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

Peter Philippen
p.philippen@fz-juelich.de

2013-04-03

- Several performance tools co-exist
  - With own measurement systems and output formats
- Complementary features and overlapping functionality
- Redundant effort for development and maintenance
  - Limited or expensive interoperability
- Complications for user experience, support, training

| Vampir | Scalasca | TAU | Periscope |
|---|---|---|---|
| VampirTrace OTF | EPILOG / CUBE | TAU native formats | Online measurement |

- Start a community effort for a common infrastructure
  - Score-P instrumentation and measurement system
  - Common data formats OTF2 and CUBE4

- Developer perspective:
  - Save manpower by sharing development resources
  - Invest in new analysis functionality and scalability
  - Save efforts for maintenance, testing, porting, support, training

- User perspective:
  - Single learning curve
  - Single installation, fewer version updates
  - Interoperability and data exchange

- SILC project funded by BMBF

- Close collaboration PRIMA project funded by DOE

- Continued in LMAC project (BMBF)

- Other projects:
  - H4H
  - HOPSA
  - Catwalk

GEFÖRDERT VOM

Bundesministerium
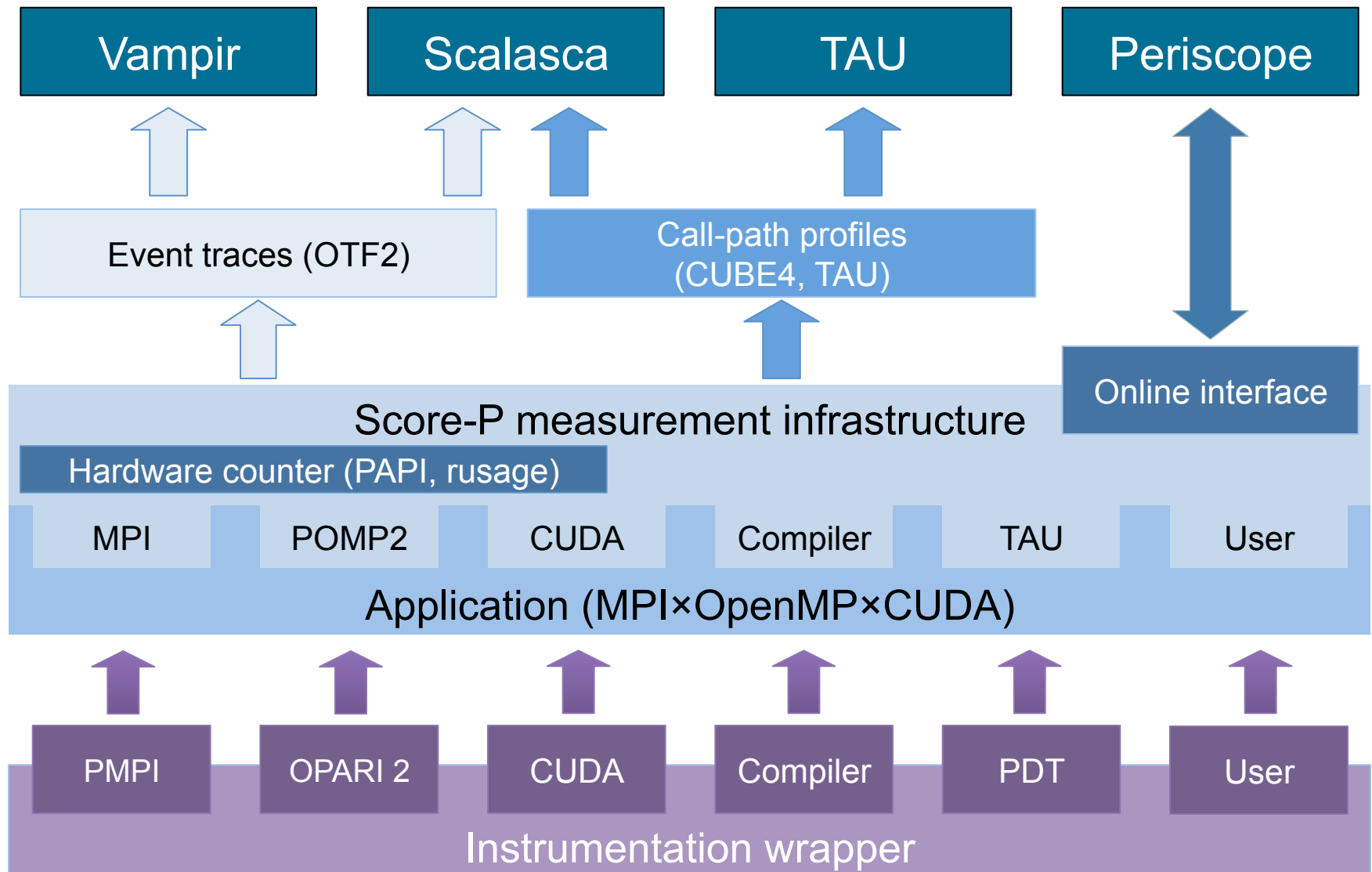für Bildung
und Forschung

- Forschungszentrum Jülich, Germany

- German Research School for Simulation Sciences, Aachen, Germany

- Gesellschaft für numerische Simulation mbH Braunschweig, Germany

- RWTH Aachen, Germany

- Technische Universität Dresden, Germany

- Technische Universität München, Germany
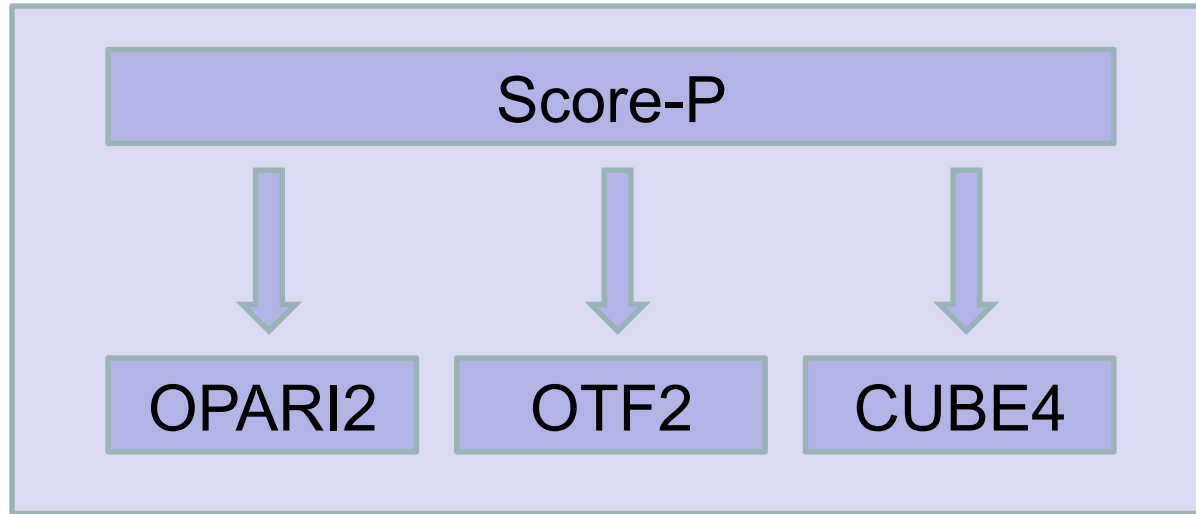
- University of Oregon, Eugene, USA

- Provide typical functionality for HPC performance tools
- Support all fundamental concepts of partner's tools

- Instrumentation (various methods)
- Flexible measurement without re-compilation:
  – Basic and advanced profile generation
  – Event trace recording
  – Online access to profiling data

- MPI, OpenMP, and hybrid parallelism (and serial)
- Enhanced functionality (OpenMP 3.0, CUDA, highly scalable I/O)

- ## Functional requirements
  - – Generation of call-path profiles and event traces
  - – Using direct instrumentation, later also sampling
  - – Recording time, visits, communication data, hardware counters
  - – Access and reconfiguration also at runtime
  - – Support for MPI, OpenMP, basic CUDA, and all combinations
    - • Later also OpenCL/HMPP/PTHREAD/…

- ## Non-functional requirements
  - – Portability: all major HPC platforms
  - – Scalability: petascale
  - – Low measurement overhead
  - – Easy and uniform installation through UNITE framework
  - – Robustness
  - – Open Source: New BSD License

Score-P

↓　　　　↓　　　　↓

OPARI2　　OTF2　　CUBE4

- Separate, stand-alone packages

- Uniform look & feel

- Common functionality factored out

- Automated builds and tests

- Instrumenter scorep

- Links application to measurement library

- Records time, visits, communication metrics, hardware counter using internal memory management
  - Store data in thread-local chunks of preallocated memory
  - Efficient utilization of available memory
  - Minimize perturbation/overhead
  - Useful for unification too
  - Variable number of threads
  - Access data during runtime

- Switch modes (tracing/profiling/online) w/o recompilation

- Layer architecture, minimal coupling, easy to extend

- Prefix compile/link commands, e.g.

```
mpicc -c foo.c    ➡    scorep mpicc -c foo.c
```
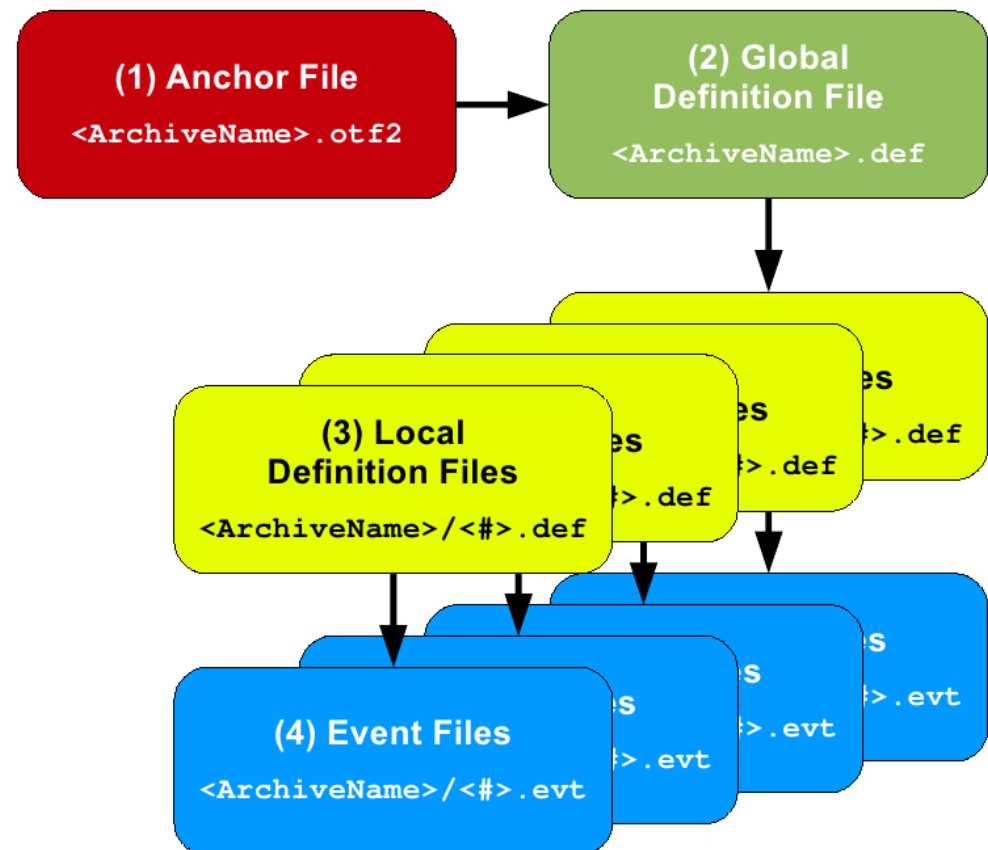
or, more convenient

```
# in Makefile
MPICC = $(PREP) mpicc

% make PREP="scorep [options]"
```
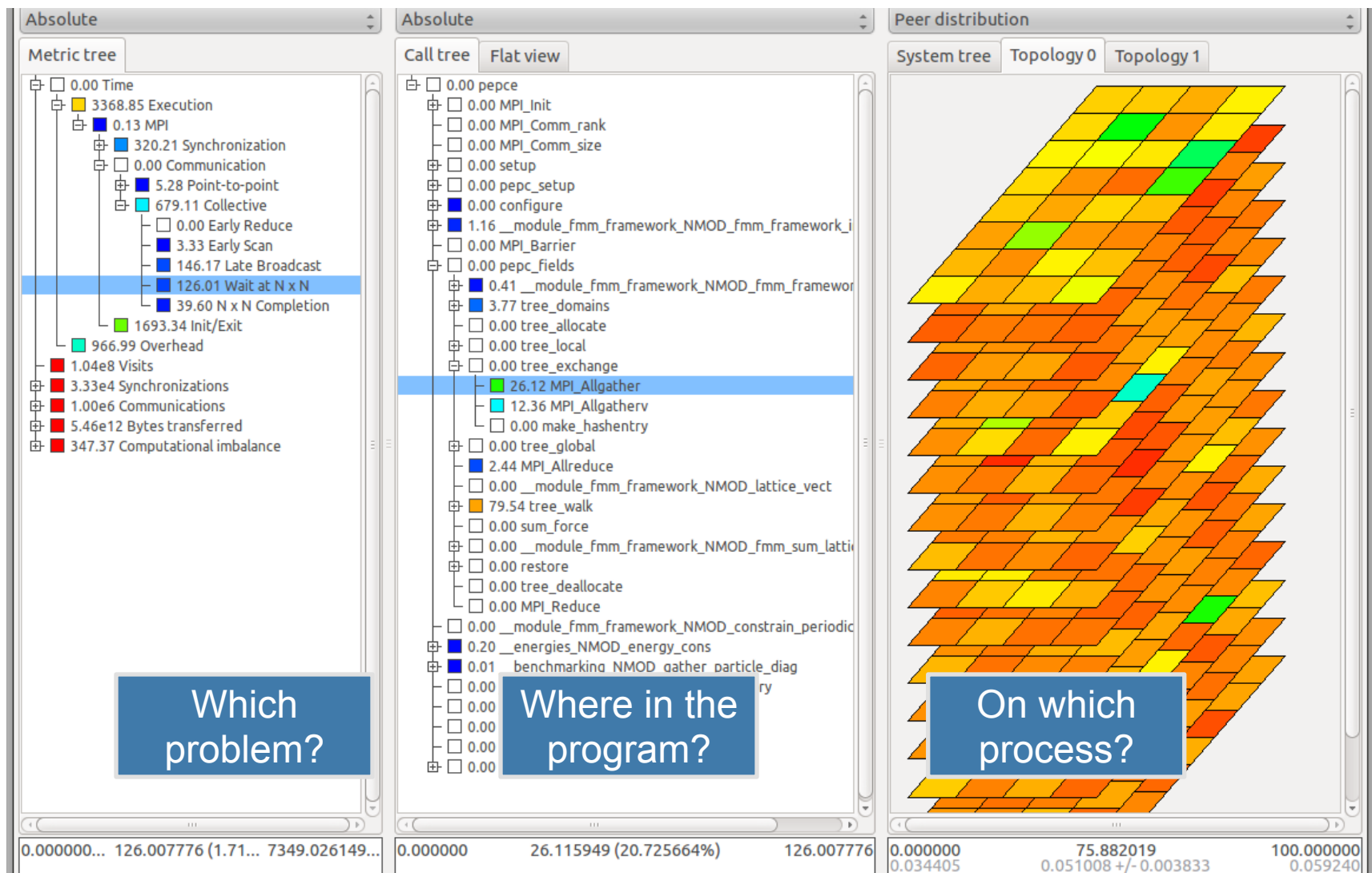
- Customization via options, e.g. --pdt --user
  % scorep --help for all options
- Automatic paradigm detection (serial/OpenMP/MPI/hybrid)
- Manual approach, get libs and includes via scorep_config

- Just run your instrumented program!

- Customize via environment variables
    - Switch between profiling/tracing
    - Select MPI groups to record, e.g. P2P, COLL, …
    - Specify total measurement memory
    - Trace output options, e.g. compression, SION
    - Hardware counter (PAPI, rusage)

- Customize via files
    - In/exclude regions by name/wildcards from recording (filtering)
    - Restrict trace recording to specified executions of a region (selective tracing)

- Data written to uniquely named directory

- Typical event trace data format
  - Event record types + definition record types
  - Stored per process/thread in temporal order
- Multi-file format
  - Anchor file
  - Global and local definitions + mappings
  - Event files
- OTF2 API + read/write library + support tools

- Generic format for call-path profiles of parallel apps, three-dimensional performance space:
  Metric hierarchy x Call-tree hierarchy x System hierarchy

- File organization, since version 4
  - Metadata stored in XML format
  - Metric values stored in binary format
    (Two files per metric: data + index for storage-efficient sparse representation)
  - Bundled into one file

- Optimized for
  - High write bandwidth
  - Fast interactive analysis through incremental loading

# CUBE4 GUI

- The Online Access Interface (OA) allows to:
  - Retrieve measured performance data
    while application is still running
  - Control application execution, interrupt between phases
  - Configure and re-configure measurement settings
  - Successive measurement iterations with refined settings
- Designed according to Periscope's analysis method
  - Multiple performance experiments within one run
  - Remote analysis with measurements acquisition over sockets
  - Faster analysis: one iteration of the time loop could be enough
- May induce new analysis methods

- Early preview release at SC'10

- Score-P BOF at SC'11

- Latest release: V 1.1.1

- Scalability to maximum available CPU core count
- Support for OpenCL, and OpenACC
- Support for sampling, binary instrumentation
- Support for new programming models, e.g. PGAS
- Support for new architectures

- Ensure a single official release version at every time which will always work with the tools
- Allow experimental versions for new features or research

- Open for new partners after SILC funding period
- Commitment to joint long-term cooperation

- Future integration in Open MPI releases

# Long-term Cooperation, Other Projects

- SILC (2009 – 2011)
  - Initial version

- PRIMA (2009 – 2012)
  - Integration with TAU

- H4H (2010 – 2013)
  - Heterogeneous architectures

- HOPSA (EU-Russia project, 2011-2012)
  - Integration with system monitoring

- LMAC (2011 – 2014)
  - Evolution of Score-P, performance dynamics

- GASPI (2011-2014)
  - PGAS, performance API, tool support

GEFÖRDERT VOM

Bundesministerium
für Bildung
und Forschung

GEFÖRDERT VOM

Bundesministerium
für Bildung
und Forschung

GEFÖRDERT VOM

Bundesministerium
für Bildung
und Forschung

- Common measurement part is community effort
  - Use released resources for analysis
- Unite file formats, open for adoption
- Online access interface, open for adoption
- Scalability/flexibility/overhead limitations addressed
- Easy to extend due to layered architecture
- Robust due to extensive and continuous testing
- Long-term commitment
  - Partners have extensive history of collaboration

- Dieter an Mey, Scott Biersdorf, Kai Diethelm, Dominic Eschweiler, Markus Geimer, Michael Gerndt, Houssam Haitof, Rene Jäkel, Koutaiba Kassem, Andreas Knüpfer, Daniel Lorenz, Suzanne Millstein, Bernd Mohr, Yury Oleynik, Peter Philippen, Christian Rössel, Pavel Saviankou, Dirk Schmidl, Sameer Shende, Wyatt Spear, Ronny Tschüter, Michael Wagner, Bert Wesarg, Felix Wolf, Brian Wylie