

From Open | SpeedShop to a Component Based Tool Framework

***Jim Galarowicz
Don Maghrak
The Krell Institute***

Open | SpeedShop
Performance Tool Status Update

Overall Agenda

- **Open | SpeedShop Project Update**
- **Component Based Tool Framework Project Introduction**

- **Open | SpeedShop Project Overview**
- **Current Features**
- **What has changed since last update**
- **What we are working on now**
- **Current Release and Status**

■ **Comprehensive open source performance analysis framework**

- Combining **Profiling and Tracing**
- Common **workflow** for all experiments
- Flexible **instrumentation** (dynamic and offline)
- **Extensibility** through plugins
- **GUI, CLI**, immediate **command** and **Python API** user interfaces

■ **Partners**

- DOE/NNSA Tri-Labs (LLNL, LANL, SNLs)
- Krell Institute
- Universities of Wisconsin and Maryland
- ORNL

■ **What can Open | SpeedShop do for the user?**

- *Give lightweight overview of where program spends time*
- *Find hot call paths in user program and libraries*
- *Give access to hardware counter event information*
- *Trace calls to POSIX I/O functions, give timing, call paths, and optional info like: bytes read, file names...*
- *Trace calls to MPI functions. give timing, call paths, and optional info like: source, destination ranks,*
- *Help pinpoint numerical problem areas by tracking FPEs*

■ **Maps the performance information back to the source and displays source annotated with the performance information.**

▣ **Platforms supported currently:**

- **Linux Clusters** with x86, IA-64, Opteron, and EM64T
- Ports to Linux: PPC, BlueGene, Cray-XT in progress

▣ Gather performance data on **unmodified application binaries**

- Where no shared library support build statically
 - ▣ Open | SpeedShop provides “**osslink**” script to help re-link our collector code into the application

▣ Concept of an *Experiment*

- *What to measure (metric) and what to analyze (appl.)*
- *Experiment chosen by user*

▣ *Experiment* consists of *Collectors* and *Views*

- *Collectors define specific performance data sources*
 - ▣ *Hardware counters*
 - ▣ *Tracing of certain routines*
- *Views specify data aggregation and presentation*
- *Multiple collectors per experiment possible*

■ PC Sampling (*pcsamp*)

- Record PC in user defined time intervals
- Low overhead overview of time distribution
- Good first step to find hot spots in program

■ User Time (*usertime*)

- PC Sampling + Call stacks for each sample
- Provides inclusive & exclusive timing data
- Find hot call paths in application

■ Hardware Counters (*hwc, hwctime*)

- Sample HWC overflow events
- Access to data like cache and TLB misses

■ I/O Tracing (*io, iot*)

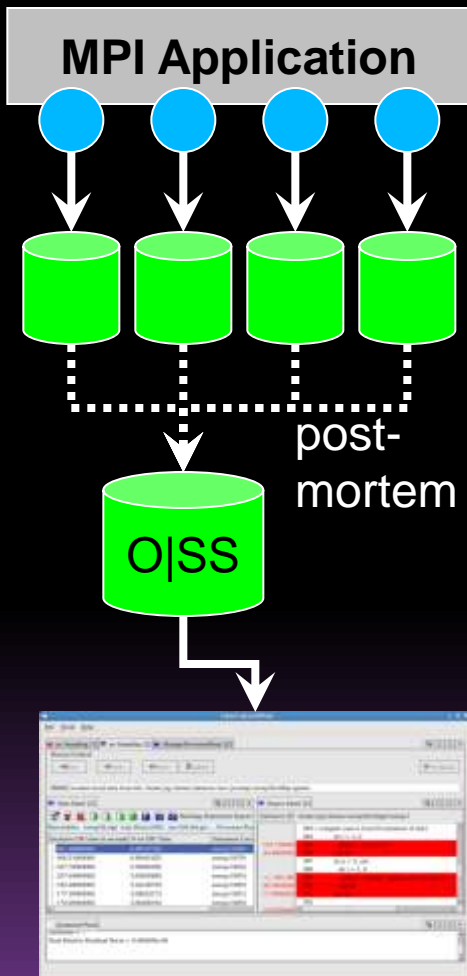
- Record invocation of all POSIX I/O events
- Provides I/O aggregate and individual timings
- *iot* – Shows bytes read/written, etc. & event by event list

■ MPI Tracing (*mpi, mpit, mpiotf*)

- Record invocation of all MPI routines
- Provides MPI aggregate and individual timings
- *mpit* – Shows bytes transferred, ranks involved, etc. & event by event list
- *mpiotf* – Writes open trace format files using vampirtrace under the hood.

- Floating Point Exception Tracing (*fpe*)
 - Triggered by any FPE caused by the code
 - Helps pinpoint numerical problem areas
 - Mapped back to source where FPE occurred

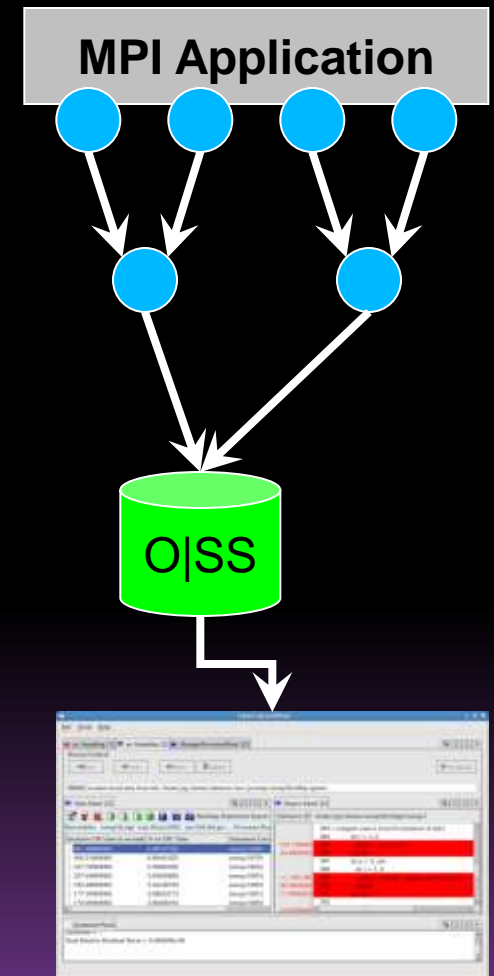
Offline



Easy setup
Low overhead
No additional resources
Higher portability

Online analysis
Intermediate updates
Attach to running code
Optional aggregation

MRNet



- **Moved away from DPCL to MRNet as online transport**
- **Developed the offline mode of operation.**
 - *Using libmonitor (Rice) to hook into application, monitor sys calls*
 - *Link our collectors into the application to gather data*
 - *Write raw data files, then create OSS database file*
 - *Transitioned to having offline the default instrumentation mode*
 - *Low start-up overhead and works well in batch environments*
- **Continued to update the open source components we use**
 - *sqlite, libdwarf, libunwind, libmonitor, Dyninst, MRNet, PAPI,...*
- **Improved installation scripts, tools**

■ **Usability Improvements**

- *Optional View window to select which metrics to be used to create the view*
 - *Ability to quickly switch to function, statement, or library view*
 - *Improvements (tool bar) for custom comparison view*
 - *Integrated offline mode support into GUI wizards*
 - *Created offline convenience scripts to hide the previous syntax*
 - ▣ *ossamps, ossusertime, osshwc, osshwcptime, ossio,*
- ***In general, tool is more robust. Has been exposed to more applications, compilers, job schedulers, MPI versions.***

- **Work on selected modularization of Open | SpeedShop**
 - *Ability to build a viewer only version*
 - *Ability to build only the runtime libraries and collectors*
 - **Refactor runtime library component to be more modular**
- **Porting Open | SpeedShop:**
 - *Linux PPC*
 - *BG/L and BG/P*
 - *CNL: Cray-XT4 and Cray-XT5*
- **Supporting current users and assisting new users**
 - *Release updates*
 - *New features and bug fixes to existing code*

- **Scalability Improvements**
- **Integrate the latest versions of MRNet and Dyninst into Open | SpeedShop (CBTF project)**
 - *Using Dyninst-6.1 and MRNet 2.2 beta for development*
 - *More on this later in the talk*
- **Component Based Tool Framework project**
 - *Subject of next half of this talk*

- **Open | SpeedShop 1.9.3.3 available**
 - *Packages and source from sourceforge.net*
 - *Tested on a variety of platforms*
- **Cray-XT, BG, and PPC versions coming soon**
- **Open | SpeedShop website:**
<http://www.openspeedshop.org/>
- **Download options:**
 - *Package with Install Script (install.sh or install-oss)*
 - *Source for tool and base libraries*

***Building a Community Infrastructure for Scalable
On-Line Performance Analysis Tools***

***Component Based Tool Framework
“CBTF”***

***Jim Galarowicz
Don Maghrak
The Krell Institute***

- **Project Origin and Team**
- **Project Rationale**
- **Project Goals/Objectives**
- **Research Challenges/Project Requirements**
- **Performance Tools Pipeline**
- **Project Results/Outcomes**
- **Current Status**

- ***Project Origin***

- ***OASCR Proposal: "Building a Community Infrastructure for Scalable On-Line Performance Analysis Tools Around Open SpeedShop" for Software Development Tools for Improved Ease-of-Use of Petascale Systems***
- ***Jointly funded by OASCR and NNSA***
- ***Three year project***

- **Project Team**
 - *The Krell Institute*
 - *University of Maryland*
 - *University of Wisconsin*
 - *Oak Ridge National Laboratory*
 - *Lawrence Livermore National Laboratory*
 - *Los Alamos National Laboratory*
 - *Sandia National Laboratories*
 - *Carnegie Mellon University*
 - *Others welcome.....*

- **Why the need for the project?**
 - *Petascale environments need tool sets that are flexible*
 - *Need to quickly create new and specialized tools*
 - *Better availability of tools across more platforms*
 - *Need to avoid creating stove pipe tools*

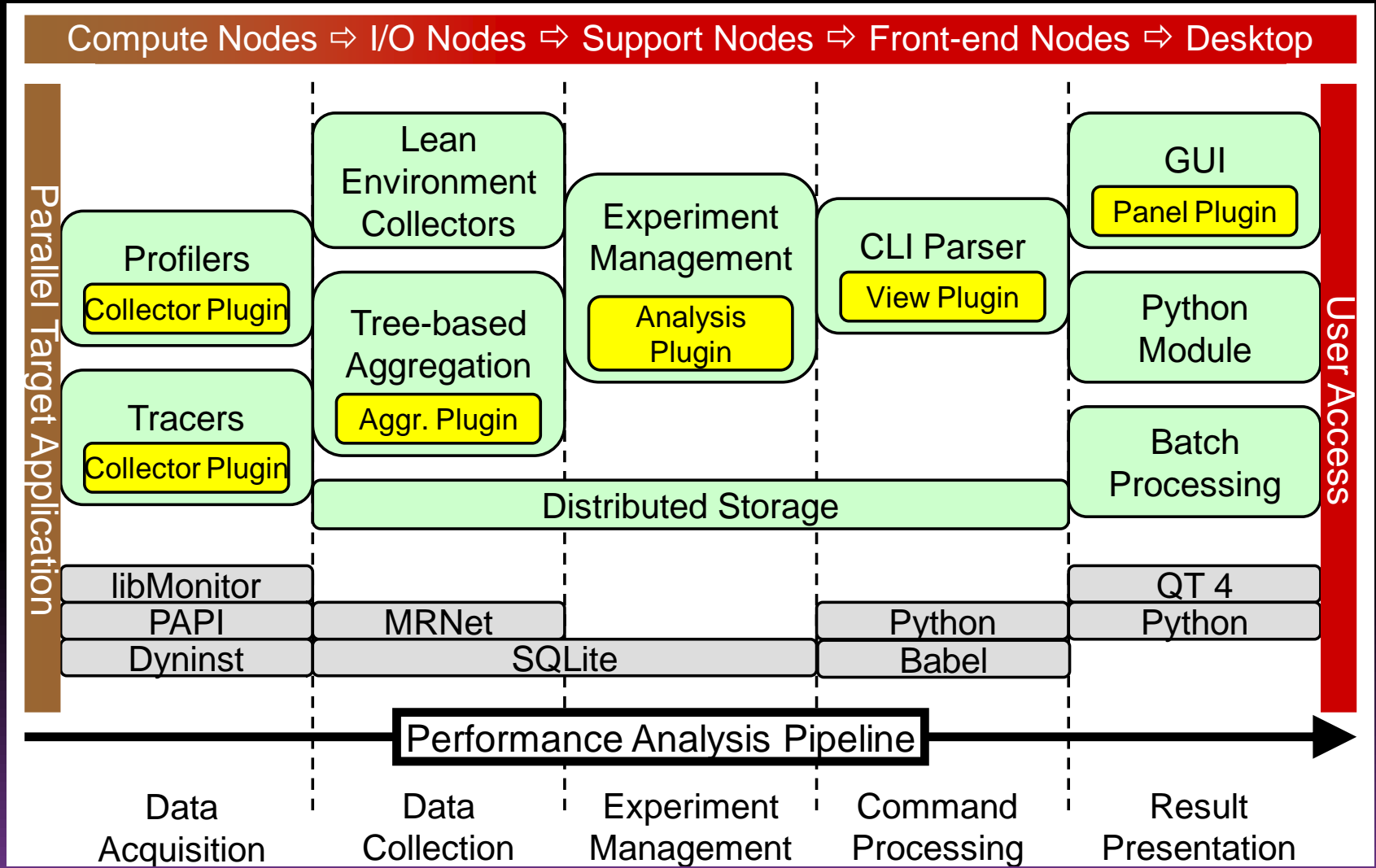
- **Project Goals/Objectives:**
 - *Create a toolbox of components for building high-level end user tools and/or quickly build tool prototypes.*
 - *Tools should be easily configurable/adjustable w/o rebuilding.*
 - *Able to mix components from several groups and/or vendors. Everyone should be able to contribute and use the new components.*
 - *We would like contributors to define the interfaces with us so that we can share components later in both directions.*

- **Project Goals/Objectives:**
 - *Research into efficient and effective online data aggregation, reduction, filtering, and data transfer*
 - *MRNet*
 - *Research lightweight data acquisition techniques*
 - *Binary rewriting*
 - *Assemble new tool components to create a more modular Open | SpeedShop performance tool*
 - *Support BlueGene and Cray-XT platforms*

- **Research Challenges/Project Requirements:**
 - *Components must be designed for scale but also have a need for generality.*
 - *Support specialized tool components intended for serial or small scale usage.*
 - *Infrastructure must support online data aggregation because of potentially high data volume at scale.*
 - *Petascale machines are likely to have limited OS capabilities requiring new and light-weight data acquisition techniques.*
 - *Must be able to efficiently store the performance data.*
 - *Must be able to map any combination of tool components to the target architecture.*

Building a Community Infrastructure for Scalable On-Line Performance Analysis Tools

Performance Analysis Pipeline



- **Creating a first Performance Tools Pipeline prototype**
 - Start with Open | SpeedShop components as one set of examples for such an infrastructure.
 - Decompose core components into general building blocks.
 - Arrange building blocks into a logical performance analysis pipeline.
 - Allows users and tool builders to select individual components for each pipeline stage.
 - Supports a flexible mapping onto the target architecture which provides efficient execution and visualization (incl. remote operation) environments.

- **Project Results/Deliverables:**
 - *Set of reusable components for creating performance tools*
 - *Modified version of gprof using reusable components*
 - *Components for online data aggregation, reduction, filtering, and transfer at high scale*
 - *Tool or Open | SpeedShop experiment based on Active Harmony*
 - *A new, more modular Open | SpeedShop performance tool*
 - *Support for BlueGene and Cray-XT platforms*
 - *Special purpose tool, based on need at ORNL*

Dyninst/MRNet Features/Requirements/Desires

- *Plan to use symtabAPI*
- *Plan to be using the MRNet lightweight library*
- *Plan to use the detach on the fly feature*
- *Plan to use the binary rewriter feature*
- *Would like a floating point register fix up feature*
- *Plan to use the “1st party” stackwalker API*
- *Plan to create an “new” OSS feature based on Active Harmony*
- *Plan to use MRNet transport mechanism*

- **Current Status**

- *Open | SpeedShop team design meetings*
- *Holding extended CBTF team meetings to discuss ideas for component interfaces*
- *Created a CBTF wiki*
- *Started prototyping the component interface design*
- *Doing a number of improvements and decompositions in Open | SpeedShop in preparation to move to CBTF*
- *Plan to focus on transport components first*

Questions?

jeg@krellinst.org

dpm@krellinst.org

Open | SpeedShop
Performance Tool Status Update

Open | SpeedShop Appendix

Open | SpeedShop Appendix

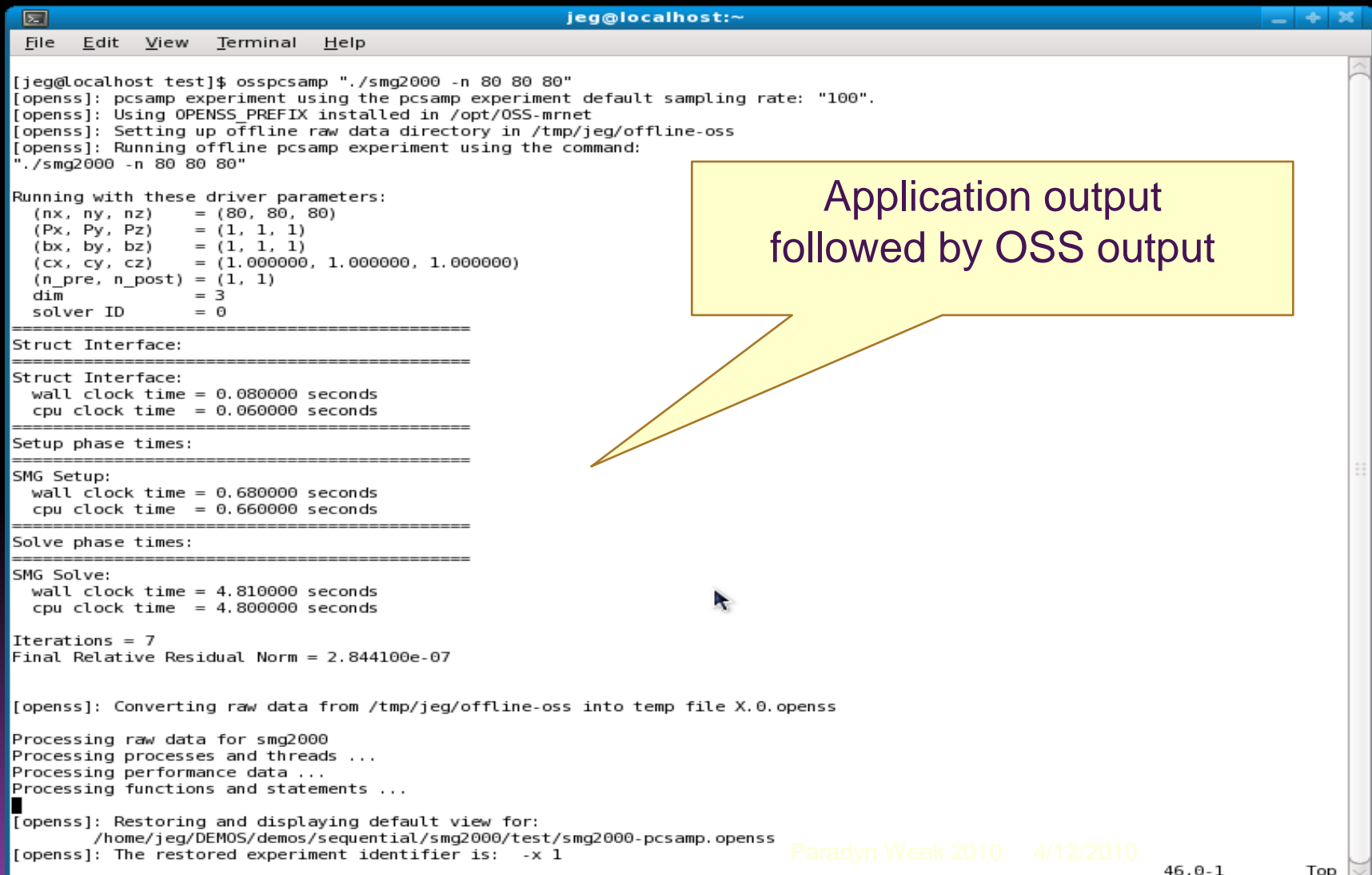
www.openspeedshop.org

- **osspcsamp “<executable> <arguments>”**
 - One line command to gather PC Sampling results
 - Note: “” around executable line
 - Run command without extra arguments for help or view man page
- **Separate command for each experiment**
 - osspcsamp, ossusertime, osshwc, osshwctime
 - ossio, ossiot, ossmpi, ossmpit, ossfpe
- **Example Sequential run: (example run in following slides)**
 - osspcsamp “./smg2000 -n 80 80 80”
- **Example multi-process run:**
 - ossmpi “mpirun -np 64 sweep3d.mpi”

Open | SpeedShop Performance Tool Status Update

Example Offline Run With Output

```
osspsamp "./smg2000 -n 80 80 80"
```



The image shows a terminal window titled 'jeg@localhost:~' with a menu bar (File, Edit, View, Terminal, Help). The terminal displays the command `osspsamp "./smg2000 -n 80 80 80"` and its output. The output includes driver parameters, timing information for various phases (Struct Interface, SMG Setup, SMG Solve), and performance data. A yellow callout box points to the output text, stating 'Application output followed by OSS output'. The terminal output is as follows:

```
[jeg@localhost test]$ osspsamp "./smg2000 -n 80 80 80"
[openss]: pcsamp experiment using the pcsamp experiment default sampling rate: "100".
[openss]: Using OPENSS_PREFIX installed in /opt/OSS-mrnet
[openss]: Setting up offline raw data directory in /tmp/jeg/offline-oss
[openss]: Running offline pcsamp experiment using the command:
"./smg2000 -n 80 80 80"

Running with these driver parameters:
(nx, ny, nz)   = (80, 80, 80)
(Px, Py, Pz)  = (1, 1, 1)
(bx, by, bz)  = (1, 1, 1)
(cx, cy, cz)  = (1.000000, 1.000000, 1.000000)
(n_pre, n_post) = (1, 1)
dim           = 3
solver ID     = 0
=====
Struct Interface:
=====
Struct Interface:
  wall clock time = 0.080000 seconds
  cpu clock time  = 0.060000 seconds
=====
Setup phase times:
=====
SMG Setup:
  wall clock time = 0.680000 seconds
  cpu clock time  = 0.660000 seconds
=====
Solve phase times:
=====
SMG Solve:
  wall clock time = 4.810000 seconds
  cpu clock time  = 4.800000 seconds

Iterations = 7
Final Relative Residual Norm = 2.844100e-07

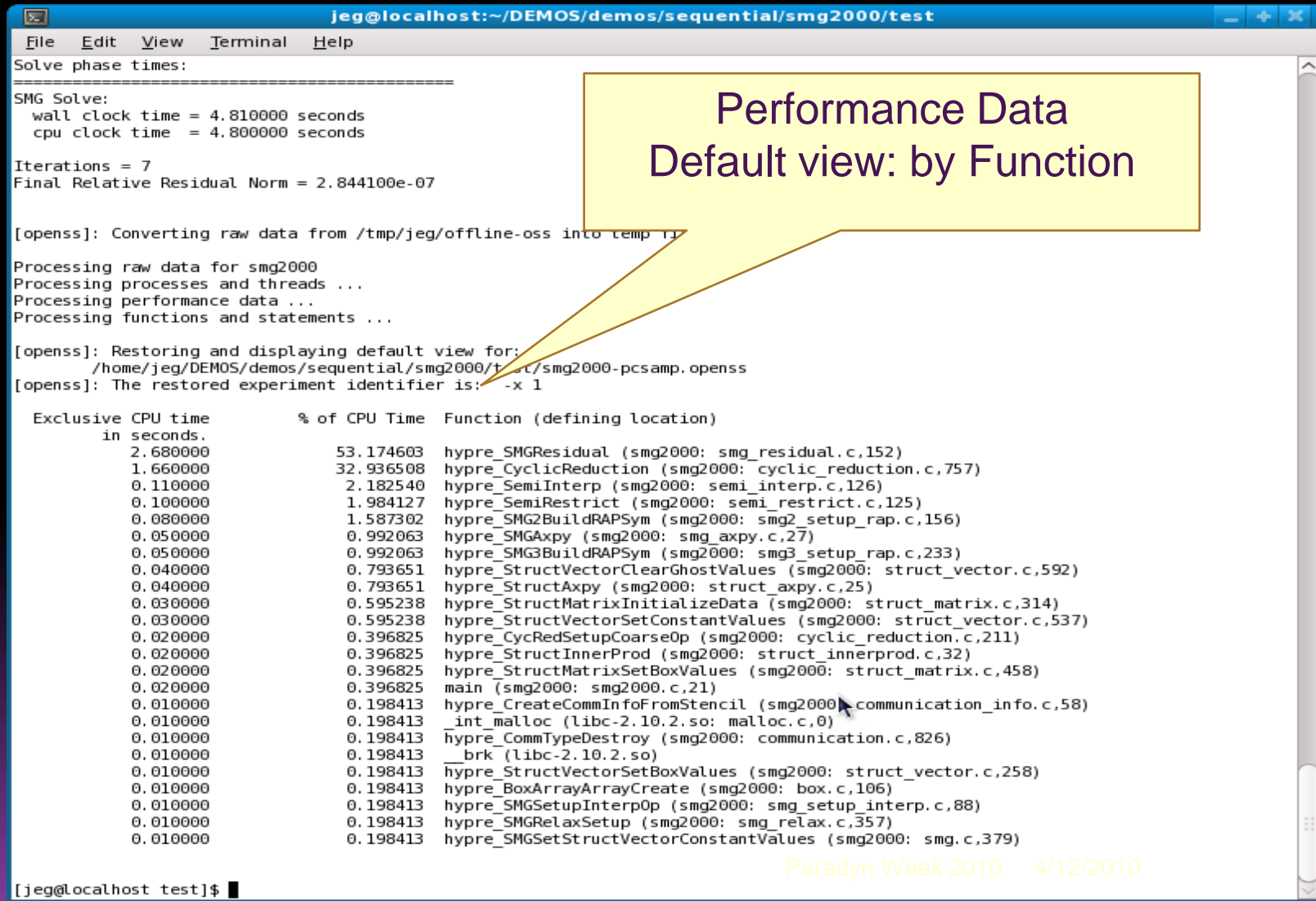
[openss]: Converting raw data from /tmp/jeg/offline-oss into temp file X.0.openss
Processing raw data for smg2000
Processing processes and threads ...
Processing performance data ...
Processing functions and statements ...
[openss]: Restoring and displaying default view for:
/home/jeg/DEMOS/demos/sequential/smg2000/test/smg2000-pcsamp.openss
[openss]: The restored experiment identifier is: -x 1
```

Application output
followed by OSS output

Open | SpeedShop Performance Tool Status Update

Example Offline Run With Output (2)

osspsamp “./smg2000 -n 80 80 80”



Solve phase times:
=====

SMG Solve:
wall clock time = 4.810000 seconds
cpu clock time = 4.800000 seconds

Iterations = 7
Final Relative Residual Norm = 2.844100e-07

[openss]: Converting raw data from /tmp/jeg/offline-oss into temp file

Processing raw data for smg2000
Processing processes and threads ...
Processing performance data ...
Processing functions and statements ...

[openss]: Restoring and displaying default view for:
/home/jeg/DEMOS/demos/sequential/smg2000/test/smg2000-pcsamp.openss
[openss]: The restored experiment identifier is: -x 1

Exclusive CPU time % of CPU Time Function (defining location)
in seconds.

Exclusive CPU time in seconds.	% of CPU Time	Function (defining location)
2.680000	53.174603	hypr_SMGResidual (smg2000: smg_residual.c,152)
1.660000	32.936508	hypr_CyclicReduction (smg2000: cyclic_reduction.c,757)
0.110000	2.182540	hypr_SemiInterp (smg2000: semi_interp.c,126)
0.100000	1.984127	hypr_SemiRestrict (smg2000: semi_restrict.c,125)
0.080000	1.587302	hypr_SMG2BuildRAPSym (smg2000: smg2_setup_rap.c,156)
0.050000	0.992063	hypr_SMGAxy (smg2000: smg_axpy.c,27)
0.050000	0.992063	hypr_SMG3BuildRAPSym (smg2000: smg3_setup_rap.c,233)
0.040000	0.793651	hypr_StructVectorClearGhostValues (smg2000: struct_vector.c,592)
0.040000	0.793651	hypr_StructAxy (smg2000: struct_axpy.c,25)
0.030000	0.595238	hypr_StructMatrixInitializeData (smg2000: struct_matrix.c,314)
0.030000	0.595238	hypr_StructVectorSetConstantValues (smg2000: struct_vector.c,537)
0.020000	0.396825	hypr_CycRedSetupCoarseOp (smg2000: cyclic_reduction.c,211)
0.020000	0.396825	hypr_StructInnerProd (smg2000: struct_innerprod.c,32)
0.020000	0.396825	hypr_StructMatrixSetBoxValues (smg2000: struct_matrix.c,458)
0.020000	0.396825	main (smg2000: smg2000.c,21)
0.010000	0.198413	hypr_CreateCommInfoFromStencil (smg2000: communication_info.c,58)
0.010000	0.198413	_int_malloc (libc-2.10.2.so: malloc.c,0)
0.010000	0.198413	hypr_CommTypeDestroy (smg2000: communication.c,826)
0.010000	0.198413	__brk (libc-2.10.2.so)
0.010000	0.198413	hypr_StructVectorSetBoxValues (smg2000: struct_vector.c,258)
0.010000	0.198413	hypr_BoxArrayCreate (smg2000: box.c,106)
0.010000	0.198413	hypr_SMGSetupInterpOp (smg2000: smg_setup_interp.c,88)
0.010000	0.198413	hypr_SMGRelaxSetup (smg2000: smg_relax.c,357)
0.010000	0.198413	hypr_SMGSetStructVectorConstantValues (smg2000: smg.c,379)

[jeg@localhost test]\$

Performance Data
Default view: by Function

▣ **Outputs from: osspcsamp “executable”**

- *Normal program output while executable is running*
- *The sorted list of performance information*
 - ▣ *A list of the functions taking the most time*
 - ▣ *The corresponding sample derived time for each function*
- *A performance information database file (.openss file)*
 - ▣ *The database file contains all the information needed to view the data at anytime in the future without the executable(s).*
 - *Symbol table information from executable(s) and system libraries*
 - *Performance data openss gathered*
 - *Time stamps for when dynamic shared libraries were loaded and unloaded*

Open | SpeedShop Performance Tool Status Update

Default GUI View

Toolbar to switch Views

Performance Data
Default view: by Function

Graphical Representation

The screenshot shows the Open | SpeedShop GUI with a performance report. The report is titled "Showing Functions Report:" and lists executables: "sweep3d.mpi Hosts:(16) yra089.yr.lanl.gov ... Processes/Ranks/Threads:(16) 0 ...". The main table displays performance data for various functions. A callout points to a toolbar with icons for switching views. Another callout points to the table, indicating it is the default view for performance data. A third callout points to a graphical representation of the data.

Function (defining location)	Exclusive CPU time in s	% of CPU Time	Function (defining location)
sweep_ (sweep3d.mpi: sweep.f,2)	14532.36000000	85.48577802	sweep_ (sweep3d.mpi: sweep.f,2)
mca_btl_sm_component_progress (libmpi: ...)	-907.24000000	5.33678750	mca_btl_sm_component_progress (libmpi.so.0.0.0)
btl_openib_component_progress (libmpi: ...)	-584.22000000	3.43664080	btl_openib_component_progress (libmpi.so.0.0.0)
pthread_spin_lock (libpthread-2.3.4.so)	-165.18000000	0.97166192	pthread_spin_lock (libpthread-2.3.4.so)
source_ (sweep3d.mpi: source.f,2)	-131.21000000	0.77183533	source_ (sweep3d.mpi: source.f,2)
mca_bml_r2_progress (libmpi.so.0.0.0)	-115.48000000	0.67930451	mca_bml_r2_progress (libmpi.so.0.0.0)
opal_progress (libopen-pal.so.0.0.0)	-95.69000000	0.56289096	opal_progress (libopen-pal.so.0.0.0)

Open | SpeedShop Performance Tool Status Update

Associate Source and Performance Data

The screenshot shows the OpenSpeedShop application window titled "Open|SpeedShop (on yra089)". The interface includes a menu bar (File, Tools, Help), a status bar ("Status: Loaded saved data from file X.0.openss."), and two main panels: "Source Panel [1]" and "Stats Panel [1]".

Source Panel [1]: Displays source code for a file named "sweep.f". The code is as follows:

```
140.4900 474 do i = 1, it
683.7400 475 flux(i,j,k,1) = flux(i,j,k,1) + w(m)*phi(i)
476 end do
20.10000 477 do n = 2, nm
451.3200 478 do i = 1, it
479 flux(i,j,k,n) = flux(i,j,k,n)
>> 3047 480 & + pn(m,n,iq)*w(m)*phi(i)
481 end do
end do
```

Stats Panel [1]: Displays a table of performance data. The table has columns for "Exclusive CPU time in s", "% of CPU Time", and "Statement Location". The data is as follows:

Exclusive CPU time in s	% of CPU Time	Statement Location
3047.46000000	17.92650946	sweep.f(480)
2785.60000000	16.38613296	sweep.f(389)
1011.96000000	5.95279693	sweep.f(492)
989.41000000	5.82014784	sweep.f(490)
683.74000000	4.02206151	sweep.f(475)
667.60000000	3.92711889	sweep.f(488)
451.32000000	2.65486413	sweep.f(478)

Annotations:

- "Use window controls to split/arrange windows" points to the window control buttons in the Source Panel.
- "Double click to open source window" points to the Source Panel title bar.
- "Selected performance data point" points to the first row in the Stats Panel table.