

# Automatic trace analysis with the Scalasca Trace Tools

Marc Schlütter Jülich Supercomputing Centre





# Scalasca Trace Tools

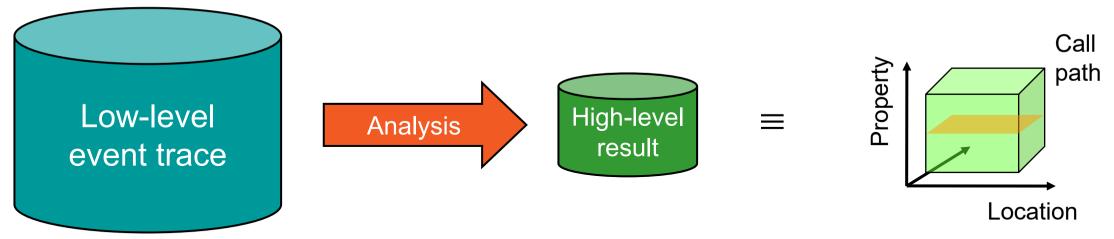
DOI 10.5281/zenodo.15125898

- Scalable trace-based performance analysis toolset for the most popular
  - parallel programming paradigms
  - Current focus: MPI, OpenMP, and (to a limited extend) POSIX threads
  - Analysis of traces including only host-side events from applications using CUDA, OpenCL, or OpenACC (also in combination with MPI and/or OpenMP) is possible, but results need to be interpreted with some care
- Specifically targeting large-scale parallel applications
  - Demonstrated scalability up to 1.8 million parallel threads
  - Of course also works at small/medium scale
- Latest release:
  - Scalasca Trace Tools v2.6.2 (April 2025)

#### **Automatic trace analysis**

#### Idea

- 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

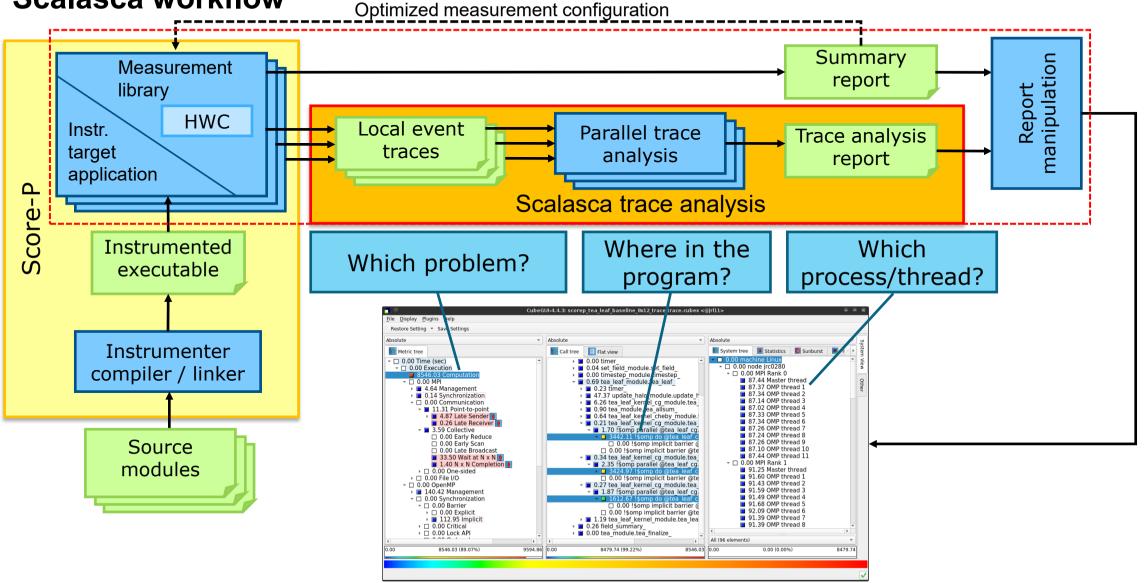
## Scalasca Trace Tools: Features

- Open source, 3-clause BSD license
- Supports all major HPC platforms
- Uses Score-P instrumenter & measurement libraries
  - Scalasca v2 core package focuses on trace-based analyses
  - Provides convenience commands for measurement, analysis, and post-processing
  - Supports common data formats
    - Reads event traces in OTF2 format
    - Writes analysis reports in CUBE4 format

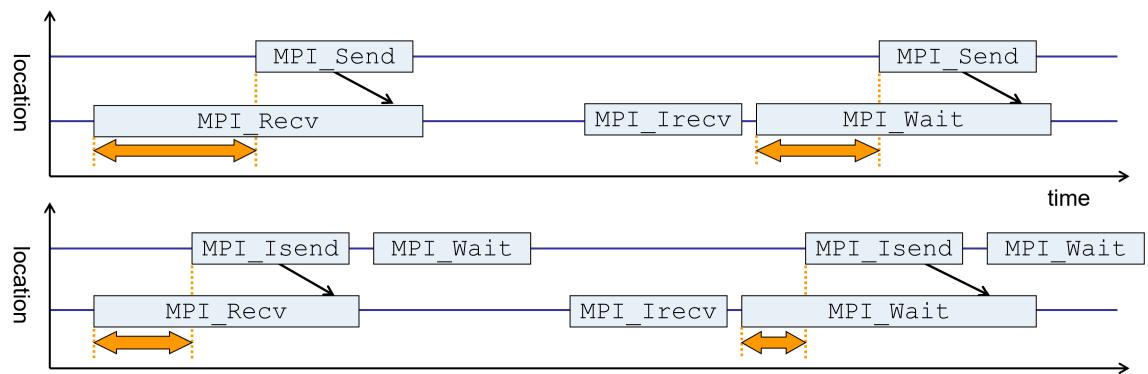
#### Current limitations:

- Unable to handle traces ...
  - with MPI thread level exceeding MPI\_THREAD\_FUNNELED
  - containing memory events, CUDA/OpenCL device events (kernel, memcpy), SHMEM, or OpenMP nested parallelism
- PAPI/rusage metrics for trace events are ignored

#### Scalasca workflow



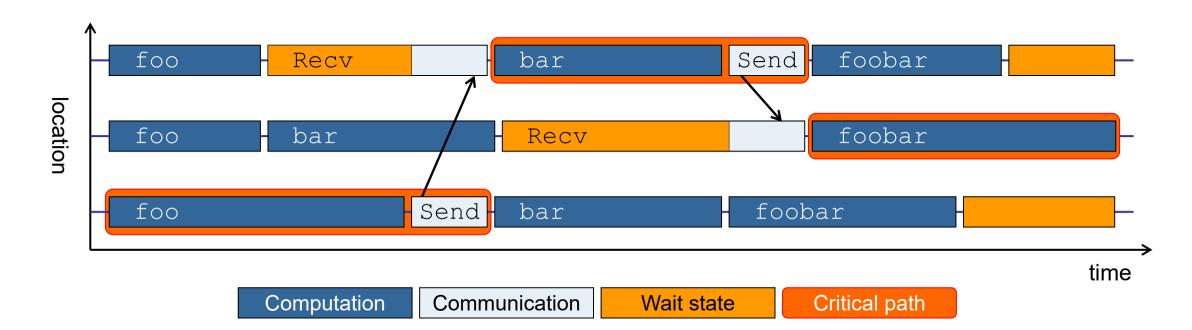
#### Example: "Late Sender" wait state



time

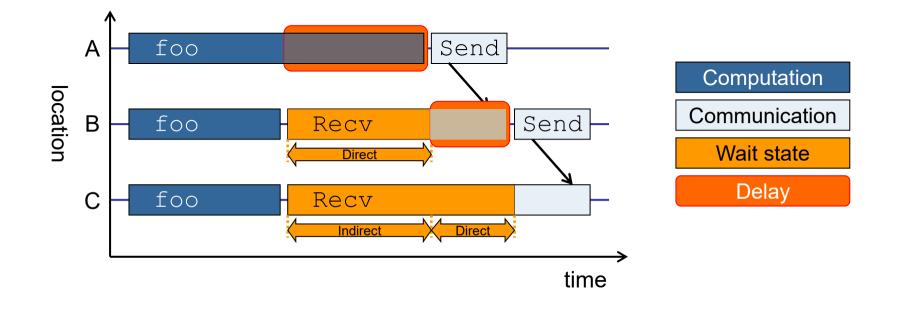
- Waiting time caused by a blocking receive operation posted earlier than the corresponding send
- Applies to blocking as well as non-blocking communication

## **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

#### **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





#### **Recap: Setup for exercises**

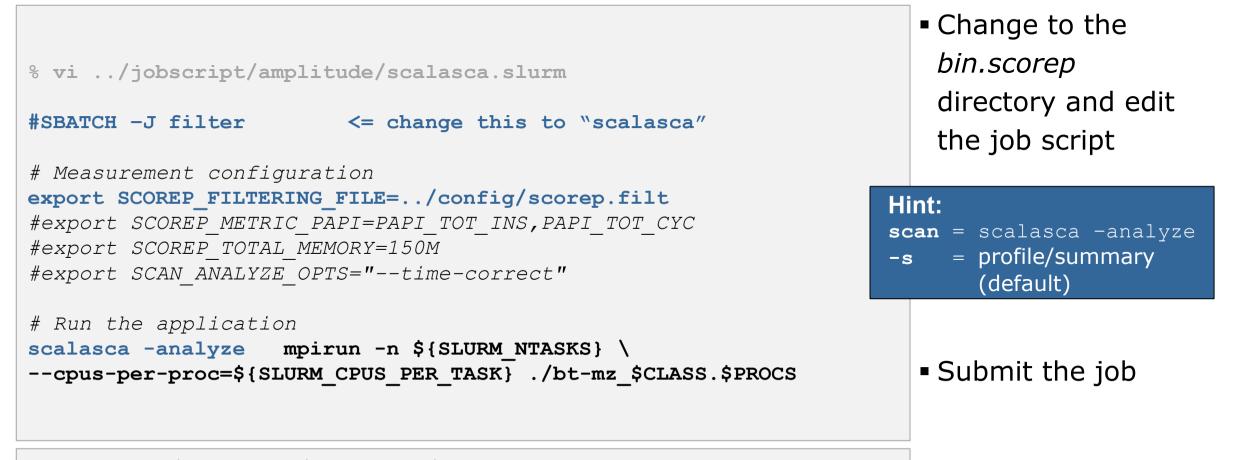
- Connect to your account on the local training system
  - Amplitude:

% ssh -Y <account>@login1.amplitude.uni-due.de

- Change to directory containing BT-MZ sources
  - Existing instrumented executable can be reused

% module use /cluster/vi-hps\_tuning\_workshop/examples/opt/modules % module load scalasca/2.6.2-gcc-openmpi scorep/9.0-gcc-openmpi-cuda cube % cd <your bt-mz path>/bin.scorep

#### **Demo: BT-MZ summary measurement collection...**



% sbatch ../jobscript/amplitude/scalasca.slurm

# scan: 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)

#### **Demo: BT-MZ summary measurement**

```
S=C=A=N: Scalasca 2.6.2 runtime summarization
S=C=A=N: scalasca/scorep_bt-mz_C_8x8_sum experiment archive
S=C=A=N: Fri Feb 23 11:54:48 2024: Collect start
srun ... bin.scorep/bt-mz_C.8
```

```
NAS Parallel Benchmarks (NPB3.3-MZ-MPI) -
BT-MZ MPI+OpenMP Benchmark
```

Number of zones: 8 x 8 Iterations: 200 dt: 0.000100 Number of active processes: 8

```
[... More application output ...]
```

S=C=A=N: Fri Feb 23 11:55:09 2024: Collect done (status=0) 21s
S=C=A=N: scalasca/scorep-8-8-summary complete.

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

 Creates experiment directory: scorep\_bt-mz\_C\_8x8\_sum

#### **Demo: BT-MZ summary analysis report examination**

#### Score summary analysis report

% square -s scorep\_bt-mz\_C\_8x8\_sum
INFO: Post-processing runtime summarization report (profile.cubex)...
INFO: Score report written to scorep bt-mz C 8x8 sum/scorep.score

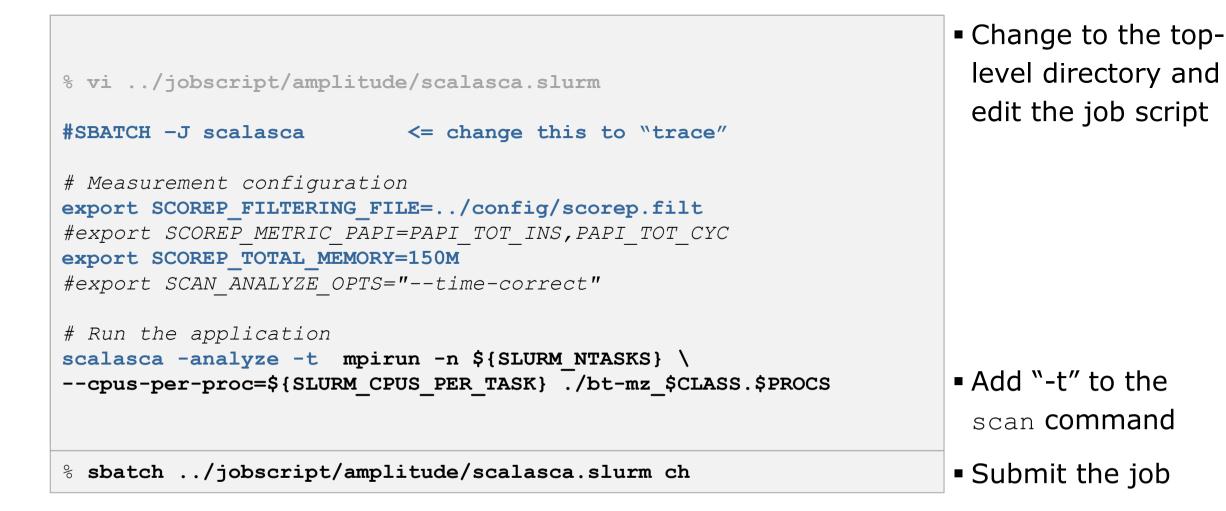
#### Post-processing and interactive exploration with Cube

% square scorep\_bt-mz\_C\_8x8\_sum INFO: Displaying scorep bt-mz C 8x8 sum/summary.cubex Hint: Copy 'summary.cubex' to local system (laptop) using 'scp' to improve responsiveness of GUI

[GUI showing summary analysis report]

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

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



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

S=C=A=N: Scalasca 2.6.2 trace collection and analysis S=C=A=N: Fri Feb 23 12:49:25 2024: Collect start srun ... bin.scorep/bt-mz\_C.8

NAS Parallel Benchmarks (NPB3.3-MZ-MPI) - BT-MZ MPI+OpenMP \ >Benchmark

Number of zones: 8 x 8 Iterations: 200 dt: 0.000100 Number of active processes: 8

[... More application output ...]

S=C=A=N: Fri Feb 23 12:49:45 2024: Collect done (status=0) 20s

 Starts measurement with collection of trace files ...

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

```
S=C=A=N: Fri Feb 23 12:49:45 2024: Analyze start
srun [...] scout.hyb --time-correct \
> scorep bt-mz C 8x8 trace/traces.otf2
SCOUT (Scalasca 2.6.1)
Analyzing experiment archive scorep bt-mz C 8x8 trace/traces.otf2
Opening experiment archive ... done (0.002s).
Reading definition data<br/>Reading event trace data... done (0.002s).Preprocessing<br/>Timestamp correction<br/>Analyzing trace data... done (1.117s).<br/>... done (0.729s).<br/>... done (4.370s).<br/>... done (23.496s).<br/>... done (0.284s).
                                                 : 3048.270MB
Max. memory usage
               # passes : 1
# violated : 0
Total processing time : 30.087s
S=C=A=N: Fri Feb 23 12:50:21 2024: Analyze done (status=0) 36s
```

 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\_8x8\_trace INFO: Post-processing runtime summarization report (profile.cubex)... INFO: Post-processing trace analysis report (scout.cubex)... INFO: Displaying scorep\_bt-mz\_C\_8x8\_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



# Case study: TeaLeaf MPI+OpenMP





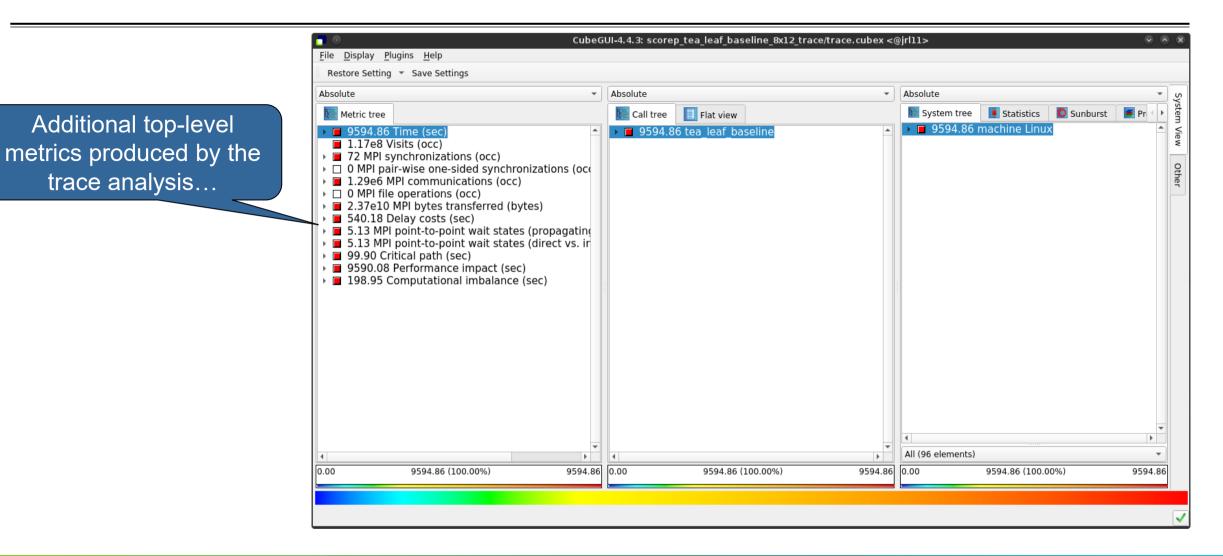
# **Case study: TeaLeaf MPI+OpenMP**

- HPC mini-app developed by the UK Mini-App Consortium
  - Solves the linear 2D heat conduction equation on a spatially decomposed regular grid using a 5 point stencil with implicit solvers
  - Part of the Mantevo 3.0 suite
  - Available on GitHub: https://uk-mac.github.io/TeaLeaf/
- Measurements of TeaLeaf reference v1.0 taken on Jureca cluster @ JSC
  - Using Intel 19.0.3 compilers, Intel MPI 2019.3, Score-P 5.0, and Scalasca 2.5
  - Run configuration
    - 8 MPI ranks with 12 OpenMP threads each
    - Distributed across 4 compute nodes (2 ranks per node)
    - Test problem "5": 4000 × 4000 cells, CG solver

% cube scorep\_tea\_leaf\_baseline\_8x12\_trace/trace.cubex

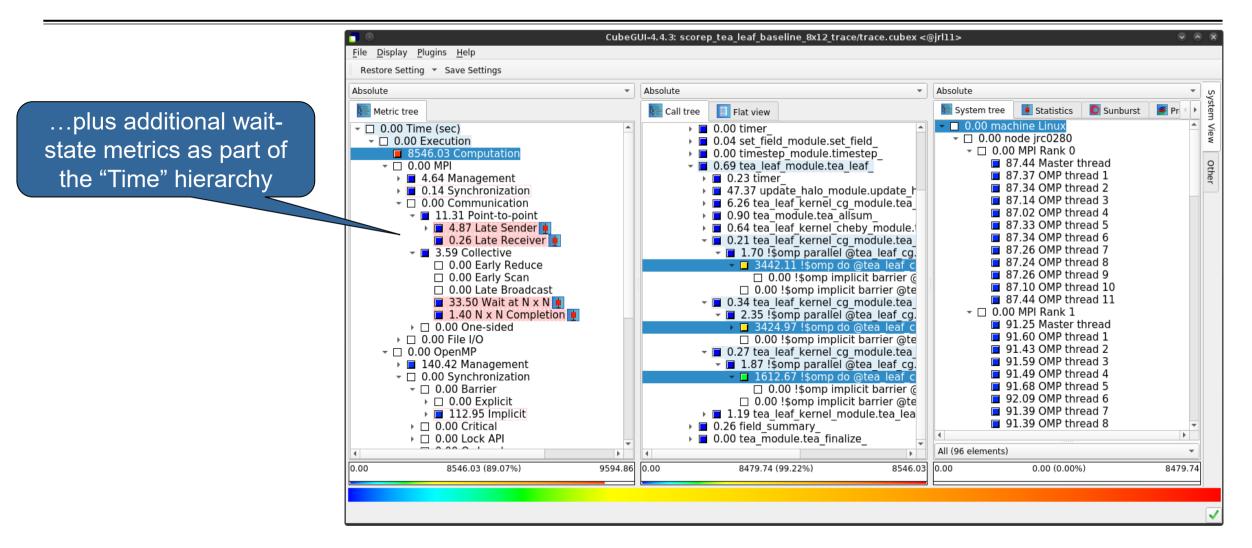
[GUI showing post-processed trace analysis report]

# Scalasca analysis report exploration (opening view)



#### **Scalasca wait-state metrics**

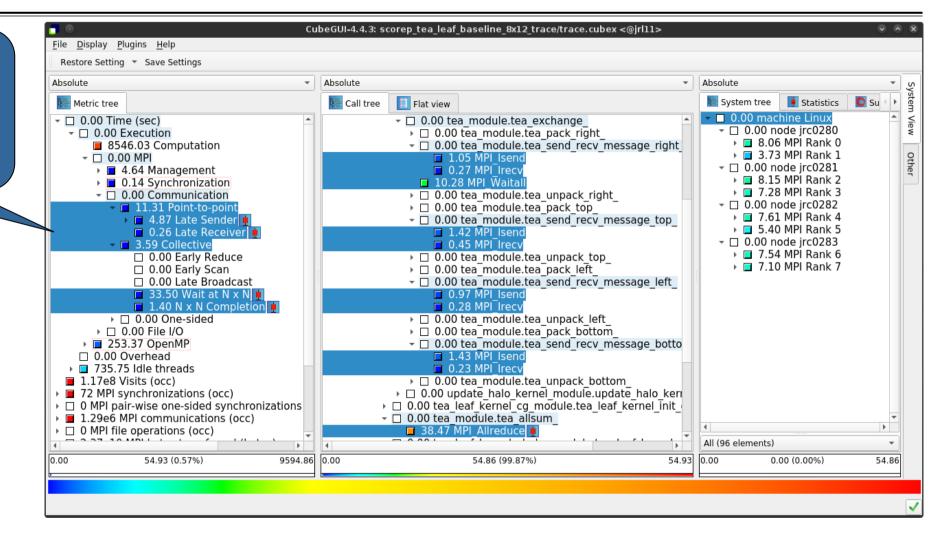




# **TeaLeaf Scalasca report analysis (I)**



While MPI communication time and wait states are small (~0.6% of the total execution time)...



# **TeaLeaf Scalasca report analysis (II)**



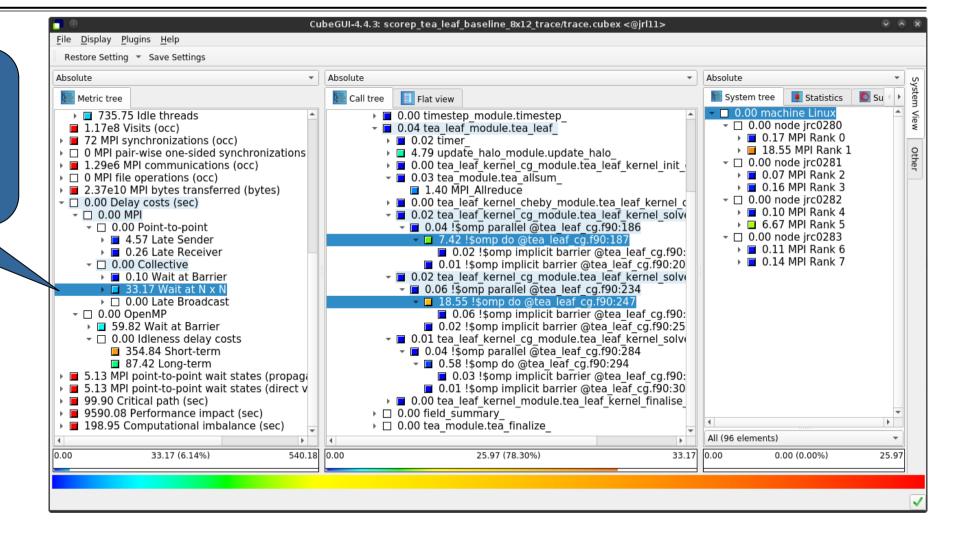
CubeGUI-4.4.3: scorep tea leaf baseline 8x12 trace/trace.cubex <@irl11> File Display Plugins Help Restore Setting \* Save Settings Absolute Absolute Absolute -Flat view Statistics O Su 🔚 Call tree 🔚 System tree Metric tree 0.00 machine Linux 735.75 Idle threads - 🗖 7.00 tea module.tea exchange ▶ ■ 4.76 tea module.tea pack right → □ 0.00 node irc0280 1.17e8 Visits (occ) 24.77 MPI Rank 0 1.98 tea module.tea send recy message right 72 MPI synchronizations (occ) 24.21 MPI Rank 1 0 MPI pair-wise one-sided synchronizations othe 11.56 MPI Isend 1.29e6 MPI communications (occ) → □ 0.00 node irc0281 2.99 MPI Trecv 20.93 MPI Rank 2 I 0 MPI file operations (occ) 56.82 MPI Waital 21.55 MPI Rank 3 4.79 tea module.tea unpack right 2.37e10 MPI bytes transferred (bytes) - □ 0.00 node irc0282 6.85 tea module.tea pack top 0.00 Delay costs (sec) 23.46 MPI Rank 4 1.25 tea module.tea send recy message top → □ 0.00 MPL 24.15 MPI Rank 5 15.65 MPI Isend - 0.00 Point-to-point - 0.00 node jrc0283 4.57 Late Sender □ 4.95 MPI Irecv 19.39 MPI Rank 6 7.10 tea module.tea unpack top 0.26 Late Receiver 20.40 MPI Rank 7 4.87 tea module.tea pack left → □ 0.00 Collective 1.92 tea module.tea send recy message left 0.10 Wait at Barrier 33.17 Wait at N x N 10.63 MPI Isend D 0.00 Late Broadcast 3.13 MPI Trecv 0.00 OpenMP 6.98 tea module.tea pack bottom 59.82 Wait at Barrier 1.34 tea module.tea send recv message botto 0.00 Idleness delay costs 354.84 Short-term 15.69 MPI Isend 87.42 Long-term 2.55 MPI Trecv 6.96 tea module.tea unpack bottom 5.13 MPI point-to-point wait states (propage 3.83 update halo kernel module.update halo keri 5.13 MPI point-to-point wait states (direct v ▶ 99.90 Critical path (sec) > 3.55 tea leaf kernel cg module.tea leaf kernel init 9590.08 Performance impact (sec) 9.87 tea module.tea allsum Þ. 54.90 MPL Allreduce 198.95 Computational imbalance (sec) All (96 elements) 4 F 540.18 0.00 0.00 354.84 (65.69%) 178.86 (50.41%) 354.84 0.00 0.00 (0.00%) 178.86

...they directly cause a significant amount of the OpenMP thread idleness

# **TeaLeaf Scalasca report analysis (III)**



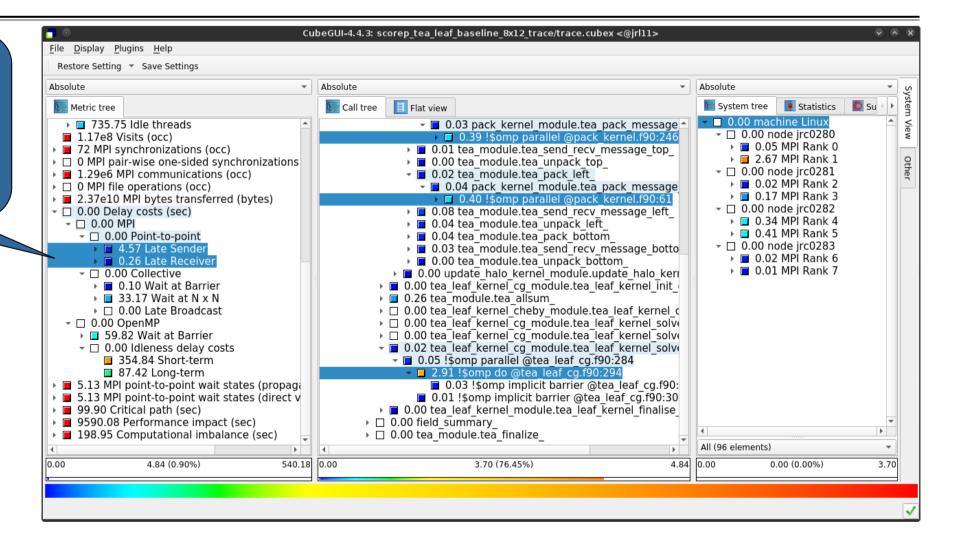
The "Wait at NxN" collective wait states are mostly caused by the first 2 OpenMP do loops of the solver (on ranks 5 & 1, resp.)...



# **TeaLeaf Scalasca report analysis (IV)**



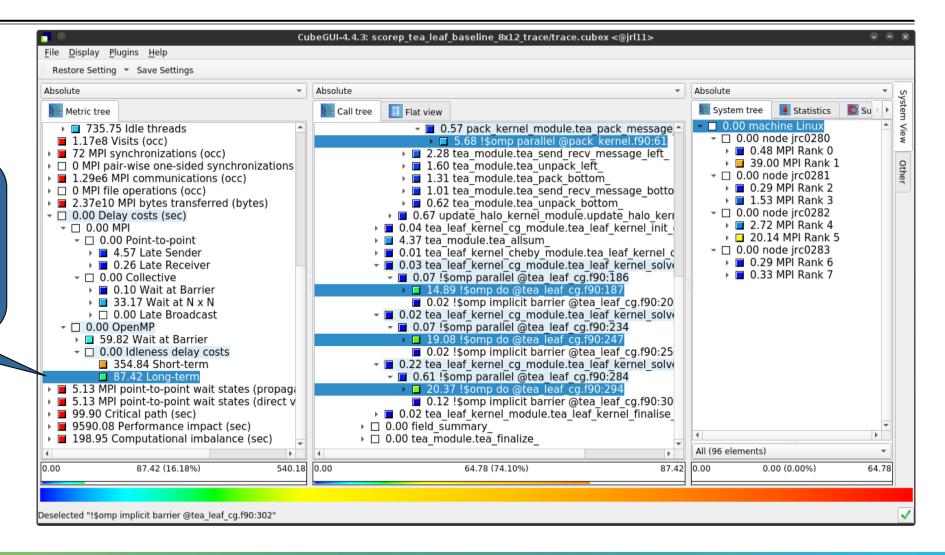
...while the MPI pointto-point wait states are caused by the 3<sup>rd</sup> solver do loop (on rank 1) and two loops in the halo exchange



# **TeaLeaf Scalasca report analysis (V)**

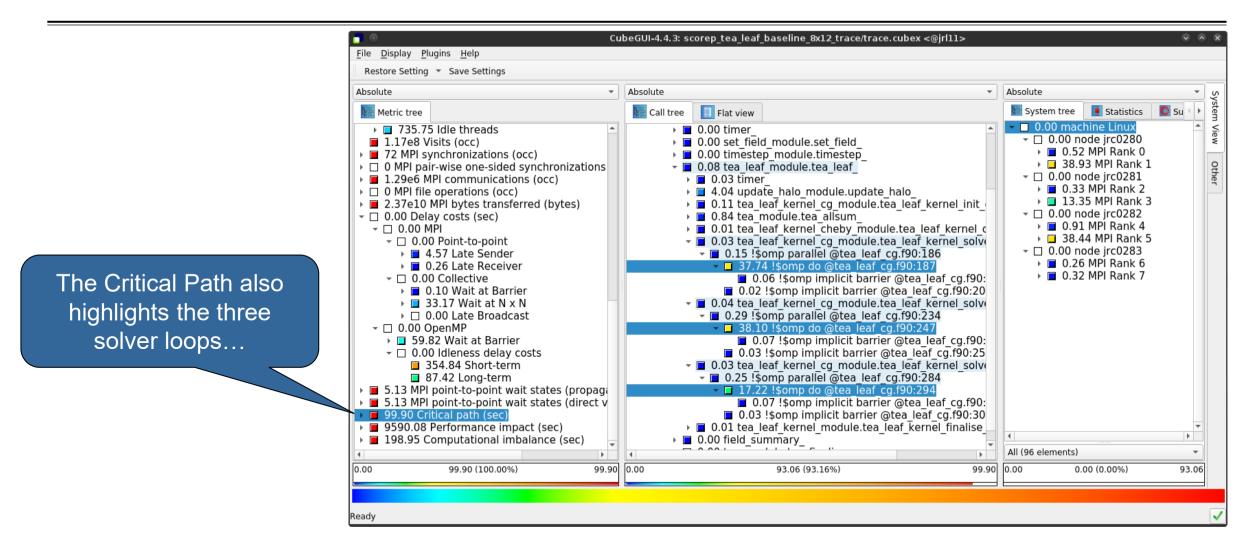


Various OpenMP do loops (incl. the solver loops) also cause OpenMP thread idleness on other ranks via propagation



# **TeaLeaf Scalasca report analysis (VI)**





# **TeaLeaf Scalasca report analysis (VII)**



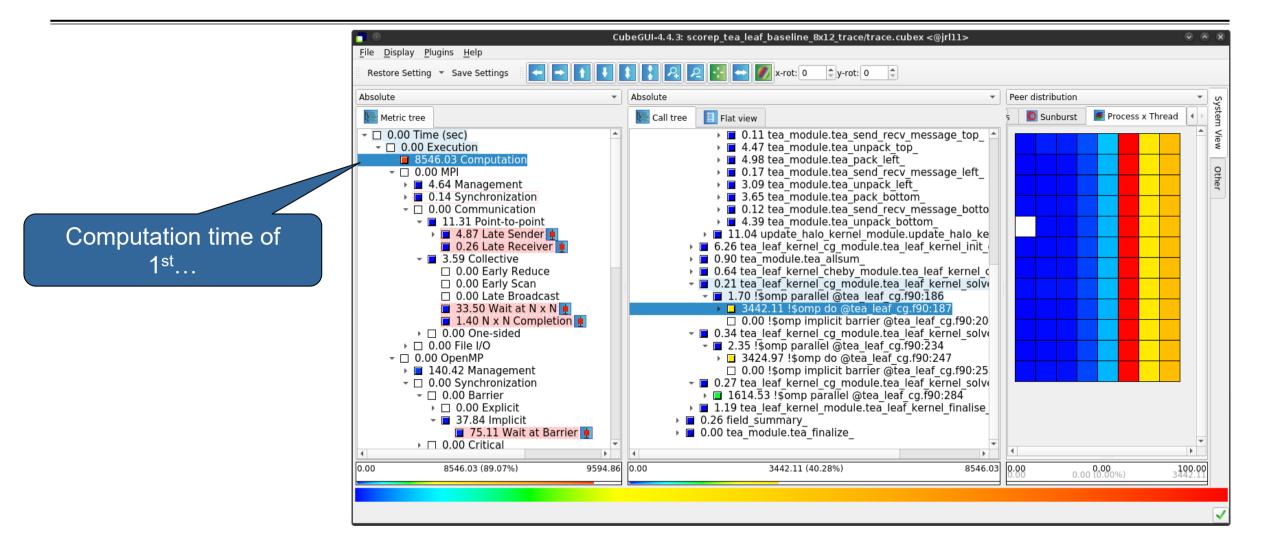
CubeGUI-4.4.3: scorep tea leaf baseline 8x12 trace/trace.cubex <@irl11> File Display Plugins Help Restore Setting \* Save Settings Absolute Absolute Absolute \* Flat view Statistics 🚺 Su 🔚 System tree Metric tree Call tree 0.00 machine Linux 735.75 Idle threads 0.72 MPI Waitall 0.00 tea module.tea unpack right ¬ □ 0.00 node irc0280 1.17e8 Visits (occ) 0.03 MPI Rank 0 0.20 tea module.tea pack top 72 MPI synchronizations (occ) 3.07 MPI Rank 1 0 MPI pair-wise one-sided synchronizations othe ▶ ■ 0.20 tea module.tea send recv message top - □ 0.00 node irc0281 I.29e6 MPI communications (occ) 0.21 tea module.tea unpack top 0.01 MPI Rank 2 I 0 MPI file operations (occ) ▶ ■ 0.28 tea module.tea pack left 0.28 MPI Rank 3 ▶ ■ 0.32 tea module.tea send recv message left 2.37e10 MPI bytes transferred (bytes) 0.23 tea module.tea unpack left → □ 0.00 node irc0282 0.02 MPI Rank 4 0.31 tea module.tea pack bottom → □ 0.00 MPI 2.30 MPI Rank 5 0.28 tea module.tea send recv message botto - 0.00 Point-to-point - 0.00 node jrc0283 I 0.18 tea module.tea unpack bottom 4.57 Late Sender 0.01 MPI Rank 6 0.12 update halo kernel module.update halo keri 0.26 Late Receiver 0.02 MPI Rank 7 0.02 tea leaf kernel cg module.tea leaf kernel init 
 □ 0.00 Collective
 - 🗖 0.09 tea module.tea allsum 0.10 Wait at Barrier 0.68 MPI Allreduce 33.17 Wait at N x N 
 0.00 tea\_leaf\_kernel\_cheby\_module.tea\_leaf\_kernel\_c

 0.00 tea\_leaf\_kernel\_cg\_module.tea\_leaf\_kernel\_solve
 → □ 0.00 Late Broadcast - 0.00 OpenMP 0.01 !somp parallel @tea leaf cg.f90:186 59.82 Wait at Barrier ▶ 🗖 1.90 !\$omp do @tea leaf cg.f90:187 → □ 0.00 Idleness delay costs 0.01 !\$omp implicit barrier @tea leaf cg.f90:20 354.84 Short-term 0.00 tea leaf kernel cg module tea leaf kernel solve 87.42 Long-term □ 0.02 !\$omp parallel @tea\_leaf\_cg.f90:234
 □ 2.45 !\$omp do @tea\_leaf\_cg.f90:247 5.13 MPI point-to-point wait states (propage 5.13 MPI point-to-point wait states (direct v 90.49 Critical path (sec) 0.01 !\$omp implicit barrier @tea leaf cg.f90:25 - 🔲 0.00 tea leaf kernel cg module tea leaf kernel solv 9.41 Imbalance Þ. ▶ ■ 0.47 !\$omp parallel @tea leaf cg.f90:284 9590.08 Performance impact (sec) All (96 elements) Ŧ b. 99.90 0.00 9.41 (9.42%) 5.75 (61.12%) 9.41 0.00 0.00 (0.00%) 5.75 0.00

...with imbalance (time on critical path above average) mostly in the first two loops and MPI communication

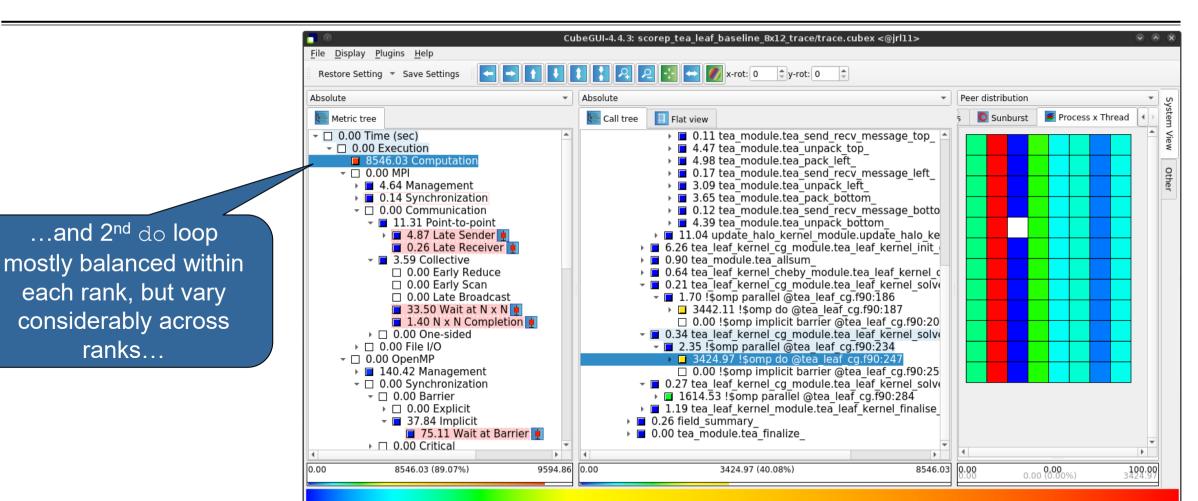
# **TeaLeaf Scalasca report analysis (VIII)**





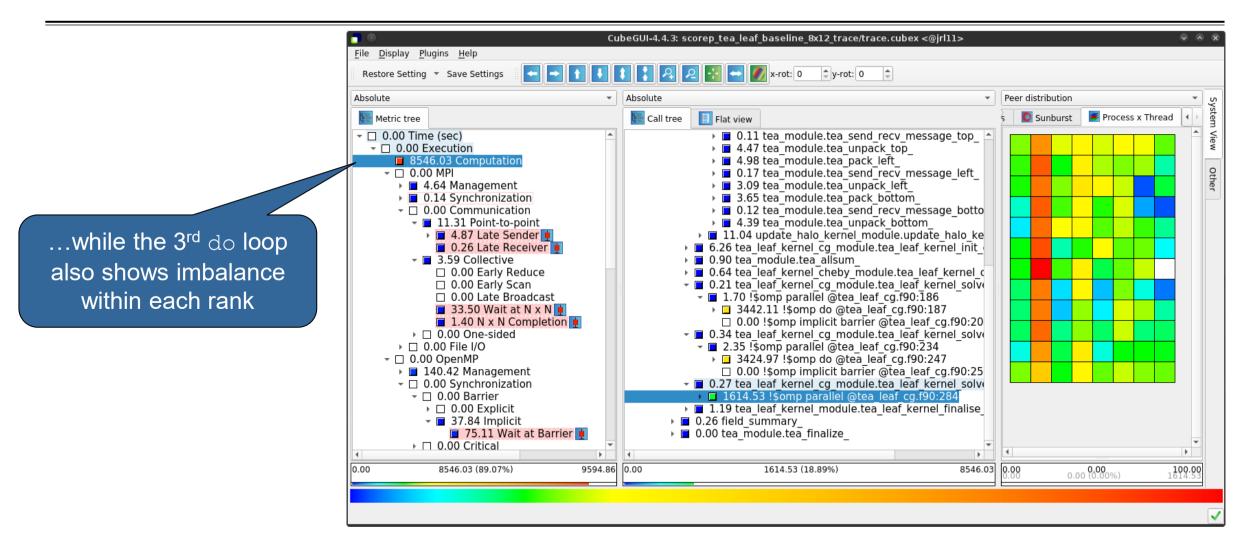
# TeaLeaf Scalasca report analysis (IX)





# **TeaLeaf Scalasca report analysis (X)**





# **TeaLeaf analysis summary**

- The first two OpenMP do loops of the solver are well balanced within a rank, but are imbalanced across ranks
  - → Requires a global load balancing strategy
- The third OpenMP do loop, however, is imbalanced within ranks,
  - causing direct "Wait at OpenMP Barrier" wait states,
  - which cause indirect MPI point-to-point wait states,
  - which in turn cause OpenMP thread idleness
  - Low-hanging fruit
- Adding a SCHEDULE (guided) clause reduced
  - the MPI point-to-point wait states by ~66%
  - the MPI collective wait states by ~50%
  - the OpenMP "Wait at Barrier" wait states by ~55%
  - the OpenMP thread idleness by ~11%
  - → Overall runtime (wall-clock) reduction by ~5%

# **Scalasca Trace Tools: 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





# **Reference material**





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



	<pre>% scalasca</pre>										
	Scalasc	lasca 2.6.2									
	Toolset	Toolset for scalable performance analysis of large-scale parallel applications									
	usage:	age: scalasca [OPTION] ACTION <argument></argument>									
	1. prepare application objects and executable for measurement:										
	scalasca -instrument <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.	3. interactively explore measurement analysis report:									
	scalasca -examine <experiment-archive report></experiment-archive report>										
Options:											
		show-config	show configuration summary and exit								
		help	show this help and exit								
		dry-run	show actions without taking them								
		quickref	show quick reference guide and exit								
		remap-specfile	show path to remapper specification file and exit								
		verbose	enable verbose commentary								
	-V,	version	show version information and exit								

■ The `scalasca -instrument' command is deprecated and will be remove in the next major release
⇒ use Score-P instrumenter directly

#### Scalasca convenience command: scan / scalasca -analyze



Scan											
Scalasca 2.6.2: measurement collection & analysis nexus											
usage: scan {options} [launchcmd [launchargs]] target [targetargs]											
where {options} may include:											
-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.											
-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 cfgfile : Specify a config file for a multi-run measurement.											
-P preset : Specify a preset for a multi-run measurement, e.g.,	'pop'.										
-L : List available multi-run presets.											
	: Check a multi-run config file for validity and dump										
: the processed configuration for comparison.											

#### Scalasca measurement collection & analysis nexus

# Scalasca convenience command: square / scalasca -examine

<pre>% square Scalasca 2.6.2: analysis report explorer</pre>										
usage: square [OPTIONS] <experiment archive="" cube="" file=""  =""></experiment>										
-C <none full="" quick=""  =""> : Level of sanity checks for newly created reports</none>										
	<u> </u>									
-c <number></number>	: Consider number of counters when doing scoring (-s)									
-F	: Force remapping of already existing reports									
-f filtfile	: Use specified filter file when doing scoring (-s)									
-s	: Skip display and output textual score report									
-v	: Enable verbose mode									
-n	: Do not include idle thread metric									
-S <mean merge=""  =""></mean>	: Aggregation method for summarization results of									
	each configuration (default: merge)									
-T <mean merge=""  =""></mean>	: Aggregation method for trace analysis results of									
	each configuration (default: merge)									
-A	: Post-process every step of a multi-run experiment									
-I	: Ignore structural sanity checks and force aggregation									
	of measurements in a multi-run experiment									
-x <scorep-score opt=""></scorep-score>	: Pass option(s) to scorep-score									

Scalasca analysis report explorer (Cube)

#### Scalasca advanced command: scout - Scalasca automatic trace analyzer



```
% scout.hvb --help
SCOUT
      (Scalasca 2.6.2)
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
                     Enables critical-path analysis [default]
  --critical-path
  --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 advanced command: clc\_synchronize



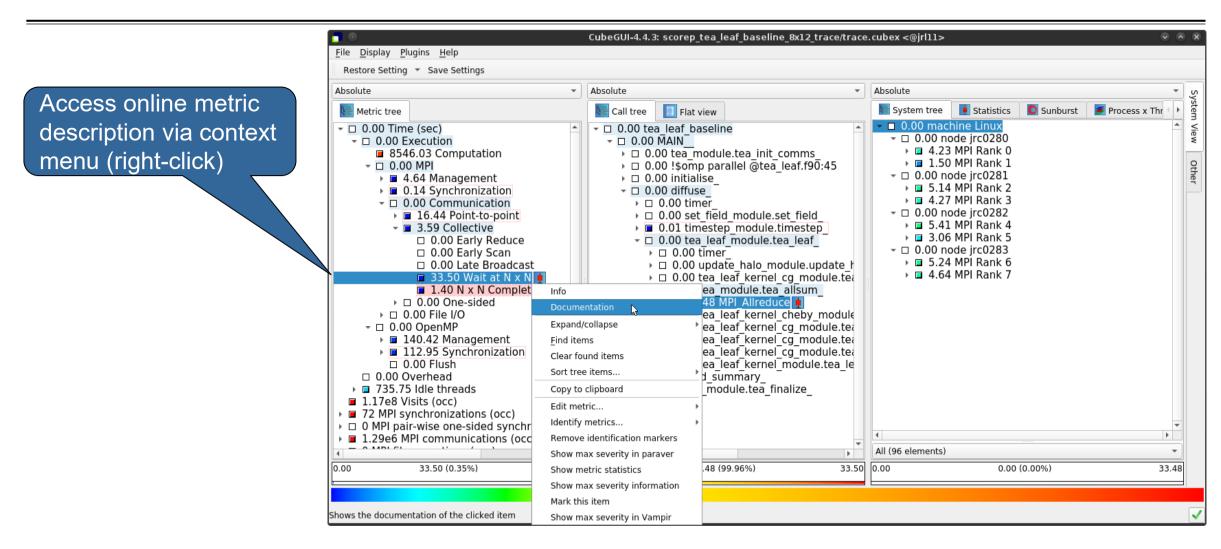
Scalasca trace event timestamp consistency correction

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

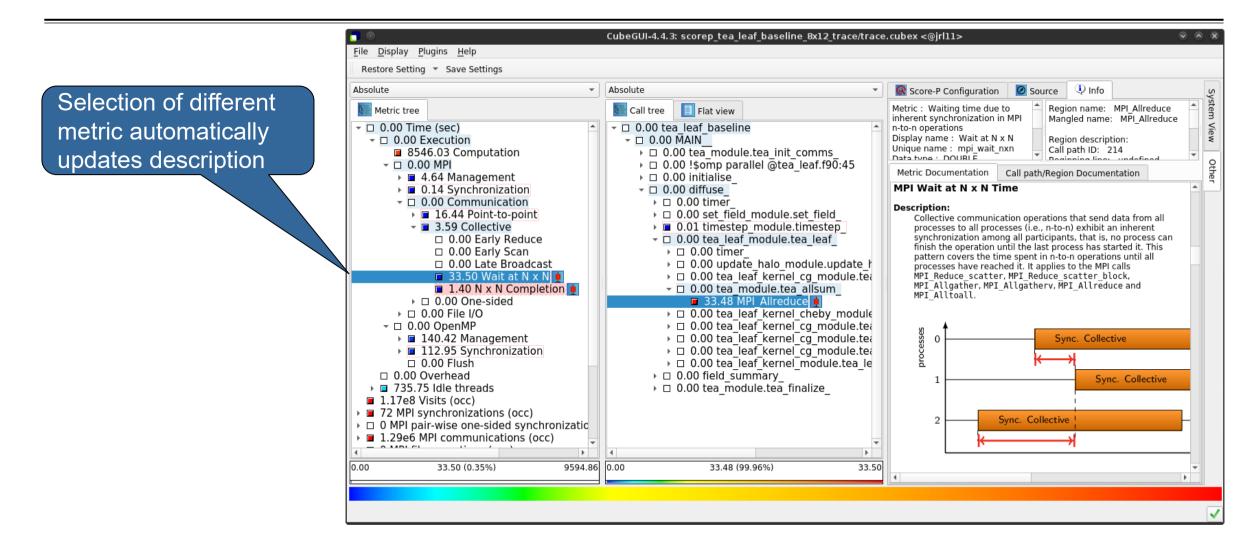
#### **Online metric description**





# **Online metric description (cont.)**





VIRTUAL INSTITUTE – HIGH PRODUCTIVITY SUPERCOMPUTING

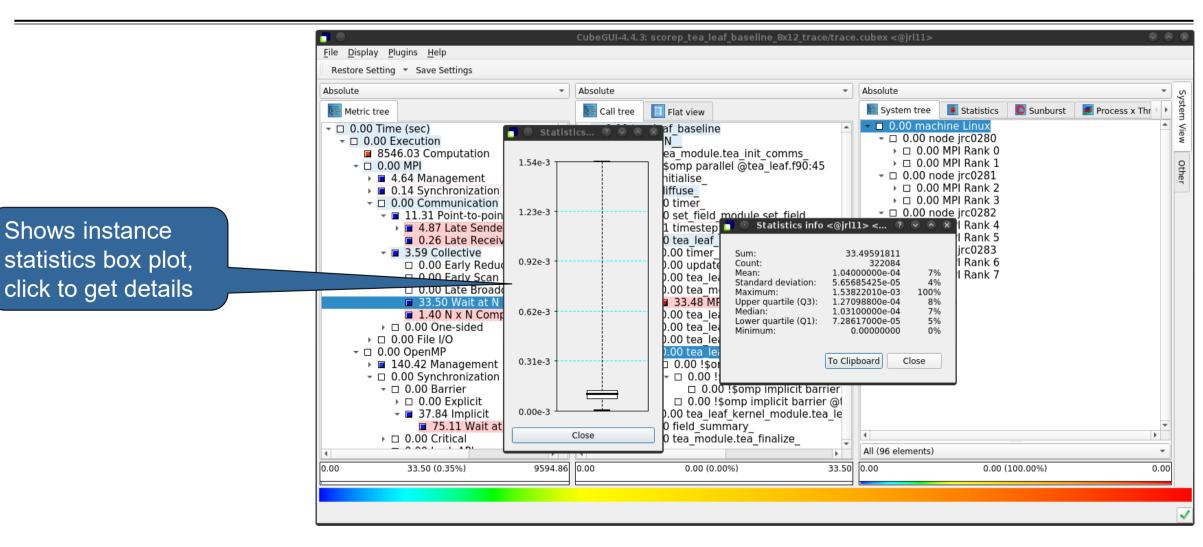
#### **Metric statistics**



	<b>0</b>		CubeGUI-4.4.3: score	ep_tea_leaf_baseline_8x12_trace/trac	e.cubex <@jrl11>		•		
	File Display Plugins Help								
	Restore Setting  Save Settings								
	Absolute	•	Absolute	·	Absolute		· · · · · · · · · · · · · · · · · · ·		
	tree Metric tree		🔚 Call tree 📗 F	lat view	🔚 System tree		👻 burst 🛛 🖉 Process x Thr		
Access metric statistics for metrics marked with box plot icon from context menu		Info Docume Expand/ Find iter Clear for Sort tree Copy to Edit mel Identify	<ul> <li>0.00 !\$c</li> <li>0.00 init</li> <li>0.00 dif</li> <li>0.00</li> <li>0.00</li> <li>0.00</li> <li>0.00</li> <li>0.01</li> <li>0.00</li> <li>0.00</li></ul>	a_module.tea_init_comms omp parallel @tea_leaf.f90:45 tialise_ fuse_ timer_ set_field_module.set_field_ timestep_module.timestep_ tea_leaf_module.tea_leaf_ 00 update_halo_module.update_f 00 update_halo_module.update_f 00 tea_leaf_kernel_cg_module.tea 00 tea_module.tea_allsum_ <b>48 MPI_Allreduce_ff</b> ea_leaf_kernel_cg_module.tea ea_leaf_kernel_cg_module.tea ea_leaf_kernel_cg_module.tea ea_leaf_kernel_cg_module.tea ea_leaf_kernel_cg_module.tea 0.00 !\$omp parallel @tea_leaf_cg.f 0.00 !\$omp implicit barrier 0.00 !\$omp implicit barrier	<ul> <li>▶ ■ 1.50</li> <li>▼ □ 0.00 no</li> <li>▶ ■ 5.14</li> <li>▶ ■ 4.27</li> <li>▼ □ 0.00 no</li> <li>▶ ■ 3.06</li> <li>▼ □ 0.00 no</li> <li>▶ ■ 3.06</li> <li>▼ □ 0.00 no</li> <li>▶ ■ 4.64</li> <li>▲ 4.64</li> <li>▲ 4.64</li> </ul>	de jrc0280 MPI Rank 0 MPI Rank 1 de jrc0281 MPI Rank 2 MPI Rank 3 de jrc0282 MPI Rank 4 MPI Rank 5			
	0.00 33.50 (0.35%)	Show m	ax severity in paraver	.48 (99.96%) 33.50	0.00	0.00 (0.00%)	33.48		
		Show m	etric statistics 🛛 📡						
			Show max severity information						
		Mark thi							
		Show m	ax severity in Vampir						

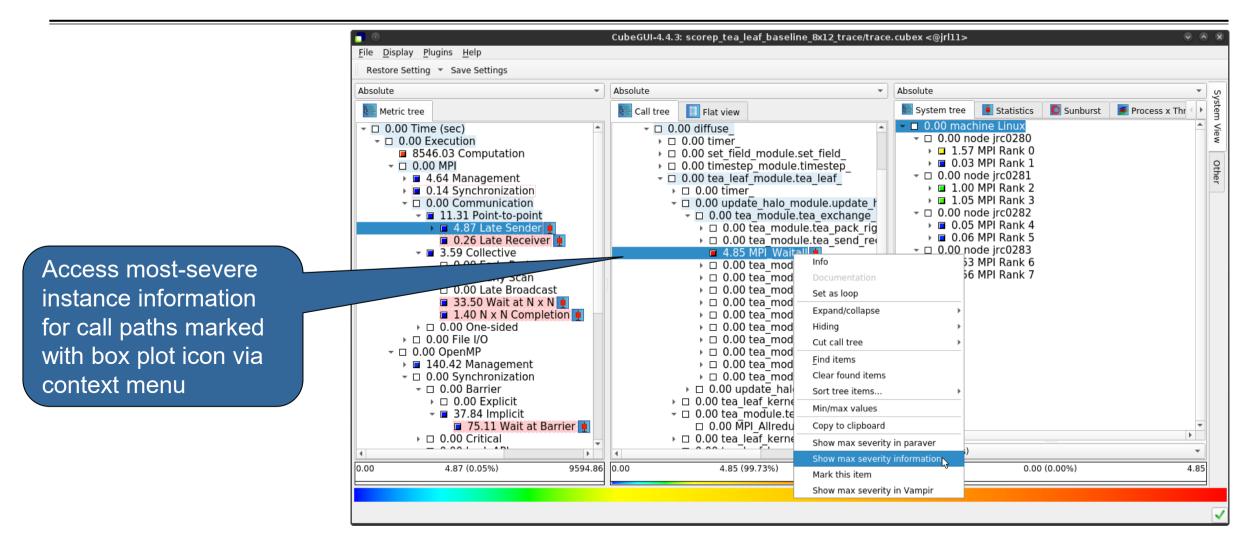
#### Metric statistics (cont.)





#### **Metric instance statistics**





# Metric instance statistics (cont.)



