

# Automatic trace analysis with Scalasca

David Boehme Lawrence Livermore National Lab

# scalasca 🗖



# Automatic trace analysis

- 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

# The Scalasca project: Objective

- 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 Blue Gene or Cray systems with one million or more processes/threads
- Latest release:
  - Scalasca v2.3.1 coordinated with Score-P v2.0.2 (May 2016)

## Scalasca features

- Open source, 3-clause BSD license
- Fairly portable
  - IBM Blue Gene, Cray XT/XE/XK/XC, SGI Altix, Fujitsu FX10/100 & K computer, Linux clusters (x86, Power, ARM), Intel Xeon Phi, ...
- Uses Score-P instrumenter & measurement libraries
  - Scalasca 2 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:
  - Unable to handle traces containing CUDA or SHMEM events, or OpenMP nested parallelism
  - PAPI/rusage metrics for trace events are ignored

# VI-HPS

# Scalasca workflow



VIRTUAL INSTITUTE – HIGH PRODUCTIVITY SUPERCOMPUTING

#### Example: "Late Sender" wait state



time

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

# **Example: Critical path**



- Shows call paths and processes/threads that are responsible for the program's wall-clock runtime
- Identifies good optimization candidates and parallelization bottlenecks

VIRTUAL INSTITUTE – HIGH PRODUCTIVITY SUPERCOMPUTING

#### **Example: Root-cause analysis**



- Classifies wait states into direct and indirect (i.e., caused by other wait states)
- Identifies delays (excess computation/communication) as root causes of wait states
- Attributes wait states as *delay costs*



# Hands-on: NPB-MZ-MPI / BT

# scalasca 🗖



## Scalasca command – One command for (almost) everything

<pre>% scalasca Scalasca 2.3.1</pre>		
Toolset for scalable performance analysis of large-scale parallel applications		
usage: scalasca [OPTION] ACTION <argument></argument>		
1.	prepare application	on objects and executable for measurement:
	scalasca -instrum	ent <compile-or-link-command> # skin (using scorep)</compile-or-link-command>
2.	. run application under control of measurement system:	
	scalasca -analyze	<application-launch-command> # scan</application-launch-command>
3.	interactively exp	lore measurement analysis report:
	scalasca -examine	<experiment-archive report></experiment-archive report>
Options:		
-C,	show-config	show configuration summary and exit
-h,	help	show this help and exit
-n,	dry-run	show actions without taking them
	quickref	show quick reference guide and exit
	remap-specfile	show path to remapper specification file and exit
-v,	verbose	enable verbose commentary
-V,	version	show version information and exit

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

#### Scalasca compatibility command: skin

- Scalasca application instrumenter
  - Provides compatibility with Scalasca 1.x
  - Deprecated! Use Score-P instrumenter directly.

#### Scalasca convenience command: scan

% scan			
Scalasca 2.3.1: measurement collection & analysis nexus			
<pre>usage: scan {options} [launchcmd [launchargs]] target [targetargs]   where {options} may include:</pre>			
-h 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.			
<pre>-e exptdir : Experiment archive to generate and/or analyze. (overrides default experiment archive title)</pre>			
-f filtfile : File specifying measurement filter.			
-l lockfile : File that blocks start of measurement.			
-m metrics : Metric specification for measurement.			

#### Scalasca measurement collection & analysis nexus

# Scalasca advanced command: scout - Scalasca automatic trace analyzer

```
% scout.hyb --help
        Copyright (c) 1998-2016 Forschungszentrum Juelich GmbH
SCOUT
        Copyright (c) 2009-2014 German Research School for Simulation
                                Sciences GmbH
Usage: <launchcmd> scout.hyb [OPTION]... <ANCHORFILE | EPIK DIRECTORY>
Options:
                    Enables instance tracking and statistics [default]
  --statistics
  --no-statistics Disables instance tracking and statistics
  --critical-path
                     Enables critical-path analysis [default]
  --no-critical-path Disables critical-path analysis
                     Enables root-cause analysis [default]
  --rootcause
                     Disables root-cause analysis
  --no-rootcause
  --single-pass
                     Single-pass forward analysis only
  --time-correct
                     Enables enhanced timestamp correction
                     Disables enhanced timestamp correction [default]
  --no-time-correct
                     Increase verbosity
  --verbose, -v
  --help
                     Display this information and exit
```

Provided in serial (.ser), OpenMP (.omp), MPI (.mpi) and MPI+OpenMP (.hyb) variants

# Scalasca advanced command: clc\_synchronize

Scalasca trace event timestamp consistency correction

Usage: <launchcmd> clc\_synchronize.hyb <ANCHORFILE | EPIK\_DIRECTORY>

- Provided in MPI (.mpi) and MPI+OpenMP (.hyb) variants
- Takes as input a trace experiment archive where the events may have timestamp inconsistencies
  E.g., multi-node measurements on systems without adequately synchronized clocks on each compute node
- Generates a new experiment archive (always called ./clc\_sync) containing a trace with event timestamp inconsistencies resolved
  - E.g., suitable for detailed examination with a time-line visualizer

#### Scalasca convenience command: square

```
% square
Scalasca 2.3.1: analysis report explorer
usage: square [-v] [-s] [-f filtfile] [-F] <experiment archive | cube file>
-c <none | quick | full> : Level of sanity checks for newly created reports
-F : Force remapping of already existing reports
-f filtfile : Use specified filter file when doing scoring
-s : Skip display and output textual score report
-v : Enable verbose mode
-n : Do not include idle thread metric
```

Scalasca analysis report explorer

# Automatic measurement configuration

- scan configures Score-P measurement by automatically setting some environment variables and exporting them
  - 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)

# **Setup environment**

- Remember load modules to add local tool installations to \$PATH
  - Must match the compiler you intend to use (Intel or GCC)

```
% use scorep-2.0.2-mvapich2-gnu-2.2
% use scalasca-2.3.1-mvapich2-gnu-2.2
```

- Change to directory containing NPB3.3-MZ-MPI sources
- Existing instrumented executable in bin.scorep/ directory can be reused

% cd \$HOME/scratch/NPB3.3-MZ-MPI

# Performance analysis steps

- 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 scoring
- 2.1 Summary measurement collection with filtering
- 2.2 Filtered summary analysis report examination
- 3.0 Event trace collection
- 3.1 Event trace examination & analysis

#### **BT-MZ trace measurement collection...**

```
% cd bin.scorep
% cp ../jobscript/darwin/scalasca.sbatch .
% vi scalasca.sbatch
[...]
export SCOREP_FILTERING_FILE=../config/scorep.filt
export SCOREP_TOTAL_MEMORY=21MB
# Scalasca configuration
export SCAN ANALYZE OPTS="--time-correct"
```

scalasca -analyze -t srun -n8 -ppdebug./bt-mz\_\${CLASS}.\${PROCS}

 Change to directory with the executable and edit the job script

Submit the job

% sbatch scalasca.sbatch

#### **BT-MZ trace measurement ... collection**

S=C=A=N: Scalasca 2.3.1 trace collection and analysis S=C=A=N: ./scorep\_bt-mz\_B\_8x4\_trace experiment archive S=C=A=N: Sat Jul 2 16:20:56 2016: Collect start \ mpirun -np 8 ./bt-mz B.8

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: 8

[... More application output ...]

S=C=A=N: Sat Jul 2 16:21:11 2016: Collect done (status=0) 15s

 Starts measurement with collection of trace files ...

#### **BT-MZ** trace measurement ... analysis

S=C=A=N: Sat Jul 2 16:21:11 2016: Analyze start  $\setminus$ mpirun -np 8 scout.hyb ./scorep bt-mz B 8x4 trace/traces.otf2 Copyright (c) 1998-2016 Forschungszentrum Juelich GmbH SCOUT Copyright (c) 2009-2014 German Research School for Simulation Sciences GmbH Analyzing experiment archive ./scorep bt-mz B 8x4 trace/traces.otf2 Opening experiment archive ... done (0.073s). Reading definition data ... done (0.053s). Reading event trace data ... done (0.106s). Preprocessing ... done (0.125s). Analyzing trace data ... done (2.778s). Writing analysis report ... done (0.260s). Max. memory usage : 171.227MB Total processing time : 3.436s S=C=A=N: Sat Jul 2 16:21:16 2016: Analyze done (status=0) 5s

 Continues with automatic (parallel) analysis of trace files

# **BT-MZ trace analysis report exploration**

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

% square scorep\_bt-mz\_B\_8x4\_trace INFO: Post-processing runtime summarization result... INFO: Post-processing trace analysis report... INFO: Displaying ./scorep\_bt-mz\_B\_8x4\_trace/trace.cubex...

[GUI showing trace analysis report]

# Post-processed trace analysis report



#### **Online metric description**



#### **Online metric description**



VIRTUAL INSTITUTE – HIGH PRODUCTIVITY SUPERCOMPUTING

# **Critical-path analysis**



VIRTUAL INSTITUTE – HIGH PRODUCTIVITY SUPERCOMPUTING

## **Critical-path analysis**



#### Pattern instance statistics



#### **Connect to Vampir trace browser**



#### Show most severe pattern instances



#### Investigate most severe instance in Vampir



VIRTUAL INSTITUTE – HIGH PRODUCTIVITY SUPERCOMPUTING

# **Further information**

# 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 3-clause BSD open-source license
- Sources, documentation & publications:
  - http://www.scalasca.org
  - mailto: scalasca@fz-juelich.de

