

Automatic trace analysis with Scalasca

Markus Geimer, Brian Wylie, David Böhme Jülich Supercomputing Centre

















- Idea
 - Automatic search for patterns of inefficient behavior
 - Classification of behavior & quantification of significance



- Guaranteed to cover the entire event trace
- Quicker than manual/visual trace analysis
- Parallel replay analysis exploits available memory & processors to deliver scalability



- Project started in 2006
 - Initial funding by Helmholtz Initiative & Networking Fund
 - Many follow-up projects
- Follow-up to pioneering KOJAK project (started 1998)
 - Automatic pattern-based trace analysis
- Now joint development of
 - Jülich Supercomputing Centre
 - German Research School for Simulation Sciences





- Development of a scalable performance analysis toolset for most popular parallel programming paradigms
- Specifically targeting large-scale parallel applications
 - such as those running on IBM BlueGene or Cray XT systems with one million or more processes/threads
- Latest release:
 - Scalasca v2.0 with Score-P support (August 2013)



- Open source, New BSD license
- Fairly portable
 - IBM Blue Gene, IBM SP & blade clusters, Cray XT, SGI Altix, Solaris & Linux clusters, ...
- Uses Score-P instrumenter & measurement libraries
 - Scalasca 2.0 core package focuses on trace-based analyses
 - Supports common data formats
 - Reads event traces in OTF2 format
 - Writes analysis reports in CUBE4 format
- Current limitations:
 - No support for nested OpenMP parallelism and tasking
 - Unable to handle OTF2 traces containing CUDA events



Example: Wait at NxN

locatior





time

- Time spent waiting in front of synchronizing collective operation until the last process reaches the operation
- Applies to: MPI_Allgather, MPI_Allgatherv, MPI_Alltoall, MPI_Reduce_scatter, MPI_Reduce_scatter_block, MPI_Allreduce

Example: Late Broadcast



- Waiting times if the destination processes of a collective 1-to-N operation enter the operation earlier than the source process (root)
- Applies to: MPI_Bcast, MPI_Scatter, MPI_Scatterv

Example: Late Sender





- Waiting time caused by a blocking receive operation posted earlier than the corresponding send
- Applies to blocking as well as non-blocking communication



Hands-on: NPB-MZ-MPI / BT

scalasca 🗖





• One command for (almost) everything...

```
% scalasca
Scalasca 2.0
Toolset for scalable performance analysis of large-scale applications
usage: scalasca [-v][-n][c] {action}
1. prepare application objects and executable for measurement:
    scalasca -instrument <compile-or-link-command> # skin (using scorep)
2. run application under control of measurement system:
    scalasca -analyze <application-launch-command> # scan
3. interactively explore measurement analysis report:
    scalasca -examine <experiment-archive|report> # square
-v, --verbose enable verbose commentary
-n, --dry-run show actions without taking them
-c, --show-config show configuration and exit
```

- The 'scalasca -instrument' command is deprecated and only provided for backwards compatibility with Scalasca 1.x.
- Recommended: use Score-P instrumenter directly



```
% skin
Scalasca 2.0: application instrumenter using scorep
usage: skin [-v] [-comp] [-pdt] [-pomp] [-user] <compile-or-link-cmd>
-comp={all|none|...}: routines to be instrumented by compiler
        (... custom instrumentation specification for compiler)
        -pdt: process source files with PDT instrumenter
        -pomp: process source files for POMP directives
        -user: enable EPIK user instrumentation API macros in source code
        -v: enable verbose commentary when instrumenting
        --*: options to pass to Score-P instrumenter
```

- Provides compatibility with Scalasca 1.x
- Recommended: use Score-P instrumenter directly



• Scalasca measurement collection & analysis nexus

% scan		
Scalasca 2.0: measurement collection & analysis nexus		
usage: scan {options} [launchcmd [launchargs]] target [targetargs] where {options} may include:		
h Holp: show this brief usage message and exit		
-II Help. Show this brief usage message and exit.		
-v Verbose: increase verbosity.		
-n Preview: show command(s) to be launched but don't execute.		
-q Quiescent: execution with neither summarization nor tracing.		
-s Summary: enable runtime summarization. [Default]		
-t Tracing: enable trace collection and analysis.		
-a Analyze: skip measurement to (re-)analyze an existing trace.		
-e exptdir : Experiment archive to generate and/or analyze.		
(overrides default experiment archive title)		
-f filtfile : File specifying measurement filter.		
-l lockfile : File that blocks start of measurement.		

VI-HPS



- scan configures Score-P measurement by setting some environment variables automatically
 - e.g., experiment title, profiling/tracing mode, filter file, ...
 - Precedence order:
 - Command-line arguments
 - Environment variables already set
 - Automatically determined values
- Also, scan includes consistency checks and prevents corrupting existing experiment directories
- For tracing experiments, after trace collection completes then automatic parallel trace analysis is initiated
 - uses identical launch configuration to that used for measurement (i.e., the same allocated compute resources)

- VI-HPS
- Run the application using the Scalasca measurement collection & analysis nexus prefixed to launch command

```
% export SCOREP_EXPERIMENT_DIRECTORY=scorep_bt-mz_W_4x4_sum
% OMP NUM THREADS=4 scan mpiexec -np 4 ./bt-mz_W.4
S=C=A=N: Scalasca 2.0 runtime summarization
S=C=A=N: ./scorep_bt-mz_W_4x4_sum experiment archive
S=C=A=N: Thu Sep 13 18:05:17 2012: Collect start
mpiexec -np 4 ./bt-mz_W.4
NAS Parallel Benchmarks (NPB3.3-MZ-MPI) - BT-MZ MPI+OpenMP Benchmark
Number of zones: 8 x 8
Iterations: 200 dt: 0.000300
Number of active processes:
                                4
 [... More application output ...]
S=C=A=N: Thu Sep 13 18:05:39 2012: Collect done (status=0) 22s
S=C=A=N: ./scorep_bt-mz_W_4x4_sum_complete.
```

Creates experiment directory ./scorep_bt-mz_W_4x4_sum



```
% square -s scorep_bt-mz_W_4x4_sum
INFO: Post-processing runtime summarization result...
INFO: Score report written to ./scorep_bt-mz_W_4x4_sum/scorep.score
```

Post-processing and interactive exploration with CUBE

% square scorep_bt-mz_W_4x4_sum INFO: Displaying ./scorep_bt-mz_W_4x4_sum/summary.cubex...

[GUI showing summary analysis report]

 The post-processing derives additional metrics and generates a structured metric hierarchy

Post-processed summary analysis report





VI-HPS

0.0 Reference preparation for validation

- 1.0 Program instrumentation
- 1.1 Summary measurement collection
- 1.2 Summary analysis report examination
- 2.0 Summary experiment scoring2.1 Summary measurement collection with filtering2.2 Filtered summary analysis report examination
- 3.0 Event trace collection
- 3.1 Event trace examination & analysis



Load modules

```
% module load UNITE
UNITE loaded
% module load scorep/1.2.1
scorep/1.2.1 loaded
% module load scalasca/2.1-alpha1
scalasca/2.1-alpha1 loaded
% module load cube/4.2
cube/4.2 loaded
```

- Change to directory containing NPB BT-MZ sources
- Existing instrumented binary in bin.scorep/ can be reused



Change to executable directory and edit job script

• Submit the job

% llsubmit scalasca2.ll



• Continues with automatic (parallel) analysis of trace files

```
S=C=A=N: Fri Sep 20 15:09:59 2013: Analyze start
Analyzing experiment archive ./scorep_bt-mz_C_2p64x8_trace/traces.otf2
Opening experiment archive ... done (0.019s).
Reading definition data ... done (0.178s).
Reading event trace data ... done (2.068s).
Preprocessing ... done (3.789s).
Analyzing trace data ...
  Wait-state detection (fwd) (1/5) ... done (2.889s).
  Wait-state detection (bwd) (2/5) ... done (1.136s).
  Synchpoint exchange (fws) (3/5) ... done (0.813s).
 Critical-path & delay analysis (4/5) ... done (0.568s).
done (5.413s).
Writing analysis report ... done (1.994s).
Max. memory usage : 181.066MB
Total processing time : 13.645s
S=C=A=N: Fri Sep 20 15:10:16 2013: Analyze done (status=0) 17s
```





 Produces trace analysis report in experiment directory containing trace-based wait-state metrics

```
% square scorep_bt-mz_C_32x8_trace
INFO: Post-processing runtime summarization result...
INFO: Post-processing trace analysis report...
INFO: Displaying ./scorep_bt-mz_C_32x8_trace/trace.cubex...
```

[GUI showing trace analysis report]

Post-processed trace analysis report











Performance properties

Late Sender Time

Description:

Refers to the time lost waiting caused by a blocking receive operation (e.g., MPI_Recv or MPI_Wait) that is posted earlier than the corresponding send operation.



If the receiving process is waiting for multiple messages to arrive (e.g., in an call to MPI_Waitall), the maximum waiting time is accounted, i.e., the waiting time due to the latest sender.

Unit:

Seconds

Diagnosis:

Try to replace MPI_Recv with a non-blocking receive MPI_Irecv that can be posted earlier, proceed concurrently with computation, and complete with a wait operation after the message is expected to have been sent. Try to post sends earlier, such that they are available when receivers need them. Note that outstanding messages (i.e., sent before the receiver is ready) will occupy internal message buffers, and that large numbers of posted receive buffers will also introduce message management overhead, therefore moderation is advisable.

Parent:

MPI Point-to-point Communication Time Children:

<u>C</u>lose







VI-HPS



VI-HPS





Connect to Vampir trace browser









cube 4.1.1 livedvd2: scorep_bt-mz_B_4x4_trace/trace.cubex File Display Topology Help		
Absolute	Absolute ~	Absolute ~
Metric tree	💽 Call tree 🔲 Flat view	System tree 🚺 Box Plot
0.00 Time 300.91 Execution 0.00 MPI 0.01 Synchronization 0.00 Communication 0.39 Point-to-point 1.38 Late Sender 0.00 Collective 0.00 Early Reduce 0.00 Early Reduce 0.00 Early Scan 0.00 Wait at N x N 0.00 File I/O 0.00 File I/O 0.00 File I/O 0.00 Flush 0.00 Flush 1.38 (0.41%) 337.45	Call site O.00 MAIN O.00 mpi_set O.00 MPI_Bca O.00 env_set O.00 zone_se O.00 zone_se O.00 zone_st O.00 zone_st O.00 set_cons O.00 exact_rt O.00	e browser 0.00 Thread 1 0.00 Thread 2 0.00 Thread 3 - MPI Rank 1 0.00 Thread 3 - MPI Rank 1 0.00 Thread 1 0.00 Thread 1 0.00 Thread 3 - MPI Rank 2 0.00 Thread 3 - MPI Rank 2 0.00 Thread 3 - MPI Rank 3 - MPI Rank 3 - 0.00 Thread 3 - 0.00 Thread 3 - MPI Rank 3 - 0.00 Thread 3 - 0.00 Th
Shows the most severe instance of pattern i	of call paths marke	ed with a









Scalable performance analysis of large-scale parallel applications

- toolset for scalable performance measurement & analysis of MPI, OpenMP & hybrid parallel applications
- supporting most popular HPC computer systems
- available under New BSD open-source license
- sources, documentation & publications:
 - http://www.scalasca.org
 - mailto: scalasca@fz-juelich.de

