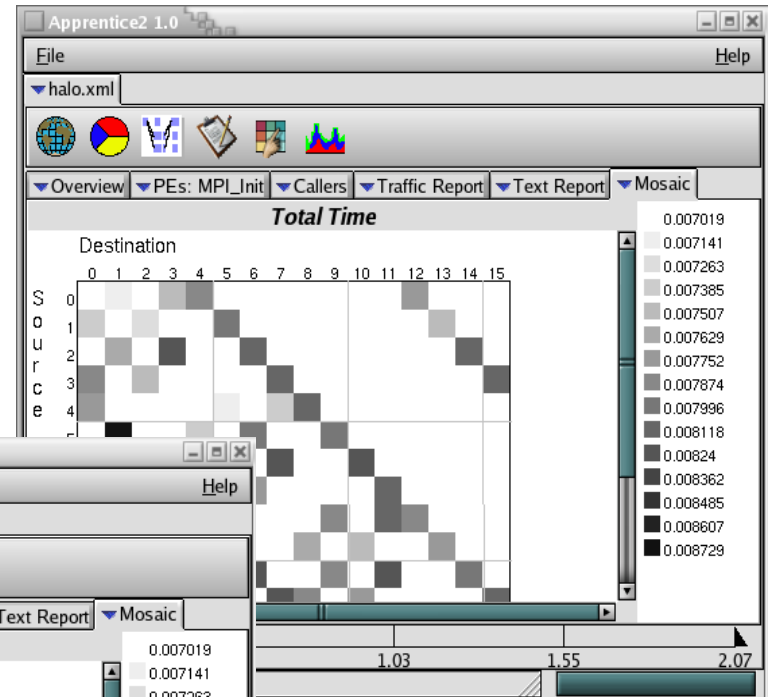
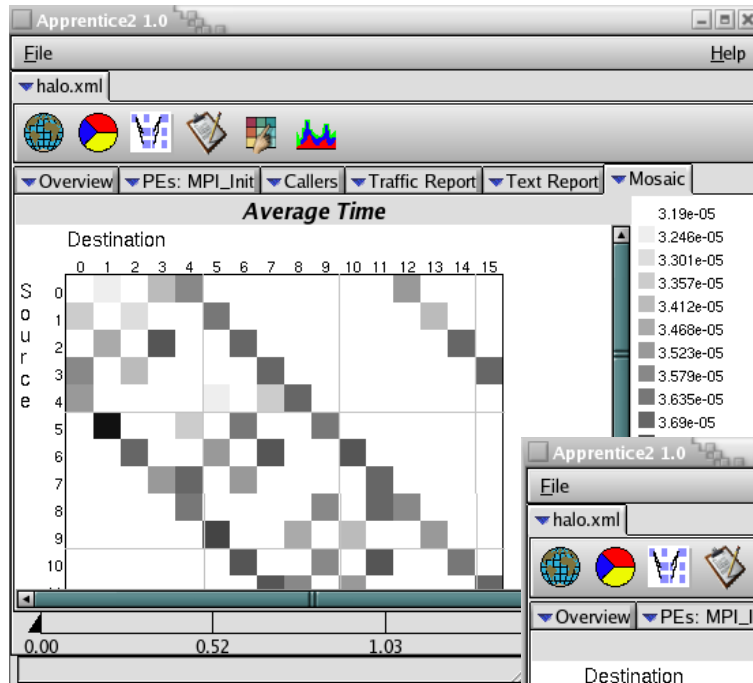
Timeline View



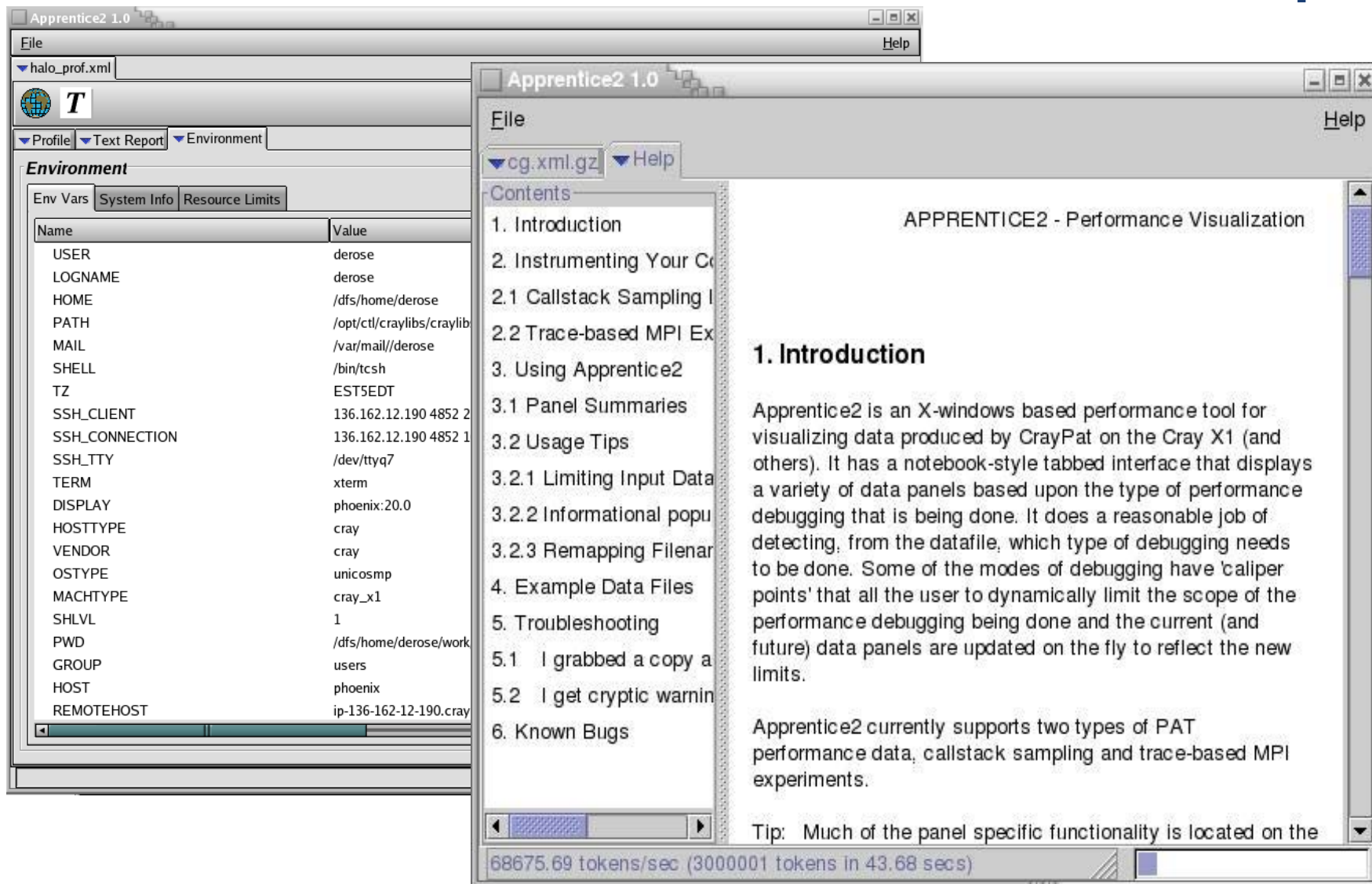
Timeline View



Pair-wise Statistics



Run Information and On-line Help



Apprentice2 1.0

File

halo_prof.xml

Profile Text Report Environment

Environment

Env Vars System Info Resource Limits

Name	Value
USER	derose
LOGNAME	derose
HOME	/dfs/home/derose
PATH	/opt/ctl/craylibs/craylib
MAIL	/var/mail//derose
SHELL	/bin/tcsh
TZ	EST5EDT
SSH_CLIENT	136.162.12.190 4852 2
SSH_CONNECTION	136.162.12.190 4852 1
SSH_TTY	/dev/ttyq7
TERM	xterm
DISPLAY	phoenix:20.0
HOSTTYPE	cray
VENDOR	cray
OSTYPE	unicosmp
MACHTYPE	cray_x1
SHLVL	1
PWD	/dfs/home/derose/work
GROUP	users
HOST	phoenix
REMOTEHOST	ip-136-162-12-190.cray

Apprentice2 1.0

File Help

cg.xml.gz Help

Contents

- 1. Introduction
- 2. Instrumenting Your Co
 - 2.1 Callstack Sampling I
 - 2.2 Trace-based MPI Ex
- 3. Using Apprentice2
 - 3.1 Panel Summaries
 - 3.2 Usage Tips
 - 3.2.1 Limiting Input Data
 - 3.2.2 Informational popu
 - 3.2.3 Remapping Filenar
- 4. Example Data Files
- 5. Troubleshooting
 - 5.1 I grabbed a copy a
 - 5.2 I get cryptic warnin
- 6. Known Bugs

1. Introduction

Apprentice2 is an X-windows based performance tool for visualizing data produced by CrayPat on the Cray X1 (and others). It has a notebook-style tabbed interface that displays a variety of data panels based upon the type of performance debugging that is being done. It does a reasonable job of detecting, from the datafile, which type of debugging needs to be done. Some of the modes of debugging have 'caliper points' that all the user to dynamically limit the scope of the performance debugging being done and the current (and future) data panels are updated on the fly to reflect the new limits.

Apprentice2 currently supports two types of PAT performance data, callstack sampling and trace-based MPI experiments.

Tip: Much of the panel specific functionality is located on the

68675.69 tokens/sec (3000001 tokens in 43.68 secs)

New Features Summary

- Improved usability of Cray PAT on the X1 Series
 - Cray Apprentice²
 - pat_run
 - Linux Cross Compiler environment support
 - Improved performance of trace generation
 - Runtime summarization
 - Better OpenMP support
- Cray XT3 and XD1
 - Port of Cray PAT to both platforms
 - Beta available by the end of the year
 - Cray Apprentice²



Conclusions

- Previous status
 - Good infrastructure
 - Lacking user interface
 - Lacking support on all platforms
- Current status
 - State of the art tools on all Cray systems
 - Functionality to answer all questions in the performance tuning and optimization cycle
- Near future
 - Best performance analysis tools in the industry





The Supercomputer Company

The New Generation of Cray Performance Tools

Questions? / Comments!

Thank You!

Luiz DeRose



ldr@cray.com

