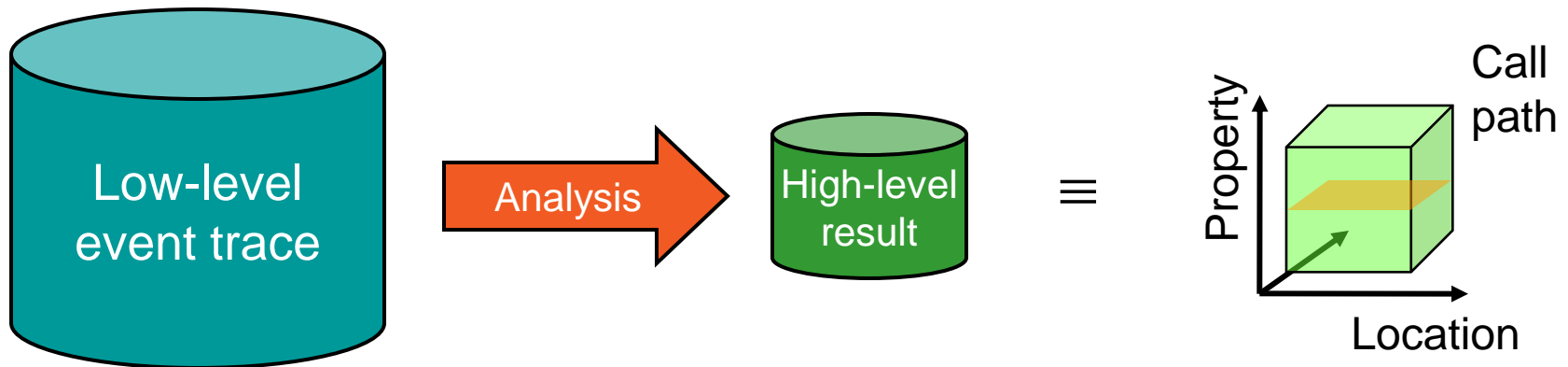# Automatic trace analysis with Scalasca

Marc Schlütter[1], Robert Dietrich[2]

Jülich Supercomputing Centre[1]

ZIH TU Dresden[2]

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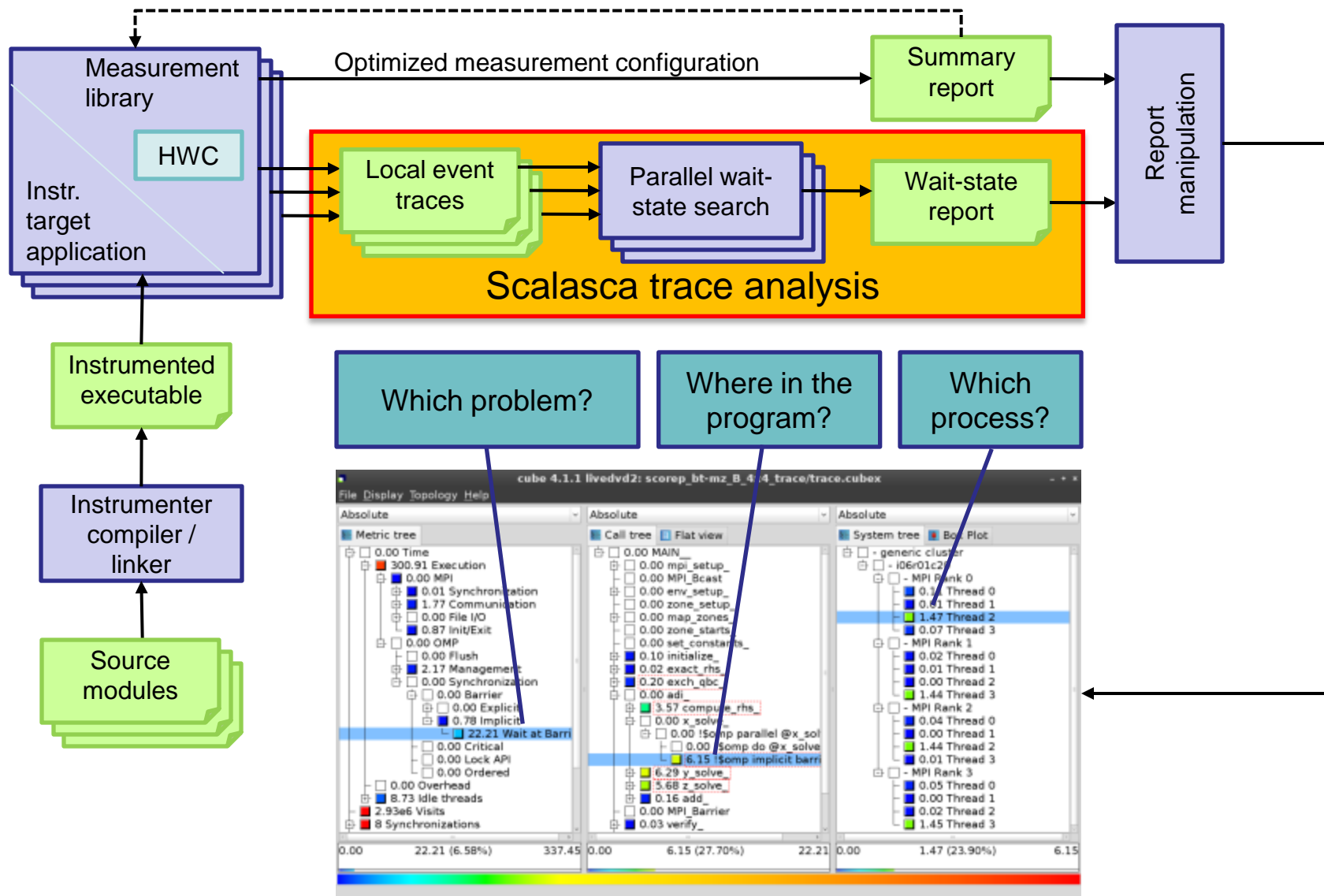


  - Guaranteed to cover the entire event trace
  - Quicker than manual/visual trace analysis
  - Parallel replay analysis exploits available memory & processors to deliver scalability
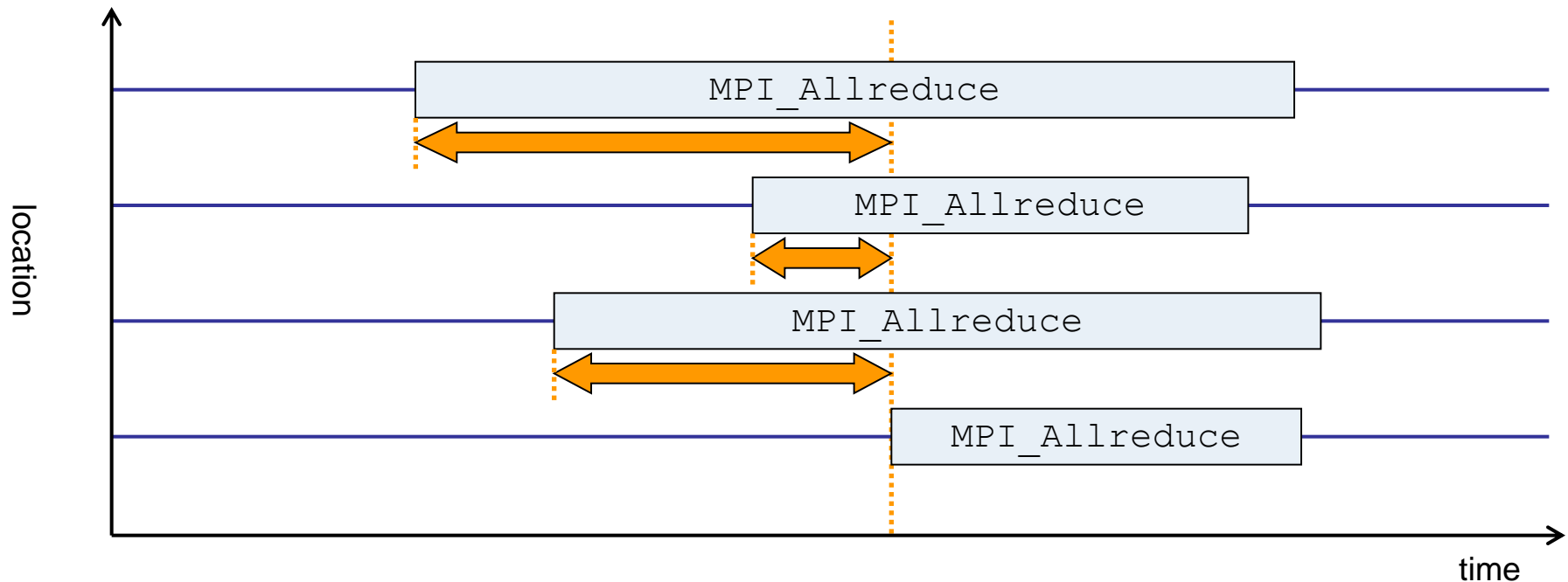
- Project started in 2006
  - Initial funding by Helmholtz Initiative & Networking Fund
  - Many follow-up projects
- Follow-up to pioneering KOJAK project (started 1998)
  - Automatic pattern-based trace analysis
- Now joint development of
  - Jülich Supercomputing Centre

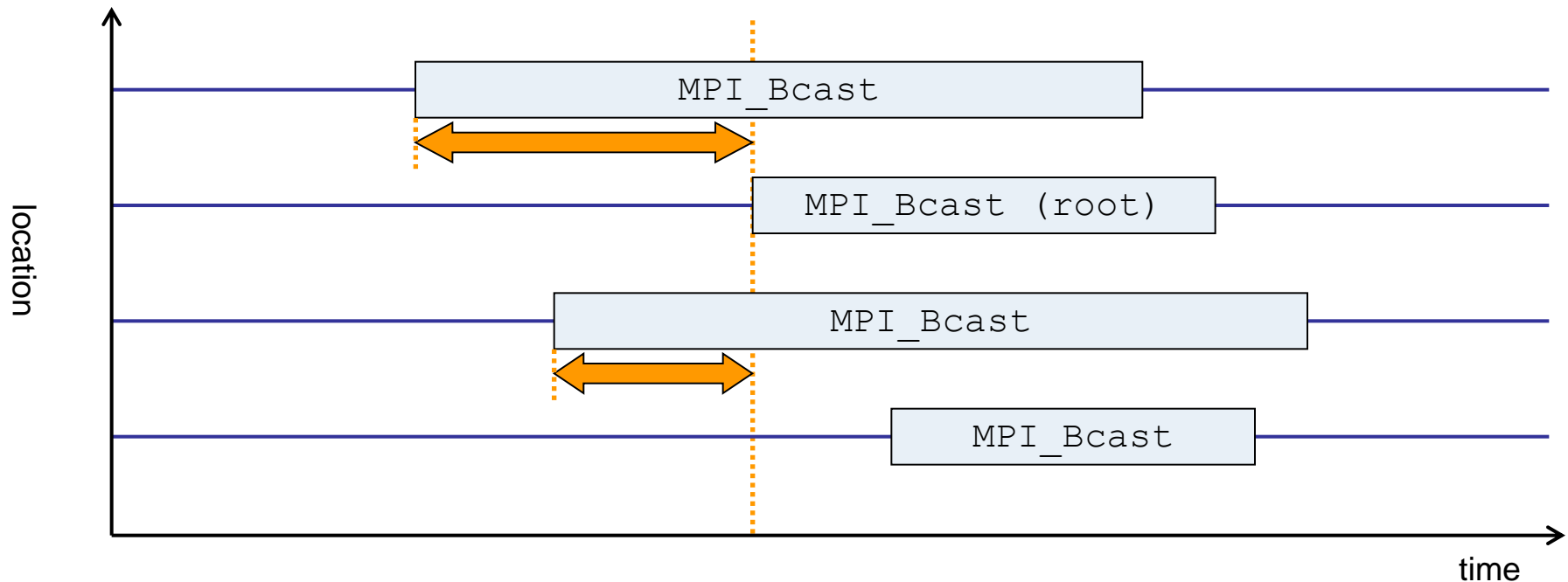  - German Research School for Simulation Sciences

- 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
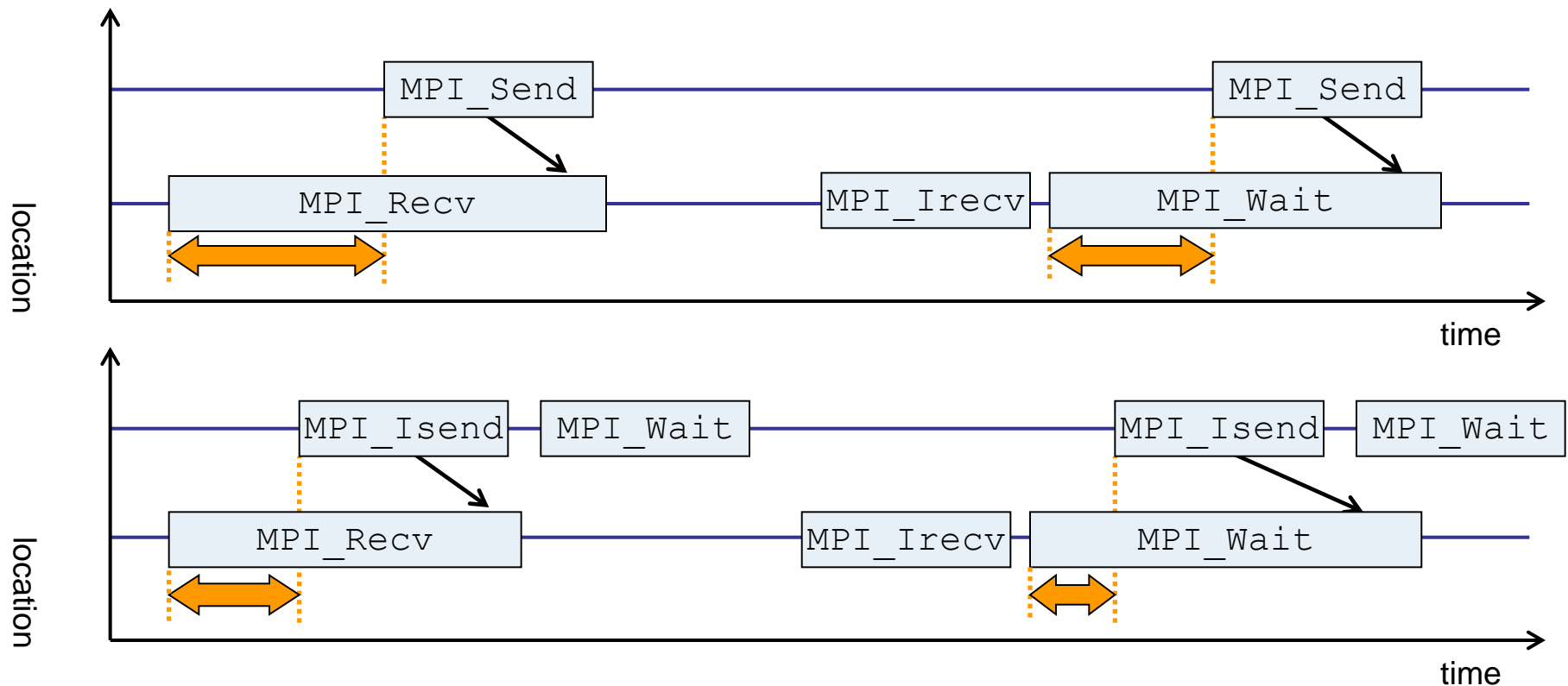  - Unable to handle OTF2 traces containing CUDA events

- 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

- 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

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

- ## One command for (almost) everything…

```
% scalasca
Scalasca 2.1-alpha2
Toolset for scalable performance analysis of large-scale parallel applications
usage: scalasca [-v][-n] {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

- ## Scalasca application instrumenter

```
% skin
Scalasca 2.1-alpha2: application instrumenter (using <PATH_TO_SCOREP>/scorep)
usage: skin [-v] [-comp] [-pdt] [-pomp] [-user] [--*] <compile-or-link-command>
  -comp={all|none|...}: routines to be instrumented by compiler [default: all]
                 (... custom instrumentation specification depends on compiler)
  -pdt:  process source files with PDT/TAU 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 <PATH_TO_SCOREP>/scorep
```

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

- ## Scalasca measurement collection & analysis nexus

```
% scan
Scalasca 2.1-alpha2: 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.
  -m metrics   : Metric specification for measurement.
```

- Scalasca analysis report explorer

```
% square
Scalasca 2.1-alpha2: analysis report explorer
usage: square [-v] [-s] [-f filtfile] [-F] <experiment archive | cube file>
   -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
```

- **`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)

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

```
% export SCOREP_EXPERIMENT_DIRECTORY=scorep_bt-mz_B_4p4x4_sum
% OMP_NUM_THREADS=4 scan mpirun –np 4 –npernode 4./bt-mz_B.4
S=C=A=N: Fri Jan 31 12:50:12 2014: Collect start
/apps/OPENMPI/1.5.4//bin/mpirun -n 4 -npernode 4 ./bt-mz_B.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

 Use the default load factors with threads
 Total number of threads:     16  (  4.0 threads/process)

 [... More application output ...]
S=C=A=N: Fri Jan 31 12:51:04 2014: Collect done (status=0) 52s
```

- Creates experiment directory ./scorep_bt-mz_B_4p4x4_sum

- ## Score summary analysis report

```
% square -s   scorep_bt-mz_B_4p4x4_sum
INFO: Post-processing runtime summarization report...
<PATH>/scorep-score -r ./scorep_bt-mz_B_4p4x4_sum/profile.cubex >
  ./scorep_bt-mz_B_4p4x4_sum/scorep.score
INFO: Score report written to ./scorep_bt-mz_B_4p4x4_sum/scorep.score
```

- ## Post-processing and interactive exploration with CUBE

```
% square   scorep_bt-mz_B_4p4x4_sum
INFO: Displaying ./scorep_bt-mz_B_4p4x4_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

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

- ## Load modules

```
% module use /gpfs/projects/nct00/nct00001/UNITE/tutorial/mf
% module load UNITE
UNITE loaded
% module load scalasca
cube4/4.2.1 loaded
papi/5.3.0 loaded
scorep/1.2.3-beta-intel-openmpi loaded
scalasca/2.1-alpha2-intel-openmpi 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

```
% cd bin.scorep
% cp ../jobscript/marenostrum3/scalasca2.lsf .
% vim scalasca2.lsf

  [...]

module load UNITE scalasca/2.1-alpha1

export SCOREP_FILTERING_FILE=../config/scorep.filt
export SCOREP_TOTAL_MEMORY=50M
export SCOREP_METRIC_PAPI=PAPI_FP_OPS
export OMP_NUM_THREADS=4

scalasca -analyze -t mpirun -n 4 -npernode 4 ./bt-mz_B.4
```

- ## Submit the job

```
% bsub < scalasca2.lsf
```

- ## 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_B_4p4x4_trace/traces.otf2

Opening experiment archive ... done (0.044s).
Reading definition data    ... done (0.007s).
Reading event trace data   ... done (0.593s).
Preprocessing              ... done (1.100s).
Analyzing trace data       ...
  Wait-state detection (fwd)      (1/5) ... done (2.464s).
  Wait-state detection (bwd)      (2/5) ... done (0.052s).
  Synchpoint exchange             (3/5) ... done (0.198s).
  Critical-path & delay analysis  (4/5) ... done (0.306s).
done (3.188s).
Writing analysis report    ... done (0.228s).

Max. memory usage          : 167.840MB

Total processing time      : 5.651s
S=C=A=N: Fri Jan 31 13:09:04 2014: Analyze done (status=0) 7s
```

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

```
% square  scorep_bt-mz_B_4p4x4_trace
INFO: Post-processing runtime summarization report...
INFO: Post-processing trace analysis report...
INFO: Displaying ./scorep_bt-mz_B_4p4x4_trace/trace.cubex...

              [GUI showing trace analysis report]
```

# Post-processed trace analysis report



Additional trace-based metrics in metric hierarchy

Access online metric description via context menu

# Online metric description



13th VI-HPS Tuning Workshop, 10-14 February 2014, BSC, Barcelona, Spain

Critical-path profile shows wall-clock time impact

Critical-path imbalance highlights inefficient parallelism

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

…and select trace file from the experiment directory

# **Sc**alable performance **a**nalysis of **la**rge-**sc**ale parallel **a**pplications

– 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

scalasca