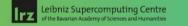
Scalasca Trace Tools

Demo/Hands-on: Automatic trace analysis

























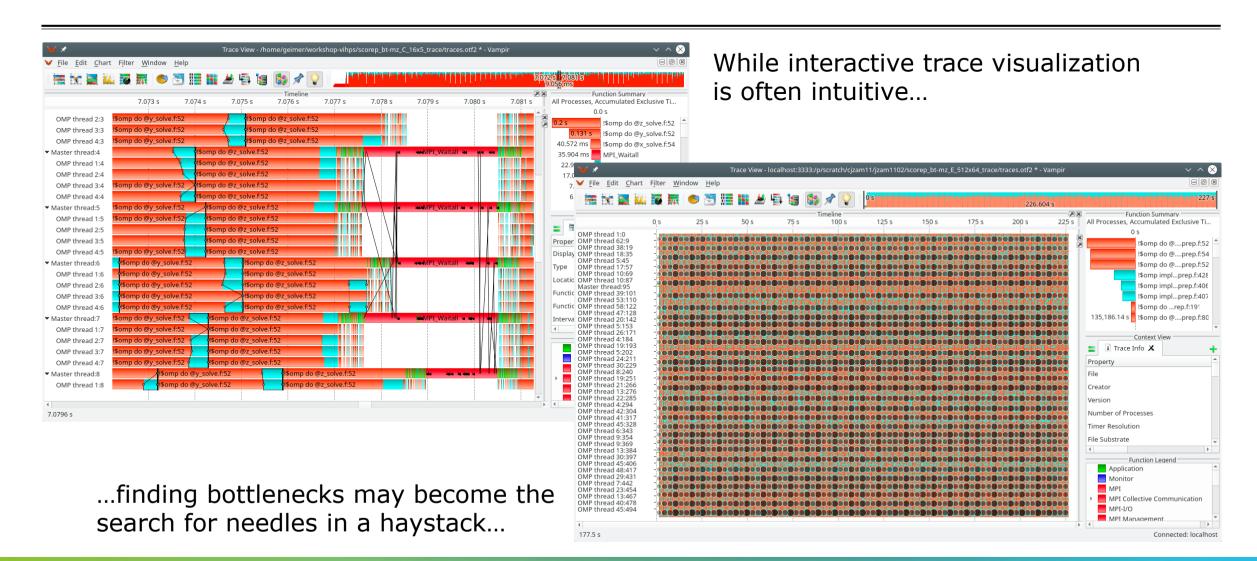






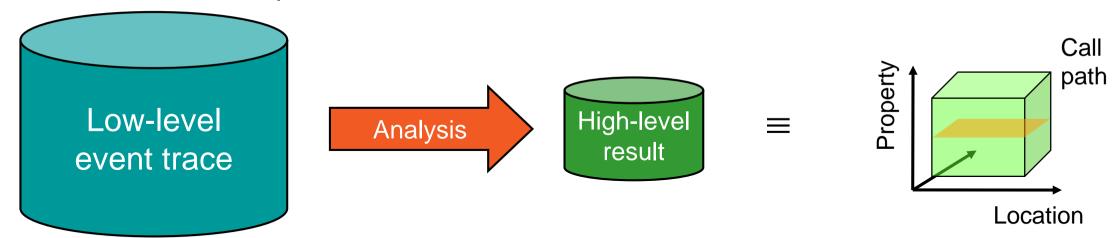


Motivation



Idea: Automatic trace analysis

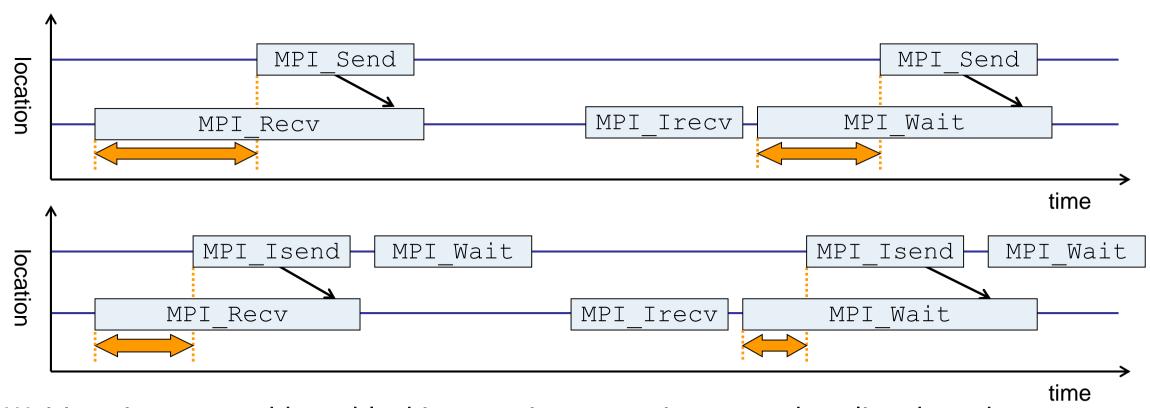
- Automatic search for patterns of inefficient behavior
- Classification of behavior & quantification of significance
- Identification of delays as root causes of inefficiencies



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



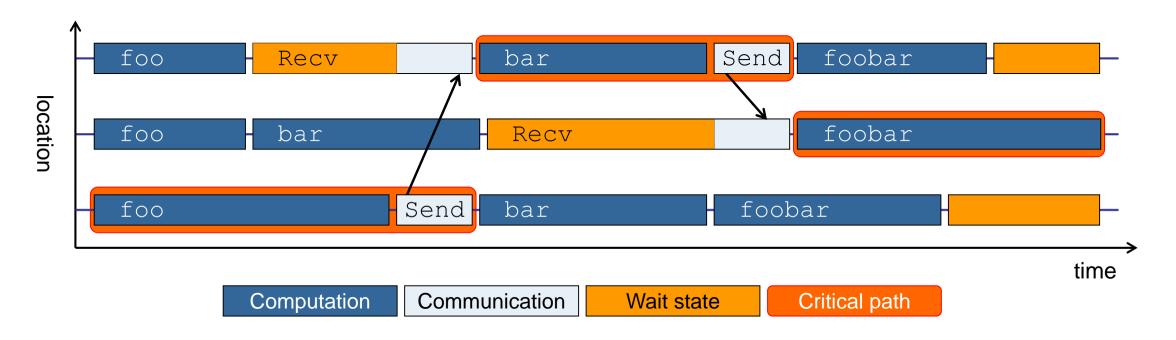
Example: "Late Sender" wait state



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



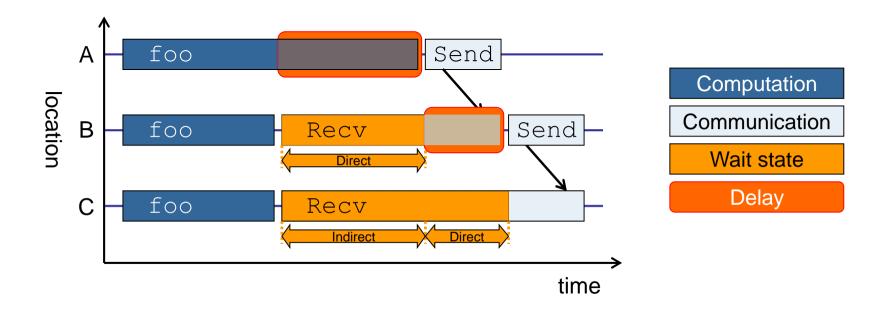
Example: Critical path



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



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



Scalasca command - One command for (almost) everything

```
% scalasca
Scalasca 2.6.1
Toolset for scalable performance analysis of large-scale parallel applications
usage: scalasca [OPTION]... ACTION <argument>...
    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
Options:
  -c, --show-config
                         show configuration summary and exit
  -h, --help
                         show this help and exit
   -n, --dry-run
                         show actions without taking them
                         show quick reference quide and exit
      --auickref
      --remap-specfile show path to remapper specification file and exit
   -v, --verbose
                         enable verbose commentary
                         show version information and exit
   -V, --version
```

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



Scalasca convenience command: scan / scalasca -analyze

```
% scan
Scalasca 2.6.1: measurement collection & analysis nexus
usage: scan {options} [launchcmd [launchargs]] target [targetargs]
      where {options} may include:
       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.
       Analvze
                  : 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.
  -R #runs
                  : Specify the number of measurement runs per config.
  -M cfafile
                  : Specify a config file for a multi-run measurement.
```

Scalasca measurement collection & analysis nexus



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)



Scalasca advanced command: scout - Scalasca automatic trace analyzer



```
% scout.hvb --help
SCOUT (Scalasca 2.6.1)
Copyright(c) 1998-2022 Forschungszentrum Juelich GmbH
Copyright(c) 2014-2021 RWTH Aachen University
Copyright(c) 2009-2014 German Research School for Simulation Sciences GmbH
Usage: <launchcmd> scout.hyb [OPTION]... <ANCHORFILE | EPIK DIRECTORY>
Options:
  --statistics
                    Enables instance tracking and statistics [default]
                     Disables instance tracking and statistics
  --no-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
                     Enables enhanced timestamp correction
  --time-correct
  --no-time-correct
                     Disables enhanced timestamp correction [default]
  --verbose, -v
                     Increase verbosity
                     Display this information and exit
  --help
```

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



Scalasca convenience command: square / scalasca -examine

```
% square
Scalasca 2.6.1: analysis report explorer
usage: square [OPTIONS] <experiment archive | cube file>
   -c <none | quick | full> : Level of sanity checks for newly created reports
                            : Force remapping of already existing reports
   -F
  -f filtfile
                            : Use specified filter file when doing scoring (-s)
                            : Skip display and output textual score report
  -s
                            : Enable verbose mode
                            : Do not include idle thread metric
   -n
                            : Aggregation method for summarization results of
   -S <mean | merge>
                              each configuration (default: merge)
                            : Aggregation method for trace analysis results of
   -T <mean | merge>
                              each configuration (default: merge)
                            : Post-process every step of a multi-run experiment
   -A
```

Scalasca analysis report explorer (Cube)



Toolchain, Score-P, and Scalasca modules (DINE)

Select modules for the Intel + IntelMPI tool chain

```
% module load intel_comp/2020-update2 intel_mpi/2020-update2
```

- Load Score-P, Scalasca and Cube modules
 - Score-P & Scalasca installations are toolchain specific!

```
% module load scorep/8.4 scalasca/2.6.1 cube/4.8.2
```

BT-MZ summary measurement collection...

```
% cd bin.scorep
% cp ../jobscript/dine/scan.sbatch .
% cat scan sbatch
# set up environment
module purge
module load intel comp/2020-update2 intel mpi/2020-update2
module load scalasca/2.6.1 scorep/8.4
# measurement configuration
export SCOREP FILTERING FILE=../config/scorep.filt
#export SCOREP TOTAL MEMORY=100M
#export SCOREP METRIC PAPI=PAPI TOT INS, PAPI TOT CYC, ...
#export SCAN ANALYZE OPTS="-time-correct"
set -x
export OMP NUM THREADS=6
scan -s mpiexec -np 8 ./bt-mz C.8
```

Change to
 directory with the
 Score-P
 instrumented
 executable and
 edit the job script

```
Hint:
scan = scalasca -analyze
-s = profile/summary (def)
```

Submit the job

% sbatch scan.sbatch

BT-MZ summary measurement

```
S=C=A=N: Scalasca 2.6.1 runtime summarization
S=C=A=N: ./scorep bt-mz C 8x6 sum experiment archive
S=C=A=N: Sat Apr 20 10:11:19 2024: Collect start
mpiexec ./bt-mz C.8
 NAS Parallel Benchmarks (NPB3.3-MZ-MPI) -
    BT-MZ MPI+OpenMP Benchmark
 Number of zones: 16 x 16
 Iterations: 200 dt: 0.000100
 Number of active processes:
 [... More application output ...]
S=C=A=N: Sat Apr 20 10:11:39 2024: Collect done (status=0) 20s
S=C=A=N: ./scorep bt-mz C 8x6 sum complete.
```

- Run the application using the Scalasca measurement collection & analysis nexus prefixed to launch command
- Creates experiment directory:scorep_bt-mz_C_8x6_sum

BT-MZ summary analysis report examination

Score summary analysis report

```
% square -s scorep_bt-mz_C_8x6_sum
INFO: Post-processing runtime summarization report (profile.cubex)...
INFO: Score report written to ./scorep_bt-mz_C_8x6_sum/scorep.score
```

Post-processing and interactive exploration with Cube

```
% square scorep_bt-mz_C_8x6_sum
INFO: Displaying ./scorep_bt-mz_C_8x6_sum/summary.cubex...

[GUI showing summary analysis report]
```

Hint:

Copy 'summary.cubex' to local system (laptop) using 'scp' to improve responsiveness of GUI

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



Post-processed summary analysis report

CubeGUI-4.4.3: scorep tea leaf baseline 8x12 sum/summary.cubex <@irl11> Display Plugins Help Restore Setting ▼ Save Settings Absolute Absolute Absolute Statistics Sunburst Flat view System tree Metric tree Call tree ¬ □ 0.00 machine Linux 0.00 tea leaf baseline → □ 0.00 Execution ■ 0.03 MAIN Split base metrics into → □ 0.00 MPÍ Rank 0 ▶ ■ 0.00 tea module.tea init comms 8478.33 Computation 35.30 Master thread ▶ ■ 0.00 !\$omp parallel @tea leaf.f90:45 ¬ □ 0 00 MPI more specific metrics 35.28 OMP thread 1 ▶ ■ 7.97 Management ▶ ■ 2.15 initialise 35.29 OMP thread 2 0.38 Synchronization ▼ ■ 0.00 diffuse 35.29 OMP thread 3 → □ 0.00 Communication 0.00 timer 35.28 OMP thread 4 ▶ ■ 0.06 set field module.set field ■ 12.08 Point-to-point 35.29 OMP thread 5 ▶ ■ 0.00 timestep module.timestep 19.82 Collective 35.28 OMP thread 6 ▼ ■ 0.75 tea leaf module.tea leaf □ 0.00 One-sided 35.29 OMP thread 7 ▶ ■ 0.26 timer → □ 0.00 File I/O 35.29 OMP thread 8 ▶ ■ 51.30 update halo module.update h 35.29 OMP thread 9 → ■ 6.22 tea leaf kernel cg module.tea ▶ ■ 1.03 tea module.tea_allsum 35.29 OMP thread 10 → ■ 0.61 tea leaf kernel cheby module. 35.30 OMP thread 11 □ 0.00 Explicit ■ 114.81 Implicit ▼ ■ 1.24 tea leaf kernel cg module.tea ■ 1.69 !\$omp parallel @tea leaf cg. 35.59 Master thread □ 0.00 Critical - ■ 3421.11 !\$omp do @tea leaf c 35.58 OMP thread 1 □ 0.00 Lock API 35.58 OMP thread 2 □ 0.00 !\$omp implicit barrier @ □ 0.00 Ordered 35.58 OMP thread 3 □ 0.00 Task Wait □ 0.00 !\$omp implicit barrier @te 35.58 OMP thread 4 - ■ 2.43 tea leaf kernel ca module.tea □ 0.00 Flush 35.58 OMP thread 5 ▼ ■ 2.01 !\$omp parallel @tea leaf cg. □ 0.00 Overhead 35.59 OMP thread 6 → ■ 3402.24 !\$omp do @tea leaf co ▶ ■ 656.11 Idle threads 35.59 OMP thread 7 1.17e8 Visits (occ) □ 0.00 !somp implicit barrier @ 35.58 OMP thread 8 □ 0.00 !\$omp implicit barrier @te 2.37e10 Bytes transferred (bytes) □ 0 MPI file operations (occ) 2.04 tea leaf kernel cg module.tea All (96 elements) 9289.50 0.00 8478.33 (91.27%) 3421.11 (40.35%) 8478.33 0.00 0.00 (0.00%) 3421.11

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 analysis & report examination



BT-MZ trace measurement collection...

```
% cd bin.scorep
% cp ../iobscript/dine/scan.sbatch .
% vim scan sbatch
# set up environment
module purge
module load intel comp/2020-update2 intel mpi/2020-update2
module load scalasca/2.6.1 scorep/8.4
# measurement configuration
export SCOREP FILTERING FILE=../config/scorep.filt
export SCOREP TOTAL MEMORY=100M
#export SCOREP METRIC PAPI=PAPI TOT INS, PAPI TOT CYC, ...
export SCAN ANALYZE OPTS="-time-correct"
set -x
export OMP NUM THREADS=6
scan -t mpiexec -np 8 ./bt-mz C.8
```

 Change to directory with the Score-P instrumented executable and edit the job script

- Use "-t" with the scan command
- Submit the job

% sbatch scan.sbatch



BT-MZ trace measurement ... collection

```
S=C=A=N: Scalasca 2.6.1 trace collection and analysis
S=C=A=N: Sat Apr 20 10:16:41 2024: Collect start
        ./bt-mz C.8
mpiexec
NAS Parallel Benchmarks (NPB3.3-MZ-MPI) - BT-MZ MPI+OpenMP \
>Benchmark
 Number of zones: 16 x 16
 Iterations: 200 dt: 0.000100
 Number of active processes:
 [... More application output ...]
S=C=A=N: Sat Apr 20 10:17:03 2024: Collect done (status=0) 22s
```

 Starts measurement with collection of trace files ...



BT-MZ trace measurement ... analysis

```
S=C=A=N: Sat Apr 20 10:17:03 2024: Analyze start mpiexec scout.hyb ./scorep bt-mz C 8x6 trace/traces.otf2
SCOUT
       (Scalasca 2.6.1)
Analyzing experiment archive ./scorep bt-mz C 8x6 trace/traces.otf2
Opening experiment archive ... done (0.007s).
Reading definition data ... done (0.009s).
Reading event trace data ... done (0.209s).
Preprocessing ... done (0.565s).
Timestamp correction ... done (1.316s).
Analyzing trace data ... done (15.125s).
Writing analysis report
                               ... done (0.227s).
                               : 996.031MB
Max. memory usage
Total processing time : 17.508s
S=C=A=N: Sat Apr 20 10:17:22 2024: Analyze done (status=0) 19s
```

 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_C_8x6_trace
INFO: Post-processing runtime summarization report (profile.cubex)...
INFO: Post-processing trace analysis report (scout.cubex)...
INFO: Displaying ./scorep_bt-mz_C_8x6_trace/trace.cubex...

[GUI showing trace analysis report]
```

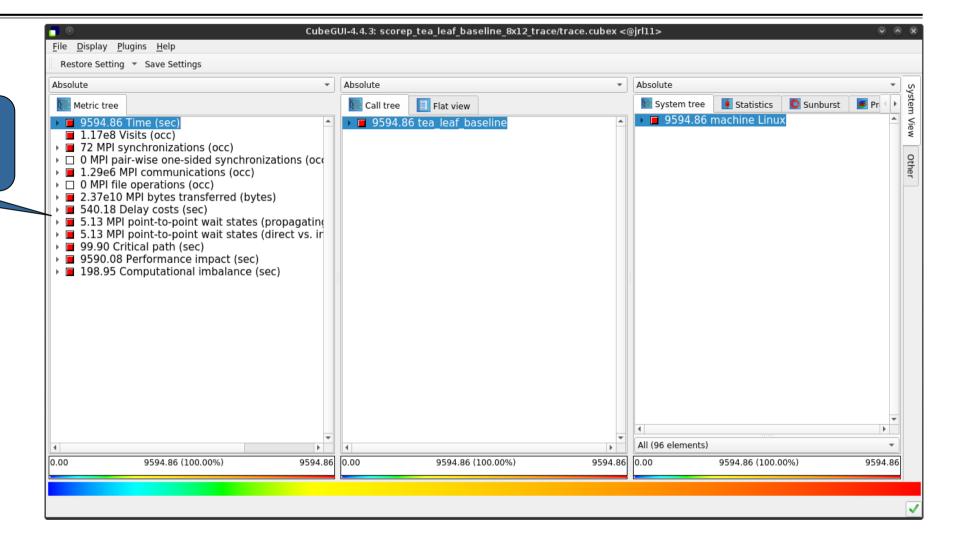
Hint:

Run 'square -s' first and then copy 'trace.cubex' to local system (laptop) using 'scp' to improve responsiveness of GUI



Scalasca analysis report exploration (opening view)

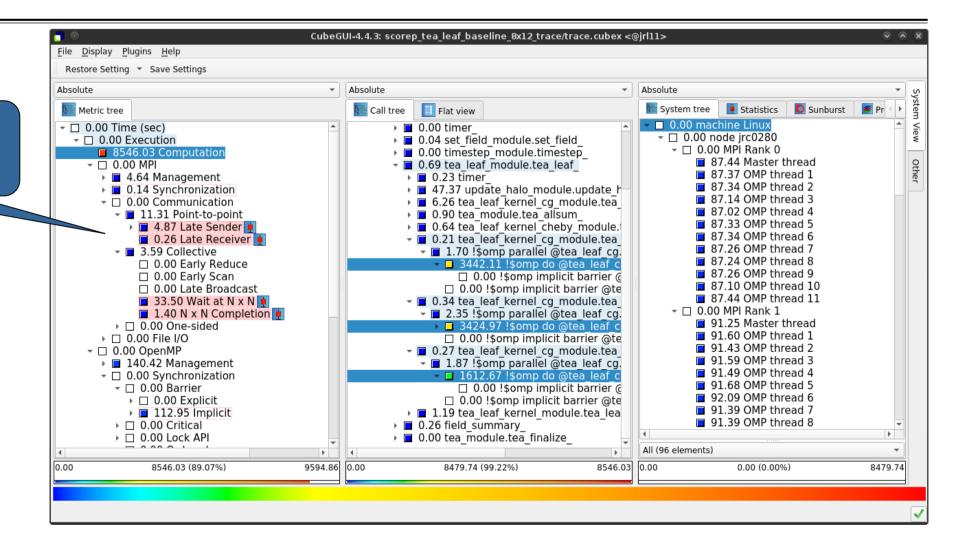
Additional top-level metrics produced by the trace analysis...





Scalasca wait-state metrics

...plus additional waitstate metrics as part of the "Time" hierarchy

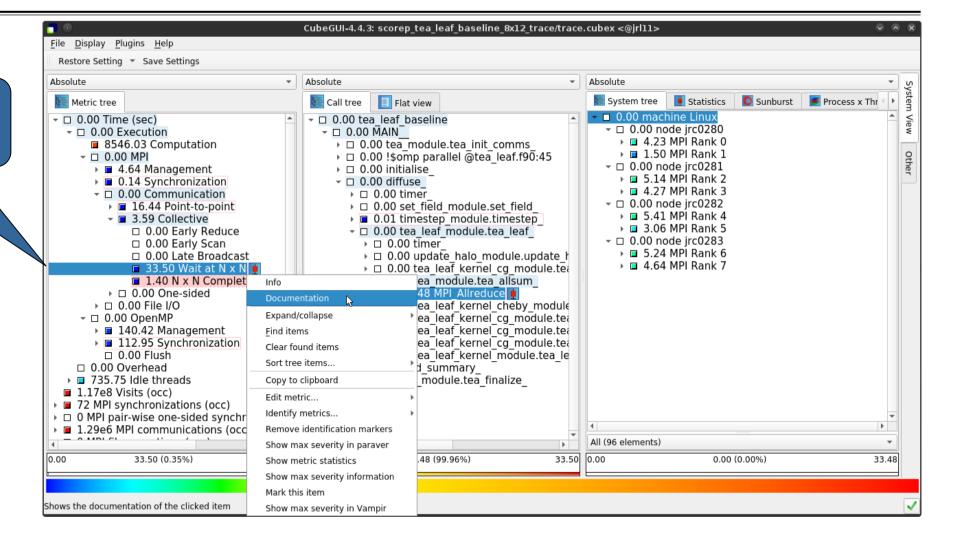




Online metric description



Access online metric description via context menu (right-click)

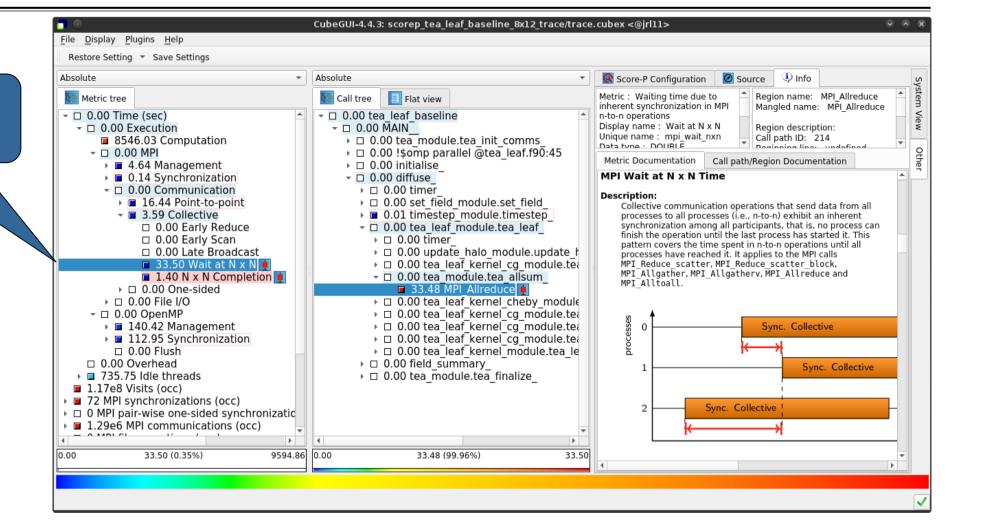




Online metric description (cont.)



Selection of different metric automatically updates description

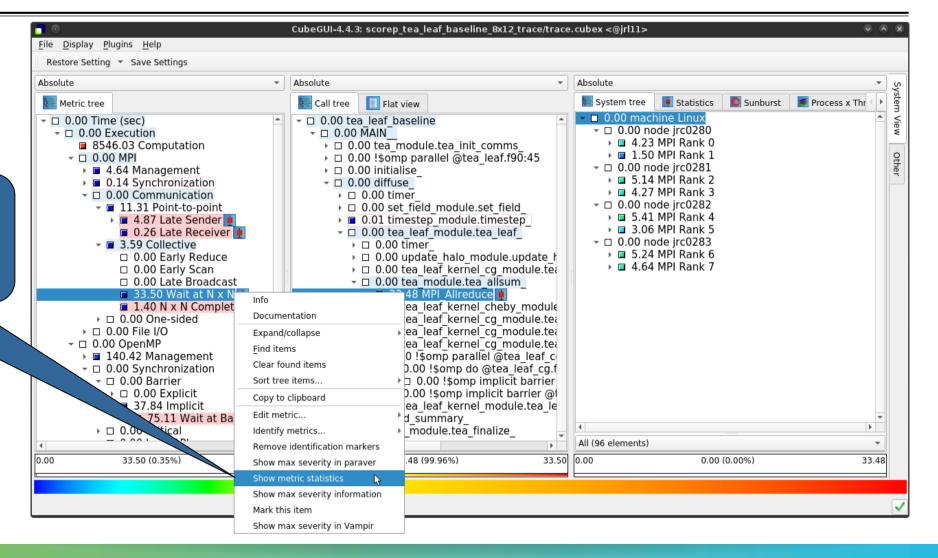




Metric statistics



Access metric statistics for metrics marked with box plot icon from context menu

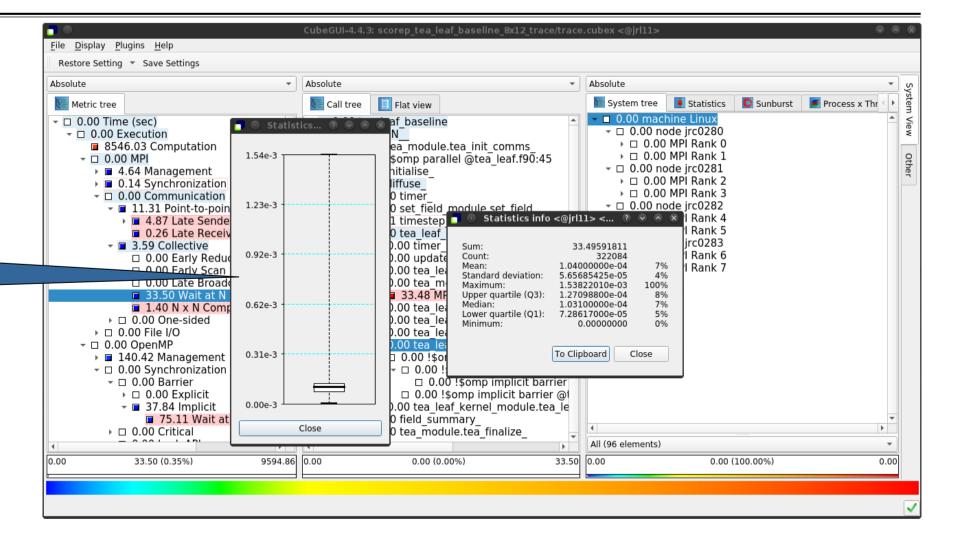


VI-HPS

Metric statistics (cont.)



Shows instance statistics box plot, click to get details

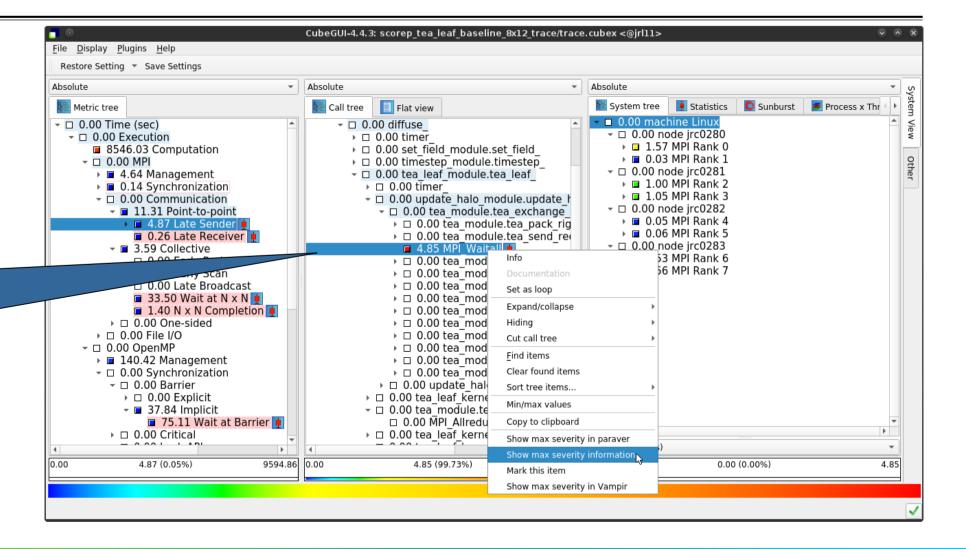




Metric instance statistics



Access most-severe instance information for call paths marked with box plot icon via context menu



Statistics

Max severity <@irl11>

Sunburst Process x Thr

To Clipboard

0.00 (0.00%)

(P) (Q) (A) (X

Close



Metric instance statistics (cont.)

File Display Plugins Help Restore Setting ▼ Save Settings

8546.03 Computation

→ ■ 4.64 Management

→ □ 0.00 One-sided

▶ ■ 140.42 Management

→ □ 0.00 Synchronization

→ □ 0.00 Explicit

4.87 (0.05%)

→ **37.84** Implicit

75.11 Wait at Barrier

→ □ 0.00 Barrier

→ □ 0.00 Critical _ ^ ^ ^ _ _ _ _ _ _

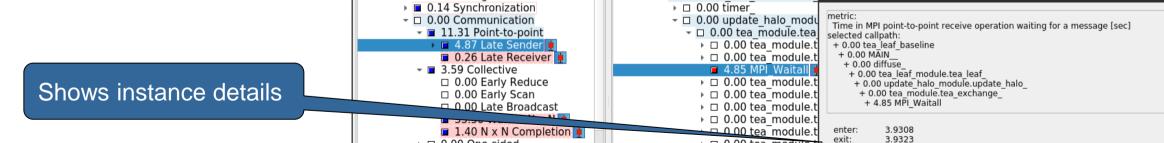
→ □ 0.00 File I/O

→ □ 0.00 OpenMP

Absolute

Metric tree





9594.86

0.00

▼ Absolute

Call tree

CubeGUI-4.4.3; scorep tea leaf baseline 8x12 trace/trace.cubex <@irl11>

Flat view

→ □ 0.00 set field module.set field

→ □ 0.00 timestep module.timestep

→ □ 0.00 tea module.t

→ □ 0.00 tea module.tea allsum

→ □ 0.00 tea leaf kernel cheby module

□ 0.00 MPI Allreduce

4.85 (99.73%)

→ □ 0.00 update halo ke → □ 0.00 tea leaf kernel cd

→ □ 0.00 timer

Absolute

0.0015

0.0014

4.87 0.00

All (96 elements)

duration:

severity:

■ ■ 0.00 machine Linux

¬ □ 0.00 node irc0280

→ □ 1.57 MPÍ Rank 0

▶ ■ 0.03 MPI Rank 1

4.85



Further information

- Collection of trace-based performance tools
 - Specifically designed for large-scale systems
 - Features an automatic trace analyzer providing wait-state, critical-path, and delay analysis
 - Supports MPI, OpenMP, POSIX threads, and hybrid MPI+OpenMP/Pthreads
- Available under 3-clause BSD open-source license
- Documentation & sources:
 - https://www.scalasca.org
- Contact:
 - mailto: scalasca@fz-juelich.de

