

# Automatic trace analysis with the Scalasca Trace Tools

---

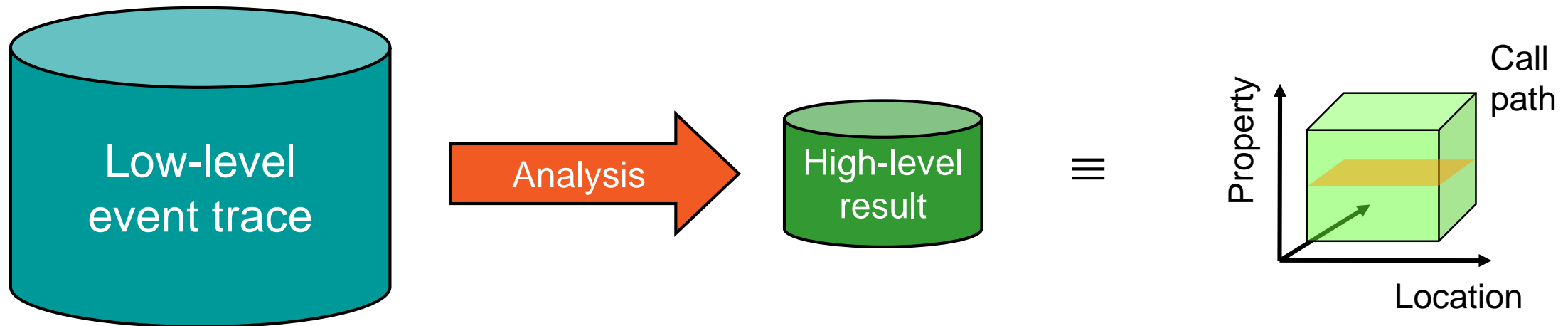
Brian Wylie  
Jülich Supercomputing Centre



# Automatic trace analysis

## ▪ Idea

- Automatic search for patterns of inefficient behaviour
- Classification of behaviour & 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: Objective

---

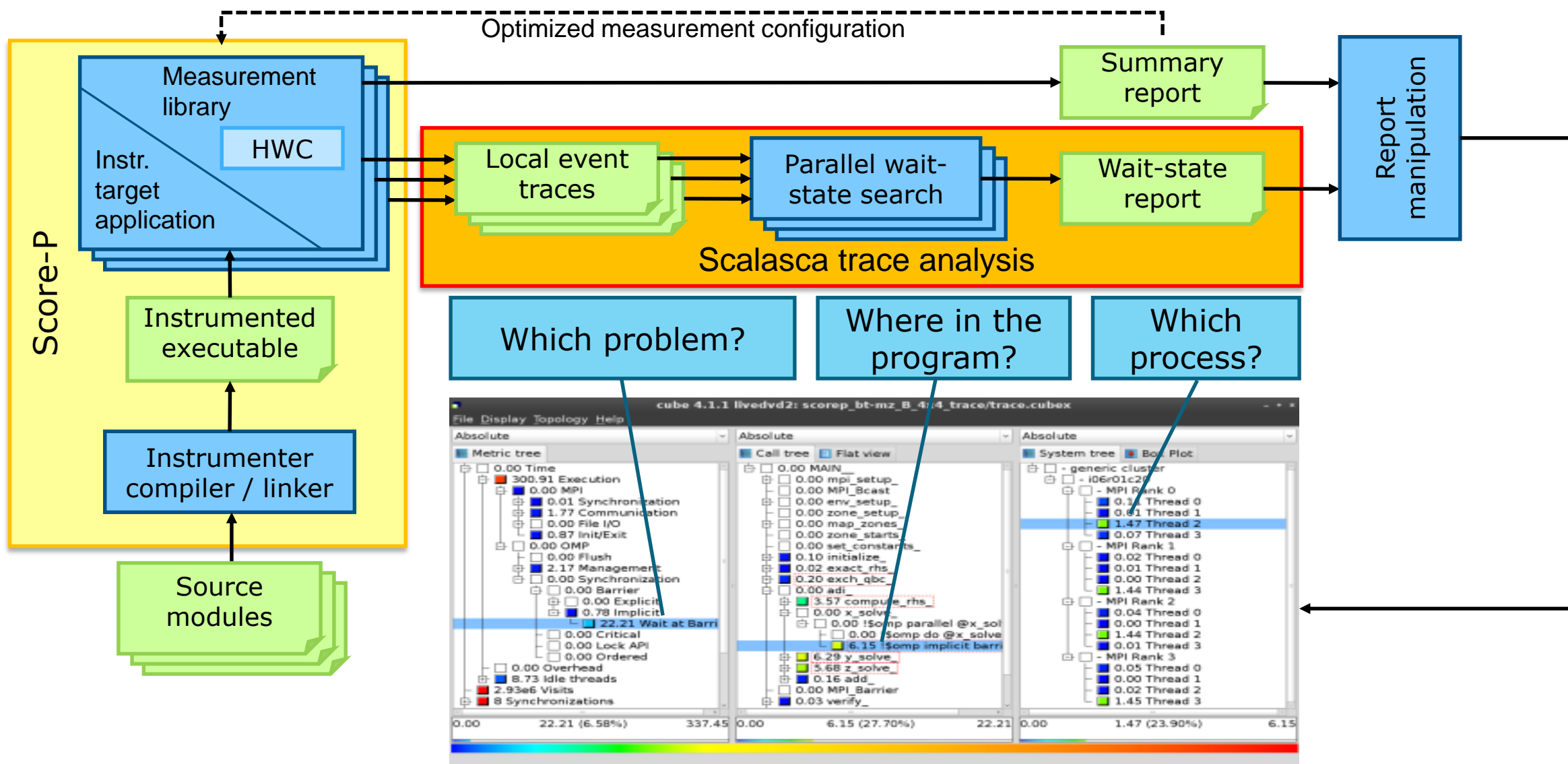
- Development of a **scalable trace-based** performance analysis toolset for the most popular parallel programming paradigms
  - Current focus: MPI, OpenMP, and POSIX threads
- Specifically targeting large-scale parallel applications
  - Such as those running on IBM Blue Gene or Cray systems with one million or more processes/threads
- Latest release:
  - Scalasca v2.5 coordinated with Score-P v5.0 (March 2019) also works with Score-P v6.0

# Scalasca Trace Tools features

---

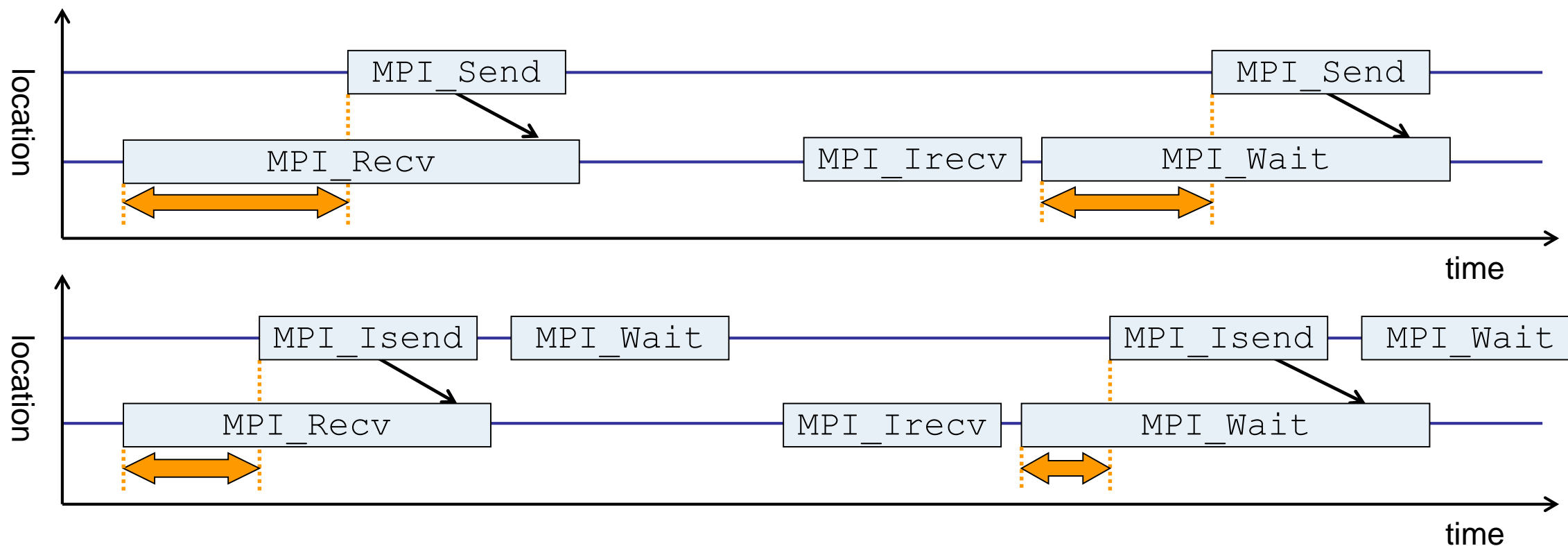
- Open source, 3-clause BSD license
- Fairly portable
  - IBM Blue Gene, Cray XT/XE/XK/XC, SGI Altix, Fujitsu FX10/100 & K computer, Linux clusters (x86, Power, ARM), Intel Xeon Phi, ...
- Uses Score-P instrumenter & measurement libraries
  - Scalasca v2 core package focuses on trace-based analyses
  - Supports common data formats
    - Reads event traces in OTF2 format
    - Writes analysis reports in CUBE4 format
- Current limitations:
  - Unable to handle traces
    - With MPI thread level exceeding MPI\_THREAD\_FUNNELED
    - Containing CUDA or SHMEM events, or OpenMP nested parallelism
  - PAPI/rusage metrics for trace events are ignored

# Scalasca workflow



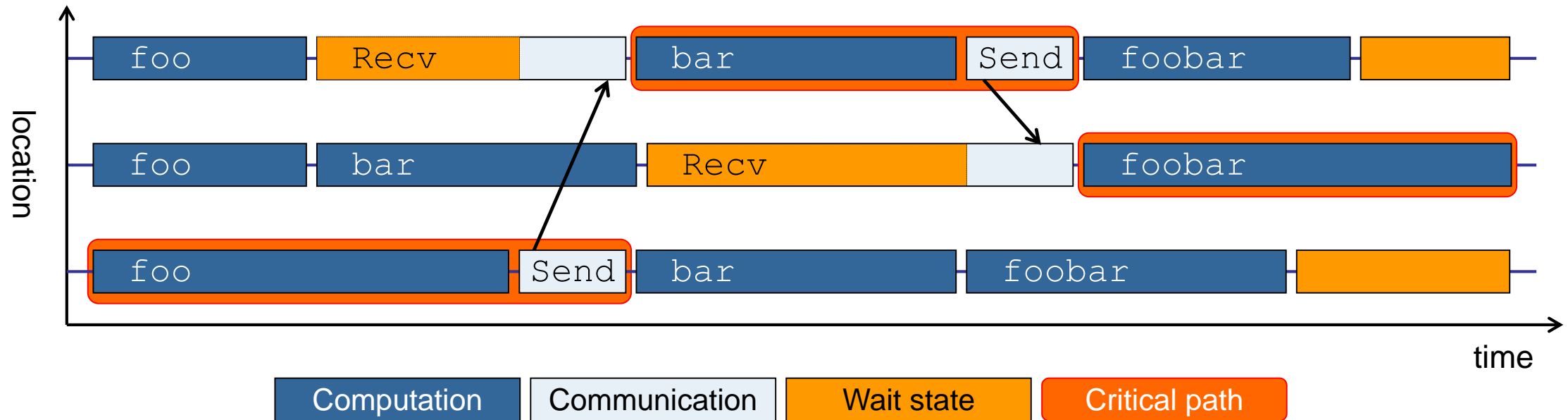


## Example: “*Late Sender*” wait state



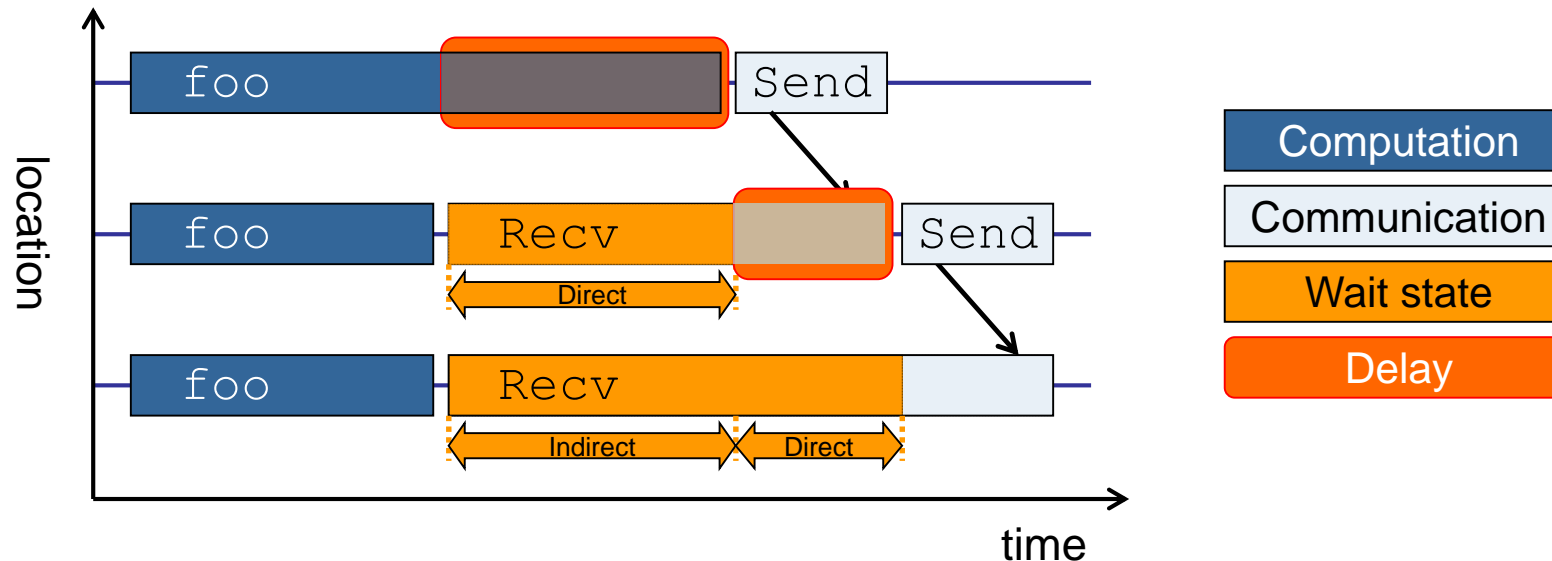
- Waiting time caused by a blocking receive operation posted earlier than the corresponding send
- Applies to blocking as well as non-blocking communication

## 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*



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

---



# 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

# Scalasca command – One command for (almost) everything

```
% scalasca
Scalasca 2.5
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
                        --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 convenience command: scan / scalasca -analyze

```
% scan
Scalasca 2.5: 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 filtfle          : 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.
```

- Scalasca measurement collection & analysis nexus

## Scalasca convenience command: square / scalasca -examine

```
% square
Scalasca 2.5: analysis report explorer
usage: square [OPTIONS] <experiment archive | cube file>
  -c <none | quick | full> : Level of sanity checks for newly created reports
  -F                        : Force remapping of already existing reports
  -f filtfiler              : 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>         : Aggregation method for summarization results of
                             each configuration (default: merge)
  -T <mean | merge>         : Aggregation method for trace analysis results of
                             each configuration (default: merge)
  -A                        : Post-process every step of a multi-run experiment
```

- Scalasca analysis report explorer (Cube)



# 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)



## Recap: Local installation (Archer)

---

- Select appropriate PrgEnv (cray, gnu, or intel: gnu recommended/assumed)

```
% module switch PrgEnv-cray PrgEnv-gnu
```

- Latest/recent versions of VI-HPS tools not yet installed system-wide

- Add extra module path
- Required for each shell session

```
% module use /home/y07/y07/scalasca/modules  
% module load scalasca # implicitly loads scorep and cube
```

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

```
% cd $WORK/NPB3.3-MZ-MPI
```

## BT-MZ summary measurement collection...

```
% cd bin.scorep
% cp ../jobscript/archer/scalasca.pbs .
% cat scalasca.pbs

# Score-P measurement configuration
export SCOREP_FILTERING_FILE=../config/scorep.filt
#export SCOREP_METRIC_PAPI=PAPI_TOT_INS,PAPI_TOT_CYC
#export SCOREP_METRIC_RUSAGE=ru_stime
#export SCOREP_METRIC_RUSAGE_PER_PROCESS=ru_maxrss
#export SCOREP_TOTAL_MEMORY=90M

# Run the application using Scalasca nexus
scalasca -analyze aprun -n $NPROCS -d $OMP_NUM_THREADS $EXE
```

```
% qsub scalasca.pbs
```

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

**Hint:**

```
scan = scalasca -analyze
```

- Submit the job

## BT-MZ summary measurement

```
S=C=A=N: Scalasca 2.5 runtime summarization
S=C=A=N: ./scorep_bt-mz_C_8x6_sum experiment archive
S=C=A=N: Thu Apr 11 13:25:35 2019: Collect start
aprun -n 8 -d 6 ./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:      8
```

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

```
S=C=A=N: Thu Apr 11 13:25:50 2019: Collect done (status=0) 15s
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 result...  
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]
```

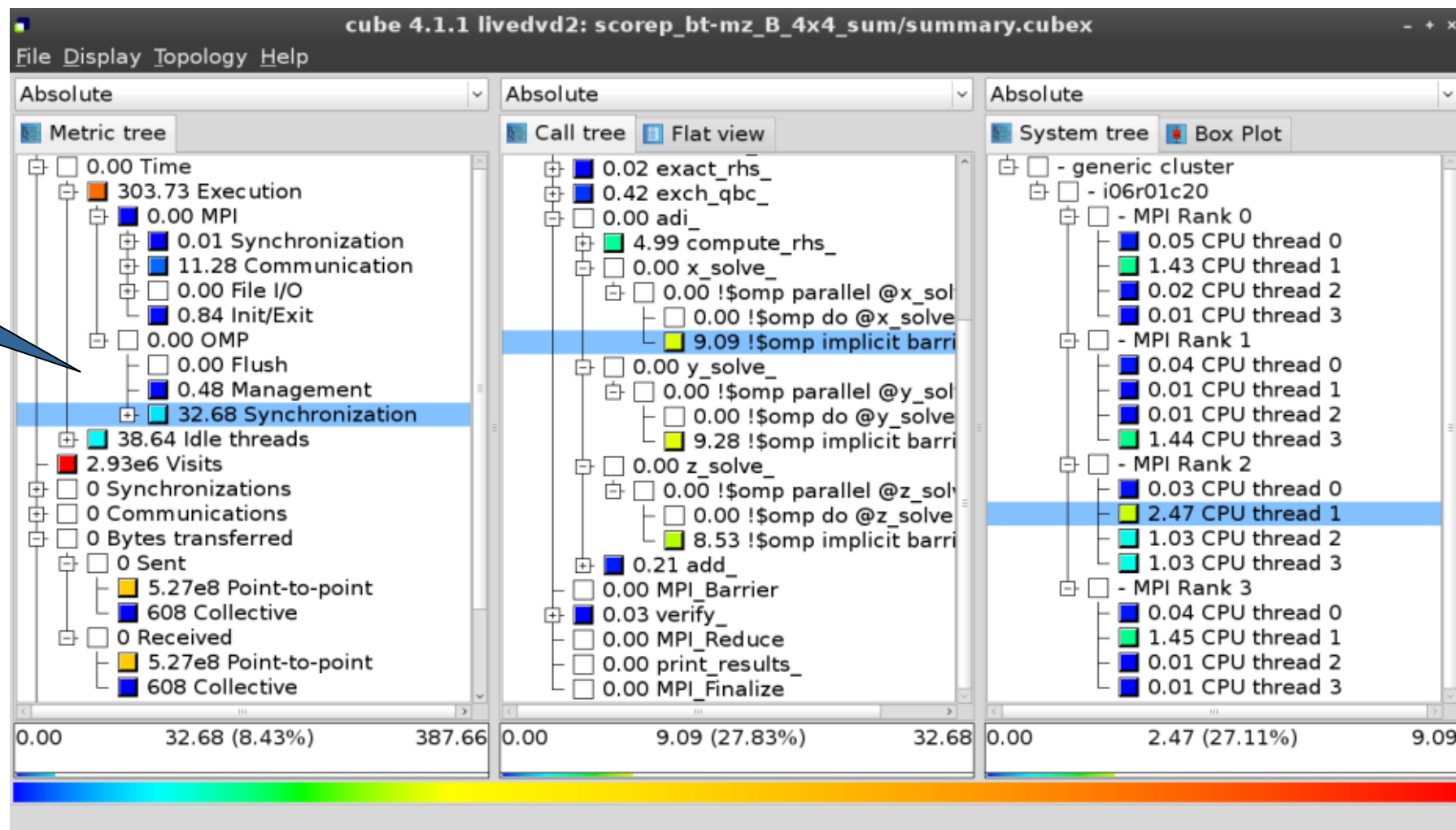
**Hint:**

Copy 'profile.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

Split base metrics into more specific metrics



# 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/archer/scalasca.pbs .
% vim scalasca.pbs

# Score-P measurement configuration
export SCOREP_FILTERING_FILE=../config/scorep.filt
#export SCOREP_METRIC_PAPI=PAPI_TOT_INS,PAPI_TOT_CYC
#export SCOREP_METRIC_RUSAGE=ru_stime
#export SCOREP_METRIC_RUSAGE_PER_PROCESS=ru_maxrss
export SCOREP_TOTAL_MEMORY=90M

# Run the application using Scalasca nexus
scalasca -analyze -t aprun -n $NPROCS -d $OMP_NUM_THREADS $EXE

% qsub scalasca.pbs
```

- Change to directory with the Score-P instrumented executable and edit the job script
- Add “-t” to the scalasca -analyze command
- Submit the job

## BT-MZ trace measurement ... collection

---

```
S=C=A=N: Scalasca 2.5 trace collection and analysis
S=C=A=N: Thu Apr 11 13:35:31 2019: Collect start
aprun -n 8 -d 6 ./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:      8

[... More application output ...]

S=C=A=N: Thu Apr 11 13:35:48 2019: Collect done (status=0) 17s
```

- Starts measurement with collection of trace files ...

## BT-MZ trace measurement ... analysis

```
...
S=C=A=N: Thu Apr 11 13:35:48 2019: Analyze start
  aprun -n 8 -d 6 scout.hyb --time-correct \
>    ./scorep_bt-mz_C_8x6_trace/traces.otf2

SCOUT    (Scalasca 2.5)

Analyzing experiment archive ./scorep_bt-mz_C_8x6_trace/traces.otf2

Opening experiment archive ... done (0.022s).
Reading definition data    ... done (0.005s).
Reading event trace data  ... done (0.365s).
Preprocessing              ... done (1.621s).
Timestamp correction       ... done (3.048s).
Analyzing trace data       ... done (26.146s).
Writing analysis report    ... done (0.343s).

Max. memory usage          : 519.883MB

      # passes              : 1
      # violated            : 0

Total processing time       : 31.604s
S=C=A=N: Thu Apr 11 13:36:26 2019: Analyze done (status=0) 38s
```

- Continues with automatic (parallel) analysis of trace files

## BT-MZ trace analysis report exploration

---

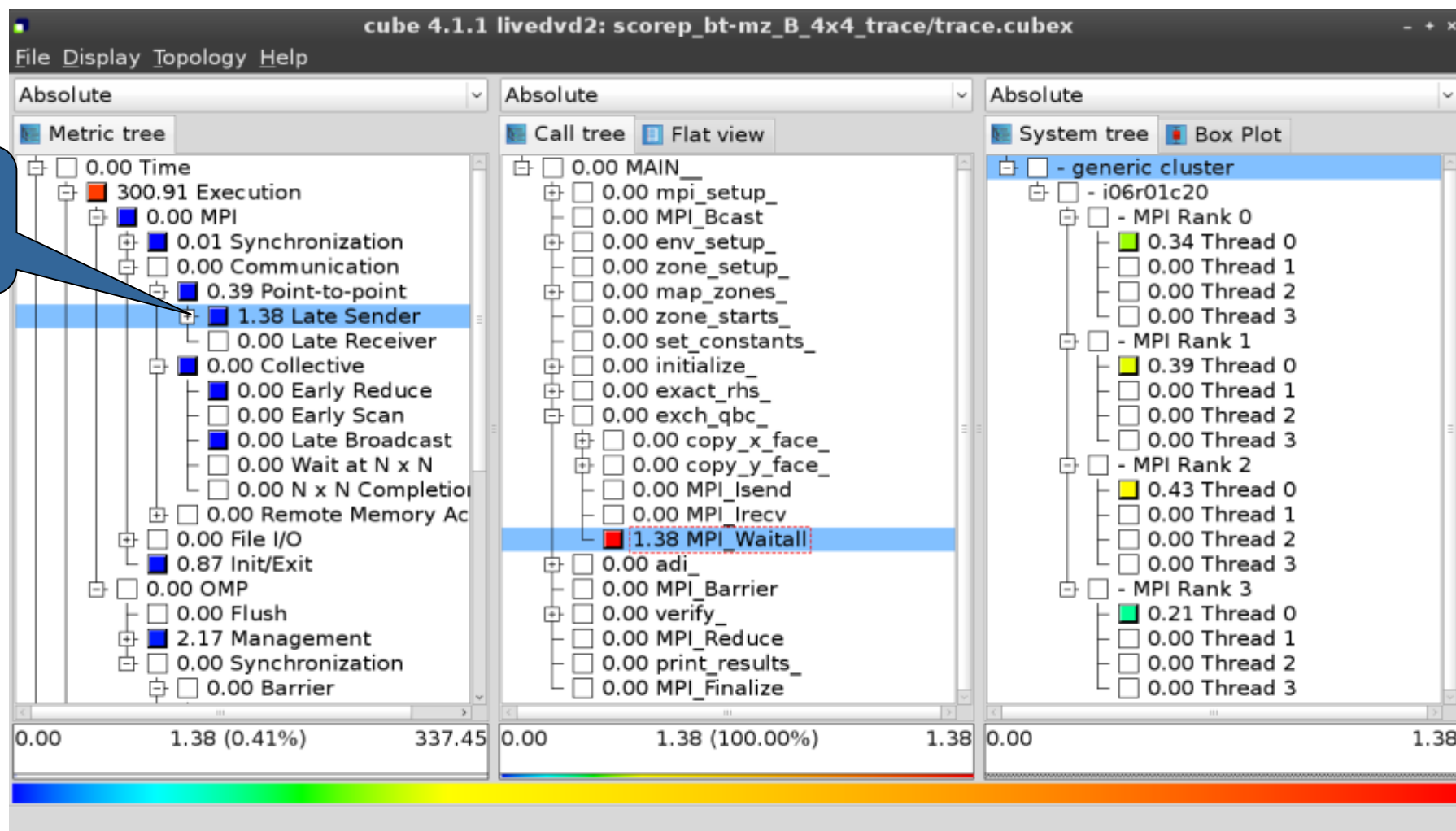
- Produces trace analysis report in the experiment directory containing trace-based wait-state metrics

```
% square scorep_bt-mz_C_8x6_trace  
INFO: Post-processing runtime summarization result...  
INFO: Post-processing trace analysis report...  
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

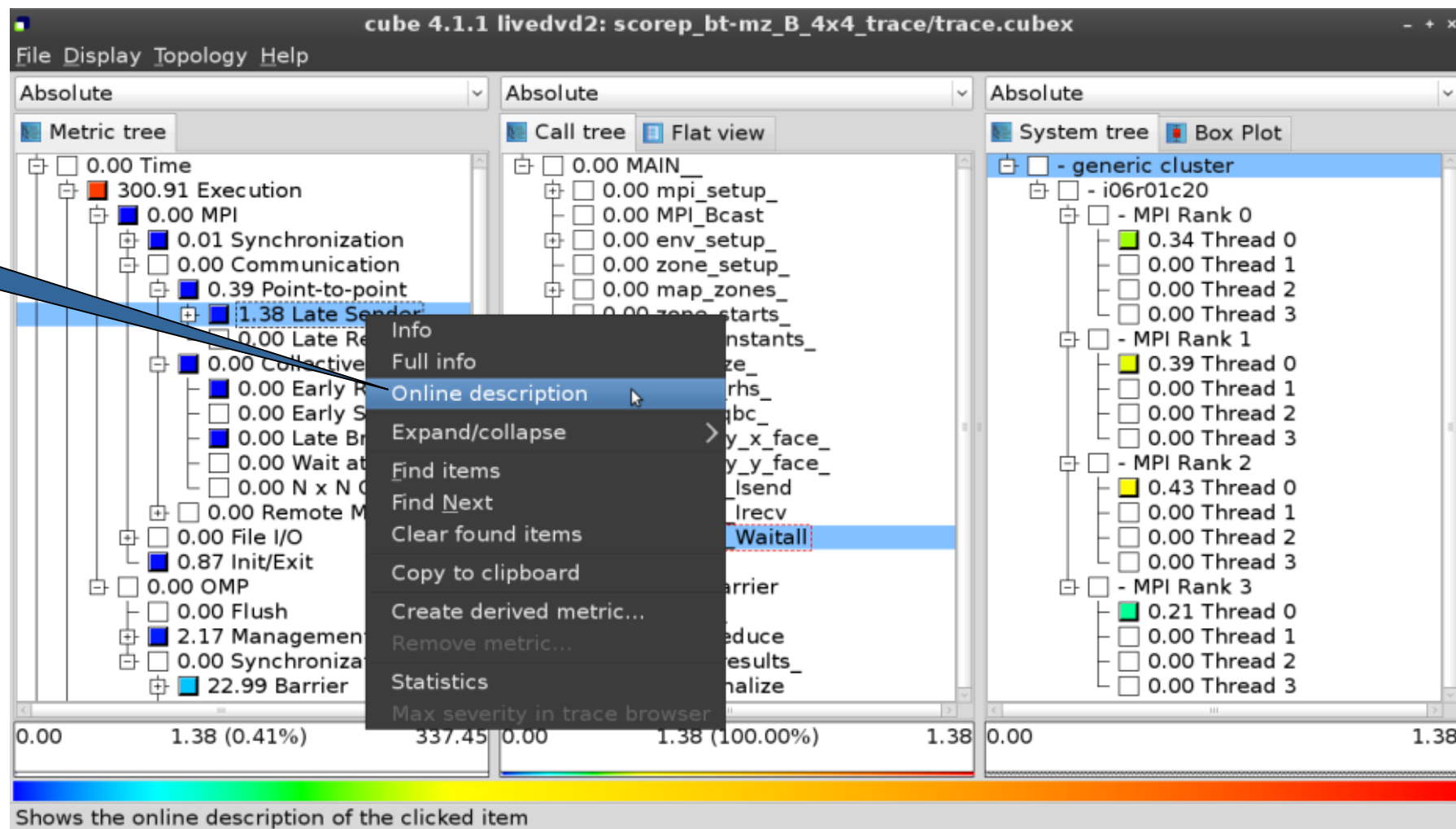
# Post-processed trace analysis report





# Online metric description

Access online metric description via context menu



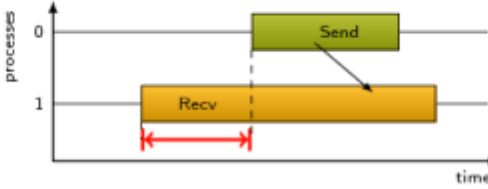


# Online metric description

**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.

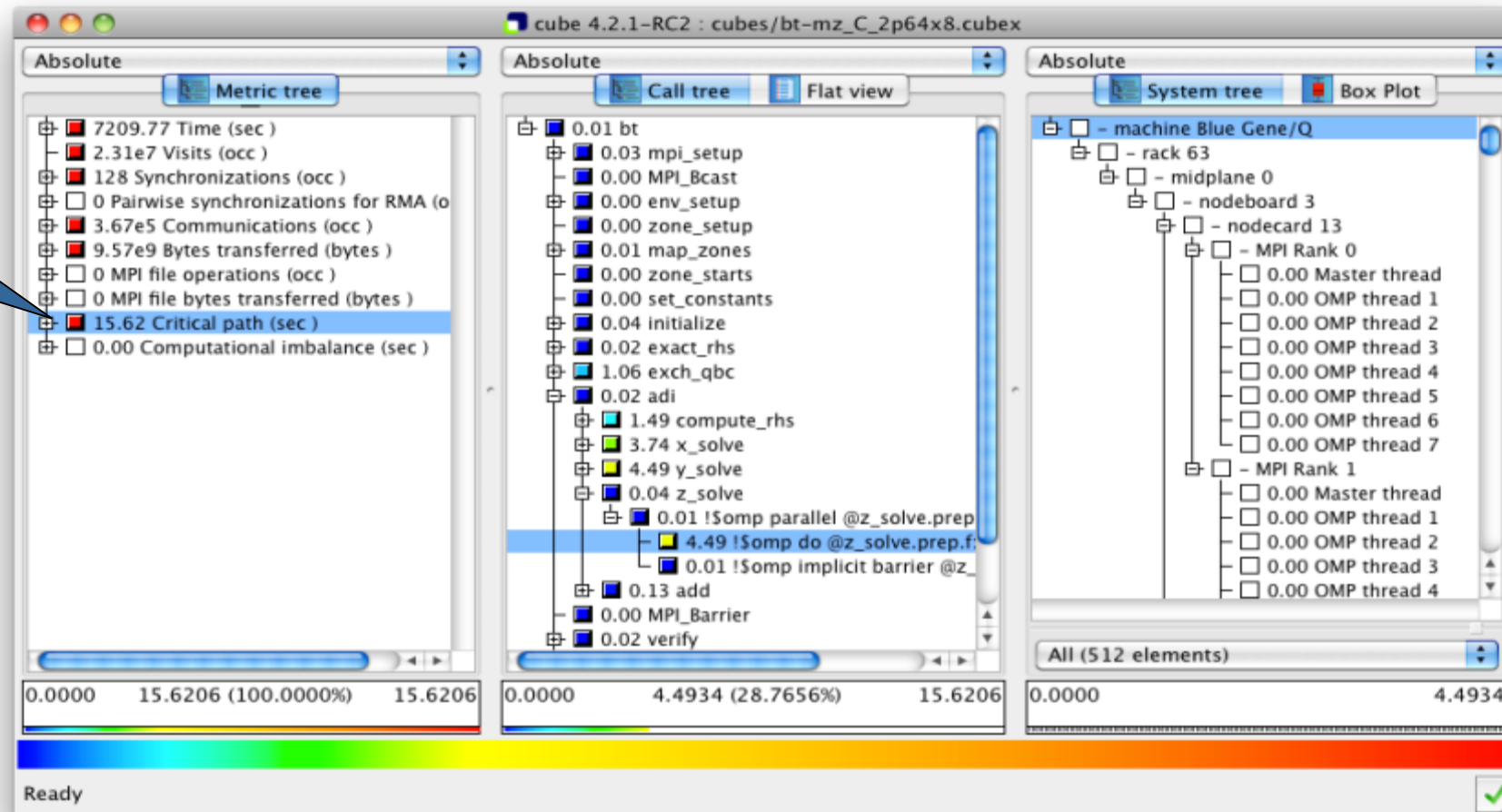
**Parent:**  
[MPI Point-to-point Communication Time](#)

**Children:**

Close

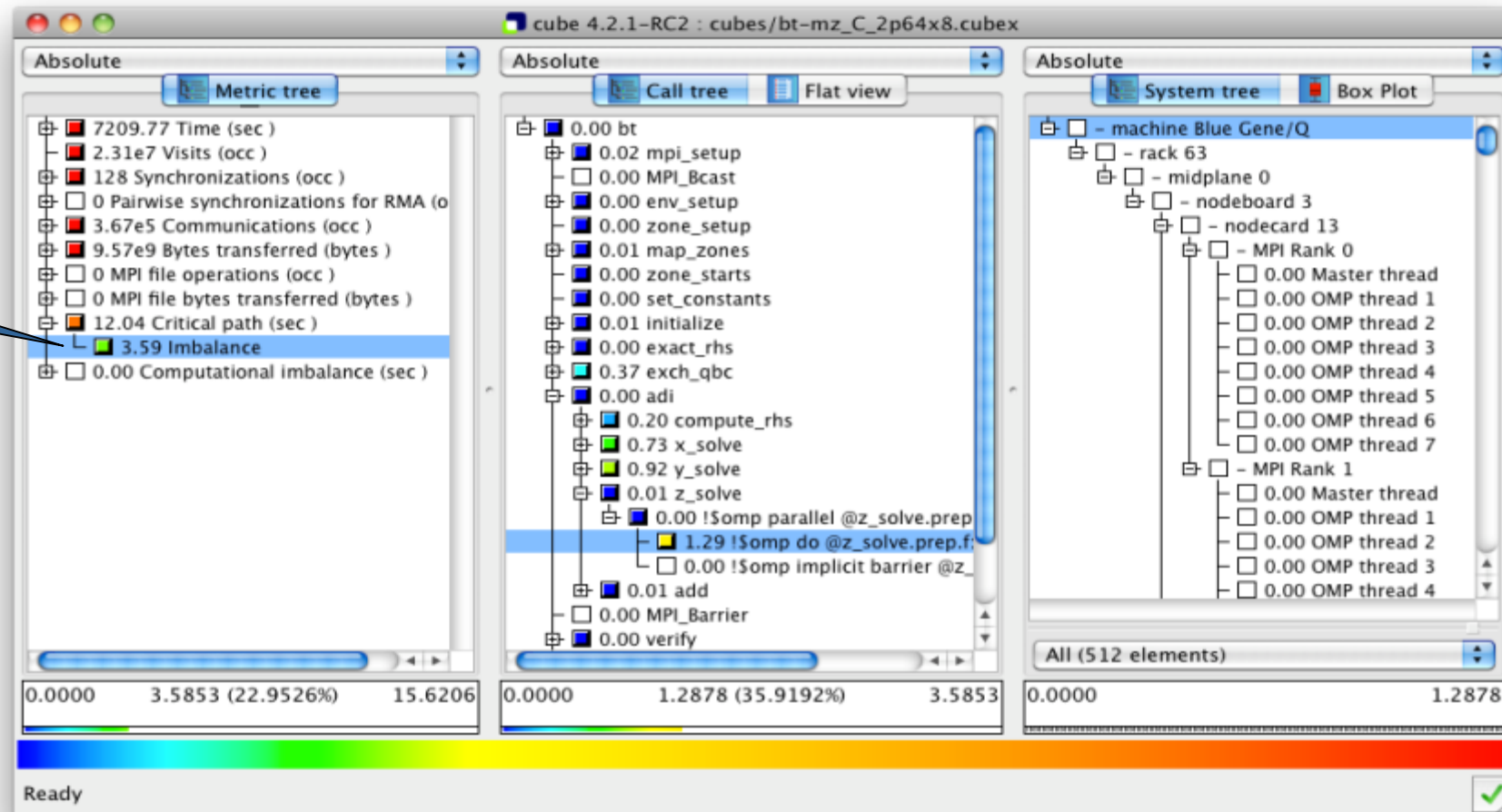
# Critical-path analysis

Critical-path profile shows wall-clock time impact

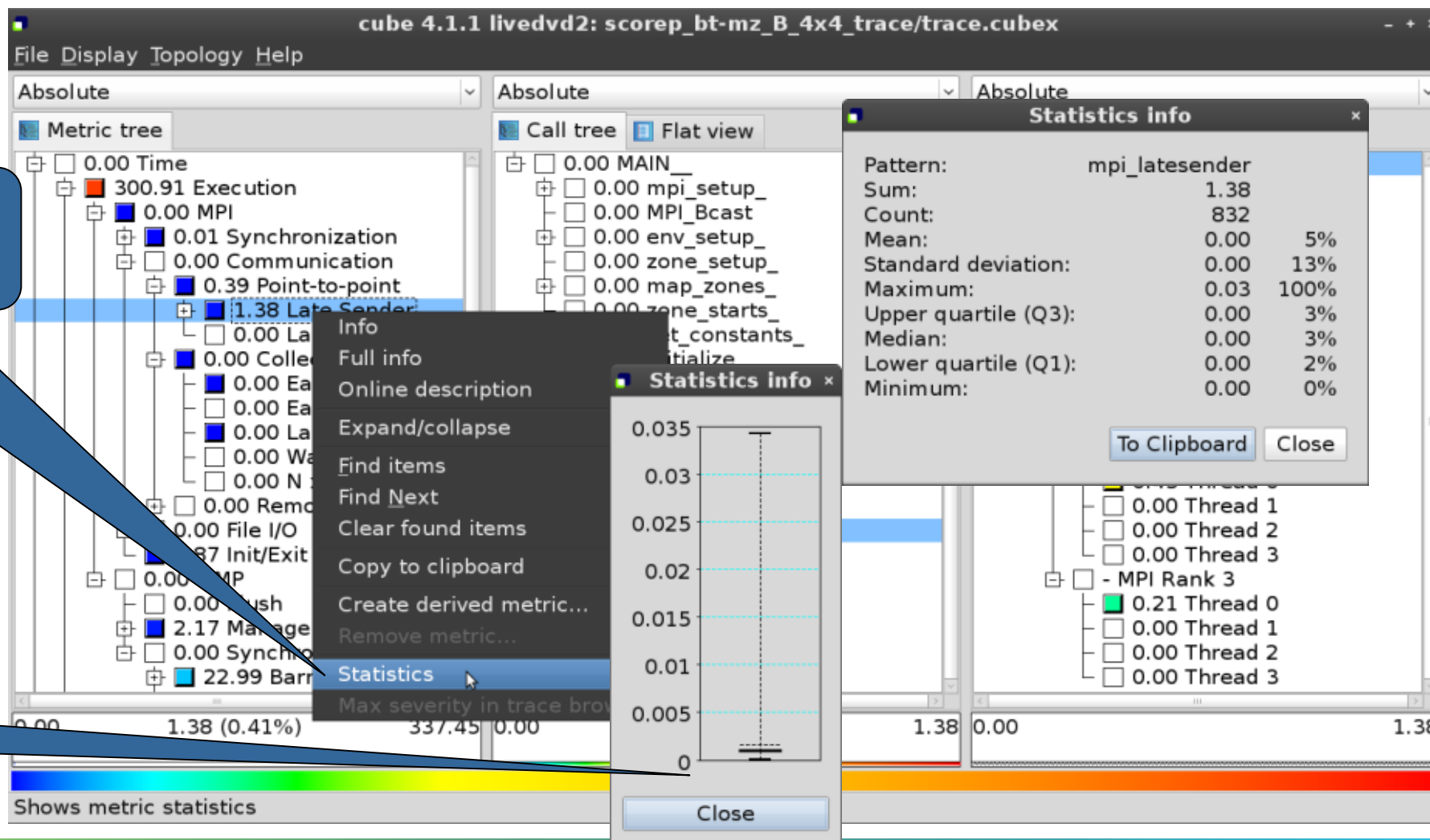


# Critical-path analysis

Critical-path imbalance highlights inefficient parallelism



# Pattern instance statistics



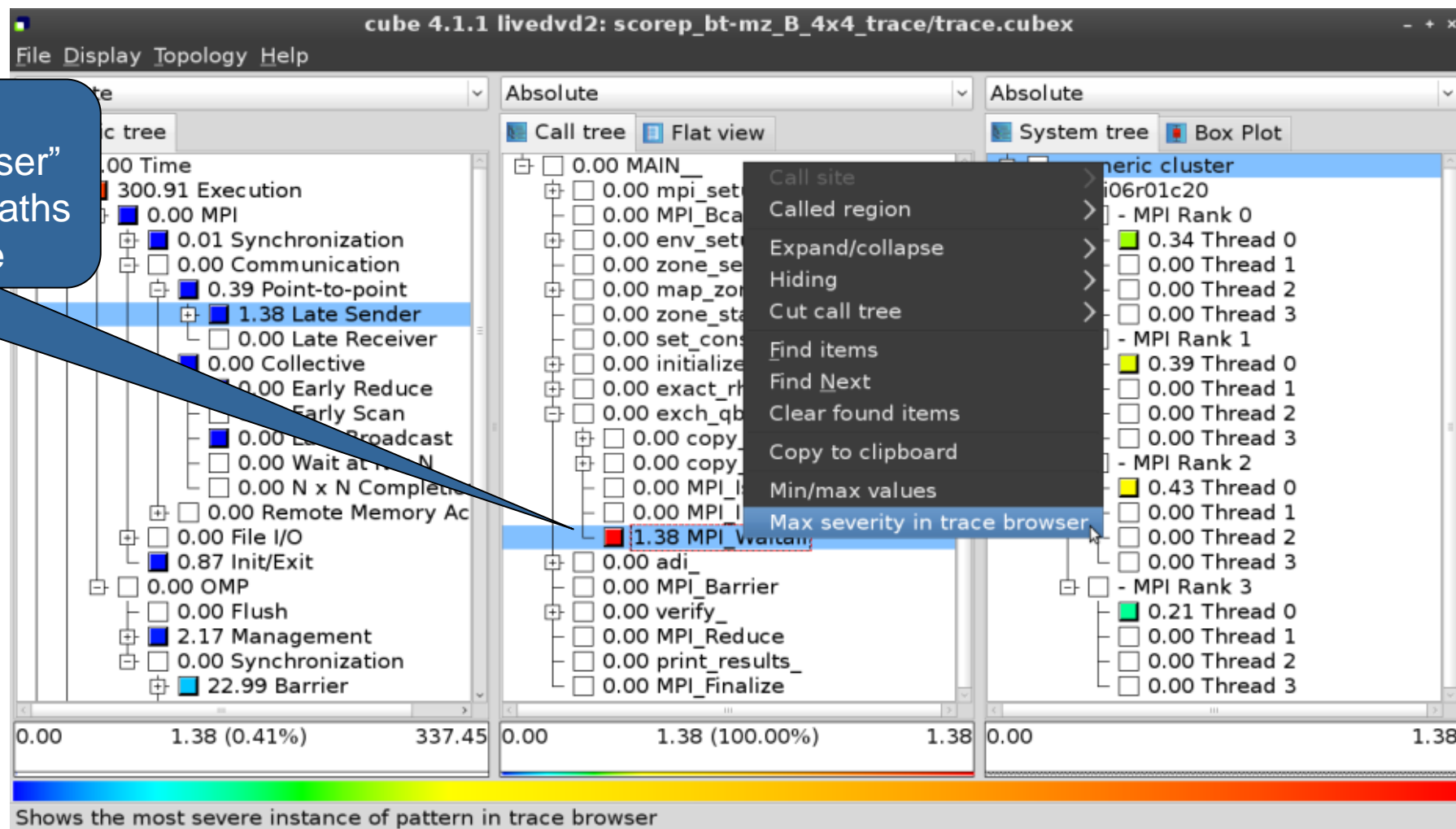


To investigate most severe pattern instances, connect to a trace browser...



# Show most severe pattern instances

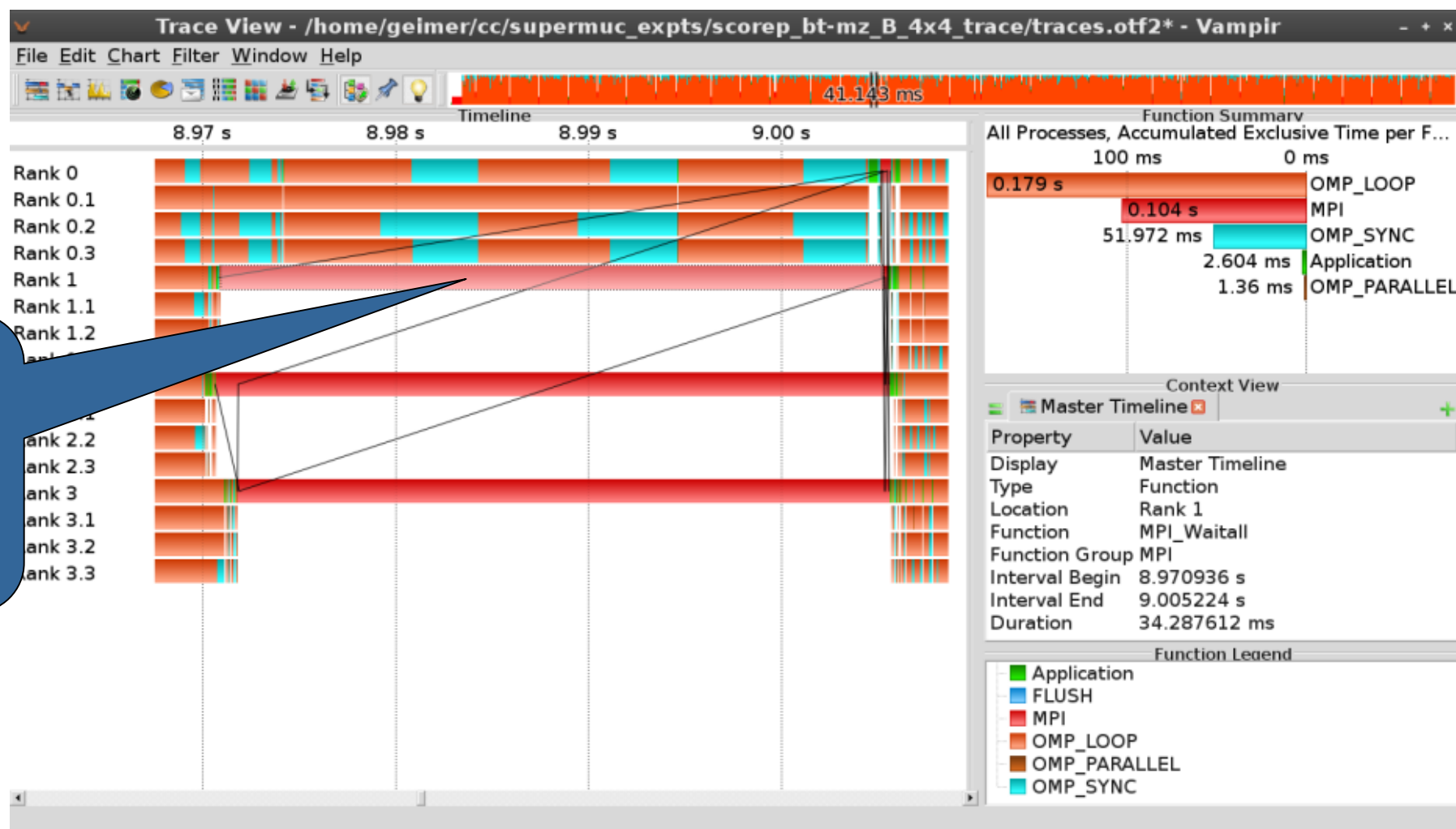
Select  
“Max severity in trace browser”  
from context menu of call paths  
marked with a red frame





# Investigate most severe instance in Vampir

Vampir will automatically zoom to the worst instance in multiple steps (i.e., undo zoom provides more context)



## 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:
  - <http://www.scalasca.org>
- Contact:
  - [mailto: scalasca@fz-juelich.de](mailto:scalasca@fz-juelich.de)

