





SC2018 Tutorial

How to Analyze the Performance of Parallel Codes 101 A case study with Open/SpeedShop

Martin Schulz: TU-München Jennifer Green: LANL Dave Montoya: LANL Don Maghrak: Krell Institute Jim Galarowicz: Krell Institute











Open | SpeedShop™

Why This Tutorial?



Performance Analysis is becoming more important

- Complex architectures and complex applications
- Mapping applications onto architectures is hard
- > Today's applications only use a fraction of the machine

Performance analysis is more than just measuring time

- > What are the critical sections in a code?
- Is a part of the code running efficiently or not?
- > Is the code using the resources well (memory, TLB, I/O, ...)?
- Where is the greatest payoff for optimization?

Often hard to know where to start

- > Which experiments to run first?
- > How to plan follow-on experiments?
- > What kind of problems can be explored?
- How to interpret the data?

2

Tutorial Goals



Sasic introduction into performance analysis

- > Typical pitfalls wrt. performance
- Wide range of types of performance tools and techniques

Provide basic guidance on ...

- > How to understand the performance of a code?
- How to answer basic performance questions?
- How to plan performance experiments?

Provide you with the ability to ...

- > Run these experiments on your own code
- Provide starting point for performance optimizations

Practical Experience: Demos and hands-on Experience

- Introduction into Open | SpeedShop as one possible tool solution
- Basic usage instructions and pointers to documentation
- Lessons and strategies apply to any tool

Open | SpeedShop™

Open | SpeedShop Tool Set



Open Source Performance Analysis Tool Framework

- > Most common performance analysis steps *all in one tool*
- Combines tracing and sampling techniques
- Extensible by plugins for data collection and representation
- Gathers and displays several types of performance information

Flexible and Easy to use

User access through:

GUI, Command Line, Python Scripting, convenience scripts

* Scalable Data Collection

- Instrumentation of *unmodified application binaries*
- New option for *hierarchical online data aggregation*

* Supports a wide range of systems

- > Extensively used and tested on a variety of *Linux clusters*
- Cray and Blue Gene support

Open | SpeedShop™

How to Analyze the Performance of Parallel Codes 101 - A Tutorial at SC'18

"Plan"/"Rules"



Staggered approach/agenda

- First session: performance analysis basics and getting ready
- Second session: Digging deeper and going parallel
- Third session: more specialized topics (HWC and I/O)
- Fourth session: new architectural challenges (memory and GPU)
- Hands-on experiments in each session

Let's keep this interactive

- Feel free to ask questions as we go along
- > Ask if you would like to see anything specific in the demos

We are interested in feedback!

- What was clear / what didn't make sense?
- What scenarios are missing?

Updated slides available before SC

- <u>https://www.openspeedshop.org/wp/category/tutorials</u>
- Then choose SC2018 Sunday Nov 11 tutorial

Presenters

- Martin Schulz: TU-München
- Im Galarowicz: Krell Institute
- **Donald Maghrak: Krell Institute** **
- Sennifer Green: LANL
- David Montoya: LANL
- Scantlen: CreativeC
- Hannes Schweiger: CreativeC

Larger team: **

- William Hachfeld, David Whitney: Krell Institute \triangleright
- Gregory Schultz: Argo Navis Technologies, LLC. \geq
- Mike Mason, David Shrader: LANL \geq
- Douglas Pase, Anthony Angelastos, Joel Stevenson: SNL \triangleright
- Matt Legendre and Chris Chambreau: LLNL \geq
- Dyninst group (Bart Miller: UW & Jeff Hollingsworth: UMD) \geq
- Phil Roth: ORNL \triangleright
- Koushik Ghosh: Engility \geq
- Mahesh Rajan: New Mexico Consortium \geq









Outline



- Welcome
- Concepts in performance analysis
- Introduction into Tools and Open | SpeedShop
- How to run basic timing experiments and what they can do?
- How to deal with parallelism (MPI and threads)?
- How to properly use hardware counters?
- Slightly more advanced targets for analysis
 - How to understand and optimize I/O activity?
 - How to evaluate memory efficiency?
 - How to analyze codes running on GPUs?
- DIY and Conclusions: DIY and Future trends
- Hands-on Exercises (after each section)
 - On site cluster available
 - We will provide exercises and test codes

Tutorial Survey



- * Tutorial surveys are entirely electronic this year
 - No paper forms
 - Tutorial attendees will receive an email reminder with the evaluation information.
- & QR code:

https://submissions.supercomputing.org/eval.png

- Evaluation site URL: <u>http://bit.ly/SC18-eval</u>
- * Thanks for attending our tutorial!







SC2018 Tutorial

How to Analyze the Performance of Parallel Codes 101 A case study with Open/SpeedShop

Section 1 Concepts in Performance Analysis











Open | SpeedShop™

Typical Development Cycle

Performance tuning is an essential part of the development cycle

- Potential impact at every stage
 - Message patterns
 - Data structure layout
 - Algorithms
- Should be done from early on in the life of a new HPC code
- Ideally continuously and automatically

Typical use

- Measure performance and store data
- > Analyze data
- Modify code and/or algorithm
- Repeat measurements
- Analyze differences

10





Algorithm

Coding

A Case for Performance Tools

First line of defense

- Full execution timings (UNIX: "time" command)
- Comparisons between input parameters
- Keep and track historical trends

Disadvantages

- Measurements are coarse grain
- Can't pin performance bottlenecks

Alternative: code integration of performance probes

- Hard to maintain
- Requirements significant a priori knowledge

Performance tools

- Enable fine grain instrumentation
- Show relation to source code
- Work universally across applications

Open | SpeedShop™

Performance Tools Overview



Basic OS tools

time, gprof, strace

Hardware counters

- > PAPI API & tool set
- hwctime (AIX)

Sampling tools

- > Typically unmodified binaries
- Callstack analysis
- > HPCToolkit (Rice U.)

Profiling/direct measurements

- MPI or OpenMP profiles
- > mpiP (LLNL&ORNL)
- ompP (LMU Munich)

Tracing tool kits

- > Capture all MPI events
- Present as timeline
- Vampir (TU-Dresden)
- Jumpshot (ANL)

✤ Trace Analysis

- Profile and trace capture
- > Automatic (parallel) trace analysis
- Kojak/Scalasca (JSC)
- Paraver (BSC)

Integrated tool kits

- Typically profiling and tracing
- Combined workflow
- > Typically GUI/some vis. support
- > Binary: Open|SpeedShop (Krell/TriLab)
- Source: TAU (U. of Oregon)

Specialized tools/techniques

- Libra (LLNL)
 Load balance analysis
- Boxfish (LLNL/Utah/Davis)
 3D visualization of torus networks
- Rubik (LLNL)
 Node mapping on torus architectures

11/11/2018

Vendor Tools

Open | SpeedShop™

How to Analyze the Performance of Parallel Codes 101 - A Tutorial at SC'18

How to Select a Tool?

A tool with the right features

- Must be easy to use
- Provides performance analysis of the code at different levels: libraries, functions, loops, statements

A tool must match the application's workflow

- Requirements from instrumentation technique
 - Access to and knowledge about source code? Recompilation time?
 - Machine environments? Supported platforms?
- Interactive and batch mode analysis options
- Support iterative tuning with ability to compare key metrics across runs

Why We Picked/Developed Open | SpeedShop?

- Sampling and tracing in a single framework
- Easy to use GUI & command line options for remote execution
 - Low learning curve for end users
- Transparent instrumentation (preloading & binary)
 - No need to recompile application

13



Next Step: Interpret Data

Tools can collect lots of data

- > At varying granularity
- > At varying cost
- At varying accuracy

Issue 1: Understand your tool and its limitations

- No tool can do everything (at least not well)
- Choose the right tool for the right task

Issue 2: Ask the right question

- Need to know basic issues to look for to get started
- Need to understand expected behavior



File Tools

π 11

2.182540

1.984127

1.587302

Click on the "Run" button to begin the experimen



₽ 🗆 🗆

ng2000: smg residual.c,152) CyclicReduction (smg2000: cyclic reduction.c.757)

Issue 1: Tool Types



Data acquisition

- Event based data: triggered by explicit events
 - Direct correlation possible, but may come in bursts
- > Sampling based data: triggered by external events like timers
 - Even distribution, but requires statistical analysis

Instrumentation

- Source code instrumentation: exact, but invasive
- Compiler instrumentation: requires source, but transparent
- Binary instrumentation: can be transparent, but still costly
- Link-level: transparent, less costly, but limited to APIs
- > Tradeoff: invasiveness vs. overhead vs. ability to correlate
- Big question: granularity

Aggregation

- No aggregation: trace
- > Aggregation over time and space: simplified profile
- Many shades of gray in between

11/11/2018

Issue 2: Asking the Right Questions



Step 1: Find where the problem actually is

- > Where is the code spending time?
 - Which code sections are even worth looking at?
- > Where should it spend time?
 - Have a (mental) model of your application

Se overview experiments

- Identify bottlenecks for your application
 - Which resource in the system is holding you back?
- Decide where to dig deeper
 - Important resource AND worth optimizing AND unexpected behavior

Pick the right tool or experiment in a tool

- Target the specific bottleneck
- Decide on instrumentation approach
- Decide on useful aggregation
- Understand impact on code perturbation

What to Look For: Sequential Runs

Step 1: Identify computational intensive parts

- > Where am I spending my time?
 - Modules/Libraries
 - Loops
 - Statements
 - Functions
- > Is the time spent in the computational kernels?
- Does this match my intuition?

Impact of memory hierarchy

- Do I have excessive cache misses?
- How is my data locality?
- Impact of TLB misses?

Sector External resources

- Is my I/O efficient?
- > Time spent in system libraries?



What to Look For: Shared Memory

K R E L L H n s t i t u t e ASC

Shared memory model

- Single shared storage
- Accessible from any CPU

Common programming models

- > Explicit threads (e.g., POSIX threads)
- > OpenMP

Typical performance issues

- False cache sharing
- Excessive Synchronization
- Limited work per thread
- Threading overhead

Complications: NUMA

- Memory locality critical
- > Thread:Memory assignments





11/11/2018

What to Look For: Message Passing

Distributed Memory Model

- Sequential/shared memory nodes coupled by a network
- > Only local memory access
- > Data exchange using message passing (e.g., MPI)

Typical performance issues

- Load imbalance; Processes waiting for data
- Large fraction of time on collective operations
- Network and I/O contention
- Non-optimal process placement & binding



19

Application

MPI Library





Overview of Open | SpeedShop

Help to understand demos and hands-on exercises

Basic questions

- > Where am I spending my time?
- How to understand the context of this information?

Hardware/Resource utilization

- > How to use hardware counters efficiently?
- > How to turn this information into actionable insight?

Next step beyond the computational core

- > How well is my I/O doing?
- How well am I utilizing memory?
- How can I understand the performance on accelerators?

20







SC2018 Tutorial

How to Analyze the Performance of Parallel Codes 101 A case study with Open/SpeedShop

Section 2 Emerging Performance Analysis models











Open | SpeedShop™

Emerging Performance Monitoring



- New approach to performance tracking
 - Light weight but with broad collection abilities
 - > Broad but not Deep
- Why?
 - Discovery looking for a place to start
 - Easy less difficult to run and understand
 - Comparison High level across architectures
 - Or compilers
 - Or System environments, File systems, etc.
 - Always on Option to add to monitoring infrastructures

Monitoring and Analysis

> Add to broader monitoring infrastructures that are collecting other data from multiple sources and aggregated for analysis

Open | SpeedShop^{*}



Open | SpeedShop[™]

How to Analyze the Performance of Parallel Codes 101 - A Tutorial at SC'18

High-level versus In-depth perf. tools



- High level performance analysis versus traditional indepth tools
 - In depth tools give per function, per statement, per loop type information
 - > High level tools give an overview, per execution view
- Use high level tools to get an understanding of application performance forensic approach
- Use In-depth tools to "home-in" on solving the issues found by using the high level tools.

Lightweight and Monitoring Tool efforts

- * TAU execution command- tau_exec
 - MPI example % mpirun -np 256 tau_exec ./a.out[<]
- * ARM Performance Reports (originally Allinea)
- * Lightweight Distributed Metric Service LDMS (SNL)
 - Provides capabilities for lightweight run-time collection of highfidelity data. Node level, system data.
- * OSS CBTF Summary details later
- * Caliper /SPOT (LLNL- in-development)
 - Instrumented calipers, included library, on-going collection and analysis via web page

Monitoring

- LLNL Sonar monitoring infrastructure
- LANL Insight monitoring infrastructure

Open | SpeedShop™

How to Analyze the Performance of Parallel Codes 101 - A Tutorial at SC'18









* High-level view valuable to:

- Analysts (End Users)
- System Administrators
- System Architects
- Code Developers

* Analysts/Users

- Show load balance (min,max,average) across processes (mpi) and/or threads (openMP)
- Hardware counters metrics can give insight into processor instruction and cache usage
- Memory metrics can show highwater memory mark (see the maximum memory resident in your program), total allocation calls and size, and total free calls. Comparing total allocations to total frees can give an indication of memory leaks.
- I/O metrics give insight to frequency of read/write calls and total size of read/write calls.



System Administrators

- Use continuous integration performance analysis to monitor applications to see which ones are performing badly
- Monitor set of applications before and after system changes

System Architects

- Use the information to extrapolate how new architecture features would impact performance of hypothetical new machine
 - How to improve scalar performance?
 - What happens if more vectorization could be done?
 - Run the survey tool on multiple architectures and make projections

Code Developers

- > Analyze an application before and after source changes to it
- Decide what areas of the code need to have further improvements

Open | SpeedShop™

High-level Performance Analysis: O|SS



Open|SpeedShop's cbtfsummary tool

- The summary experiment currently gathers high-level performance metrics, such as:
 - Time spent in MPI routines
 - Time spent in OpenMP (idle time, barrier time, task time)
 - Hardware counters (multiplexes HW counters)
 - Time spent in I/O (breaks down read and write times and byte totals)
 - Memory information
 - Allocation calls, bytes, and time
 - Free calls and time
 - PAPI dmem and statistics including high water mark
 - rusage max rss
 - rusage utime and stime

High level Performance Analysis: O|SS



Open|SpeedShop's cbtfsummary tool

- Usage: cbtfsummary "normal app run script"
- Outputs:
 - > Human readable report to stdout
 - > Human readable report to text file
 - Composite csv file
 - Directory structure containing per-thread of execution csv files

High-level Performance Analysis: O|SS

cbtfsummary tool example:

setenv CBTF_MPI_IMPLEMENTATION mvapich2

setenv CBTF_CSVDATA_DIR ./rzgenie_lulesh2_n64_omp2_csvdata

cbtfsummary "srun -n 64 ./lulesh2.0 -p -i 90"

Processing csv files in ./rzgenie_lulesh2_n64_omp2_csvdata/lulesh2-overview-csvdata-1 Metrics for thread 0 in 64 ranks

metric name	max	min	avg
implicit_task_time_seconds	s 2.290756	1.707665	1.970446
serial_time_seconds	3.323400	2.740352	3.070070
PAPI TOT INS	32141739510	26835045450	29373460884
PAPI DP OPS	6851778300	5937733304	6162486573
PAPITLDTINS	10179972804	8494295326	9310108000
PAPI VEC DP	6064913156	5325306542	5511166403
PAPI TOT CYC	17585366336	16103836394	16816938221
allocation calls	20208	19435	19795
allocation bytes	2393954217	2281932664	2308348507
allocation time seconds	0.025247	0.009769	0.011933
total time seconds	5.372663	5.214180	5.265309
free calls	20044	19288	19642
free time seconds	0.510833	0.271222	0.399044
io total time seconds	0.005593	0.002904	0.004288
write time seconds	0.005573	0.002877	0.004256
write bytes	450	180	337
dmem size kB	405080	345592	366302
dmem heap kB	159272	158968	161593
dmem high water mark k	B 74532	72180	73242
dmem shared kB	10520	8428	9122
dmem resident kB	34132	31816	32665
stime seconds	1.782871	1.247015	1.476539
utime_seconds	3.603105	3.281256	3.457614
maxrss kB	74532	72180	73242
total_mpi_time_seconds	2.928412	2.251567	2.596778

...







SC2018 Tutorial

How to Analyze the Performance of Parallel Codes 101 A case study with Open/SpeedShop

Section 3 Introduction into Tools and Open|SpeedShop











Open | SpeedShop™

Open | SpeedShop Tool Set



Open Source Performance Analysis Tool Framework

- > Most common performance analysis steps *all in one tool*
- Combines tracing and sampling techniques
- > **Extensible** by plugins for data collection and representation
- Gathers and displays several types of performance information

* Flexible and Easy to use

> User access through:

GUI, Command Line, Python Scripting, convenience scripts

Scalable Data Collection

- Instrumentation of *unmodified application binaries*
- New option for *hierarchical online data aggregation*

* Supports a wide range of systems

- > Extensively used and tested on a variety of *Linux clusters*
- > Cray, Blue Gene, ARM, Power 8, Intel Phi, GPU support

Open | SpeedShop^{*}

Classifying Open | SpeedShop



Offers both sampling and direct instrumentation

- Sampling for overview and hardware counter experiments
 - Even and low overhead, overview information
- Direct instrumentation for more detailed experiments
 - More in-depth information, but potentially bursty
- All instrumentation at link-time of runtime

Multiple direct instrumentation options

- > API level instrumentation (e.g., I/O or memory)
- Loop analysis based on binary instrumentation techniques
- Programming model specific instrumentation (e.g., MPI or OpenMP)

Aggregation

- > By default: aggregate profile data over time
 - Example: intervals, functions, ...
 - Full traces possible for some experiments (e.g.. MPI), but costly
- For parallel experiments: by default aggregation over threads, processes, ...
 - However, users can query per process/thread data

Open | SpeedShop Workflow



srun -n4 -N1 smg2000 -n 65 65 65

osspcsamp "srun -n4 -N1 smg2000 -n 65 65 65"





http://www.openspeedshop.org/

Open | SpeedShop Workflow



srun -n4 -N1 smg2000 -n 65 65 65

osspcsamp -- offline "srun -- n4 -- N1 smg2000 -- n 65 65 65"



MPI Application

http://www.openspeedshop.org/

Open | SpeedShop^{*}

Open | SpeedShop™

How to Analyze the Performance of Parallel Codes 101 - A Tutorial at SC'18

36

11/11/2018

Alternative Interfaces

Scripting language

- Immediate command interface
- > O|SS interactive command line (CLI)
 - openss -cli

Experiment Commands expView expCompare expStatus

List Commands list -v exp list -v hosts

Python module

```
import openss
my_filename=openss.FileList("myprog.a.out")
my_exptype=openss.ExpTypeList("pcsamp")
my_id=openss.expCreate(my_filename,my_exptype)
openss.expGo()
My_metric_list = openss.MetricList("exclusive")
my_viewtype = openss.ViewTypeList("pcsamp")
result = openss.expView(my_id,my_viewtype,my_metric_list)
```


Central Concept: Experiments

K R E L L H n s t i t u t e ASC

***** Users pick experiments:

- What to measure and from which sources?
- > How to select, view, and analyze the resulting data?

Two main classes of performance data collection:

- Statistical Sampling
 - Periodically interrupt execution and record location
 - Useful to get an overview
 - Low and uniform overhead
- Event Tracing
 - Gather and store individual application events
 - Provides detailed per event information
 - Can lead to huge data volumes

O|SS can be extended with additional experiments

Sampling Experiments in O|SS



PC Sampling (pcsamp)

- Record PC repeatedly at user defined time interval
- Low overhead overview of time distribution
- Good first step, lightweight overview

Call Path Profiling (usertime)

- PC Sampling and Call stacks for each sample
- Provides inclusive and exclusive timing data
- Use to find hot call paths, caller and callee relationships

Hardware Counters (hwc, hwctime, hwcsamp)

- Provides profile of hardware counter events like cache & TLB misses
- hwcsamp:
 - Periodically sample to capture profile of the code against the chosen counter
 - Default events are PAPI_TOT_INS and PAPI_TOT_CYC
- hwc, hwctime:
 - Sample a hardware counter till a certain number of events (called threshold) is recorded and get Call Stack
 - Default event is PAPI_TOT_CYC

Tracing Experiments in O|SS



Input/Output Tracing (io, iot, iop)

- Record invocation of all POSIX I/O events
- Provides aggregate and individual timings
- Store function arguments and return code for each call (iot)
- Lightweight I/O profiling because not tracking individual call details (iop)

MPI Tracing (mpi, mpit, mpip)

- Record invocation of all MPI routines
- Provides aggregate and individual timings
- Store function arguments and return code for each call (mpit)
- Lightweight MPI profiling because not tracking individual call details (mpip)

Tracing Experiments in O|SS



Memory Tracing (mem)

- Tracks potential memory allocation call that is not later destroyed (leak).
- Records any memory allocation event that set a new high-water of allocated memory current thread or process.
- Creates an event for each unique call path to a traced memory call and records:
 - The total number of times this call path was followed
 - The max allocation size
 - The min allocation size
 - The total allocation
 - The total time spent in the call path
 - The start time for the first call

40

Additional Experiments in OSS/CBTF



CUDA NVIDIA GPU Event Tracing (cuda)

- Record CUDA events, provides timeline and event timings
- Traces all NVIDIA CUDA kernel executions and the data transfers between main memory and the GPU.
- Records the call sites, time spent, and data transfer sizes.

POSIX thread tracing (pthreads)

- Record invocation of all POSIX thread events
- Provides aggregate and individual rank, thread, or process timings

OpenMP specific profiling/tracing (omptp)

Report task idle, barrier, and barrier wait times per OpenMP thread and attribute those times to the OpenMP parallel regions.

Open | SpeedShop^{*}

Performance Analysis in Parallel



* How to deal with concurrency?

- > Any experiment can be applied to parallel application
 - Important step: aggregation or selection of data
- Special experiments targeting parallelism/synchronization

***** O|SS supports MPI and threaded codes

- > Automatically applied to all tasks/threads
- > Default views aggregate across all tasks/threads
- Data from individual tasks/threads available
- Thread support (incl. OpenMP) based on POSIX threads

Specific parallel experiments (e.g., MPI)

- Wraps MPI calls and reports
 - MPI routine time
 - MPI routine parameter information
- The mpit experiment also stores function arguments and return codes for each call

Open | SpeedShop™

How to Run a First Experiment in O|SS?

1. Picking the experiment

- What do I want to measure?
- We will start with pcsamp to get a first overview

2. Launching the application

- How do I control my application under O|SS?
- > Enclose how you normally run your application in quotes
- osspcsamp "mpirun –np 4 smg2000 –n 50 50 50"

3. Storing the results

- > O/SS will create a database
- Name: smg2000-pcsamp-0.openss

4. Exploring the gathered data

- How do I interpret the data?
- O|SS will print a default report
- Open the GUI to analyze data in detail (run: "openss")

Example Run with Output



osspcsamp "mpirun –np 4 smg2000 –n 50 50 50" (1/2)

Bash>osspcsamp "mpirun -np 4 ./smg2000 -n 50 50" [openss]: pcsamp experiment using the default sampling rate: "100". Creating topology file for frontend host localhost Generated topology file: ./cbtfAutoTopology Running pcsamp collector. Program: mpirun -np 4 ./smg2000 -n 50 50 50 Number of mrnet backends: 4 Topology file used: ./cbtfAutoTopology executing mpi program: mpirun -np 4 cbtfrun --mpi --mrnet -c pcsamp ./smg2000 -n 50 50 50 Running with these driver parameters: (nx, ny, nz) = (65, 65, 65) ...

<SMG native output>

```
Final Relative Residual Norm = 1.774415e-07
All Threads are finished.
default view for ./smg2000-pcsamp-0.openss
[openss]: The restored experiment identifier is: -x 1
Performance data spans 2.257689 seconds from 2016/11/09 13:33:33 to 2016/11/09 13:33:35
```

Example Run with Output



osspcsamp "mpirun –np 4 smg2000 –n 50 50 50" (2/2)

```
Exclusive % of CPU Function (defining location)
CPU time
   in
seconds.
2.850000 36.821705 hypre SMGResidual (smg2000: smg residual.c, 152)
1.740000 22.480620 hypre CyclicReduction (smg2000: cyclic reduction.c,757)
0.410000 5.297158 mca btl vader check fboxes (libmpi.so.12.0.2: btl vader fbox.h,184)
0.250000 3.229974 opal progress (libopen-pal.so.13.0.2: opal progress.c,151)
0.250000 3.229974 hypre_SemiInterp (smg2000: semi_interp.c,126)
0.190000 2.454780 pack predefined data (libopen-pal.so.13.0.2: opal datatype pack.h,35)
0.190000 2.454780 unpack predefined data (libopen-pal.so.13.0.2: opal datatype unpack.h,34)
0.120000 1.550388 int malloc (libc-2.17.so)
0.100000 1.291990 hypre SemiRestrict (smg2000: semi restrict.c,125)
0.100000 1.291990 opal_generic_simple_pack (libopen-pal.so.13.0.2: opal_datatype_pack.c,274)
0.090000 1.162791 memcpy ssse3 back (libc-2.17.so)
0.080000 1.033592 int free (libc-2.17.so)
0.080000 1.033592 opal_generic_simple_unpack (libopen-pal.so.13.0.2:
opal datatype unpack.c,263)
```

View with GUI: openss –f smg2000-pcsamp-0.openss

Open | SpeedShop[™]

45

Default Output Report View





Open | SpeedShop^{*}

Statement Report Output View





Associate Source & Performance Data





Open | SpeedShop[™]

Library (LinkedObject) View



		Open SpeedShop			×
<u>F</u> ile <u>T</u> ools <u>H</u> elp				-	
▼ pc Sampling [1]		– Select Li – View typ	nkedObje e and Cli	ct ck	⊑ 🗆 = ×
→ Run → Cont →	Pause 🗗 Update	on I	D-icon		Terminate
Status: Process Loaded: Click on the "Run"	button to begin the experiment.				
Source Panel [1] Stats Panel [1]					
			Г	View/Display	Choice
T T CL I STOV LE CA CC	Showing Linked Objects Report			○ Functions	◯ Statements ● Linked Objects ◯ Loops
Executables: smg2000 Host: localhost Pids:	4 Ranks: 4 Threads: 4				
% of CPU Time	Exclusive CPU time in seconds.	% of CPU Time	LinkedObject		
78.143877	- 14.230000	78.143877	smg2000		
	1.960000	10.763317	libmpi.so.1.4.0		
10.763317	1.300000	7.138935	libopen-pal.so.6.1.1		
	-0.700000	3.844042	libc-2.17.so		
7.138935	0.010000	0.054915	pcsamp-rt-offline.so		
	····· 0.010000	0.054915	libmonitor.so.0.0.0		
3.844042					
Command Panel					№ 🗆 🗆 ×
openss>> Shows	time spent in				
libroriog	Con indicato				
noraries.	Can indicate	;	- 11 - 1	•	
' <u> im</u>	balance.		-Librarie	s in th	e application

Loop View



Open SpeedShop 2						
<u>F</u> ile <u>T</u> ools <u>H</u> elp		— Sele	ect Loops			
💌 pc Sampling [1]		View ty	ne and Click	⊑ 🗆 🗆 ×		
Process Control						
➡ Run 🕩 Cont ➡	Pause 🍯 Update	On	D-1con	Terminate		
Status: Process Loaded: Click on the "Run" button to begin the experiment.						
Source Panel [1] Stats Panel [1]			View/Dierley Ch			
View/Display Choice						
	U CL L Functions () Statements () Linked Objects () Loops					
Executables: smg2000 Host: localhost Pids: 4 Ranks: 4 Threads: 4						
% of CPU Time	Exclusive CPU time in seconds.	% of CPU Time	Loop Start Location (Line Number			
28.754110	7.870000	28.754110	smg_residual.c(205)			
7.672634	2.100000	7.672634	i72634 cyclic_reduction.c(882)			
	1.980000	7.234198	cyclic_reduction.c(1022)			
7.234198	1.050000	3.836317	smg_residual.c(237)			
3.836317	1.050000	3.836317	smg_residual.c(220)			
	1.050000	3.836317	smg_residual.c(237)			
<mark>3.</mark> 836317	· 1.040000	3.799781	smg_residual.c(237)			
other	- 1.010000	3.690172	btl_vader_fbox.h(117)			
	-0.440000	1.607600	opal_datatype_unpack.h(65)			
Command Panel				·····································		
openss>>						
Shows time spent in loops.						
	1	1	Statement nu	mber of start		
			— of lo	<u>000.</u>		
				S. P.		

Open | SpeedShop[™]

How to Analyze the Performance of Parallel Codes 101 - A Tutorial at SC'18

Open | SpeedShop Basics



- Place the way you run your application normally in quotes and pass it as an argument to osspcsamp, or any of the other experiment convenience scripts: ossio, ossmpi, etc.
 > osspcsamp "srun –N 8 –n 64 ./mpi_application app_args"
- Open|SpeedShop sends a summary profile to stdout
- Open | SpeedShop creates a database file
- Display alternative views of the data with the GUI via:
 > openss –f <database file>
- Display alternative views of the data with the CLI via:
 > openss -cli -f <database file>
- Start with pcsamp for overview of performance
- Then, focus on performance issues with other experiments

51



How to log into the tutorial computer system

- > Login information will be distributed at this time.
- > The "exercises" directory will be in your \$HOME directory.
- Also can find these exercises at:
 - www.openspeedshop.org/downloads
- Top-level directory has file: EXERCISES that lists all the tutorial exercises and README file has general information.
- A "docs" directory in your \$HOME has OpenSpeedShop documentation and the updated tutorial slides.

***** Exercise is in the exercise directory:

- \$HOME/exercises/loop_check
- Consult README file in each of the directories for the instructions/guidance

IBM Power system from CreativeC (Greg Scantlen)

Open | SpeedShop™

52







SC2018 Tutorial

How to Analyze the Performance of Parallel Codes 101 A case study with Open/SpeedShop

Section 4

Basic timing experiments and their Pros/Cons











Open | SpeedShop™



Flat Profile Overview

Profiles show computationally intensive code regions

First views: Time spent per functions or per statements

* Questions:

- > Are those functions/statements expected?
- Do they match the computational kernels?
- Any runtime functions taking a lot of time?

Identify bottleneck components

- View the profile aggregated by shared objects
- Correct/expected modules?
- Impact of support and runtime libraries

54

Adding Context through Stack Traces





- Missing information in flat profiles
 - Distinguish routines called from multiple callers
 - Understand the call invocation history
 - Context for performance data

Critical technique: Stack traces

- Gather stack trace for each performance sample
- Aggregate only samples with equal trace

Ser perspective:

- Butterfly views (caller/callee relationships)
- Hot call paths
 - Paths through application that take most time

Open | SpeedShop™

Inclusive vs. Exclusive Timing





- Stack traces enable calculation of inclusive/exclusive times
 - Time spent inside a function only (exclusive)
 - See: Function B
 - Time spent inside a function and its children (inclusive)
 - See Function C and children

Implementation similar to flat profiles

- Sample PC information
- Additionally collect call stack information at every sample

* Tradeoffs

- Pro: Obtain additional context information
- Con: Higher overhead/lower sampling rate

Open | SpeedShop™

How to Analyze the Performance of Parallel Codes 101 - A Tutorial at SC'18

Call path profiling & Comparisons



* Call Path Profiling

- > Take a sample: address inside a function
- > Call stack: series of program counter addresses (PCs)
- Unwinding the stack is walking through those address and recording that information for symbol resolution later.
- Leaf function is at the end of the call stack list

Open|SpeedShop: experiment called usertime

- > Time spent inside a routine vs. its children
- Key view: butterfly

* Comparisons

- > Between experiments to study improvements/changes
- Between ranks/threads to understand differences/outliers

57

Interpreting Call Context Data



Inclusive versus exclusive times

- If similar: child executions are insignificant
 - May not be useful to profile below this layer
- > If inclusive time significantly greater than exclusive time:
 - Focus attention to the execution times of the children

Hotpath analysis

- Which paths takes the most time?
- > Path time might be ok/expected, but could point to a problem

Sutterfly analysis (similar to gprof)

- Should be done on "suspicious" functions
 - Functions with large execution time
 - Functions with large difference between inclusive and exclusive time
 - Functions of interest
 - Functions that "take unexpectedly long"

•

Shows split of time in callees and callers



Basic syntax:

ossusertime "how you run your executable normally"

Examples:

ossusertime "smg2000 –n 50 50 50" ossusertime "smg2000 –n 50 50 50" low

Parameters

Sampling frequency (samples per second) Alternative parameter: high (70) | low (18) | default (35)

Recommendation: compile code with –g to get statements!

Open | SpeedShop™

Reading Inclusive/Exclusive Timings



Default View

- Similar to pcsamp view from first example
- Calculates inclusive versus exclusive times



60

Stack Trace Views: Hot Call Path



	Open S	seedShop
File Tools Help ▼ User Time [1] Process Control ⇒ Run ♦ Cont ♦] Pause ♥ User	īpdate	Hot Call Path
Status: Process Loaded: Click on the "Run" b Stats Panel [1] ManageProcessesPanel T T CL D S C L HC B TS OV SA Executables: smg 100 nost: localhost Pids: :	uteon to begin the experiment. 1 [1] 13 CA CC Showing Hot Callpath Report: 2 Ranks: 2 Threads: 2	
Exclusive CPU tiple in seconds. Inclusive	e CPU time in seconds. <mark>% of Total Exclus</mark> 9 1.754386	ve CPU Call Stack Function (defining location) start (smg2000) @ 556 inlibc_start_main (libmonitor.so.0.0.0) libc_start_main (life-2.14.90.so) @ 517 in monitor_main (libmonitor.so.0.0.0) @ 510 in main (mg2000: smg2000.c,21) @ 65 in HYPPF_StructSMGSolve (smg2000: HYPRE_struct_smg.c,64) @ 168 in hypre_SMGSolve (smg2000: smg_solve.c,57) @ 289 in hypre_SMGResidual (smg2000: smg_residual.c,152)
ccess to call paths All call paths (C+)	•	_start (smg2000)
All call naths for		

• All call paths for selected function (C♥)

Stack Trace Views: Butterfly View



Similar to well known "gprof" tool

Op	en SpeedShop ×
<u>F</u> ile <u>T</u> ools <u>H</u> elp	Callers of
User Time [1]	"hypre SMGSolve" • *
Process Control	
➡ Run 🕨 Cont ➡ Pause 5 Update	Terminate
Status: Process Loaded: Click on the "Run" button to begin the experiment.	
▼ Stats Panel [1] ▼ManageProcessesPanel [1]	G ₈ □ - ×
📅 🗊 💼 🗊 🕵 💕 📢 🕼 🔹 🐨 🐼 🧏 🍱 CA CC Showing Butterfly Repor	:
Executables: smg2000 Host: localhost Pids: 2 Ranks: 2 Threads: 2	
Inclusive CPU time in sec % of Total Inclusive CPU Call Stack Function (definit	ng lobaion)
-48.971428 98.336202 HYPRE_StructSMGSolve (s	ng2000: HYPRE_struct_smg.c,64)
-0.828571 1.663798 hypre_SMGRelax (smg2000	: smg_relax.c,228)
tar 49.799999 100.000000 hypre_SMGSolve (smg200): smg_solve.c,57)
-48.657142 97.705106 hypre_SMGRelax (smg:	.000: smg_relax.c,228)
-0.771429 1.549053 hypre_SMGResidual (st	ng2000: smg_residual.c,152)
-0.114286 0.229489 hypre_Semiliterp (smg	2000: semi_interp.c,126)
-0.1/1429 0.344234 nypre_StructionerProd	(smg, shi struct_innerprod.c,32)
-0.03/143 0.114/45 nypre_SemiRestrict (sn	g2000: semi_t_t_title.c,125)
Pivot routine	Callees of "hypre_SMGSolve"
"hypre_SMGSolve"	
<u> </u>	

Open | SpeedShop™

Comparing Performance Data



* Key functionality for any performance analysis

- Absolute numbers often don't help
- Need some kind of baseline / number to compare against

Typical examples

- > Before/after optimization
- Different configurations or inputs
- Different ranks, processes or threads

Very limited support in most tools

- Manual operation after multiple runs
- Requires lining up profile data
- Even harder for traces

Open | SpeedShop has support to line up profiles

- Perform multiple experiments and create multiple databases
- Script to load all experiments and create multiple columns

Open | SpeedShop™

Comparing Performance Data in O|SS



Convenience Script: osscompare

- Compares Open | SpeedShop up to 8 databases to each other
 - Syntax: osscompare "db1.openss,db2.openss,..." [options]
 - osscompare man page has more details
- Produces side-by-side comparison listing
- Data metric option parameter:
 - Compare based on: time, percent, a hwc counter, etc.
- Limit the number of lines by "rows=nn" option
- Specify the: viewtype=[functions|statements|linkedobjects]
 - Control the view granularity.
 - Compare based on the function, statement, or library level.
 - By default the compare will be done comparing the performance of functions in each of the databases.
 - If statements option is specified then all the comparisons will be made by looking at the performance of each statement in all the databases that are specified.
 - Similar for libraries, if linkedobject is selected as the viewtype parameter.

64

Comparison Report in O|SS



osscompare "smg2000-pcsamp.openss,smg2000-pcsamp-1.openss"

openss]: Legend: -c 2 represents smg2000-pcsamp.openss

[openss]: Legend: -c 4 represents smg2000-pcsamp-1.openss

-c 2, Exclusive CPU -c 4, Exclusive CPU Function (defining location)

time in seconds. time in seconds.

- 3.870000000 3.630000000 hypre_SMGResidual (smg2000: smg_residual.c,152)
- 2.610000000 2.860000000 hypre_CyclicReduction (smg2000: cyclic_reduction.c,757)
- 2.03000000 0.15000000 opal_progress (libopen-pal.so.0.0.)
- 1.33000000 0.10000000 mca_btl_sm_component_progress (libmpi.so.0.0.2: topo_unity_component.c,0)
 - 0.280000000 0.210000000 hypre_SemiInterp (smg2000: semi_interp.c,126)

0.28000000 0.04000000 mca_pml_ob1_progress (libmpi.so.0.0.2:

```
topo_unity_component.c,0)
```



Typical starting point:

- Flat profile
- Aggregated information on where time is spent in a code
- Low and uniform overhead when implemented as sampling

Adding context

- > From where was a routine called, which routine did it call
- Enables the calculation of exclusive and inclusive timing
- Technique: stack traces combined with sampling

Key analysis options

- Hot call paths that contains most execution time
- Butterfly view to show relations to parents/children

Comparative analysis

- > Absolute numbers often carry little meaning
- Need the correct baseline, then compare against that

11/11/2018

66



- Sasic sampling application exercise
 - Also comparing runs to each other
- ***** Exercises are in the exercise directory:
 - \$HOME/exercises/seq_lulesh/test
- Consult README file in each of the directories for the instructions/guidance

67







SC2018 Tutorial

How to Analyze the Performance of Parallel Codes 101 A case study with Open/SpeedShop

Section 5 Analysis of parallel codes: MPI, OpenMP, POSIX threads











Open | SpeedShop™

Parallel Application Performance Challenges



Architectures are Complex and Evolving Rapidly

- Changing multicore processor designs
- Emergence of accelerators (GPGPU, MIC, etc.)
- Multi-level memory hierarchy
- I/O storage sub-systems
- Increasing scale: number of processors, accelerators

Parallel processing adds more performance factors

- MPI communication time versus computation time
- Threading synchronization time versus computation time
- CPU time versus accelerator transfer and startup time tradeoffs
- I/O device multi-process contention issues
- Efficient memory referencing across processes/threads
- Changes in application performance due to adapting to new architectures

Parallel Execution Goals





Ideal scenario

- Efficient threading when using pthreads or OpenMP
 - All threads are assigned work that can execute concurrently
 - Synchronization times are low.
- Load balance for parallel jobs using MPI
 - All MPI ranks doing same amount of work, so no MPI rank waits
- > Hybrid application with both MPI and threads
 - Limited amount of serial work per MPI process

Open | SpeedShop^{*}

Parallel Execution Goals



***** What causes the ideal goal to fail?

- > For MPI:
 - Equal work was not given to each rank
 - There is an out of balance communication pattern occurring
 - The application can't scale with the number of ranks being used
- For threaded applications:
 - One or more threads doing more work than others and subsequently causing other threads to wait.
- For hybrid applications:
 - Too much time spent between parallel/threaded regions
- For multicore processors:
 - Remote memory references from the non-uniform access shared memory can cause sub-par performance
- > For accelerators:
 - Data transfers to the accelerator kernel might take more time than the speed-up for the accelerator operations on that data also is the CPU fully utilized?

Open | SpeedShop™

71

Parallel Application Analysis Techniques



What steps can we take to analyze parallel jobs?

- > Get an overview of where the time is being spent.
 - Use sampling to get a low overhead overview of time spent
 - Program counter, call stack, hardware counter
- > Examine overview information for all ranks, threads, ...
 - Analyze load balance information:
 - Min, max, and average values across the ranks and/or threads
 - Look at this information per library as well
 - $_{\odot}\,$ Too much time in MPI could indicate load balance issue.
 - Use above info to determine if the program is well balanced
 - Are the minimum, maximum values widely different? If so:
 - Probably have load imbalance and need to look for the cause of performance lost because of the imbalance.
 - $_{\circ}~$ Not all ranks or threads doing the same amount of work
 - Too much waiting at barriers or synchronous global operations like MPI_Allreduce
pcsamp Default View: NPB: LU



Default Aggregated pcsamp Experiment View



Open | SpeedShop^{*}

How to Analyze the Performance of Parallel Codes 101 - A Tutorial at SC'18

11/11/2018

Load Balance View: NPB: LU



* Load Balance View based on functions (pcsamp)

		Оре	n SpeedShop		×			
<u>F</u> ile <u>T</u> ools <u>H</u> elp								
pc Sampling [1] Process Control Run Cont	➡ Pause ■	Update	MPI library showing up high in the list Max time in rank 255					
Status: Process Loaded: Click or	1 the "Run" button to beg	in the experiment.						
 Stats Panel [1] Manage Mana	eProcessesPanel [1]	ad Balance (min,max,ave) Report: : 256 Ranks: 256 Thuras, 3		View/Displ Functio	a Choice			
Max CPU Time Across Ranks(s)	Rank of Max	Min CPU Time Across Ranks(s)	Rank of Min	Average CPU Time Acro	oss Ralks(s Function (defining location)			
4.980000	18	3.330000	15	4.020742	rhs_ (lu.C.256: rhs.f,5)			
4.550000	255	1.440000	17	2.833555	smpi_net_lookup (libmpich.so.1.0: mp			
1.600000	17	0.850000	254	1.204844	buts_ (lu.C.256: buts.f,4)			
1.530000	17	0.920000	143	1.218828	blts_ (lu.C.256: blts.f,4)			
1.390000	145	0.740000	242	1.059414	jacld_ (lu.C.256: jacld.f,5)			
1.190000	25	0.670000	243	0.918945	jacu_ (lu.C.256: jacu.f,5)			
0.750000	64	0.160000	255	0.476289	ssor_ (lu.C.256: ssor.f,4)			
0.530000	94	0.110.200	0	0.21.568	GI memcpy (libc-2.5.so)			
	189	0.090000	255	0.294453	exchange_3_ (lu.C.256: exchange_3.f,			
0.300000	111	0.050000	55	0.131406	pthread spin lock (libpthread-2.5.so)			
Command Panel				With load b	palance view we are			
openss>>	_			looking for out of norm such large differen	performance number of what is expected, as relatively ces between min ma			
				and/or	average values			

Open | SpeedShop™

Default Linked Object View: NPB: LU

Default Aggregated View based on Linked Objects (libraries)

		Open SpeedShop		×	
File Tools Help	e 5 Update n to begin the experiment.	Linked Object View (library view) Select "Linked Objects" Click D-icon			
Stats Panel [1] ManageProcessesPanel Image: The constraint of the constraint o	[1] ving Linked Objects Report v Pids: 256 Ranks: 256 Threads:	3	View/Display Choice	s O Loops	
% of CPU Time 71.339523 24.325538 1.924810	Exclusive CPU time in seconds. - 2382.790000 - 812.490000 - 64.290000 - 45.700000 - 34.800000	% of CPU Time 71.339523 24.325538 1.924810 1.368235 1.041894	LinkedObject lu.C.256 libmpich.so.1.0 libc-2.5.so mlx4-rdmav2.so libethread-2.5.so		
1.368235 1.041894	[THE SECOND SECOND	NOTE: Look at library time to get the MPI over	the MPI an idea of head.	
© <u>C</u> ommand Panel openss>>					

Parallel Execution Analysis Techniques



If imbalance detected, then what? How do you find the cause?

- Look at library time distribution across all the ranks, threads
 - Is the MPI library taking a disproportionate amount of time?
- If threaded (e.g. OpenMP), then look at the balance of time across worker threads.
 - For OpenMP look at idleness, barrier time, in addition to task times
- If MPI application, use a tool that provides per MPI function call timings
 - Can look at MPI function time distributions
 - In particular, MPI_Waitall
 - Then look at the call path to MPI_Waitall
 - Also, can look source code relative to
 - MPI rank or particular pthread that is involved.
 - Is there any special processing for the particular rank or thread
 - Examine the call paths and check code along path
- > Use Cluster Analysis type feature, if tool has this capability
 - Cluster analysis can categorize threads or ranks that have similar performance into groups identifying the outlier rank or thread

Hot Call Paths View (CLI): NPB: LU

Hot Call Paths for MPI_Wait for rank 255 only



Open | SpeedShop™

77



Identifying Load Imbalance With O|SS

Get overview of application

- > Run a lightweight experiment to verify performance expectations
 - pcsamp, usertime, hwc

Subset Use load balance view on pcsamp, usertime, hwc

- Look for performance values outside of norm
 - Somewhat large difference for the min, max, average values
 - If the MPI libraries are showing up in the load balance for pcsamp, then do an MPI specific experiment

Seload balance view on MPI experiment

- Look for performance values outside of norm
 - Somewhat large difference for the min, max, average values
- Focus on the MPI_Functions to find potential problems

Subset Set In the set of the s

Can also use expcompare across OpenMP threads

Open | SpeedShop™

Link. Obj. Load Balance: Using NPB: LU

* Load Balance View based on Linked Objects (libraries)

Open SpeedShop	_ + X
<u>F</u> ile <u>T</u> ools	<u>H</u> elp
Process Control Ann → Cont → Pause 5 Update Status: Process Loaded: Click on the "Run" button to begin the experiment.	Rank 255 has maximum MPI library time value & minimum LU time
Stats Panel [1] ManageProcessesPanel [1]	БаП ⊟ X
Image: Totel State (1) Image: Totel State (View/Display Choice
May Evolutions, Reals of May, Missionly in Processes/Raiks/ Freedo. (250) 0	
Max Exclusive Kalk of Max Dua Exclusive Kalk of Max Average Exclusive Enkedobject -10.8800 17 7.3700 255 9.3078 lu.C.256 -4.9400 255 1.7600 17 3.1738 libmpich.so.1 -0.5400 94 0.1100 0 0.2511 libc-2.5.so -0.3500 255 0.0500 219 0.1785 libmlx4-rdma -0.3000 190 0.0500 55 0.1359 libpthread-2.3	.0 IV2.so 5.so
⊡ <u>C</u> ommand Panel	
openss>>	

Open SpeedShop[™]

How to Analyze the Performance of Parallel Codes 101 - A Tutorial at SC'18

Using Cluster Analysis in O|SS



Can use with pcsamp, usertime, hwc

- > Will group like performing ranks/threads into groups
- Groups may identify outlier groups of ranks/threads
- > Can examine the performance of a member of the outlier group
- Can compare that member with member of acceptable performing group

Can use with mpi, mpit, mpip

- Same functionality as above
- But, now focuses on the performance of individual MPI_Functions.
- Key functions are MPI_Wait, MPI_WaitAll
- Can look at call paths to the key functions to analyze why they are being called to find performance issues

Link. Obj. Cluster Analysis: NPB: LU



Cluster Analysis View based on Linked Objects (libraries)

		Оре	n SpeedShop		×				
<u>F</u> ile <u>T</u> ools <u>H</u> elp									
▼ pc Sampling [1] Process Control			In	Cluster Ana	lysis results				
➡ Run 🕨	Cont 🌖 Pause 🗗	Update		outlier.					
Status: Process Loaded: O	Click on the "Run" button to begi	n the experiment.							
Stats Panel [1]	ManageProcessesPanel [1]			- View/Displan/hoise					
i i ci i 8 (CA CC Showing Con	nparative analysis Report:		Functions Statements	s 💿 Linked Objects 🔾 Loops				
Executables: lu.C.256 View consists of compariso	on columns click optime metadat	a icon "I" for details.	ost. Inclato.Init.200, Inclatt.Init.2	. Incraitz. Intractor Incraito. Intractor					
for performance data type Column(s) labeled -c 4: E for performance data type Column(s) labeled -c 5: E for performance data type	e: pcsament linkedobjects using xperiment 1 Database Name: In e: pcsamp -v linkedobjects using xperiment 1 Database Name: In e: pcsamp -v linkedobjects using	display option: ThreadAverage -v li a.C.256-pcsamp.openss : Showing H display option: ThreadAverage -v li a.C.256-pcsamp.openss : Showing H display option: ThreadAverage -v li	nkedobjects ost: hera19.1lnl.gov -p 24° 24 -t 469 nkedobjects ost: hera5.1lnl.gov p 29797, 29798 nkedobjects	12529633216 Rank: 255 , 29799, 29800, 29801, 29802, 29	803, 29804, 29805, 29806, 2980				
-c 2, Average CPU Time A	Across Ra -c 3, Average CPU Tin	ne Across Ra -c 4, Average CPU Tip	Across Ra -c 5, Average CPU Tim	e Across Ra LinkedObject 🖉					
- 9.344375	8.247368	7.370000	10.465000	lu.C.256					
- 0.132054	0.184737	0.230000	0.123333	libpthread-2.5.so					
	4.185789	4.940000	2.083333	libmpich.so.1.0					
- 0.176429	0.224211	0.350000	0.130833	libmlx4-rdmav2.so					
····· 0.255714	0.204737	0.160000	0.246667	libc-2.5.so					
Command Panel					⊑ 🗆 ×				
openss>>									

MPI/OpenMP Specific Experiments



MPI specific experiments

Record all MPI call invocations

> MPI functions are profiled (ossmpip)

- Show call paths for each MPI unique call path, but individual call information is not recorded.
- Less overhead than mpi, mpit.
- MPI functions are traced (ossmpi)
 - Record call times and call paths for each event
- > MPI functions are traced with details (ossmpit)
 - Record call times, call paths and argument info for each event

OpenMP specific experiment (ossomptp)

- Uses OMPT API to record task time, idleness, barrier, and wait barrier per OpenMP parallel region
 - Shows load balance for time
 - expcompare time across all threads

MPI Tracing Results: Default View



* Default Aggregated MPI Experiment View

	OpenISpeed	Shon		_ + X						
<u>File Tools</u>	Information	lcon		<u>H</u> elp						
MPI [1] Process Control	Displays Expe	eriment a								
Run Cont Pause	🗗 Update			Terminate						
Status: Process Leaded: Click on the "Run" button to begin the experiment.										
Stat Panel [1] ManageProcessesPanel [1] Aggregated Results lay Choice -										
Image: Comparison of the second se										
Metadata for Experiment 1: Application command: Executables: smg2000 Experiment type: mpi Host(s): hyperion583.llnl.gov hyperion584.llnl.gov hyperion585.llnl.gov hyperion586.llnl.gov hyperion588.llnl.gov hyperion588										
Minimum MPI Call Time(ms) Maximum MPI	Call Time(ms) Average Time(ms	Number of Calls	s Function (defining location)							
-555.306000 1276.275000 -151.147000 167.504000 -0.152000 0.474000 -0.043000 0.212000 -0.031000 2.034000 -0.013000 10.322000 -0.000001 611.617000 -0.000001 0.600000 -0.000001 0.069000	755.289027 163.231094 0.334205 0.133098 1.312102 0.717578 0.977852 0.001156 0.000665	512 512 512 512 512 6144 4667648 5403936 5403936	PMPI_Init (libmonitor.so.0.0: pmpi.c,94) PMPI_Finalize (libmonitor.so.0.0: pmpi.c,223) MPI_Allgatherv (libmpich.so.1.0: allgatherv.c,73) MPI_Allgather (libmpich.so.1.0: allgather.c,70) MPI_Barrier (libmpich.so.1.0: barrier.c,56) MPI_Allreduce (libmpich.so.1.0: allreduce.c,59) MPI_Waitall (libmpich.so.1.0: waitall.c,57) MPI_Isend (libmpich.so.1.0: isend.c,58) MPI_Irecv (libmpich.so.1.0: irecv.c,48)							
				М						

Open | SpeedShop™

Using OMPTP experiment in O|SS



- The following three CLI examples show the most important ways to view OMPTP experiment data.
- Default view shows the timing of the parallel regions, idle, barrier, and wait barrier as an aggregate across all threads

openss -cli -f ./matmult-omptp-0.openss openss>>expview

Exclusive	Inclusive	% of Funct	ion (defining location)
times in	times in 7	Fotal	
seconds.	seconds. E	xclusive	
	CPU Time	•	
44.638794	45.255843	93.499987	computeomp_fn.1 (matmult: matmult.c,68)
1.744841	1.775104	3.654726	compute_interchangeomp_fn.3 (matmult: matmult.c,118)
0.701720	0.701726	1.469817	compute_triangularomp_fn.2 (matmult: matmult.c,95)
0.652438	0.652438	1.366591	IDLE (omptp-collector-monitor-mrnet.so: collector.c,573)
0.004206	0.009359	0.008810	initializeomp_fn.0 (matmult: matmult.c,32)
0.000032	0.000032	0.000068	BARRIER (omptp-collector-monitor-mrnet.so: collector.c,587)
0.000000	0.000000	0.000001	WAIT_BARRIER (omptp-collector-monitor-mrnet.so: collector.c,602)

Using OMPTP experiment in O|SS



This example shows the comparison of exclusive time across all threads for the parallel regions, idle, barrier, and wait barrier

openss>>expcompare -mtime -t0:4

-t 2, -t 3, -t 4, Function (defining location) -t 0. **Exclusive Exclusive Exclusive** times in times in times in times in seconds. seconds. seconds. seconds. 11.313892 11.081346 11.313889 10.929668 compute._omp_fn.1 (matmult: matmult.c,68) 0.443713 0.430553 0.429635 0.440940 compute interchange. omp fn.3 (matmult: matmult.c, 118) 0.069975 compute triangular. omp fn.2 (matmult: matmult.c,95) 0.253632 0.213238 0.164875 0.001047 0.001100 0.001095 0.000964 initialize. omp fn.0 (matmult: matmult.c,32) 0.000008 0.000008 0.000006 0.000010 BARRIER (omptp-collector-monitor-mrnet.so: collector.c,587) 0.000000 0.000000 0.000000 0.000000 WAIT BARRIER (omptp-collector-monitor-mrnet.so: collector.c,602) 0.388890 IDLE (omptp-collector-monitor-mrnet.so: collector.c,573) 0.000000 0.247592 0.015956

Using OMPTP experiment in O|SS



This example shows the load balance of time across all threads for the parallel regions, idle, barrier, and wait barrier

openss>>expview -mloadbalance

Max Fxclusive	OpenN Thread	Mp Min IId Fxclusive	Open Threa	Mp Average Function (defining location)
Time Across	of Ma	ax Time Acro	ss of N	Vin Time Across
OpenMp		OpenMp		
ThreadIds(s))	ThreadIds(s)		ThreadIds(s)
11.313892	0	10.929668	4	11.159699 computeomp_fn.1 (matmult: matmult.c,68)
0.443713	0	0.429635	3	0.436210 compute_interchangeomp_fn.3 (matmult: matmult.c,118)
0.388890	4	0.015956	3	0.217479 IDLE (omptp-collector-monitor-mrnet.so:
collector.c,573	3)			
0.253632	0	0.069975	4	0.175430 compute_triangularomp_fn.2 (matmult: matmult.c,95)
0.001100	2	0.000964	4	0.001052 initializeomp_fn.0 (matmult: matmult.c,32)
0.000010	4	0.000006	3	0.000008 BARRIER (omptp-collector-monitor-mrnet.so:
collector.c,5	87)			
0.000000	0	0.000000	0	0.000000 WAIT_BARRIER (omptp-collector-monitor-mrnet.so:
collector.c,6	02)			

Open | SpeedShop™

How to Analyze the Performance of Parallel Codes 101 - A Tutorial at SC'18

Summary / Parallel Bottlenecks



Open | SpeedShop supports MPI, OpenMP, and threaded applications (including hybrid)

> Works with multiple MPI implementations

* Parallel experiments

- > Apply the sequential O|SS collectors to all nodes
- Specialized MPI profiling and tracing experiments
- Specialized OpenMP profiling experiment

* Result Viewing

- Results are aggregated across ranks/processes/threads
- > Optionally: select individual ranks/threads or groups
- Specialized views:
 - Load balance view
 - Cluster analysis

Solution State Sections of Problem Code

Open | SpeedShop™

87



- Parallel related application exercise (MPI)
 - More information at the tutorial

* Exercises are in the exercise directory:

- \$HOME/exercises/mpi_nbody
- > Supplemental:
 - \$HOME/exercises/smg2000/test

Consult README file in each of the directories for the instructions/guidance

Open | SpeedShop™

How to Analyze the Performance of Parallel Codes 101 - A Tutorial at SC'18



- * Parallel related parallel application exercise (threading)
- ***** Exercises are in the exercise directory:
 - \$HOME/exercises/matmul
 - > Supplemental:
 - \$HOME/exercises/lulesh2.0.3
- Parallel related application exercise (MPI)
- ***** Exercises are in the exercise directory:
 - \$HOME/exercises/mpi_nbody
- Consult README file in each of the directories for the instructions/guidance







SC2018 Tutorial

How to Analyze the Performance of Parallel Codes 101 A case study with Open/SpeedShop

Section 6

Advanced analysis: Hardware Counter Experiments











Open | SpeedShop™



* Timing information shows where you spend your time

- > Hot functions / statements / libraries
- Hot call paths

BUT: It doesn't show you why

- > Are the computationally intensive parts efficient?
- > Are the processor architectural components working optimally?

Answer can be very platform dependent

- Bottlenecks may differ
- Cause of missing performance portability
- Need to tune to architectural parameters

Next: Investigate hardware/application interaction

- Efficient use of hardware resources or Micro-architectural tuning
- > Architectural units (on/off chip) that are stressed

Open | SpeedShop™

How to Analyze the Performance of Parallel Codes 101 - A Tutorial at SC'18

91



Modern memory systems are complex

- Deep hierarchies
- Explicitly managed memory
- NUMA behavior
- Streaming/Prefetching

Data Location	Access Latency, ns (Sandy Bridge, 2.6GHZ)
L1	1.2
L2	3.5
L3	6.5
DRAM	28

* Key to performance: Data locality and Concurrency

- Accessing the same data repeatedly(Temporal)
- Accessing neighboring data(Spatial)
- Effective/parallel use of cores

Information to look for

- Load/Store Latencies
- Prefetch efficiency
- Cache miss rate at all levels
- TLB miss rates
- NUMA overheads



11/11/2018

Open | SpeedShop^{*}



* Newer processors have wide vector registers

- Intel Xeon 2670,Sandy Bridge: 256 bits floating point registers, AVX (8 Real / 4 Double)
- Intel Xeon Phi, Knights Corner: 512 bits (16 Real / 8 Double)
- Intel Haswell 256 bits Integer Registers, AVX2 : FMA (2X the peak flops)

✤ Key to performance: Vectorization

- Compiler Vectorization
- Use of 'intrinsics'
- > Use of Pragmas to help the compiler
- Assembly code

Analysis Options

- Compiler vectorization report
- > Look at assembly code
- Measure performance with PAPI counters



Going from Scalar to Intel[®] AVX can provide up to 8x faster performance

Open | SpeedShop™

How to Analyze the Performance of Parallel Codes 101 - A Tutorial at SC'18

Hardware Performance Counters

* Architectural Features

- > Typically/Mostly packaged inside the CPU
- Count hardware events transparently without overhead

Newer platforms also provide system counters

- Network cards and switches
- Environmental sensors

Drawbacks

- > Availability differs between platform & processors
- Slight semantic differences between platforms
- > In some cases : requires privileged access & kernel patches

Recommended: Access through PAPI

- > API for tools + simple runtime tools
- > Abstractions for system specific layers
- More information: http://icl.cs.utk.edu/papi/

Open | SpeedShop™

How to Analyze the Performance of Parallel Codes 101 - A Tutorial at SC'18

94

The O|SS HWC Experiments



Provides access to hardware counters

- Implemented on top of PAPI
- Access to PAPI and native counters
- Examples: cache misses, TLB misses, bus accesses

Sasic model 1: Timer Based Sampling: HWCsamp

- Samples at set sampling rate for the chosen event
- Supports multiple counters
- Lower statistical accuracy
- Can be used to estimate good threshold for hwc/hwctime

Sasic model 2: Thresholding: HWC and HWCtime

- User selects one counter
- Run until a fixed number of events have been reached
- Take PC sample at that location
 - HWCtime also records stacktrace
- Reset number of events
- Ideal number of events (threshold) depends on application

Examples of Typical Counters (Xeon E5-2670)



PAPI Name	Description	Threshold
PAPI_L1_DCM	L1 data cache misses	high
PAPI_L2_DCM	L2 data cache misses	high/medium
PAPI_L3_TCM	L3 cache misses	high
PAPI_TOT_INS	Instructions completed	high
PAPI_STL_ICY	Cycles with no instruction issue	high/medium
PAPI_BR_MSP	Miss-predicted branches	medium/low
PAPI_DP_OPS	Number of 64-Bit floating point Vector OPS	high
PAPI_LD_INS	Number of load instructions	high
PAPI_VEC_DP	Number of vector/SIMD instructions – 64Bit	high
PAPI_BR_INS	Number of branch instructions	low
PAPI_TLB_TL	Number of TLB misses	low

Note: Threshold indications are just rough guidance and depend on the application. Note: counters platform dependent (use papi_avail& papi_native_avail)

Open | SpeedShop[™]

How to Analyze the Performance of Parallel Codes 101 - A Tutorial at SC'18

Suggestions to Manage Complexity



- The number of PAPI counters and their use can be overwhelming; Some guidance here with a few "Metric-Ratios".
 - Ratios derived from a combination of hardware events can sometimes provide more useful information than raw metrics
- Develop the ability to interpret Metric-Ratios with a focus on understanding:
 - > Instructions per cycle or cycles per instruction
 - Floating point / Vectorization efficiency
 - Cache behaviors; Long latency instruction impact
 - Branch mispredictions
 - Memory and resource access patterns
 - Pipeline stalls

This presentation will illustrate with some examples of the use of Metric-Ratios

Open | SpeedShop™

How to use OSS HWCsamp experiment

- osshwcsamp "<command>< args>" [default |<PAPI event list>|<sampling rate>]
 - Sequential job example:
 - osshwcsamp "smg2000"
 - Parallel job example:
 - osshwcsamp "mpirun –np 128 smg2000 –n 50 50 50" PAPI_L1_DCM,PAPI_L1_TCA 50
- * default events: PAPI_TOT_CYC and PAPI_TOT_INS
- * default sampling_rate: 100
- <PAPI_event_list>: Comma separated PAPI event list
 (Maximum of 6 events that can be combined)
- <sampling_rate>:Integer value sampling rate
- Use event count values to guide selection of thresholds for HWC, HWCtime experiments for deeper analysis

Selecting the Counters & Sampling Rate



For osshwcsamp, Open | SpeedShop supports ...

- Derived and Non derived PAPI presets
 - All derived and non derived events reported by "papi_avail"
 - Also reported by running "osshwcsamp" with no arguments
 - Ability to sample up to six (6) counters at one time; before use test with – papi_event_chooser PRESET <list of events>
 - If a counter does not appear in the output, there may be a conflict in the hardware counters
- All native events
 - Architecture specific (incl. naming)
 - Names listed in the PAPI documentation
 - Native events reported by "papi_native_avail"

* Sampling rate depends on application

- > Overhead vs. Accuracy
 - Lower sampling rate causes less samples

Open | SpeedShop™

Useful Metric-Ratio 1: IPC



- Instructions Per Cycle(IPC) also referred to as Computational Intensity
 - IPC= PAPI_TOT_INS/PAPI_TOT_CYCLES
- ✤ Data from single-core Xeon E5-2670, Sandy Bridge
- In the table below compiler optimization -O1 used to bring out differences in IPC based on stride used with different loop order;
- If you use -O2 for this simple case compiler does the right transformations, permuting loop order and vectorizing to yield IPC = 3.594 (jki order); This improves access to memory through cache.
- Importance of stride through the data is illustrated with this simple example; Compiler may not always do the needed optimization. Use IPC values from functions and loops to understand efficiency of data access through your data structures.

- Example matrix multiply;
 Triple do loop;
 (n1=n2=n3=1000)
- code for loop order 'ijk'; All vectors 'double'

```
do i = 1, n1
do j = 1, n3
do k = 1, n2
a(i,j) = a(i,j) + b(i,k) * c(k,j)
end do
end do
end do
```

Metric	IJK	IKJ	JIK	JKI	KIJ	KJI	MATMUL	DGEMM
PAPI_TOT_INS	8.012E+09	9.011E+09	8.011E+09	9.01E+09	9.01E+09	9.011E+09	9.016E+09	7.405E+08
PAPI_TOT_CYC	2.42E+10	5.615E+10	2.423E+10	2.507E+09	5.612E+10	2.61E+09	2.601E+09	2.859E+08
IPC	0.331	0.160	0.331	3.594	0.161	3.452	3.466	2.590
MFLOPS	272	117	271	2625	117	2525	2532	19233 (93% peak)



Level	Operation	# Memory Refs or Ops	# Flops	Flops/Ops	Comments on Flops/Ops
1	y = k x + y	3n	2n	2/3	Achieved in Benchmarks
2	y = Ax + y	n ²	2n ²	2	Achieved in Benchmarks
3	C = AB + C	4n ²	2n ³	n/2	Exceeds HW MAX

Use these Flops/Ops to understand how sections of your code relate to simple memory access patterns as typified by these BLAS operations

Open | SpeedShop[™]

How to Analyze the Performance of Parallel Codes 101 - A Tutorial at SC'18 11/11/2018

Useful Metric-Ratio 2: FloatOps/Cycle



Traditionally PAPI_FP_INS/PAPI_TOT_CYC used to evaluate relative floating point density

- For a number of reasons measuring and analyzing floating point performance on Intel Sandy Bridge and Ivy bridge must be done with care. See PAPI web site for full discussion. The reasons are: instruction mix - scalar instructions + vector (AVX, SSE) packed instructions, hyperthreading, turbo-mode and speculative execution.
- > The floating point counters have been disabled in the newer Intel Haswell cpu architecture
- > On Sandy Bridge and Ivy Bridge PAPI_FP_INS is no longer an appropriate counter if loops are vectorized
- No single PAPI metric captures all floating point operations
- We provide some guidance with useful PAPI Preset counters. Data from single-core Xeon E5-2670, Sandy Bridge. Double precision array operations for Blas1(daxpy), Blas2(dgemv) and Blas3(dgemm) are benchmarked. Matrix size=nxn; vector size=nx1. Data array sizes are picked to force operations from DRAM memory
- Table below shows measured PAPI counter data for a few counters and compares the measured FLOP/Ops against theoretical expectations.
- PAPI_DP_OPS and PAPI_VEC_DP give similar values and these counter values correlate well with expected floating point operation counts for double precision.

Blas Operation	n	Thererical mem refs or Ops	Theoretical FLOP	Theoretical FLOP/Ops	wall time, secs	тот сус	TOT INS	FP INS	LD INS	SR INS	DP OPS	PAPI GFLOPS	PAPI FLOP/Ops
dayay				0.67	0.03	1 0/5+09	- 5 205 07	11 50	2 505+07	1 255,07	- E 01E+07	1 56	0.668
dgemy	1.00F+04	1.0F+08	2.0F+08	2	0.06073	2.16F+08	1.69F+08	29.12	6.25F+07	1.25E+07	2.36F+08	3.89	1.57557985
dgemm	1.00E+04	4.00E+08	2E+12	5000.00	80.937	2.67E+11	7.33E+11	7.2	1.12E+11	1.38E+09	2.01E+12	24.80	8.83518225

Open | SpeedShop™

How to Analyze the Performance of Parallel Codes 101 - A Tutorial at SC'18

102



- We again provide some guidance with data from a single-core of a Haswell Processor (Intel(R) Xeon(R) CPU E5-2698 v3 @ 2.30GHz)
- Blas1, Blas2 and Blas3 kernels as in the previous slide are benchmarked. Matrix size=nxn; vector size=nx1. Data array sizes are picked to force operations from DRAM memory
- Table below shows measured PAPI counter data for a few counters and metric ratio IPC
- ***** When operating at peak performance, Haswell can retire 4 micro-ops/cycle

	Thererical mem refs	Theoretica	Theoretical	wall time,									
n	or Ops	I FLOP	FLOP/Ops	secs	TOT_CYC	TOT_INS	IPC	CPI		LD_INS	SR_INS	GFLOPS	FLOP/mem-Ops
2 50F+07	7 50F+07	5 00F+07	0.67	3 24F-02	1 17F+08	6 25F+07	0.5	1	1 87	3 13F+07	1 25F+07	1 53932	0.57
2.502.07	7.502107	J.002+07	0.07	J.24L 02	1.172.00	0.232107	0.5		1.07	5.152.07	1.232.07	1.55552	0.57
1.00E+04	1.00E+08	2.00E+08	2	6.11E-02	2.2E+08	2.06E+08	0.9	4	1.06	7.81E+07	1.25E+07	3.272	1.10
1.00E+04	4.00E+08	2.00E+12	5000	41.8546	1.38E+11	4.65E+11	3.3	5	0.30	1.9E+11	1.23E+09	47.7655	5.23

hwcsamp with miniFE (see mantevo.org)



openss –f miniFE.x-hwcsamp.openss

			Open SpeedShop	o (on uno-login1)				_ _ _ X
<u>F</u> ile <u>T</u> ools								<u>H</u> elp
HWCSamp Pa	Also have po	csamp			Up to six event can be			
Process Control	informati		displayed Here we have 4					
⇒Run i⇒C	momu	UII		u	ispiayed		C 114 VC 1.	Terminate
Status: Process Loaded: Click o	n the "Run" butt in to be, in the ex	cperiment.			1			
Stats Panel [1] ■ ManageF	ProcessesPane [1]							∿: □ = ×
						View/Display C	Ch <mark>r</mark> ce	
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	CA CC Showing Functions Repo	t.				♦ Functions	Statements 💠 Linked C	bjects 💠 Loops
Executables: /gpfs1/mrajan/oss_	sc14_miniFE-2.0_mkl/src	e/m -E.x Hosts:(9	9) uno159 Pids: 7	2 Ranks: 72 T	ads: 1			
% of CPU Time	Exclusive CPU time in seco	% of CPU Time	papi_dp_ops	papi_l1_dcm	papi_tot_cyc	papi_tot_ins	Function (defining I	ocation)
47.873611	-2253.430000	47.873611	2824672149050	291551937219	6989975219262	8952967031967	mkl_spblas_lp64_d	csr0ng_c_mv
	-447.660000	9.510435	87267523822	7738661859	903705759149	1039254165063	GI_epoll_wait (/lil	b64/libc-2.12.so
<mark>9.51</mark> 0435	-421.040000	8.944900	831606939512	3626363789	1341678035775	2311136732678	void miniFE::perforr	n_element_loop
	-374.450000	7.955106	472834291021	30726897002	1094184657537	532891650221	mkl_blas_daxpby (/	opt/intel-12.1/mk
8 944900	-241.230000	5.124877		4361118141	768295120029	1170840620681	MatrixInitOp <minif8< th=""><th>E::CSRMatrix<do< th=""></do<></th></minif8<>	E::CSRMatrix <do< th=""></do<>
	-230.740000	4.902019	456917366559	1985851369	735787515693	1267123401365	void miniFE::Hex8::	diffusionMatrix_
	-202.540000	4.302916	1647949388	589250201	645977293841	875225391031	void miniFE::impose	e_dirichlet <minil< th=""></minil<>
<mark>7.9</mark> 55106	-127.260000	2.703610	184233405196	15545928811	400531304861	310351177239	mkl_blas_xddot (/op	ot/intel-12.1/mkl/l
	-78.900000	1.676213	155934713579	713587396	251575572613	432357079682	void miniFE::get_el	em_nodes_and
<mark>5.1</mark> 24877	-46.800000	0.994255	8479608341	748205444	94333323140	109592763431	opal_progress (/opt	openmpi-1.6-int
	-42.960000	0.912675	8832226710	776687330	87613448245	100107292999	opal_event_base_lo	oop (/opt/openm
other	-37.680000	0.800503	17557210051	477783784	120071360670	194746931361	std::_Rb_tree_decre	ement(std::_Rb_
	-32.640000	0.693429	6349358453	564732359	66230897286	76323494900	epoll_dispatch (/opt	/openmpi-1.6-int
		0 529632		2/100102	216/2195/97	2910731701	etd:wantor>doubla	etd∵allocator>dc⊠ ⊳
□ <u>C</u> ommand Panel								5: □ = ×
openss>>								



Viewing hwcsamp Data in CLI



openss -cli -f miniFE.x-hwcsamp.openss

mrajan@uno-login1:/gpfs1/mrajan/oss_sc14_miniFE/miniFE-2.0_mkl/run_oss_hwcsamp										
<u>F</u> ile <u>E</u> dit <u>V</u> iew <u>S</u> earch <u>T</u> erminal <u>H</u> elp										
[mrajan@uno	-login1 run	oss hwcsamp]\$ o	penss -cli -f	miniFE.x-hwcsa	mp.openss		^			
openss>>[ope	enss]: The ī	estored experim	ent identifie	ris: -x 1						
openss>>exp	openss>>expview -v statements hwcsamp30									
Exclusive	% of CPU	papi_dp_ops	papi_l1_dcm	papi_tot_cyc	papi_tot_ins	Statement Location (Line Number)				
CPU time in	Time									
seconds.	22 604112	07202006600	7741240402	000000000000000000000000000000000000000	1000416470601					
447.770000	23.604112	8/303880000	7741340402	903892322030	10394164/0601	/usr/src/debug///////glibc-2.12-2-gC4cctt1//misc//sysdeps/unix/syscall-template.s(82	.)			
327.880000	10 522474	048522/30/02	202/0/33//	1040430490144	180000000381304	/gpfs1/mrajan/oss_sc14_miniFe/miniFe/2.0 mkl/csc/./perform_etement_toop.npp(100)				
199.020000	0 514/4	046245017	512592005	575701/27606	761/02150152	/gpfs1/mrajal/oss_sc14_minife/minife2.0%mk1/csc/./Matrixinitop.hpp1400/ /gpfs1/mrajan/csc/sc14_minife/_0/mk1/csc//SparseMatrix/functions_bop(475)				
72 920000	3 8/396/	1//552896087	6225/0031	232540036283	/01402139132	/gpfs/mrajai/oss_sci4_minife/minife2.0_mkl/crc//jparsematif/_lunctions.npp(4/5/				
67.220000	3.543490	133137191569	576640758	214346202757	369021735348	/gp/s1/mrajan/oss_sc14_miniFE/miniFE-2.0_mkl/src/_/fem/Hex8_hon(302)				
57.890000	3.051661	114320520328	536923044	184576819581	316915984454	/gp/si/mrajan/oss_sc14_miniFE/miniFE-2.0_mkl/src/./Hex8_box_utils.hpp(136)				
35.630000	1.878229	11.020020020	12642935	27982648181	11492892526	/usr/include/c++/4.4.7/bits/vector.tcc(397)				
25.740000	1.356879	50961487331	222435633	82098246515	141360199733	/gpfs1/mraian/oss sc14 miniFE/miniFE-2.0 mkl/src//fem/Hex8.hpp(335)				
23.720000	1.250395	46993941906	201176970	75638030490	130325024838	/gpfs1/mrajan/oss_sc14_miniFE/miniFE-2.0_mkl/src//fem/Hex8.hpp(330)				
20.090000	1.059041	39815708667	172829164	64060100337	110338757571	/gpfs1/mrajan/oss_sc14_miniFE/miniFE-2.0_mkl/src//fem/Hex8.hpp(260)				
19.840000	1.045862	3395917479	296952208	39763200164	46443969250	/builddir/build/BUILD/openmpi-1.6.4/opal/runtime/opal_progress.c(185)				
18.520000	0.976278	36633304484	159683865	59047484493	101653749874	/gpfs1/mrajan/oss sc14 miniFE/miniFE-2.0 mkl/src/./perform element loop.hpp(82)				
18.150000	0.956774	35947896420	158886460	57884207500	99655374262	/gpfs1/mrajan/oss_sc14_miniFE/miniFE-2.0_mkl/src//fem/Hex8.hpp(338)				
16.110000	0.849236	31863068315	140255815	51366606737	88465155762	/gpfs1/mrajan/oss_sc14_miniFE/miniFE-2.0_mkl/src//fem/Hex8.hpp(339)				
12.120000	0.638904	2140016350	192110608	24394027282	28522826298	/builddir/build/BUILD/openmpi-1.6.4/opal/event/epoll.c(279)				
11.820000	0.623089		162934927	37648219524	58301814345	/gpfs1/mrajan/oss_sc14_miniFE/miniFE-2.0_mkl/src/./MatrixInitOp.hpp(137)				
11.280000	0.594623	22319927660	97114952	35960407509	61927707941	/gpfs1/mrajan/oss_sc14_miniFE/miniFE-2.0_mkl/src//fem/Hex8.hpp(336)				
11.170000	0.588824	197905093	34646748	35630595444	54672411406	/gpfs1/mrajan/oss_sc14_miniFE/miniFE-2.0_mkl/src/./SparseMatrix_functions.hpp(466)				
10.560000	0.556668	14511257762	964669997	32659133238	12804049091	/gpfs1/mrajan/oss_sc14_miniFE/miniFE-2.0_mkl/src/./cg_solve.hpp(179)				
9.390000	0.494992	1960742963	175257744	19206353960	21857221066	/bullddir/bulld/BULD/openmpi-1.6.4/opal/event/epoll.c(215)				
9.290000	0.489721	2017022172	255848742	27797129336	63963403261	/gpts1/mrajan/oss_scl4_min1FE/min1FE-2.0_mk(/src//Utils/box_utils.npp(297)				
9.270000	0.488000	201/0331/2	1/5//1/45	19021518723	21408/10322	/builddir/build/Build/build/openmpi-1.0.4/opa/veveni.c(81)				
7 020000	0.435424	1500756010	145174100	20201012341	41307700004	/gpist/midjal/oss_sci4_minire/_minire/2.0 min//sic/./matrixinic0p.npp142)				
7.930000	0.410020	1390/30210	61254855	2/027/01228	50552268201	/anfs//maine/oss.cc/d/miniE//anfs/				
7 780000	0.410121		36998244	24551069808	59333194907	/gprs/mrajan/oss_scl4_minife/minife2.0_mkl/src//generate_matrix_structure.hpp(10)				
7,210000	0.380074	14227594738	65160484	22988528053	39617523842	/gp/s1/mrajan/oss_scl4_miniF/miniF2.0_mkl/src/_/fm/Hex8_hpn(300)				
6.430000	0.338956	12745141419	53417855	20499184114	35356808102	/gp/s1/mrajan/oss_sc14_minifE/minifE-2.0_mkl/src/./Hex8_hox_utils_hop(147)				
6.210000	0.327359	1319041811	117512799	12639734182	14411185863	/build/ir/build/BUILD/openmoi-1.6.4/opa/event/event.c(838)				
1668.380000	87.948339	1447432312151	20626804289	4631125982124	6974582554912	Report Summary				
openss>>						······································				
l										



Selections of CLI commands used to view the data:

- expview -v linkedobjects
- expview –m loadbalance
- * expview -v statements hwcsamp<number>
 - Example to show top 10 statements:
 - expview –v statements hwcsamp10
- * expview -v calltrees,fullstack usertime<number>
- * expcompare r 1 –r 2 –m time (compares rank 1 to rank 2 for metric equal time)

Deeper Analysis with HWC and HWCtime



- osshwc[time] "<command> < args>" [default |
 <PAPI_event> | <PAPI threshold> | <PAPI_event><PAPI
 threshold>]
 - Sequential job example:
 - osshwc[time] "smg2000 –n 50 50 50" PAPI_FP_OPS 50000
 - Parallel job example:
 - osshwc[time] "mpirun –np 128 smg2000 –n 50 50 50"
- default: event (PAPI_TOT_CYC), threshold (10000000)
- <PAPI_event>: PAPI event name
- PAPI threshold>: PAPI integer threshold
- NOTE: If the output is empty, try lowering the <threshold> value. There may not have been enough PAPI event occurrences to record and present

Viewing hwc Data



* hwc default view: Counter = Instruction Cache Misses

Open SpeedShop ×									
ile Tools Help									
HW Counter [1] Process Control	5 Update		Flat hardware counter profile of a single hardware counter event.						
Executables: sweep3d.mpi Hosts:(16) ys1328 Pids:	Functions Report: 64 Ranks: 64 Threads: 64		View/Display Choice Functions O Statements O Linked Objects O Loops						
% of Total PAPI_L1_ICM Counts	Exclusive PAPI_L1_ICM Counts	% of Total PAPI_L1_ICM Coun	ts Function (defining location)						
15.819540 13.946259 6.301776 4.055792 2.918152 other	- 97200000 - 85690000 - 38720000 - 24920000 - 17930000 - 16870000 - 16170000 - 15040000 - 13150000	15.819540 13.946259 6.301776 4.055792 2.918152 2.745634 2.631707 2.447797 2.140195	sweep (sweep3d.mpi: sweep.f,2) _lapi_dispatcher <false> (libpami.so) udp_read_callback (libpamiudp.so) memcpy_ssse3 (libc-2.12.so) intel_ssse3_rep_memcpy (libirc.so) MPI_Recv (libmpich.so.3.3) PAMI::Interface::Context<pami::context>::advance (libpami.so: ContextInterface.h,158) Sam::SendContig (libpami.so) PMPI_Send (libmpich.so.3.3)</pami::context></false>						
	···· 10410000	1.694253	MPID Recv (libmpich.so.3.3)						
Command Panel									
openss>>									

Open | SpeedShop[™]
Viewing hwctime Data



hwctime default view: Counter = L1 Data Cache Misses

		Open Speed	IShop		×
File Tools Help ▼ HWCTime Panel [1] Process Control → Run For any Pause Status: Process Loaded: Click on the "Run" button to	Update		Calling con counter pro hardware c Exclusive/Ir	G □ ► ×	
Stats Panel [1] ▼ ManageProcessesPanel [1] T T CL D S C G HC B TS ON Executables: smg2000 Host: localhost Pids: 4 Ranks:	4 Threads: 4	inctions Report:		View/Display Choice	ta □ = × d Objects ♀ Loops
% of Total Exclusive PAPI_L1_DCM Counts 53.208556 30.695187 2.299465 2.139037 1.818182 other	Exclusive PAPI_L1_DCM 	Inclusive PAPI_L1_DC 77400000 501750000 3300000 3000000 29250000 3000000 1350000 9750000 8250000 12750000	M C 115 % of Total Exclusive 53.208556 30.695187 2.299465 2.139037 1.818182 1.711230 0.962567 0.695187 0.588235 0.481283	PAPI_L1_DC Function (defining location) hypre_SMGResidual (smg2000: sr hypre_CyclicReduction (smg2000) hypre_SemiRestrict (smg2000: sem unpack_predefined_data (liboper pack_predefined_data (liboper-p hypre_StructAxpy (smg2000: stru memcpy_ssse3_back (libc-2.17 hypre_SMGAxpy (smg2000: smg_ hypre_CycRedSetupCoarseOp (sm	mg_residual.c,152) : cyclic_reduction.c,; mi_restrict.c,125) i_interp.c,126) -pal.so.6.2.0: opal_c al.so.6.2.0: opal_dat ct_axpy.c,25) .so) axpy.c,27) g2000: cyclic_reduct
<u>Command Panel</u> openss>>					



Major reasons on-node scaling limitations *

- Memory Bandwidth \triangleright
- Shared L3 Cache \triangleright

L3 cache miss for 1,2,4 Pes matches * expectation for strong scaling

- Reduced data per PE \triangleright
- L3 misses decreasing up to 4 PEs linearly. \geq







- On the other hand L3 Evictions for 1,2,4 PEs * similarly decrease 'near-perfect' but dramatically increases to 100x at 8PEs and 170x at 16 PEs
- L3 evictions are a good measure of memory * bandwidth limited performance bottleneck at a node
- **General Memory BW limitation Remedies** **
 - Blocking \triangleright
 - Remove false sharing for threaded codes

AMG Intra Node Scaling

How to Analyze the Performance of Parallel Codes 101 - A Tutorial at SC'18

L3 EVICTIONS:ALL



```
! Cache line UnAligned
real*4, dimension(100,100)::c,d
!$OMP PARALLEL DO
do i=1,100
  do j=2, 100
     c(i,j) = c(i, j-1) + d(i,j)
     enddo
enddo
!$OMP END PARALLEL DO
```

! Cache line Aligned real*4, dimension(112,100)::c,d !OMP DO SCHEDULE(STATIC, 16) do i=1,100 do j=2, 100 c(i,j) = c(i, j-1) + d(i,j) enddo enddo !OMP END DO

Same computation, but careful attention to alignment and independent OMP parallel cache-line chunks can have big impact; L3_EVICTIONS a good measure;

	Run Time	L3_EVICTIONS:ALL	L3_EVICTIONS:MODIFIED
Aligned	6.5e-03	9	3
UnAligned	2.4e-02	1583	1422
Perf. Penalty	3.7	175	474

Open | SpeedShop[™]

How to Analyze the Performance of Parallel Codes 101 - A Tutorial at SC'18



113

- * Hardware counter experiments related exercises
- ***** Exercises are in the exercise directory:
 - \$HOME/exercises/soa_aos
 - > \$HOME/exercises/matrix_multiply
 - Supplemental exercises:
 - \$HOME/exercises/HPCCG-0.5
 - \$HOME/exercises/HPCCG-0.5_from_snl (no run just view)

Consult README file in each of the directories for the instructions/guidance

Open | SpeedShop^{**}

How to Analyze the Performance of Parallel Codes 101 - A Tutorial at SC'18 11/11/2018







SC2018 Tutorial

How to Analyze the Performance of Parallel Codes 101 A case study with Open/SpeedShop

Section 7 Analysis of I/O











11/11/2018

Open | SpeedShop™

Need for Understanding I/O



- * I/O could be significant percentage of execution time dependent upon:
 - > Checkpoint, analysis output, visualization & I/O frequencies
 - I/O pattern in the application:
 N-to-1, N-to-N; simultaneous writes or requests
 - Nature of application: data intensive, traditional HPC, out-of-core
 - File system and Striping: NFS, Lustre, Panasas, and # of Object Storage Targets (OSTs)
 - > I/O libraries: MPI-IO, hdf5, PLFS,...
 - Other jobs stressing the I/O sub-systems

***** Obvious candidates to explore first while tuning:

- > Use parallel file system
- > Optimize for I/O pattern
- Match checkpoint I/O frequency to MTBI of the system
- Use appropriate libraries

I/O Performance Example

Application: OOCORE benchmark from DOD HPCMO

- Out-of-core SCALAPACK benchmark from UTK
- Can be configured to be disk I/O intensive
- Characterizes a very important class of HPC application involving the use of Method of Moments (MOM) formulation for investigating electromagnetics (e.g. Radar Cross Section, Antenna design)
- Solves dense matrix equations by LU, QR or Cholesky factorization
- *"Benchmarking OOCORE, an Out-of-Core Matrix Solver,"* Cable, S.B., D'Avezedo, E. SCALAPACK Team, University of Tennessee at Knoxville/U.S. Army Engineering and Development Center



- Sed by HPCMO to evaluate I/O system scalability
- Out-of-core dense solver benchmarks demonstrate the importance of the following in performance analysis:
 - I/O overhead minimization
 - Matrix Multiply kernel possible to achieve close to peak performance of the machine if tuned well
 - "Blocking" very important to tune for deep memory hierarchies

Use O|SS to measure and tune for I/O



INPUT: testdriver.in

ScaLAPACK out-of-core LU,QR,LL
factorization input file

testdriver.out

6	device out
1	number of factorizations
LU	factorization methods QR, LU, or LT
1	number of problem sizes
31000	values of M
31000	values of N
1	values of nrhs
9200000	values of Asize
1	number of MB's and NB's
16	values of MB
16	values of NB
1	number of process grids
4	values of P
4	values of Q

Run on 16 cores on an SNL Quad-Core, Quad-Socket Opteron IB Cluster

Investigate File system impact with OpenSpeedShop: Compare Lustre I/O with striping to NFS I/O

run cmd: ossio "srun -N 1-n 16 ./testzdriver-std"

Sample Output from Lustre run:

TIME M N MB NB NRHS P Q Fact/SolveTime Error Residual

WALL 31000 31000 16 16 1 4 4 1842.20 1611.59 4.51E+15 1.45E+11

DEPS = 1.110223024625157E-016

sum(xsol_i) = (30999.9999999873,0.00000000000000E+000)

sum |xsol_i - x_i| = (3.332285336962339E-006,0.00000000000000E+000)

sum |xsol_i - x_i|/(M*eps) = (968211.548505533,0.00000000000000E+000)

From output of two separate runs using Lustre and NFS: LU Fact time with Lustre= 1842 secs; LU Fact time with NFS = 2655 secs 813 sec penalty (more than 30%) if you do not use parallel file system like Lustre!

Open | SpeedShop™

NFS and Lustre O|SS Analysis (screen shot from NFS)



I/O to Lustre instead of NFS reduces runtime 25%: (1360 + 99) – (847 + 7) = 605 secs

NFS RUN

LUSTRE RUN

Min t (secs)	Max t (secs)	Avg t (secs)	call Function	Min t (secs)	Max t (secs)	Avg t (secs)	call Function
			libc_read(/lib64/libpthread-				libc_read(/lib64/libpthread-
1102.380076	1360.727283	1261.310157	2.5.so)	368.898283	847.919127	508.658604	2.5.so)
31.19218	99.444468	49.01867	libc_write(/lib64/libpthread- 2.5.so)	6.27036	7.896153	6.850897	libc_write(/lib64/libpthread- 2.5.so)

		Open SpeedS	hop (on chama-login1)		×
<u>F</u> ile <u>T</u> ools					Help
💌 IO [1]					
Process Control					
➡ Run II Cont	➡I Pause ■ Updat	e			Terminate
Status: Process Loaded	d: Click on the "Run" bu	tton to begin the experiment.			
Stats Panel [1]	ManageProcessesPan	əl [1]			
					View/Display Choice
	HC LIS THE OV LE CA	CC Showing Load Balance (mir	n,max,ave) Report:		Functions
Executables: /projects/a	ppperf/mrajan/oocore_	mvapich/app/oocore/bin/testzdriv	ver-std Host: taco10 Pids:	16 Ranks: 16 Threads: 1	
Max Exclusive I/O call ti	me Rank of Max	Min Exclusive I/O call	time i Rank of Min	Average Exclusive I/O	call t Function (defining location)
-1360727.283000	0	1102380.076000	6	1261310.156875	libc read (/lib64/libpthread-2.5
285339.633000	15	217995.565000	5	236541.875375	libc_write (/lib64/libpthread-2.5
-301.410000	15	162.285000	10	241.361625	llseek (/lib64/libpthread-2.5.so)
-29.835000	14	1.505000	0	20.403312	libc_open (/lib64/libpthread-2.5
-2.367000	2	0.393000	8	1.373875	libc_close (/lib64/libpthread-2.5
0.500000	0	0.500000	0	0.500000	libc_open64 (/lib64/libpthread-
-					
Command Panel					
openss>>					
<u></u>					
					1

Open | SpeedShop[™]

How to Analyze the Performance of Parallel Codes 101 - A Tutorial at SC'18



Lustre File System (lfs) commands:

Ifs setstripe –s (size bytes; k, M, G) –c (count; -1 all) –I (index; -1 round robin) <file | directory> Typical defaults: -s 1M -c 4 –i -1 (usually good to try first) File striping is set upon file creation

Ifs getstripe <file | directory> Example: Ifs getstripe --verbose ./oss_Ifs_stripe_16 | grep stripe_count stripe_count: 16 stripe_size: 1048576 stripe_offset: -1



Open | SpeedShop^{*}

How to Analyze the Performance of Parallel Codes 101 - A Tutorial at SC'18

OpenSpeedShop IO-experiment used to identify optimal *lfs striping* (from load balance view (max, min & avg) for 16 way parallel run)



OOCORE I/O performance libc_read time from OpenSpeedShop



Open SpeedShop™

How to Analyze the Performance of Parallel Codes 101 - A Tutorial at SC'18 11/11/2018 1

Additional I/O analysis with O|SS



122

*** Extended I/O Tracing (iot experiment)**

- Records each event in chronological order
- Collects Additional Information
 - Function Parameters
 - Function Return Value
- > When to use extended I/O tracing?
 - When you want to trace the exact order of events
 - When you want to see the return values or bytes read or written.
 - When you want to see the parameters of the IO call

Open | SpeedShop™

Beware of Serial I/O in applications: Encountered in VOSS, code LeP: Simple code here illustrates (acknowledgment: Mike Davis, Cray, Inc.)



#include <stdio.h>
#include <stdlib.h>
#include <mpi.h>
#include <MPI.h>

```
/* Write a single restart file from many MPI processes */
int write_restart (
MPI_Comm comm
                             /// MPI communicator
, int num_cells
                             /// number of cells on this process
, double *cellv )
                             /// cell vector
                        // rank of this process within comm
int rank;
int size;
                        // size of comm
                        // for MPI_Send, MPI_Recv
int tag;
                        // for serializing I/O
int baton;
FILE *f;
                        // file handle for restart file
 /* Procedure: Get MPI parameters */
 MPI_Comm_rank (comm, &rank);
 MPI_Comm_size (comm, &size);
 tag = 4747;
 if (rank == 0) {
  /* Rank 0 create a fresh restart file,
   * and start the serial I/O;
   * write cell data, then pass the baton to rank 1 */
  f = fopen ("restart.dat", "wb");
  fwrite (cellv, num_cells, VARS_PER_CELL * sizeof (double), f);
   fclose (f);
  MPI_Send (&baton, 1, MPI_INT, 1, tag, comm);
 } else {
```

```
/* Ranks 1 and higher wait for previous rank to complete I/O,
   * then append its cell data to the restart file,
   * then pass the baton to the next rank */
 MPI Recv (&baton, 1, MPI INT, rank - 1, tag, comm, MPI STATUS IGNORE);
  f = fopen ("restart.dat", "ab");
  fwrite (cellv, num_cells, VARS_PER_CELL * sizeof (double), f);
  fclose (f);
  if (rank < size - 1) {
   MPI Send (&baton, 1, MPI INT, rank + 1, tag, comm);
 /* All ranks have posted to the restart file; return to called */
return 0;
int main (int argc, char *argv[]) {
 MPI_Comm comm;
 int comm rank;
 int comm size;
 int num_cells;
 double *cellv;
 int i;
 MPI_Init (&argc, &argv);
 MPI_Comm_dup (MPI_COMM_WORLD, &comm);
 MPI Comm rank (comm, &comm rank);
 MPI Comm size (comm, &comm size);
  /**
  * Make the cells be distributed somewhat evenly across ranks
  */
  num cells = 5000000 + 2000 * (comm size / 2 - comm rank);
  cellv = (double *) malloc (num cells * VARS PER CELL * sizeof (double));
  for (i = 0; i < num_cells * VARS_PER_CELL; i++) {</pre>
    cellv[i] = comm rank;
  }
  write_restart (comm, num_cells, cellv);
  MPI Finalize ();
return 0;
```

Open | SpeedShop^{*}

IOT O|SS Experiment of Serial I/O Example



X Open SpeedShop File Tools Help Custom Experiment [1] Process Control Run I⇒Cont	Pause 5 Upd	late		- C	× SI × C gi	HOWS EVENT BY VENT LIST: licking on this ives each call to a
Status: Process Loaded:	Click on the "Run" button	to begin the experiment.			tr	aced as shown.
Stats Panel [1] ■ Man Image: Imag	nageProcessesPanel [1]	CA CC Showing Per E	vent Report:	View/Display Choice • Functions	× B	elow is a raphical trace
Executables: (none) Host	: glory238 Processes/Ra	nks/Threads:(4) 0			vi d	ew of the same ata showing
Start Time 🔺	I/O Call Time(ms)	% of Total Time	Call Stack F	Function (defining location)	- SF	erialization of
-2010/09/08 13:22:54	0.029000	6.682028	>libc_rea	d (/lib64/libpthread-2.5.so)		
-2010/09/08 13:22:54	0.026000	5.990783	>libc_writ	e (/lib64/libpthread-2.5.so)	TV	write() (THE RED
2010/09/08 13:22:54	0.008000	1.843318	>libc_rea	d (/lib64/libpthread-2.5.so)	B	ARS for each PE)
2010/09/08 13:22:54	0.058000	13.364055	>libc_writ	e (/lib64/libpthread-2.5.so)		vith another tool
2010/09/08 13:22:54	0.061000	14.055300	>libc_writ	e (/lib64/libpthread-2.5.so)		
-2010/09/08 13:22:54	0.010000	2.304147	>libc_rea	d (/lib64/libpthread-2.5.so)		ntice2 5.0
-2010/09/08 13:22:54	0.016000	3.686636	>libc_rea	d (/lib64/libpthread-2.5.so)		
2010/09/08 13:22:54	0.015000	3.456221	>libc_rea	d (/lib64/libpthread-2.5.so)	0.000 PE #0 PE #1 PE #2	1.410 2.436 4.254 5.472 7.409 0.507 0.425 11.345 12.761
-2010/09/08 13:22:54	0.025000	5.760369	>libc_rea	d (/lib64/libpthread-2.5.so)	PE #3	
-2010/09/08 13:22:54	0.021000	4.838710	>libc_rea	d (/lib64/libpthread-2.5.so)		
	0.015000	3.456221	> libc writ	e (/lib64/libothread-2.5.so) ▶	Symc Sen Boast Reco	d Housekeeping AlfoAt MPI Tie I/0 ever Reduce Comm Offer scale = 24 1% Q Zoom In Q Zoom Qu 0 352 7 05 1057

124



Offline io/iop/iot experiment on sweep3d application

Convenience script basic syntax:

ossio[p][t] "executable" [default | <list of I/O func>]

- > Parameters
 - I/O Function list to sample(default is all)
 - creat, creat64, dup, dup2, lseek, lseek64, open, open64, pipe, pread, pread64, pwrite, pwrite64, read, readv, write, writev

Examples:

ossio "mpirun –np 256 sweep3d.mpi"

ossiop "mpirun –np 256 sweep3d.mpi" read,readv,write **ossiot** "mpirun –np 256 sweep3d.mpi" read,readv,write

Open | SpeedShop™

I/O output via GUI



* I/O Default View for IOR application "io" experiment

Open SpeedShop ×								
File Tools Help ▼ IO [1] Process Control → Run For any Pause	e 5 Update	Shows the aggregated time spent in the I/O functions traced during the application.						
Status: Process Loaded: Click on the "Run" buttor Status Panel [1] ManageProcessesPanel C Executables: IOR Host: localhost Pids: 4 Ranks: 4	 to begin the experiment. [1] CA CC Showing F Threads: 4 	unctions Report:		F₂ □ ⊨ × View/Display Choice ● Functions				
% of Total 96.427264 2.636284	Exclusive I/O Call Time 	(ms) A % of Total 96.427264 2.636284 0.897688 0.024672 0.014092	Number of Calls 8 16 16 32 8	Function (defining location) open64 (libpthread-2.17.so) write (libpthread-2.17.so) read (libpthread-2.17.so) lseek64 (libpthread-2.17.so) close (libpthread-2.17.so)				
Openss>>								

Open | SpeedShop™

I/O output via GUI



* I/O Call Path View for IOR application "io" experiment

		Open SpeedSho	p	×
File Tools Help ▼ IO [1] Process Control → Run [→ Cont	Pause 5 Updare		Shows the call paths to the I/O functions traced and the time spent along the paths.	≥ □ = × Terminate
Status: Process Loaded: Click on the "Run" Stats Panel [1] ManageProcesses Contemported Stats Panel [1] Executables: IOR Host: localhost Pids: 4 Ra	Panel [1] CA CC Showing Hot anks: 4 Threads: 4	Callpath Report:		G
Exclusive I/O Call Time(ms)	% of Total	Number of Cons	Call Stack Function (defining location) _start (IOR) @ 562 inlibc_start_main (libmonitor.se00.00: main.c,541) libc_start_main (libc-2.17.so) @ 517 in monitor_main (libmonitor.se00.0.0: main.c,492) @ 153 in main (IOR: IOR.c,108) @ 2004 in TestloSys (IOR: IOR.c,148) @ 104 in IOR Create POSIX (IOX: aiori-POSIX.c,74)	
	96.411648	4	open64 (libpthread-2.17.so) _start (IOR)	×
Command Panel				₠ 🗆 🖶 ×
openss>>				

I/O output via CLI (equivalent of HC in GUI)

128

openss>>expview -vcalltrees,fullstack iot1

I/O Call Time(ms) % of Total Time Number of Calls Call Stack Function (defining location) start (sweep3d.mpi) > @ 470 in libc start main (libmonitor.so.0.0.0: main.c,450) >> libc start main (libc-2.10.2.so) >>> @ 428 in monitor main (libmonitor.so.0.0.0: main.c,412) >>>main (sweep3d.mpi) >>>> @ 58 in MAIN (sweep3d.mpi: driver.f,1) >>>>> @ 25 in task init (sweep3d.mpi: mpi stuff.f,1) >>>>>> gfortran ftell i2 sub (libgfortran.so.3.0.0) >>>>>>>>>>_gfortran_st_read (libgfortran.so.3.0.0)

Section Summary - I/O Tradeoffs



- * Avoid writing to one file from all MPI tasks
 - If you need to, be sure to distinguish offsets for each PE at a stripe boundary, and use Buffered I/O
- If each process writes its own file, then the parallel file system attempts to load balance the Object Storage Targets (OSTs), taking advantage of the stripe characteristics
- Metadata server overhead can often create severe I/O problems
 - Minimize number of files accessed per PE and minimize each PE doing operations like seek, open, close, stat that involve inode information
- I/O time is usually not measured, even in applications that keep some function profile
 - Open|SpeedShop can shed light on time spent in I/O using io, iot experiments

Hands-on Section 6: I/O Performance



- ***** I/O experiments related application exercise
- ***** Exercises are in the exercise directory:
 - \$HOME/exercises/IOR
 - \$HOME/exercises/ser_par_io
- Consult README file in each of the directories for the instructions/guidance







SC2018 Tutorial

How to Analyze the Performance of Parallel Codes 101 A case study with Open/SpeedShop

Section 8 Analysis of Memory Usage











Open | SpeedShop™

Memory Hierarchy



132

Memory Hierarchy

- > CPU registers and cache
- > System RAM
- > Online memory, such as disks, etc.
- > Offline memory not physically connected to system
- https://en.wikipedia.org/wiki/Memory hierarchy

What do we mean by memory?

- > Memory an application requires from the system RAM
- Memory allocated on the heap by system calls, such as malloc and friends



Memory Leaks

Is the application releasing memory back to the system?

Memory Footprint

- > How much memory is the application using?
- Finding the High Water Mark (HWM) of memory allocated
- > Out Of Memory (OOM) potential
- Swap and paging issues

Memory allocation patterns

- > Memory allocations longer than expected
- > Allocations that consume large amounts of heap space
- Short lived allocations

Example Memory Heap Analysis Tools



MemP is a parallel heap profiling library

- Requires mpi
- http://sourceforge.net/projects/memp

ValGrind provides two heap profilers.

- Massif is a heap profiler
 - http://valgrind.org/docs/manual/ms-manual.html
- > DHAT is a dynamic heap analysis tool
 - http://valgrind.org/docs/manual/dh-manual.html

Dmalloc - Debug Malloc Library

<u>http://dmalloc.com/</u>

Soogle PerfTools heap analysis and leak detection.

https://github.com/gperftools/gperftools

Open | SpeedShop™

Supports sequential, mpi and threaded applications.

- > No instrumentation needed in application.
- Traces system calls via wrappers
 - malloc
 - calloc
 - realloc
 - free
 - memalign and posix_memalign

Provides metrics for

- > Timeline of events that set an new high-water mark.
- > List of event allocations (with calling context) to leaks.
- Overview of all unique callpaths to traced memory calls that provides max and min allocation and count of calls on this path.

Example Usage

- > ossmem "./lulesh2.0"
- > ossmem "srun -N4 -n 64 ./sweep3d.mpi"





* expview -vunique

Show times, call counts per path, min,max bytes allocation, total allocation to all unique paths to memory calls that the mem collector saw

* expview -vleaked

Show function view of allocations that were not released while the mem collector was active

* expview -vtrace,leaked

> Will show a timeline of any allocation calls that were not released

* expview -vfullstack,leaked

Display a full callpath to each unique leaked allocation

* expview -v trace,highwater

- > Is a timeline of mem calls that set a new high-water
- The last entry is the allocation call that the set the high-water for the complete run
- Investigate the last calls in the timeline and look at allocations that have the largest allocation size (size1,size2,etc) if your application is consuming lots of system ram



Shows the last 8 allocation events that set the high water mark

openss>>expview -vtrace, highwater

Start Time(d:h:m:s)	Event Id:	t Size s Arg1	Size L Arg2	Ptr Arg	Return Value New Call Stack Function (defining location) Highwater
*** trimmed all but th	e last	8 event	ts of 61 **	**	
2016/11/10 09:56:50.8 2.18.so)	824 1	1877:0	2080	0	0x7760e0 19758988 >>>>>GIlibc_malloc (libc-
2016/11/10 09:56:50.8	826 1	1877:0	1728000	0	0x11783d0 21484908 >>>>GIlibc_malloc (libc-
2.18.so) 2016/11/10 09:56:50.8	827 1	1877:0	1728000	0	0x131e1e0 23212908 >>>>GIlibc_malloc (libc-
2.18.so) 2016/11/10 09:56:50.8	827 1	1877:0	1728000	0	0x14c3ff0 24940908 >>>>GIlibc_malloc(libc-
2.18.so) 2016/11/10 09:56:50.8	827 1	1877:0	2080	0	0x776a90 24942988 >>>>> GI libc malloc (libc-
2.18.so)		4077.0	470000	•	
2016/11/10 09:56:50.9 2.18.so)	919 1	18//:0	1728000	0	0x1654030 25286604 >>>>GIlibc_malloc (libc-
2016/11/10 09:56:50.9	919 1	1877:0	1728000	0	0x17f9e40 27014604 >>>>GIlibc_malloc (libc-
2016/11/10 09:56:50.9 2.18.so)	919 1	1877:0	2080	0	0xabc6a0

O|SS Memory Experiment



- The next slide shows the default view of all unique memory calls seen while the mem collector was active. This is an overview of the memory activity. The default is display is aggregated across all processes and threads. Can view specific processes or threads.
- For all memory calls the following are displayed:
 - > The exclusive time and percent of exclusive time
 - > The number of times this memory function was called.
 - > The traced memory function name.

For allocation calls (e.g. malloc) the follow:

- > The max and min allocation size seen.
- The number of times the that max or min was seen are displayed.
- > The total allocation size of all allocations.



openss>>expview -vunique

Exclusive	% of	Number	Min	Min	Max	Max	Total	Function (defining location)
(ms)	Total	of	Request	Requested	Request	Request	ed Bytes	
	Time	Calls	Count	Bytes	Count	Bytes	Requested	
0.024847	89.028629	1546	1	192	6	4096	6316416	GIlibc_malloc (libc-2.18.so)
0.002371	8.495467	5						GIlibc_free (libc-2.18.so)
0.000369	1.322154	1	1	40	1	40	40	realloc (libc-2.18.so)
0.000322	1.153750	3	1	368	1	368	1104	calloc (libc-2.18.so)

NOTE: Number of Calls means the number of unique paths to the memory function call.

To see the paths use the CLI command: expview -vunique,fullstack

Open | SpeedShop[™]

How to Analyze the Performance of Parallel Codes 101 - A Tutorial at SC'18 11/11/2018

O|SS Memory Experiment (Leaked Calls)



In this example the sequential OpenMP version of lulesh was run under ossmem. The initial run detected 69 potential leaks of memory. Examining the calltrees using the cli command "**expview -vfullstack,leaked -mtot_bytes**"

revealed that allocations from the Domain::Domain constructor where not later released in the Domain::~Domain destructor. After adding appropriate delete's in the destructor and rerunning ossmem, we observed a resolution of the leaks detected

in the Domain class. The remaining leaks where minor and from system libraries.

Using the exprestore command to load in the initial database and the database from the second run, we can use the expcompare cli command to see the improvements. Below, database -x1 shows the initial run and -x2 shows the results from the run with the changes to address the leaks detected in the Domain class.

openss>>exprestore -f lulesh-mem-initial.openss openss>>exprestore -f lulesh-mem-improved.openss openss>>expcompare -vleaked -mtot_bytes -mcalls -x1 -x2

-x 2, Function (defining location) -x 1, -x 1, -x 2, Number Total Number Total of of **Bytes Bytes Requested Calls Requested Calls** 10599396 69 3332 8 ___GI___libc_malloc (libc-2.17.so) 72 72 realloc (libc-2.17.so) 1 1



***** Benefits of Memory Heap Analysis

- Detect leaks
- Inefficient use of system memory
- Find potential OOM, paging, swapping conditions
- Determine memory footprint over lifetime of application run

Observations of Memory Analysis Tools

- > Less concerned with the time spent in memory calls
- Emphasis is placed on the relationship of allocation calls to free calls.
- Can slow down and impact application while running



- Memory experiment related application exercise
 - More information provided at the tutorial
- Exercises are in the exercise directory in
 - \$HOME/exercises/matmul
 - \$HOME/exercises/lulesh2.0.3
 - \$HOME/exercises/lulesh2.0.3-fixed

***** Look for the README file for instructions.

Open | SpeedShop™

How to Analyze the Performance of Parallel Codes 101 - A Tutorial at SC'18 11/11/2018







SC2018 Tutorial

How to Analyze the Performance of Parallel Codes 101 A case study with Open/SpeedShop

Section 9 Analysis of heterogeneous codes











Open | SpeedShop™



- Heterogeneous computing refers to systems that use more than one kind of processor.
- What led to increased heterogeneous processing in HPC?
 - > Limits on ability to continue to scale processor frequencies
 - Power consumption hitting realistic upper bound
 - Programmability advances lead to more wide-spread, general usage of graphics processing unit (GPU).
 - Advances in manycore, multi-core hardware technology (MIC)

Heterogeneous accelerator processing: (GPU, MIC)

- Data level parallelism (GPU)
 - Vector units, SIMD execution
 - Single instruction operates on multiple data items
- > Thread level parallelism (MIC)
 - Multithreading, multi-core, manycore



***** GPU (Graphics Processing Unit)

- General-purpose computing on graphics processing units (GPGPU)
- Solve problems of type: Single-instruction, multiple thread (SIMT) model
- Vectors of data where each element of the vector can be treated independently
- > Offload model where data is transferred into/out-of the GPU
- Program using CUDA/OpenCL language or use directive based OpenACC

Intel MIC (Many Integrated Cores)

- > Has a less specialized architecture than a GPU
- > Can execute parallel code written for:
 - Traditional programming models including POSIX threads, OpenMP
- Initially offload based (transfer data to and from co-processor)
- Now/future: programs to run natively


GPU versus CPU comparison

Different goals produce different designs

- GPU assumes work load is highly parallel
- > CPU must be good at everything, parallel or not

CPU: minimize latency experienced by 1 thread

- Big on-chip caches
- Sophisticated control logic

***** GPU: maximize throughput of all threads

- # threads in flight limited by resources => lots of resources (registers, bandwidth, etc.)
- Multi-threading can hide latency => skip the big caches
- Shared control logic across many threads

*based on NVIDIA presentation

Open | SpeedShop™



Mixing GPU and CPU usage in applications

Multicore CPU

Manycore GPU



Data must be transferred to/from the CPU to the GPU in order for the GPU to operate on it and return the new values.

*NVIDIA image

Open | SpeedShop™

11/11/2018

Heterogeneous Programming



***** There are four main ways to use an accelerator

- > Explicit programming:
 - The programmer writes explicit instructions for the accelerator device to execute as well as instructions for transferring data to and from the device (e.g. CUDA-C for GPUs or OpenMP+Cilk Plus for Phis). This method requires to most effort and knowledge from programmers because algorithms must be ported and optimized on the accelerator device.
- > Accelerator-specific pragmas/directives:
 - Accelerator code is automatically generated from your serial code by a compiler (e.g. OpenACC, OpenMP 4.0). For many applications, adding a few lines of code (pragmas/directives) can result in good performance gains on the accelerator.

Accelerator-enabled libraries:

- Only requires the use of the library, no explicit accelerator programming is necessary once the library has been written. The programmer effort is similar to using a non-accelerator enabled scientific library.
- > Accelerator-aware applications:
 - These software packages have been programed by other scientists/engineers/software developers to use accelerators and may require little or no programming for the end-user.

Credit: http://www.hpc.mcgill.ca/index.php/starthere/81-doc-pages/255-accelerator-overview

Open | SpeedShop™

How to Analyze the Performance of Parallel Codes 101 - A Tutorial at SC'18 11/11/2018

Programming for GPGPU



Prominent models for programming the GPGPU

Augment current languages to access GPU strengths

NVIDIA CUDA

- Scalable parallel programming model
- Extensions to familiar C/C++ environment
- Heterogeneous serial-parallel computing
- Supports NVIDIA only
- OpenCL (Open Computing Language)
 - > Open source, royalty-free
 - Portable, can run on different types of devices
 - Runs on AMD, Intel, and NVIDIA

OpenACC

- Provides directives (hint commands inserted into source)
- Directives tell the compiler where to create acceleration (GPU) code without the user modifying or adapting the code.

Optimal Heterogeneous Execution



GPGPU considerations for best performance?

- How is the parallel scaling for the application overall?
- Can you balance the GPU and CPU workload?
 - > Keep both the GPU and CPU busy for best performance
- ✤ Is it profitable to send a piece of work to the GPU?
 - > What is the cost of the transfer of data to and from the GPU?
- How much work is there to be done inside the GPU?
 - Will the work to be done fully populate and keep the GPU processors busy
 - Are there opportunities to chain together operations so the data can stay in the GPU for multiple operations?
- Is there a vectorization opportunity?

Intel MIC considerations for best performance?

- Program should be heavily threaded
- Parallel scaling should be high with an OpenMP version

Open | SpeedShop™

How to Analyze the Performance of Parallel Codes 101 - A Tutorial at SC'18 11/11/2018

Accelerator Performance Monitoring



How can performance tools help optimize code?

- ✤ Is profitable to send a piece of work to the GPU?
 - > Can tell you this by measuring the costs:
 - Transferring data to and from the GPU
 - How much time is spent in the GPU versus the CPU
- Is there a vectorization opportunity?
 - Could measure the mathematical operations versus the vector operations occurring in the application
 - Experiment with compiler optimization levels, re-measure operations and compare
- How is the parallel scaling for the application overall?
 - Use performance tool to get idea of real performance versus expected parallel speed-up
- Provide OpenMP programming model to source code insights
 - Use OpenMP performance analysis to map performance issues to source code

Open | SpeedShop accelerator support



What performance info does Open | SpeedShop provide?

***** For GPGPU it reports information to help understand:

- > Time spent in the GPU device
- Cost and size of data transferred to/from the GPU
- Balance of CPU versus GPU utilization
- Transfer of data between the host and device memory versus the execution of computational kernels
- Performance of the internal computational kernel code running on the GPU device
- Open | SpeedShop is able to monitor CUDA scientific libraries because it operates on application binaries.
- Support for CUDA based applications is provided by tracing actual CUDA events
- ✤ OpenACC support is conditional on the CUDA RT.

Open | SpeedShop accelerator support



What performance info does Open | SpeedShop provide?

- For Intel MIC (non-offload model):
 - Reports the same range of performance information that it does for CPU based applications
 - Open|SpeedShop will operate on MIC (co-processor KNC) similar to targeted platforms where the compute node processer is different than the front-end node processor
 - Only non-offload support is in our current plans
 - A specific OpenMP profiling experiment (omptp) has been developed. Initial version is available now.
 - Will help to better support analysis of MIC based applications
 - OpenMP performance analysis key to understanding performance

CUDA GUI View: Default CUDA view



Note: The left pane shows the executable and the nodes it ran on. In future, will effect views. Internal GPU activity is shown in thread t1 (GPU) graphic (shaded area) Red lines indicate data transfers, Green lines indication GPU kernel executions Source panel displays source for metrics clicked on in the Metric pane.



Open | SpeedShop^{*}

How to Analyze the Performance of Parallel Codes 101 - A Tutorial at SC'18

CUDA GUI View: All Events Trace



Note: The chronological list of data transfers and kernel executions in bottom pane. Duration of kernel execution and data transfer available. Internal GPU activity is shown thread t1 (GPU) graphic (shaded area) Red lines indicate data transfers, Green lines indication GPU kernel executions



CUDA GUI View: Kernel Trace



Note: The chronological list of kernel executions with details is in bottom pane. Internal GPU activity is shown in thread t1 (GPU) graphic (shaded area) Red lines indicate data transfers, Green lines indication GPU kernel executions



CUDA GUI View: Transfers Trace



Note: The chronological list of data transfers with details is in bottom pane. Internal GPU activity is shown in thread t1 (GPU) graphic (shaded area) Red lines indicate data transfers, Green lines indication GPU kernel executions



CUDA GUI View: Timeline Zoom



Note: Here is a zoomed in view of the data transfer and kernel execution timeline Red lines indicate data transfers, Green lines indication GPU kernel executions The metric view is dependent on what is active in the timeline view.



Open | SpeedShop CUDA CLI Views



openss>>expview [-vExec]

Exclusive % of Exclu Time (ms) Total Co Exclusive Time	usive Function ount	n (defining locati	on)
14.810702 52.042113	300 void	RunTest <float>(s</float>	td::string, ResultDatabase&, OptionParser&) (GEMM: GEMM.cpp,19)
13.648369 47.957887	300 void	RunTest <double< th=""><th>>(std::string, ResultDatabase&, OptionParser&) (GEMM: GEMM.cpp,19)</th></double<>	>(std::string, ResultDatabase&, OptionParser&) (GEMM: GEMM.cpp,19)
openss>>expview -vXfe	er		
Exclusive % of Exclu	usive Function	n (defining locati	on)
Time (ms) Total C	ount		,
Exclusive			
Time			
1.774178 75.232917	69 void R	unTest <float>(st</float>	d::string, ResultDatabase&, OptionParser&) (GEMM: GEMM.cpp,19)
0.584069 24.767083	69 void R	unTest <double></double>	std::string, ResultDatabase&, OptionParser&) (GEMM: GEMM.cpp,19)
openss>>expview -v tra	ace, Xfer		
Start Time (d:h:m:s)	Exclusive	e % of Si	ze Kind Call Stack Function (defining location)
	Time (ms	s) Total	
		Exclusive	
		Time	
2016/08/24 10:01:03.84	5 0.001217	0.051606 112	HostToDevice >>void RunTest <float>(std::string, ResultDatabase&, OptionParser&) (GEMM: GEMM.cpp,19)</float>
2016/08/24 10:01:03.85	0 0.027392	1.161541 262144	HostToDevice >>void RunTest <float>(std::string, ResultDatabase&, OptionParser&) (GEMM: GEMM.cpp,19)</float>
2016/08/24 10:01:03.85	0 0.027553	1.168368 262144	HostToDevice >>void RunTest <float>(std::string, ResultDatabase&, OptionParser&) (GEMM: GEMM.cpp,19)</float>
2016/08/24 10:01:03.85	1 0.001217	0.051606 112	HostToDevice >>void RunTest <float>(std::string, ResultDatabase&, OptionParser&) (GEMM: GEMM.cpp,19)</float>
2016/08/24 10:01:03.85	1 0.027425	1.162940 262144	DeviceToHost >>void RunTest <float>(std::string, ResultDatabase&, OptionParser&) (GEMM: GEMM.cpp,19)</float>
2016/08/24 10:01:03.85	2 0.026721	1.133087 262144	DeviceToHost >>void RunTest <float>(std::string, ResultDatabase&, OptionParser&) (GEMM: GEMM.cpp,19)</float>
2016/08/24 10:01:03.85	2 0.026753	1.134444 262144	DeviceToHost >>void RunTest <float>(std::string, ResultDatabase&, OptionParser&) (GEMM: GEMM.cpp,19)</float>

Open | SpeedShop™

.....

Open | SpeedShop CUDA CLI Views



openss>>expview -v trace,Exec

Start Time (d:h:m:s)	Exclusive	% of	Grid	Block	Call Stack Function (defining location)
	Time (ms)	Total	Dims		
		Exclusive	9		
		Time			
2016/08/24 10:01:03.851	0.055585	0.195316	4,4,1	16,16,1	>>void RunTest <float>(std::string, ResultDatabase&, OptionParser&) (GEMM: GEMM.cpp,19)</float>
2016/08/24 10:01:03.851	0.048705	0.171141	4,4,1	16,16,1	>>void RunTest <float>(std::string, ResultDatabase&, OptionParser&) (GEMM: GEMM.cpp,19)</float>
2016/08/24 10:01:03.851	0.049761	0.174851	4,4,1	16,16,1	>>void RunTest <float>(std::string, ResultDatabase&, OptionParser&) (GEMM: GEMM.cpp,19)</float>
2016/08/24 10:01:03.851	0.051617	0.181373	4,4,1	16,16,1	>>void RunTest <float>(std::string, ResultDatabase&, OptionParser&) (GEMM: GEMM.cpp,19)</float>
2016/08/24 10:01:03.851	0.051648	0.181482	4,4,1	16,16,1	>>void RunTest <float>(std::string, ResultDatabase&, OptionParser&) (GEMM: GEMM.cpp, 19)</float>
2016/08/24 10:01:03.851	0.050817	0.178562	4,4,1	16,16,1	>>void RunTest <float>(std::string, ResultDatabase&, OptionParser&) (GEMM: GEMM.cpp, 19)</float>
2016/08/24 10:01:03.851	0.046496	0.163378	4,4,1	16,16,1	>>void RunTest <float>(std::string, ResultDatabase&, OptionParser&) (GEMM: GEMM.cpp,19)</float>
2016/08/24 10:01:03.851	0.048193	0.169341	4,4,1	16,16,1	>>void RunTest <float>(std::string, ResultDatabase&, OptionParser&) (GEMM: GEMM.cpp,19)</float>
2016/08/24 10:01:03.852	0.049633	0.174401	4,4,1	16,16,1	>>void RunTest <float>(std::string, ResultDatabase&, OptionParser&) (GEMM: GEMM.cpp,19)</float>

....

openss>>expview -vcalltrees,fullstack

	Exclusive	% of	Exclusive Call Stack Function (defining location)	
	Time (ms)	Total	Count	
		Exclusiv	e	
		Time		
			main (GEMM: main.cpp,135)	
			> @ 130 in RunBenchmark(ResultDatabase&, OptionParser&) (GEMM: GEMM.cpp,122)	
	11.818358	41.52756	1 240 >> @ 240 in void RunTest <float>(std::string, ResultDatabase&, OptionParser&) (GEMM: GEMM.cpp,19) main (GEMM: main.cpp,135)</float>	
			> @ 137 in RunBenchmark(ResultDatabase&, OptionParser&) (GEMM: GEMM.cpp,122)	
	10.894840	38.28248	6 240 >> @ 240 in void RunTest <double>(std::string, ResultDatabase&, OptionParser&) (GEMM: GEMM.cpp,19) main (GEMM: main.cpp,135)</double>	
			> @ 130 in RunBenchmark(ResultDatabase&, OptionParser&) (GEMM: GEMM.cpp,122)	
	2.992344	10.514553	60 >> @ 231 in void RunTest <float>(std::string, ResultDatabase&, OptionParser&) (GEMM: GEMM.cpp,19) main (GEMM: main.cpp,135)</float>	
			> @ 137 in RunBenchmark(ResultDatabase&, OptionParser&) (GEMM: GEMM.cpp,122)	
	2.753529	9.675400	60 >> @ 231 in void RunTest <double>(std::string, ResultDatabase&, OptionParser&) (GEMM: GEMM.cpp,19)</double>	
Оре	n Speed	dShop™	How to Analyze the Performance of Parallel Codes 101 - A Tutorial at SC'18 11/11/202	18

Open | SpeedShop CUDA CLI Views



pfe27-433>openss -cli -f GEMM-cuda-4.openss openss>>[openss]: The restored experiment identifier is: -x 1 openss>>expview

% of Exclusive Function (defining location) Exclusive

Time (ms) Total Count Exclusive

Time

14.810702 52.042113

13.648369 47.957887

300 void RunTest<float>(std::string, ResultDatabase&, OptionParser&) (GEMM: GEMM.cpp,19) 300 void RunTest<double>(std::string, ResultDatabase&, OptionParser&) (GEMM: GEMM.cpp,19)

openss>>expview -vhwpc				242	203576	0	1	1		
Time CPU All GPU All <cpu< td=""><td>253</td><td>730303</td><td>0</td><td></td><td>ĺ</td><td></td></cpu<>				253	730303	0		ĺ		
(ms) 0 15868757	0	***			264	4937670	0	*		
11 5336880 22 5205442	0	*	Í		275	24977312	0	*****		
33 5410977 44 3780335	0	*			286	53366059	0	*********		
55 2794120 66 5031483	0	*			297	75579534	0	******		
77 3289826 88 2243716	0				308	79920340	0	*****		
99 1628496 110 670313	0				319	76604975	0	******		
121 105549 132 125052	0				330	77356196	0	*****		
143 134162 154 143953	0				341	78801255	0	******		
165 146363 176 155874	0				352	68318322	0	**********		
187 182306 198 194074	0				363	66937166	0	**********		
209 176671	0				374	69401858	0	**********		
231 196431	0				385	73239976	0	*****		
253 730303	0	*			396	71365211	544298	***********		
275 24977312	Ö	*****	í l		407	70238071	3554730	**********	****	
297 75579534	0	***************************************			418	70172897	10504920	**********	*****	
319 76604975	0	*************			429	82853194	11857290	*****	*****	
341 78801255	0	*****			440	68740879	5299162	**********	****	
363 66937166	0	**********			451	20665073	0	****		
374 69401858	0	**************								
396 /1365211 407 70238071	544298 3554730	**********	****							
418 70172897 429 82853194	10504920 11857290	*****	******							
440 68740879 451 20665073	5299162 0	***************************************	******							

Hands-on Section 8: GPU Performance



- GPU related application exercises
- Second Second
 - \$HOME/exercises/cuda/matrixMul
 - > \$HOME/exercises/cuda/shoc/bindir/bin/EP/CUDA

Consult README for exercise instructions/guidance

- Run matrixMul exercise
- Run shoc benchmarks: GEMM and FFT

Open SpeedShop™

How to Analyze the Performance of Parallel Codes 101 - A Tutorial at SC'18 11/11/2018







SC2018 Tutorial

How to Analyze the Performance of Parallel Codes 101 A case study with Open/SpeedShop

Section 10 DIY & Conclusions











11/11/2018

Open | SpeedShop™



164

OpenSpeedShop booth: 2840

- On-demand Demos, discussion, new GUI feedback, etc.
- Reminder: Tutorial surveys are entirely electronic this year
- & QR code:

https://submissions.supercomputing.org/eval.png

- Evaluation site URL: <u>http://bit.ly/SC18-eval</u>
- Thanks for attending our tutorial!

How to Take This Experience Home?



Seneral questions should apply to ...

- > ... all systems
- … all applications

* Prerequisite

- Know what to expect from your application
- Know the basic architecture of your system

* Ask the right questions

- Start with simple overview questions
- Dig deeper after that

Pick the right tool for the task

- May need more than one tool
- Will depend on the question you are asking
- May depend on what is supported on your system

Open | SpeedShop™

If You Want to Give O|SS a Try?



Available on the these system architectures

- > AMD x86-64
- Intel x86, x86-64, MIC/Phi
- > IBM PowerPC, PowerPC64. Power8
- > ARM: AArch64/A64 and AArch32/A32

Work with these operating system

- Tested on Many Popular Linux Distributions
 - SLES, SUSE
 - RHEL, Fedora, CentOS
 - Debian, Ubuntu

Tested on some large scale platforms

- IBM Blue Gene and Cray
- GPU and Intel Phi support available
- > Available on many DOE/DOD systems in shared locations
- Ask your system administrator

Open | SpeedShop™

How to Install Open | SpeedShop?



Most tools are complex pieces of software

- Low-level, platform specific pieces
- Complex dependencies
- > Need for multiple versions, e.g., based on MPIs and compilers
- > Open | SpeedShop is no exception
 - In many cases even harder because of its transparency

Installation support

- Traditional installation mechanism
 - Three parts of the installation
 - Krell Root base packages
 - CBTF Component based tool framework
 - O|SS client itself
 - Install script
- Support for "spack" now available
 - https://github.com/spack/spack

When in doubt, don't hesitate, ask us:

oss-contact@openspeedshop.org

Open | SpeedShop™



- Current version: 2.4.0 has been released
- Open | SpeedShop Website
 - <u>https://www.openspeedshop.org/</u>
- Open | SpeedShop help and bug reporting
 - Direct email: oss-contact@openspeedshop.org
 - Forum/Group: oss-questions@openspeedshop.org

* Feedback

- > Bug tracking available from website
- Feel free to contact presenters directly

***** Support contracts and onsite training available

- We are working with users to develop a support contract process through the Trenza Synergy Center.
- > Stop at booth 2840 to discuss options, if interested.

Open | SpeedShop^{*}

Download options:

- Package with install script (install-tool)
- Source for tool and base libraries

Project Wiki:

https://github.com/OpenSpeedShop/openspeedshop/wiki

Repositories access

https://github.com/OpenSpeedShop

Release Information

Release Tarball and Packages are accessible from

www.openspeedshop.org



Open | SpeedShop Documentation



- Suild and Installation Instructions
 - <u>https://www.openspeedshop.org/documentation</u>
 - Look for: Open | SpeedShop Version 2.4 Build/Install Guide

Open | SpeedShop User Guide Documentation

- <u>https://www.openspeedshop.org/documentation</u>
 - Look for Open | SpeedShop Version 2.4 Users Guide

Man pages: OpenSpeedShop, osspcsamp, ossmpi,

...

Quick start guide downloadable from web site

- <u>https://www.openspeedshop.org</u>
- Click on "Download Quick Start Guide" button

Open | SpeedShop^{*}

Tutorial Summary



Performance analysis critical on modern systems

- Complex architectures vs. complex applications
- Need to break black box behavior at multiple levels
- Lots of performance left on the table by default

Performance tools can help

- > Open|SpeedShop as one comprehensive option
- Scalability of tools is important
 - Performance problems often appear only at scale
 - We will see more and more online aggregation approaches
 - CBTF as one generic framework to implement such tools

Critical:

- Asking the right questions
- Comparing answers with good baselines or intuition
- Starting at a high level and iteratively digging deeper

Open | SpeedShop™

Questions vs. Experiments



Where do I spend my time?

- Flat profiles (pcsamp)
- Getting inclusive/exclusive timings with callstacks (usertime)
- Identifying hot callpaths (usertime + HP analysis)

How do I analyze cache performance?

- > Measure memory performance using hardware counters (hwc)
- > Compare to flat profiles (custom comparison)
- Compare multiple hardware counters (N x hwc, hwcsamp)

How to identify I/O problems?

- Study time spent in I/O routines (io)
- Compare runs under different scenarios (custom comparisons)

How do I find parallel inefficiencies?

- Study time spent in MPI routines (mpi)
- Look for load imbalance (LB view) and outliers (CA view)

Open | SpeedShop™