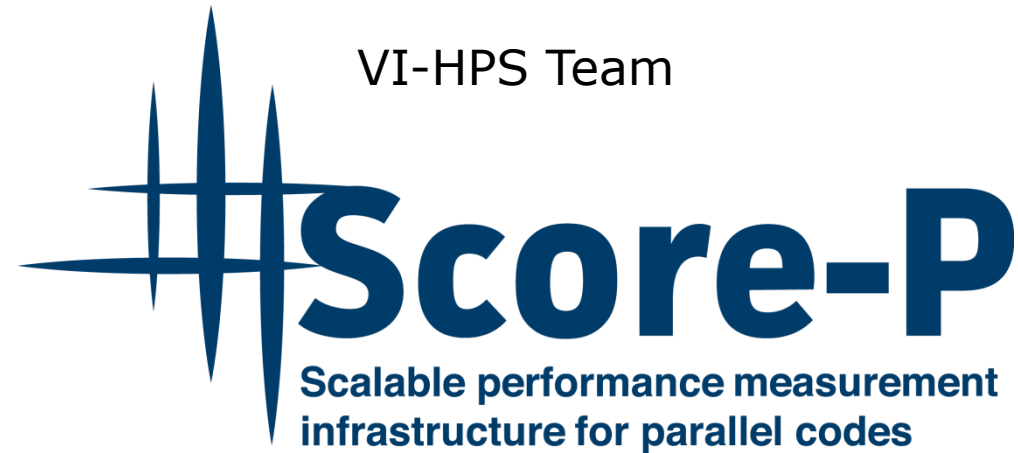


Score-P – A Joint Performance Measurement Run-Time Infrastructure for Periscope, Scalasca, TAU, and Vampir

VI-HPS Team



Congratulations!?

- If you made it this far, you successfully used Score-P to
 - instrument the application
 - analyze its execution with a summary measurement, and
 - examine it with one the interactive analysis report explorer GUIs
- ... revealing the call-path profile annotated with
 - the “Time” metric
 - Visit counts
 - MPI message statistics (bytes sent/received)
- ... but how **good** was the measurement?
 - The measured execution produced the desired valid result
 - however, the execution took rather longer than expected!
 - even when ignoring measurement start-up/completion, therefore
 - it was probably dilated by instrumentation/measurement overhead

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 summary analysis result scoring

```
% scorep-score scorep_bt-mz_sum/profile.cubex
```

Estimated aggregate size of event trace:

Estimated requirements for largest trace buffer (max_buf):

Estimated memory requirements (SCOREP_TOTAL_MEMORY):

(warning: The memory requirements cannot be satisfied by Score-P to avoid intermediate flushes when tracing. Set SCOREP_TOTAL_MEMORY=4G to get the maximum supported memory or reduce requirements using USR regions filters.)

| flt | type | max_buf[B] | visits | time[s] | time[%] | time/visit[us] | region |
|-----|------|----------------|---------------|---------|---------|----------------|--------|
| | ALL | 21,518,477,680 | 6,591,910,441 | 2825.52 | 100.0 | 0.43 | ALL |
| | USR | 21,431,996,118 | 6,574,793,529 | 1166.25 | 41.3 | 0.18 | USR |
| | OMP | 83,841,856 | 16,359,424 | 1533.00 | 54.3 | 93.71 | OMP |
| | COM | 2,351,570 | 723,560 | 2.33 | 0.1 | 3.22 | COM |
| | MPI | 288,136 | 33,928 | 123.94 | 4.4 | 3653.01 | MPI |

160 GB

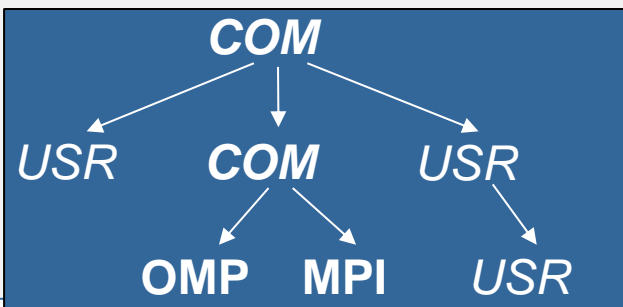
21 GB

21 GB

- Report scoring as textual output

150 GB total memory
21 GB per rank!

- Region/callpath classification
 - **MPI** pure MPI functions
 - **OMP** pure OpenMP regions
 - **USR** user-level computation
 - **COM** "combined" USR+OpenMP/MPI
 - **ANY/ALL** aggregate of all region types



BT-MZ summary analysis report breakdown

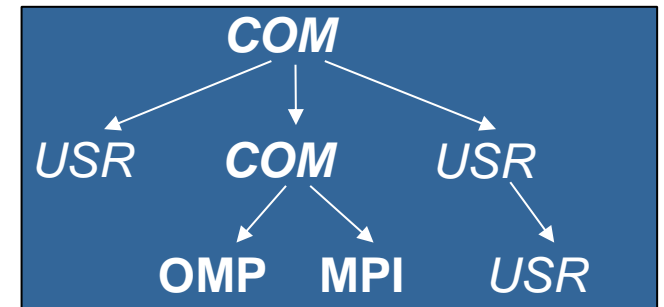
```
% scorep-score -r scorep_bt-mz_sum/profile.cubex
```

```
[...]
```

```
[...]
```

| flt | type | max_buf[B] | visits | time[s] | time[%] | time/visit[us] | region |
|-----|------|----------------|---------------|---------|---------|----------------|--------|
| | ALL | 21,518,477,680 | 6,591,910,441 | 2825.52 | 100.0 | 0.43 | ALL |
| | USR | 21,431,996,118 | 6,574,793,529 | 1166.25 | 41.3 | 0.18 | USR |
| | OMP | 83,841,856 | 16,359,424 | 1533.00 | 54.3 | 93.71 | OMP |
| | COM | 2,351,570 | 723,560 | 2.33 | 0.1 | 3.22 | COM |
| | MPI | 288,136 | 33,928 | 123.94 | 4.4 | 3653.01 | MPI |

| | | | | | | |
|-----|---------------|---------------|--------|------|------|-----------------|
| USR | 6,883,222,086 | 2,110,313,472 | 359.08 | 12.7 | 0.17 | matmul_sub_ |
| USR | 6,883,222,086 | 2,110,313,472 | 263.79 | 9.3 | 0.12 | matvec_sub_ |
| USR | 6,883,222,086 | 2,110,313,472 | 501.53 | 17.8 | 0.24 | binvcrhs_ |
| USR | 293,617,584 | 87,475,200 | 11.75 | 0.4 | 0.13 | binvrhs_ |
| USR | 293,617,584 | 87,475,200 | 21.75 | 0.8 | 0.25 | lhsinit_ |
| USR | 224,028,792 | 68,892,672 | 8.36 | 0.3 | 0.12 | exact_solution_ |



More than
20 GB just for these 6
regions

BT-MZ summary analysis score

- Summary measurement analysis score reveals
 - Total size of event trace would be ~160 GB
 - Maximum trace buffer size would be ~21 GB per rank
 - smaller buffer would require flushes to disk during measurement resulting in substantial perturbation
 - 99.9% of the trace requirements are for USR regions
 - purely computational routines never found on COM call-paths common to communication routines or OpenMP parallel regions
 - These USR regions contribute around 39% of total time
 - however, much of that is very likely to be measurement overhead for frequently-executed small routines
- Advisable to tune measurement configuration
 - Specify an adequate trace buffer size
 - Specify a filter file listing (USR) regions not to be measured

BT-MZ summary analysis report filtering

```
% cat ../config/scorep.filt
```

```
SCOREP_REGION_NAMES_BEGIN
```

```
EXCLUDE
```

```
  binvcrhs*
```

```
  matmul_sub*
```

```
  matvec_sub*
```

```
  exact_solution*
```

```
  binvrhs*
```

```
  lhs*init*
```

```
  timer_*
```

```
SCOREP_REGION_NAMES_END
```

```
% scorep-score -f ../config/scorep.filt -c 2 \  
  scorep_bt-mz_sum/profile.cubex
```

```
Estimated aggregate size of event trace:
```

```
Estimated requirements for largest trace buffer (max_buf):
```

```
Estimated memory requirements (SCOREP_TOTAL_MEMORY):
```

```
(hint: When tracing set SCOREP_TOTAL_MEMORY=215MB to avoid intermediate flushes  
or reduce requirements using USR regions filters.)
```

1624MB

203MB

215MB

- Report scoring with prospective filter listing 6 USR regions

1.6 GB of memory in total,
203 MB per rank!

(Including 2 metric values)

BT-MZ summary analysis report filtering

```
% scorep-score -r -f ../config/scorep.filt \
  scorep_bt-mz_sum/profile.cubex
```

| flt | type | max_buf[B] | visits | time[s] | time[%] | time/visit[us] | region |
|-----|------|----------------|---------------|---------|---------|----------------|---------------------------------|
| - | ALL | 21,518,477,680 | 6,591,910,441 | 2825.52 | 100.0 | 0.43 | ALL |
| - | USR | 21,431,996,118 | 6,574,793,529 | 1166.25 | 41.3 | 0.18 | USR |
| - | OMP | 83,841,856 | 16,359,424 | 1533.00 | 54.3 | 93.71 | OMP |
| - | COM | 2,351,570 | 723,560 | 2.33 | 0.1 | 3.22 | COM |
| - | MPI | 288,136 | 33,928 | 123.94 | 4.4 | 3653.01 | MPI |
| | | | | | | | |
| * | ALL | 86,513,568 | 17,126,753 | 1659.27 | 58.7 | 96.88 | ALL-FLT |
| + | FLT | 21,431,964,112 | 6,574,783,688 | 1166.25 | 41.3 | 0.18 | FLT |
| - | OMP | 83,841,856 | 16,359,424 | 1533.00 | 54.3 | 93.71 | OMP-FLT |
| * | COM | 2,351,570 | 723,560 | 2.33 | 0.1 | 3.22 | COM-FLT |
| - | MPI | 288,136 | 33,928 | 123.94 | 4.4 | 3653.01 | MPI-FLT |
| * | USR | 32,006 | 9,841 | 0.00 | 0.0 | 0.27 | USR-FLT |
| | | | | | | | |
| + | USR | 6,883,222,086 | 2,110,313,472 | 359.08 | 12.7 | 0.17 | matmul_sub |
| + | USR | 6,883,222,086 | 2,110,313,472 | 263.79 | 9.3 | 0.12 | matvec_sub |
| + | USR | 6,883,222,086 | 2,110,313,472 | 501.53 | 17.8 | 0.24 | binvcrhs |
| + | USR | 293,617,584 | 87,475,200 | 11.75 | 0.4 | 0.13 | binvrhs |
| + | USR | 293,617,584 | 87,475,200 | 21.75 | 0.8 | 0.25 | lhsinit |
| + | USR | 224,028,792 | 68,892,672 | 8.36 | 0.3 | 0.12 | exact_solution |
| - | OMP | 6,715,008 | 617,472 | 0.15 | 0.0 | 0.24 | !\$omp parallel @exch_qbc.f:215 |

- Score report breakdown by region

Filtered routines marked with '+'

BT-MZ filtered summary measurement

```
% cd bin.scorep
% cp ../jobscript/inti/scorep.msub .
% vim scorep.msub

PROCS=8
CLASS=C
export SCOREP_EXPERIMENT_DIRECTORY=scorep_bt-mz_sum_filter
export SCOREP_FILTERING_FILE=../config/scorep.filt
#export SCOREP_TOTAL_MEMORY=100M
#export SCOREP_METRIC_PAPI=PAPI_TOT_INS,PAPI_TOT_CYC
#export SCOREP_ENABLE_TRACING=true

# launch
EXE=./bt-mz_${CLASS}.${PROCS}
ccc_mprun -n $PROCS $EXE

% ccc_msub ./scorep.msub
```

- Set new experiment directory and re-run measurement with new filter configuration
- Submit job

Score-P filtering

```
% cat ../config/scorep.filt
SCOREP_REGION_NAMES_BEGIN
EXCLUDE
    binvcrhs*
    matmul_sub*
    matvec_sub*
    exact_solution*
    binvrhs*
    lhs*init*
    timer_*
SCOREP_REGION_NAMES_END

% export SCOREP_FILTERING_FILE=\
../config/scorep.filt
```

Region name
filter block
using wildcards

Apply filter

- Filtering by source file name
 - All regions in files that are excluded by the filter are ignored
- Filtering by region name
 - All regions that are excluded by the filter are ignored
 - Overruled by source file filter for excluded files
- Apply filter by
 - exporting `SCOREP_FILTERING_FILE` environment variable
- Apply filter at
 - Run-time
 - Compile-time (GCC-plugin only)
 - Add cmd-line option `--instrument-filter`
 - No overhead for filtered regions but recompilation

Source file name filter block

- Keywords
 - Case-sensitive
 - SCOREP_FILE_NAMES_BEGIN, SCOREP_FILE_NAMES_END
 - Define the source file name filter block
 - Block contains EXCLUDE, INCLUDE rules
 - EXCLUDE, INCLUDE rules
 - Followed by one or multiple white-space separated source file names
 - Names can contain bash-like wildcards *, ?, []
 - Unlike bash, * may match a string that contains slashes
 - EXCLUDE, INCLUDE rules are applied in sequential order
 - Regions in source files that are excluded after all rules are evaluated, get filtered

```
# This is a comment
SCOREP_FILE_NAMES_BEGIN
    # by default, everything is included
    EXCLUDE */foo/bar*
    INCLUDE */filter_test.c
SCOREP_FILE_NAMES_END
```

Region name filter block

- Keywords
 - Case-sensitive
 - SCOREP_REGION_NAMES_BEGIN,
SCOREP_REGION_NAMES_END
 - Define the region name filter block
 - Block contains EXCLUDE, INCLUDE rules
 - EXCLUDE, INCLUDE rules
 - Followed by one or multiple white-space separated region names
 - Names can contain bash-like wildcards *, ?, []
 - EXCLUDE, INCLUDE rules are applied in sequential order
 - Regions that are excluded after all rules are evaluated, get filtered

```
# This is a comment
SCOREP_REGION_NAMES_BEGIN
    # by default, everything is included
    EXCLUDE *
    INCLUDE bar foo
           baz
           main
SCOREP_REGION_NAMES_END
```

Region name filter block, mangling

- Name mangling
 - Filtering based on names seen by the measurement system
 - Dependent on compiler
 - Actual name may be mangled
- scorep-score names as starting point
(e.g. `matvec_sub_`)
 - Use `*` for Fortran trailing underscore(s) for portability
 - Use `?` and `*` as needed for full signatures or overloading

```
void bar(int* a) {  
    *a++;  
}  
int main() {  
    int i = 42;  
    bar(&i);  
    return 0;  
}
```

```
# filter bar:  
# for gcc-plugin, scorep-score  
# displays 'void bar(int*)',  
# other compilers may differ
```

```
SCOREP_REGION_NAMES_BEGIN  
    EXCLUDE void?bar(int?)  
SCOREP_REGION_NAMES_END
```

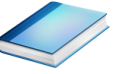

Further information

- Community instrumentation & measurement infrastructure
 - Instrumentation (various methods)
 - Basic and advanced profile generation
 - Event trace recording
 - Online access to profiling data
- Available under New BSD open-source license
- Documentation & Sources:
 - <http://www.score-p.org>
- User guide also part of installation:
 - `<prefix>/share/doc/scorep/{pdf,html}/`
- Support and feedback: support@score-p.org
- Subscribe to news@score-p.org, to be up to date

Score-P: Specialized Measurements and Analyses



Mastering build systems



- Hooking up the Score-P instrumenter `scorep` into complex build environments like *Autotools* or *CMake* was always challenging
- Score-P provides new convenience wrapper scripts to simplify this (since Score-P 2.0)
- *Autotools* and *CMake* need the used compiler already in the *configure step*, but instrumentation should not happen in this step, only in the *build step*

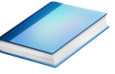
```
% SCOREP_WRAPPER=off \  
> cmake .. \  
> -DCMAKE_C_COMPILER=scorep-icc \  
> -DCMAKE_CXX_COMPILER=scorep-icpc
```

Disable instrumentation in the *configure step*

Specify the wrapper scripts as the compiler to use

- Allows to pass addition options to the Score-P instrumenter and the compiler via environment variables without modifying the *Makefiles*
- Run `scorep-wrapper --help` for a detailed description and the available wrapper scripts of the Score-P installation

Mastering C++ applications



- Automatic compiler instrumentation greatly disturbs C++ applications because of frequent/short function calls => Use sampling instead
- Novel combination of sampling events and instrumentation of MPI, OpenMP, ...
 - Sampling replaces compiler instrumentation (instrument with `--nocompiler` to further reduce overhead) => Filtering not needed anymore
 - Instrumentation is used to get accurate times for parallel activities to still be able to identify patterns of inefficiencies
- Supports profile and trace generation

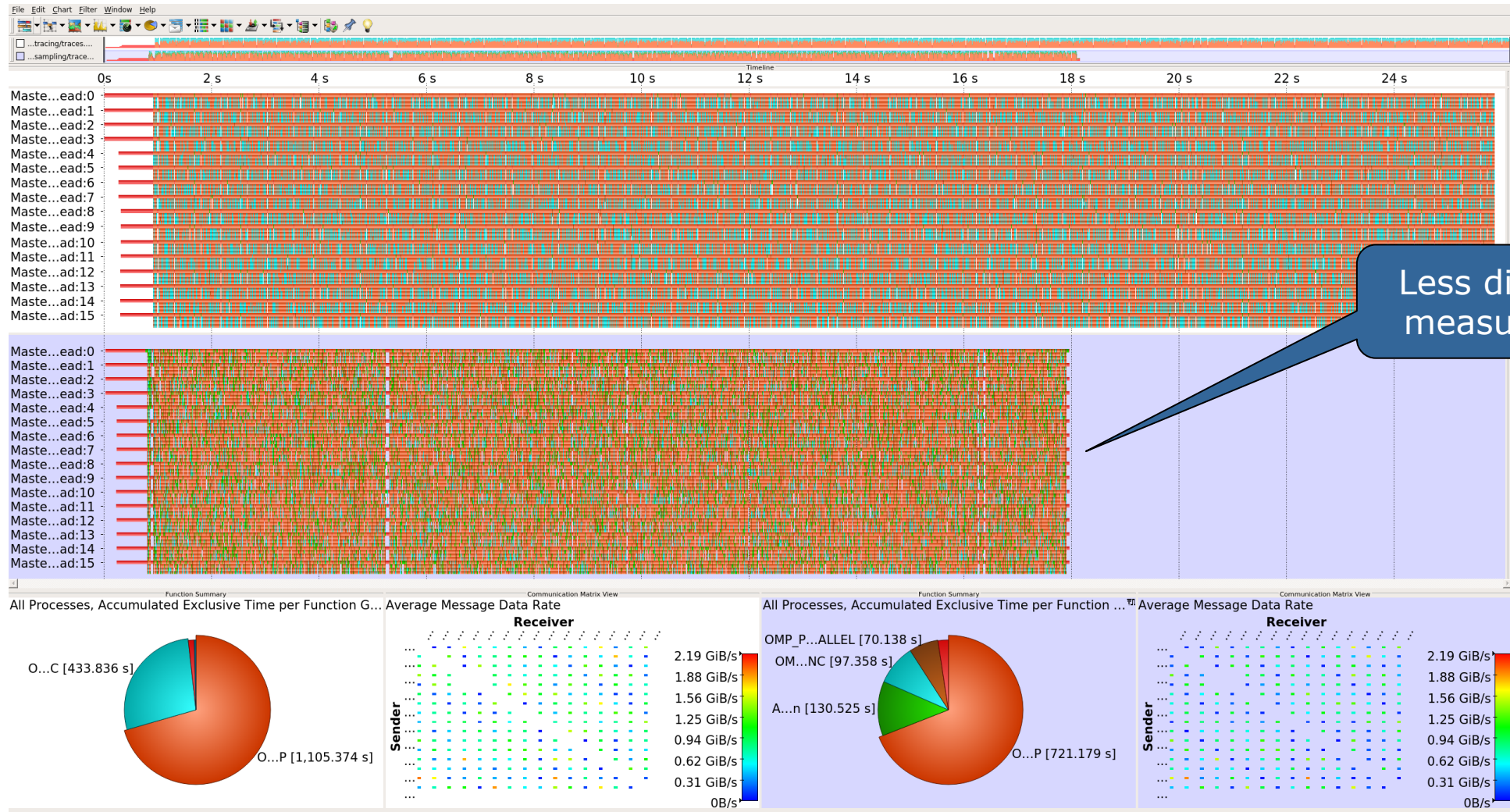
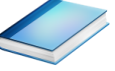
```
% export SCOREP_ENABLE_UNWINDING=true
% # use the default sampling frequency
% #export SCOREP_SAMPLING_EVENTS=perf_cycles@2000000

% OMP_NUM_THREADS=4 mpiexec -np 4 ./bt-mz_W.4
```

- Set new configuration variable to enable sampling

- Available since Score-P 2.0, only x86-64 supported currently

Mastering C++ applications



Mastering application memory usage



- Determine the maximum heap usage per process
- Find high frequent small allocation patterns
- Find memory leaks
- Support for:
 - C, C++, MPI, and SHMEM (Fortran only for GNU Compilers)
 - Profile and trace generation (profile recommended)
 - Memory leaks are recorded only in the profile
 - Resulting traces are not supported by Scalasca yet

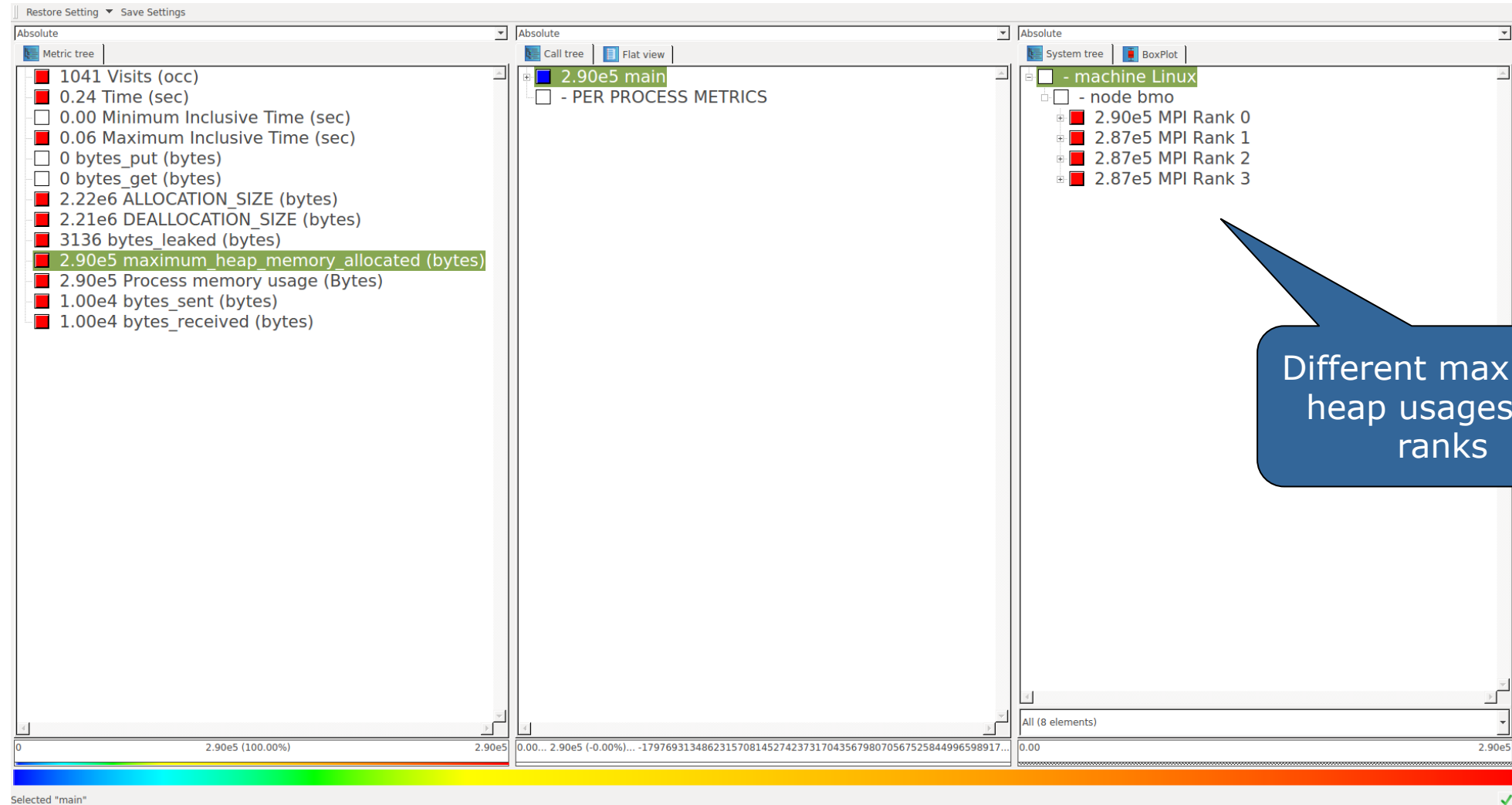
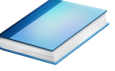
```
% export SCOREP_MEMORY_RECORDING=true
% export SCOREP_MPI_MEMORY_RECORDING=true

% OMP_NUM_THREADS=4 mpiexec -np 4 ./bt-mz_W.4
```

- Set new configuration variable to enable memory recording

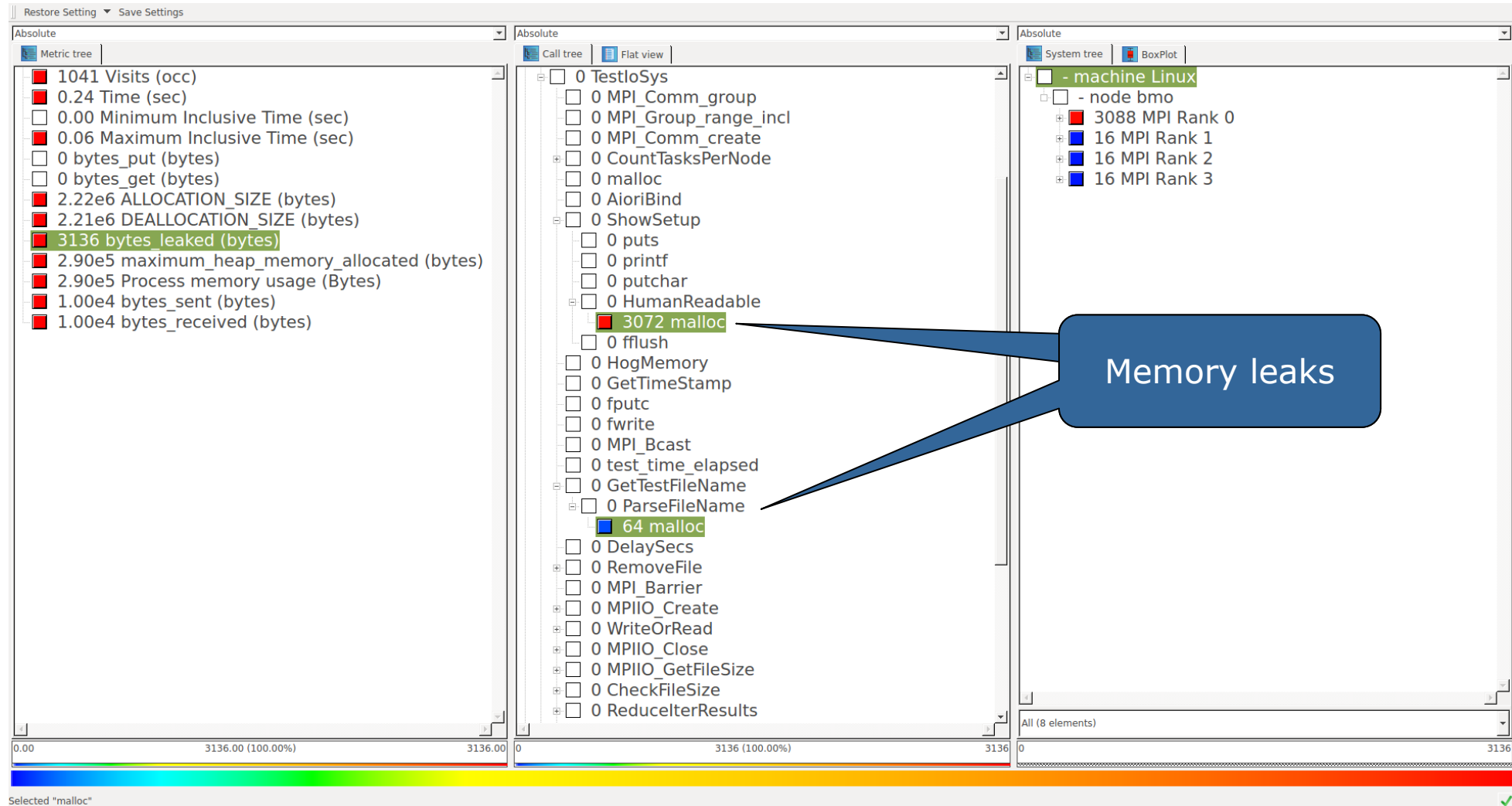
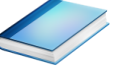
- Available since Score-P 2.0

Mastering application memory usage

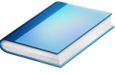


Different maximum heap usages per ranks

Mastering application memory usage



Mastering heterogeneous applications



- Record CUDA applications and device activities

```
% export SCOREP_CUDA_ENABLE=gpu,kernel,idle
```

- Record OpenCL applications and device activities

```
% export SCOREP_OPENCL_ENABLE=api,kernel
```

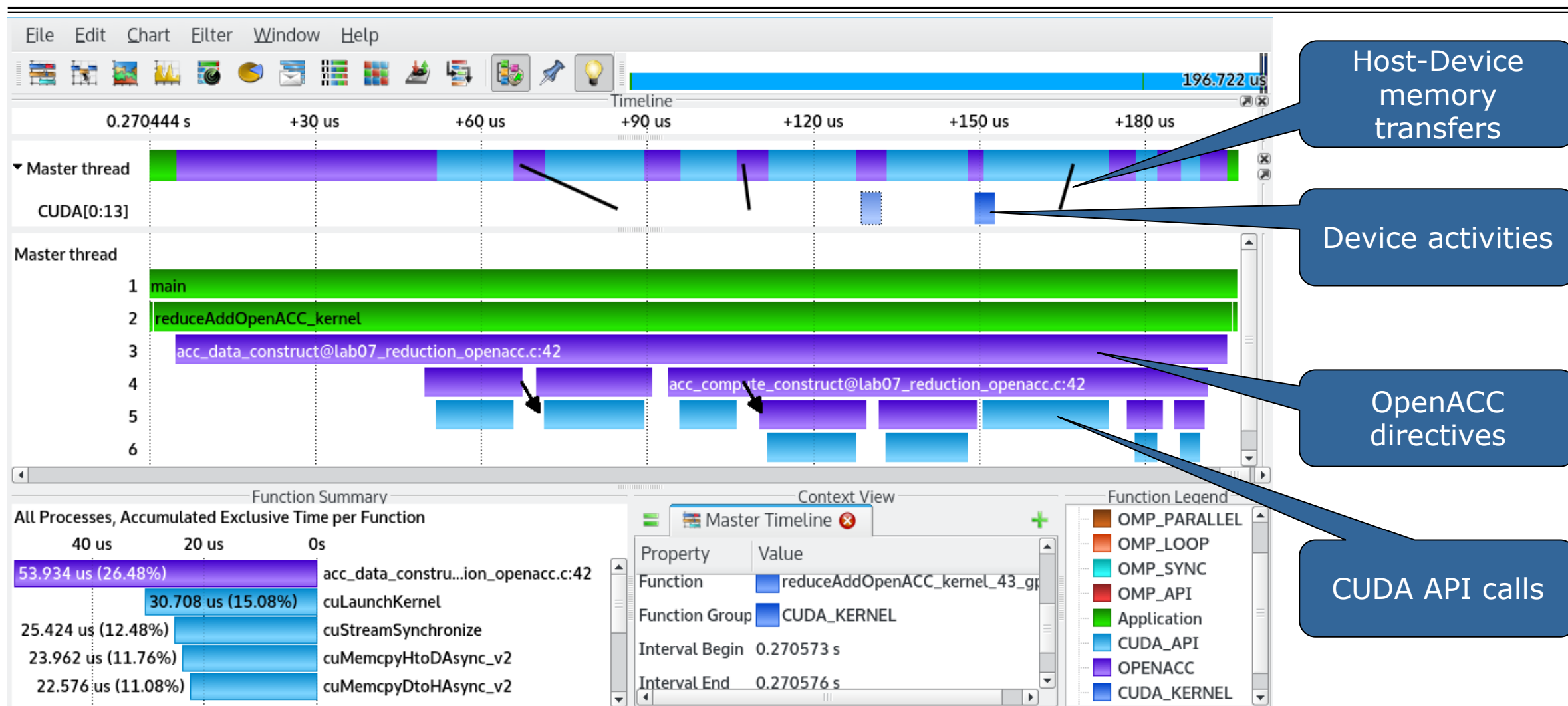
- Record OpenACC applications

```
% export SCOREP_OPENACC_ENABLE=yes
```

- Can be combined with CUDA if it is a NVIDIA device

```
% export SCOREP_CUDA_ENABLE=kernel
```

Mastering heterogeneous applications



Enriching measurements with performance counters



- Record metrics from PAPI:

```
% export SCOREP_METRIC_PAPI=PAPI_TOT_CYC
% export SCOREP_METRIC_PAPI_PER_PROCESS=PAPI_L3_TCM
```

- Use PAPI tools to get available metrics and valid combinations:

```
% papi_avail
% papi_native_avail
```

- Record metrics from Linux perf:

```
% export SCOREP_METRIC_PERF=cpu-cycles
% export SCOREP_METRIC_PERF_PER_PROCESS=LLC-load-misses
```

- Use the `perf` tool to get available metrics and valid combinations:

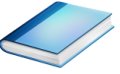
```
% perf list
```

- Write your own metric plugin

- Repository of available plugins: <https://github.com/score-p>

Only the master thread records the metric (assuming all threads of the process access the same L3 cache)

Score-P user instrumentation API



- No replacement for automatic compiler instrumentation
- Can be used to further subdivide functions
 - E.g., multiple loops inside a function
- Can be used to partition application into coarse grain phases
 - E.g., initialization, solver, & finalization
- Enabled with `--user` flag to Score-P instrumenter
- Available for Fortran / C / C++

Score-P user instrumentation API (Fortran)



```
#include "scorep/SCOREP_User.inc"

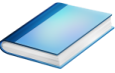
subroutine foo(...)
  ! Declarations
  SCOREP_USER_REGION_DEFINE( solve )

  ! Some code...
  SCOREP_USER_REGION_BEGIN( solve, "<solver>", \
                           SCOREP_USER_REGION_TYPE_LOOP )

  do i=1,100
    [...]
  end do
  SCOREP_USER_REGION_END( solve )
  ! Some more code...
end subroutine
```

- Requires processing by the C preprocessor
 - For most compilers, this can be automatically achieved by having an uppercase file extension, e.g., `main.F` or `main.F90`

Score-P user instrumentation API (C/C++)

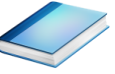


```
#include "scorep/SCOREP_User.h"

void foo()
{
    /* Declarations */
    SCOREP_USER_REGION_DEFINE( solve )

    /* Some code... */
    SCOREP_USER_REGION_BEGIN( solve, "<solver>",
                             SCOREP_USER_REGION_TYPE_LOOP )
    for (i = 0; i < 100; i++)
    {
        [...]
    }
    SCOREP_USER_REGION_END( solve )
    /* Some more code... */
}
```

Score-P user instrumentation API (C++)



```
#include "scorep/SCOREP_User.h"

void foo()
{
    // Declarations

    // Some code...
    {
        SCOREP_USER_REGION( "<solver>",
                           SCOREP_USER_REGION_TYPE_LOOP )
        for (i = 0; i < 100; i++)
        {
            [...]
        }
    }
    // Some more code...
}
```


Score-P measurement control API



- Can be used to temporarily disable measurement for certain intervals
 - Annotation macros ignored by default
 - Enabled with `--user` flag

```
#include "scorep/SCOREP_User.inc"
```

```
subroutine foo(...)
  ! Some code...
  SCOREP_RECORDING_OFF()
  ! Loop will not be measured
  do i=1,100
    [...]
  end do
  SCOREP_RECORDING_ON()
  ! Some more code...
end subroutine
```

Fortran (requires C preprocessor)

```
#include "scorep/SCOREP_User.h"
```

```
void foo(...) {
  /* Some code... */
  SCOREP_RECORDING_OFF()
  /* Loop will not be measured */
  for (i = 0; i < 100; i++) {
    [...]
  }
  SCOREP_RECORDING_ON()
  /* Some more code... */
}
```

C / C++

Score-P: Conclusion and Outlook



Project management

- Ensure a single official release version at all times which will always work with the tools
- Allow experimental versions for new features or research
- Commitment to joint long-term cooperation
 - Development based on meritocratic governance model
 - Open for contributions and new partners

Future features

- Scalability to maximum available CPU core count
- Support for emerging architectures and new programming models
- Features currently worked on:
 - User provided wrappers to 3rd party libraries
 - Hardware and MPI topologies
 - Basic support of measurements without re-compiling/-linking
 - I/O recording
 - Java recording
 - Persistent memory recording (e.g., PMEM, NVRAM, ...)

