# Extending the Functionality of Score-P through Plugins: Interfaces and Examples

## 10th International Parallel Tools Workshop, Oct 5th 2016

Robert Schöne
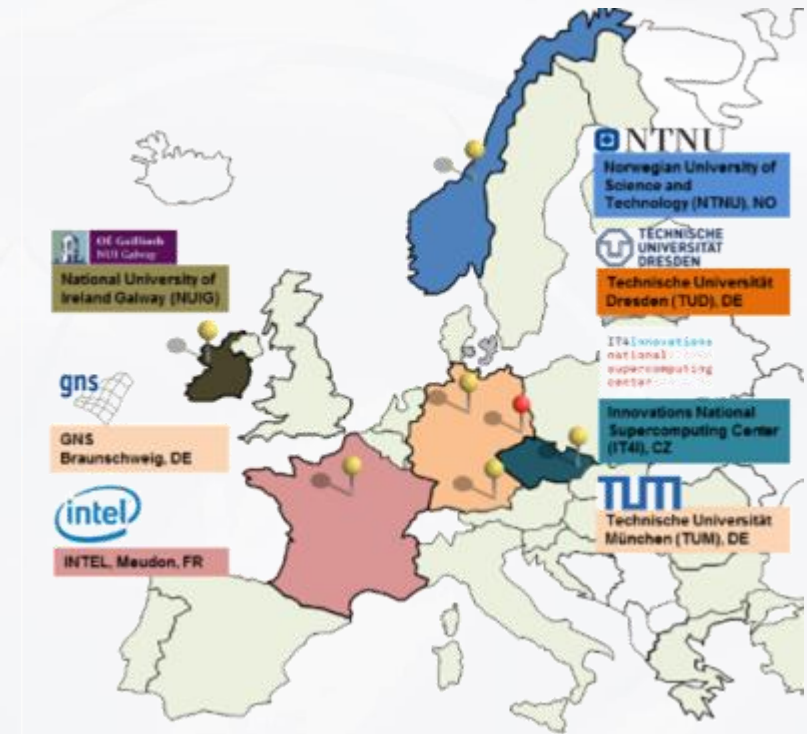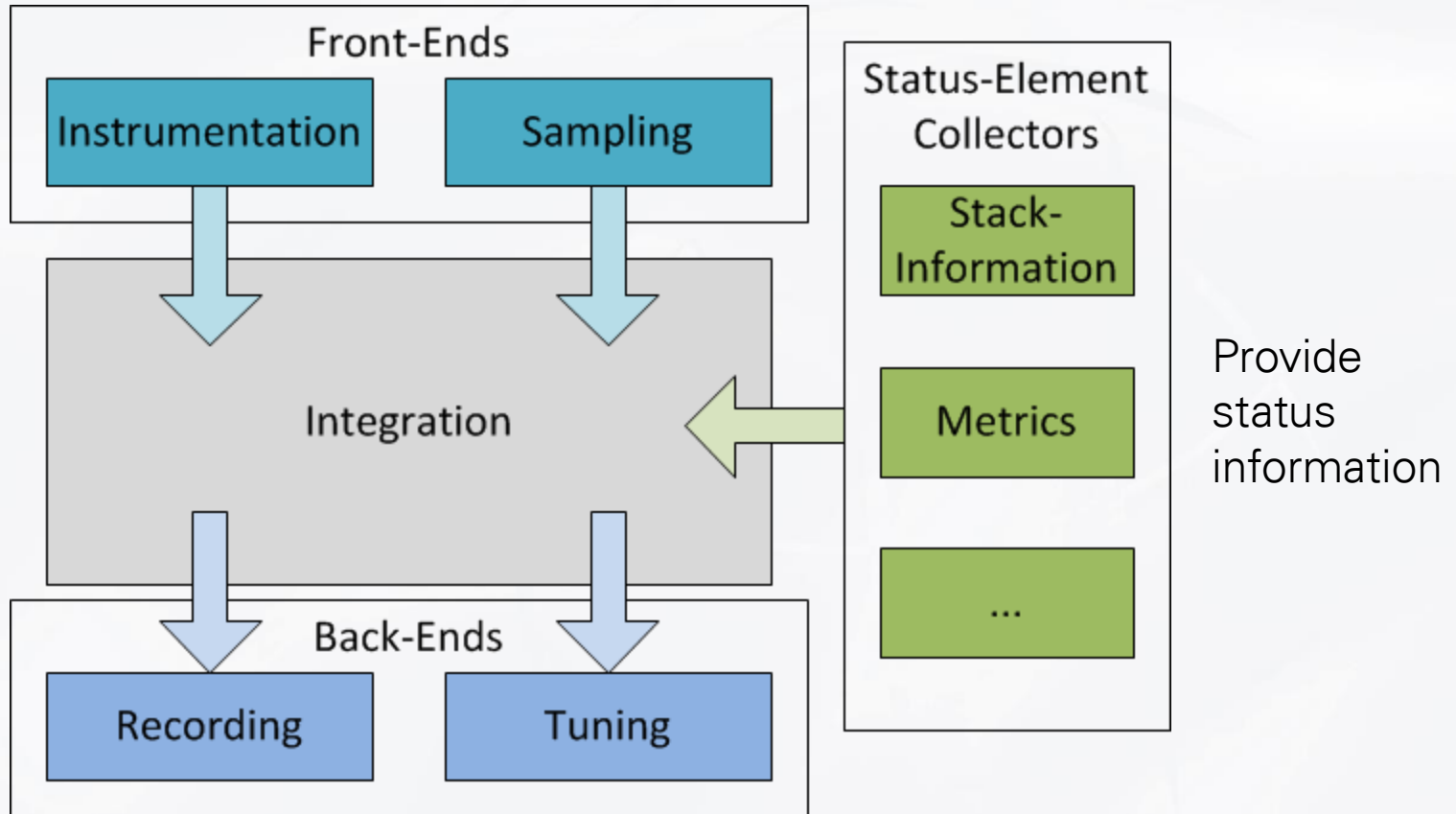
Ronny Tschüter

Thomas Ilsche

Joseph Schuchart

Daniel Hackenberg

**ZIH**
Center for Information Services &
High Performance Computing

- Technische Universität Dresden (Coordinator), Germany

- Norwegian University of Science and Technology, Norway

- IT4Innovations, Czech Republic

- Technische Universität München, Germany

- Intel European Exascale Labs, France

- GNS Braunschweig, Germany

- National University of Ireland Galway, Ireland

Robert Schöne

# Demands (Metrics)
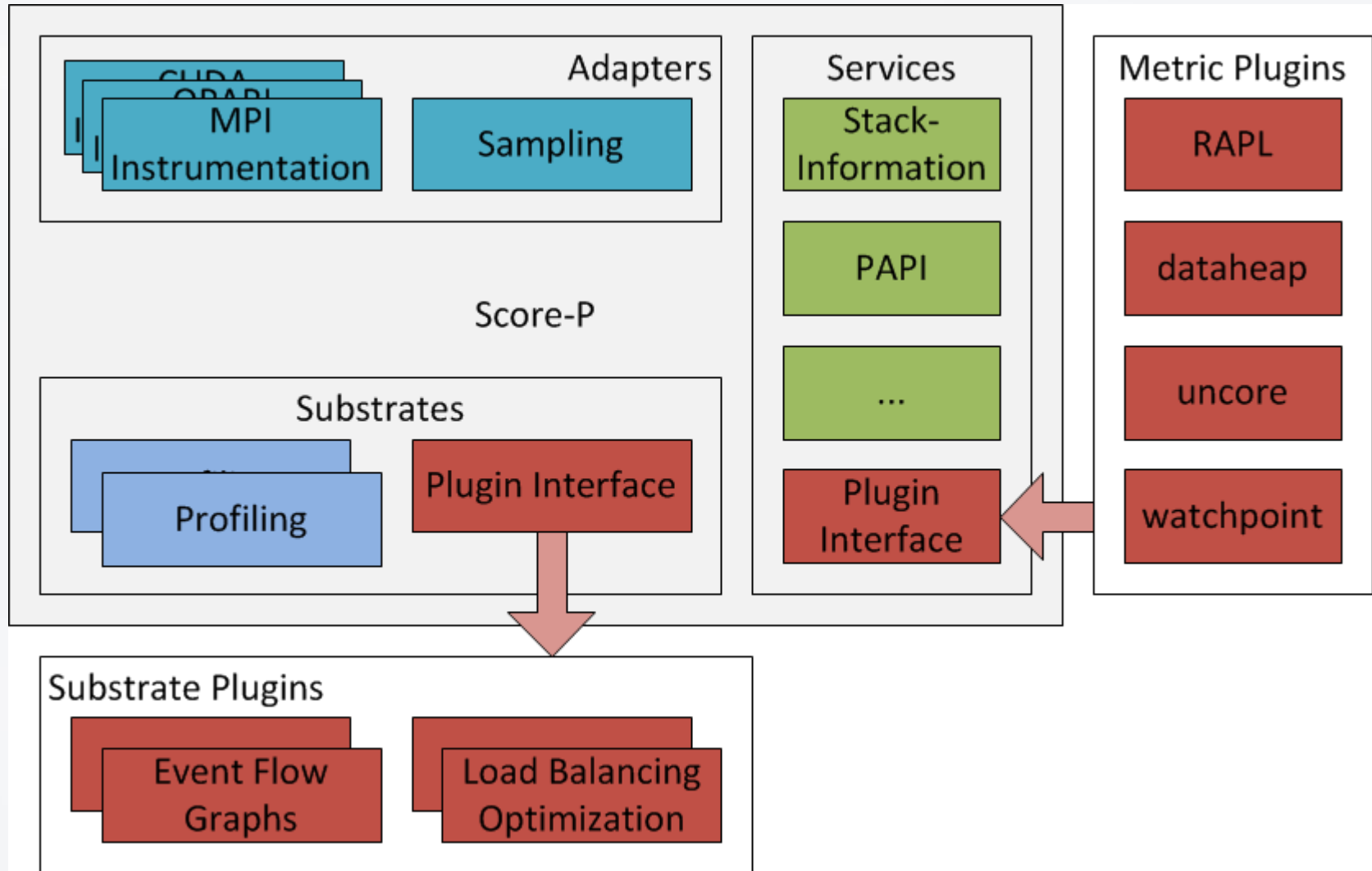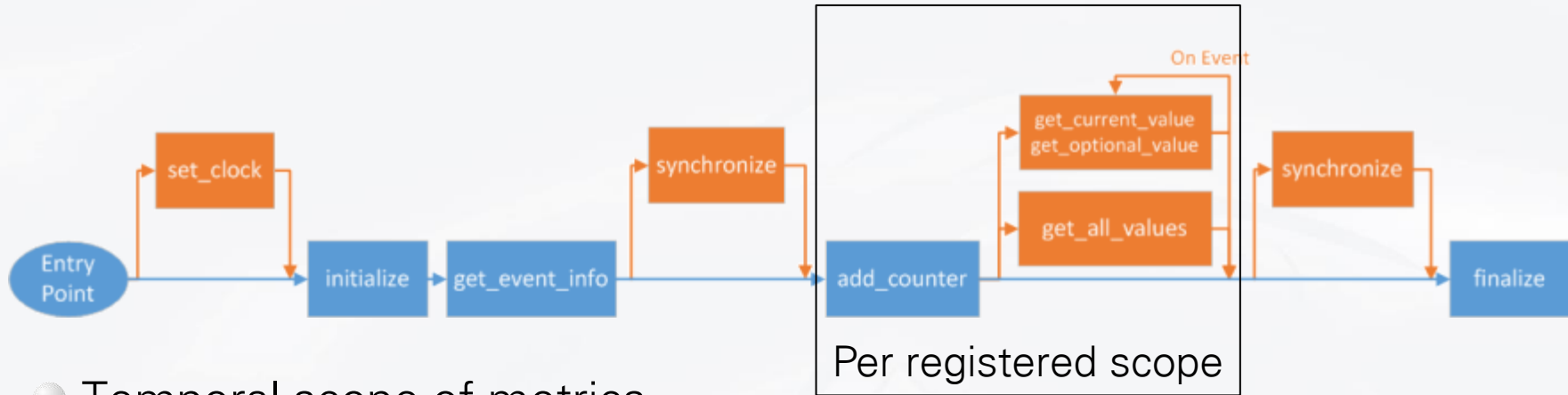
- Performance counter metrics

  – Get information on performance bottlenecks and dynamism

  – Part of Score-P

  – PAPI

  – perf

- Non traditional information

  – Open interface for different metrics

  – Power consumption information

  – Different scopes

  – Score-P Metric Plugins interface already available

  – Comparable to VampirTrace plugins [STHI11]

# Demands (Back-Ends)

- Consume program events

  - Act according to specific **regions**

  - Act according to **metrics**

  - Act depending on **location**

  - Act on instructions from **Online-Access** server

  - READEX Runtime Library

- Goal: Adapt HW/SW environment to increase energy efficiency

  - Independent of READEX

  - Different for different architectures → exchangeable

Robert Schöne

# Score-P overview

# Metric Plugins



- Temporal scope of metrics

  – synchronous/strictly synchronous/asynchronous

  – Instantaneous/backward-related/forward-related

- Spatial scope of metrics

  – Per thread/process/computing node/system

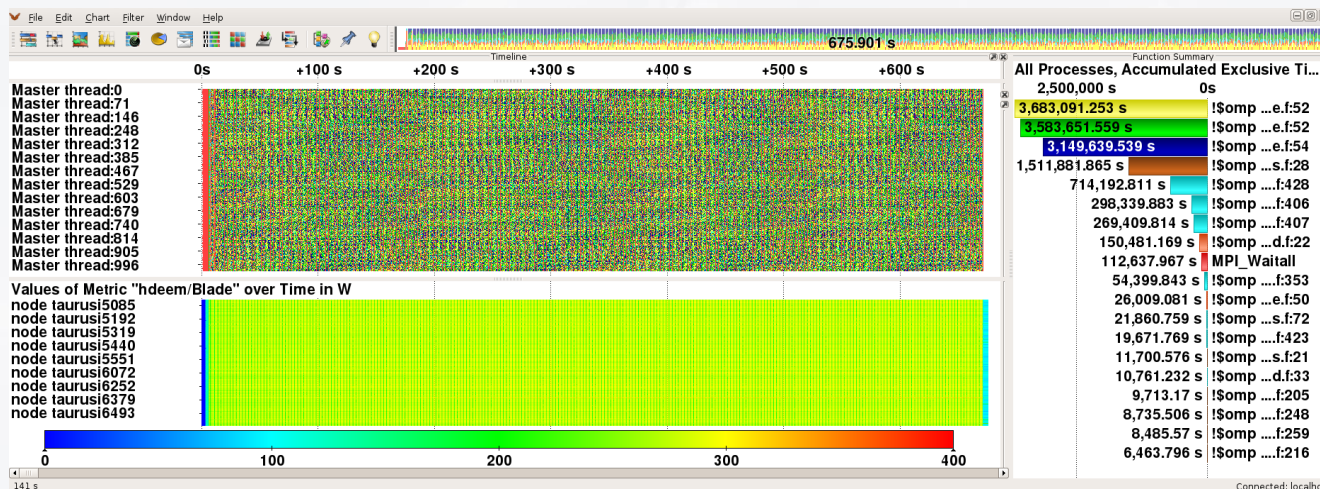- `SCOREP_METRIC_PLUGINS=foo # libfoo.so`

- `SCOREP_METRIC_FOO=bar # event bar in libfoo`

- Plugins on https://github.com/score-p/

# Examples Power Consumption



SPEC OMP applu, measurements with ZES-ZIMMER LMG 450/dataheap, RAPL, and HAEC infrastructure [HIS+13]



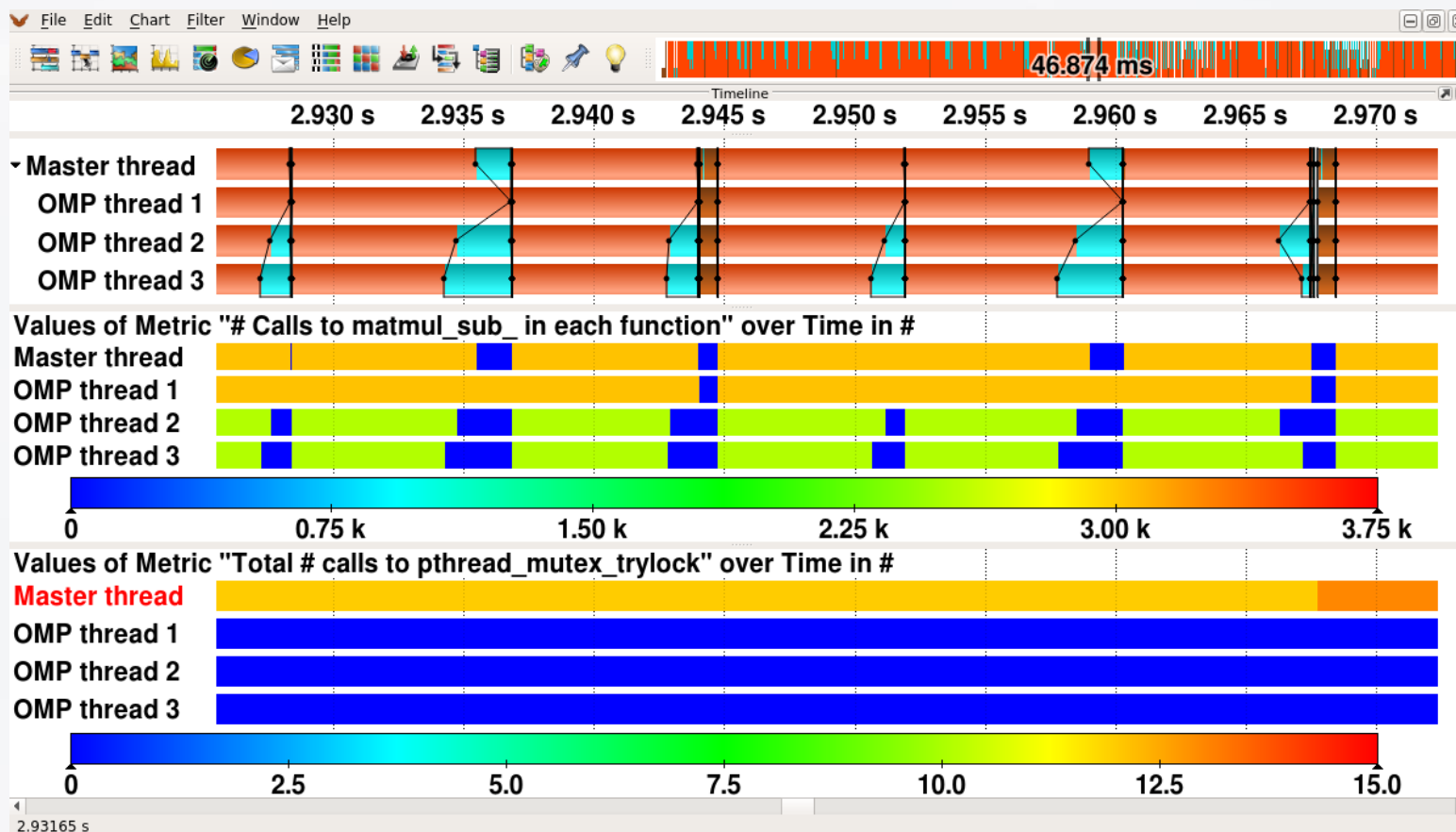NPB-MZ benchmark BT, Class F, 1024 nodes, measurement with HDEEM [HIS+14]

Robert Schöne

# Example Watchpoints

- Synchronous

  - Register variables or functions that shall be watched

  - Watch for number of read/write/execute access to variable

  - E.g.,
    `SCOREP_METRIC_PERFBREAKPOINT_PLUGIN=x_pthread_mutex_trylock`
    to watch the number of calls of `pthread_mutex_trylock`

  - Statistic information on accesses to global variables

- Asynchronous

  - Register variables whose content will be watched

  - Parse binary with libbfd, find local & global variables

  - E.g., `SCOREP_METRIC_WATCHPOINT_PLUGIN=foo:uint64_t:bar` registers the content of the local variable `bar`, defined in function `foo`

  - High overhead (from kernel infrastructure), for debugging purposes

# Example Watchpoints (Synchronous)



NPB-OpenMP benchmark BT, CLASS W, 4 Threads

Robert Schöne

# Example Watchpoints (Asynchronous)

- OpenMP parallel program, 2 Threads

- Threads concurrently write their loop iteration to a global variable

- Depending on scheduling strategy, the content of the global variable changes over time
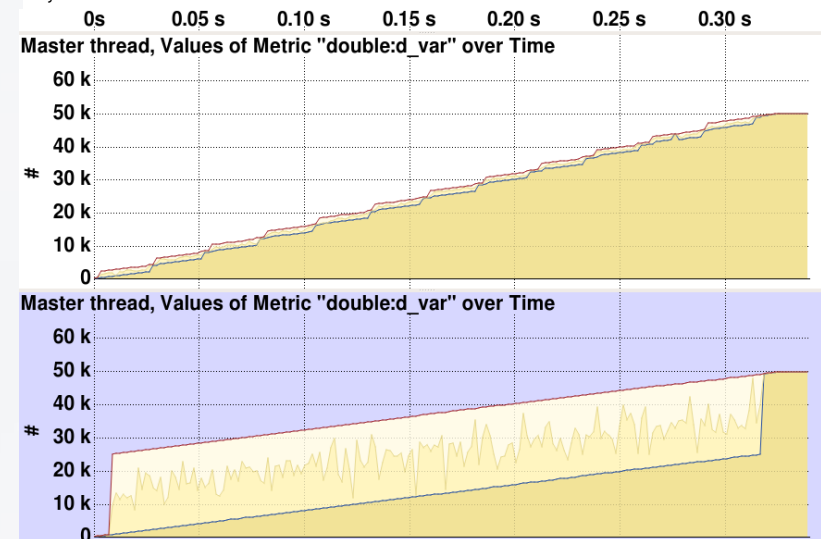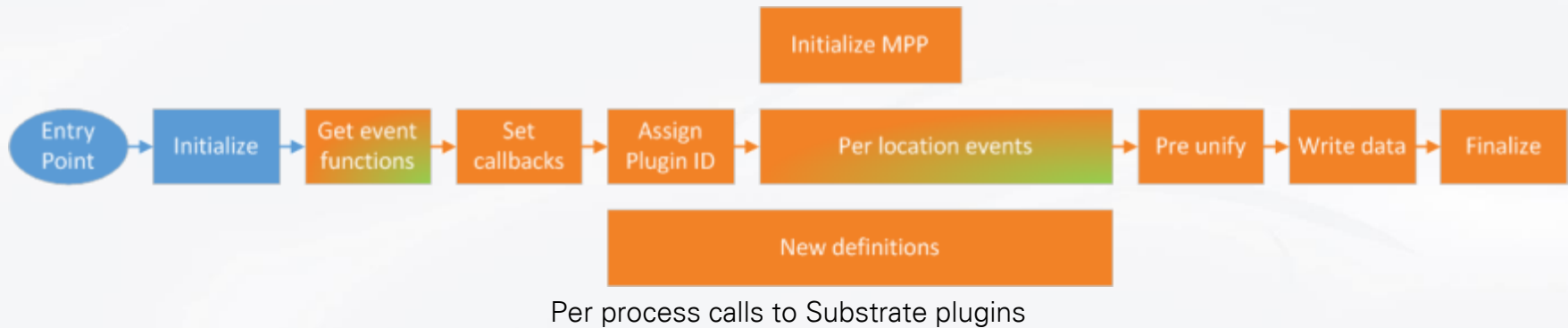
- White:

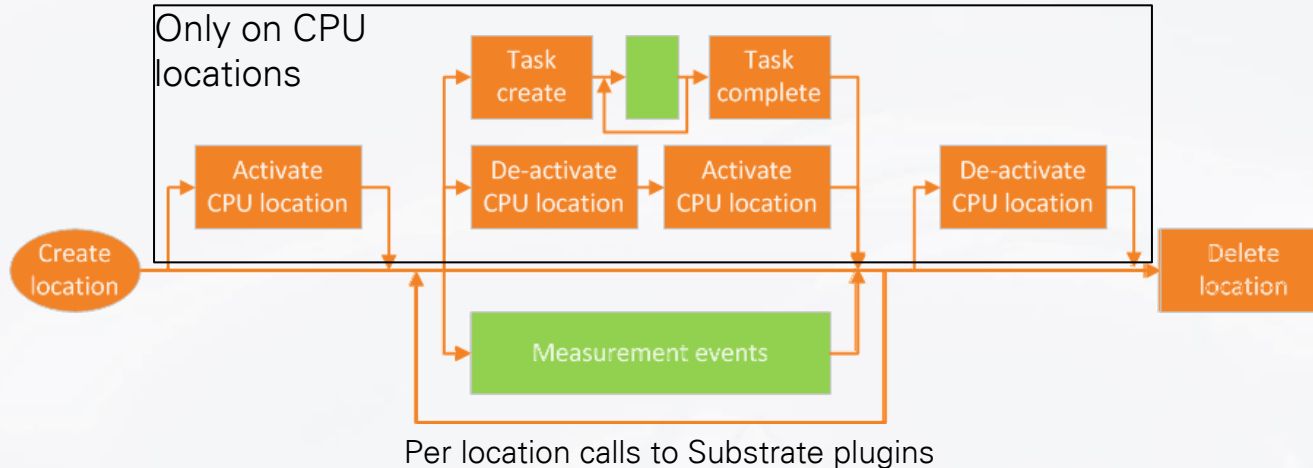  `OMP_SCHEDULE=dynamic,4096`

- Purple:

  `OMP_SCHEDULE=static`

```
static double d_var=0;
void func( int i ){
    d_var=0.5*i;
}
int main( int argc, char** argv ){
    int i=0;
#pragma omp parallel for schedule(runtime)
    for(i=0;i<100000;i++){
        func(i);
    }
    return 0;
}
```



Master thread, Values of Metric "double:d_var" over Time



Master thread, Values of Metric "double:d_var" over Time

TECHNISCHE UNIVERSITÄT DRESDEN

European Commission

Horizon 2020
European Union funding
for Research & Innovation

ZIH
Center for Information Services &
High Performance Computing

# Substrate Plugins

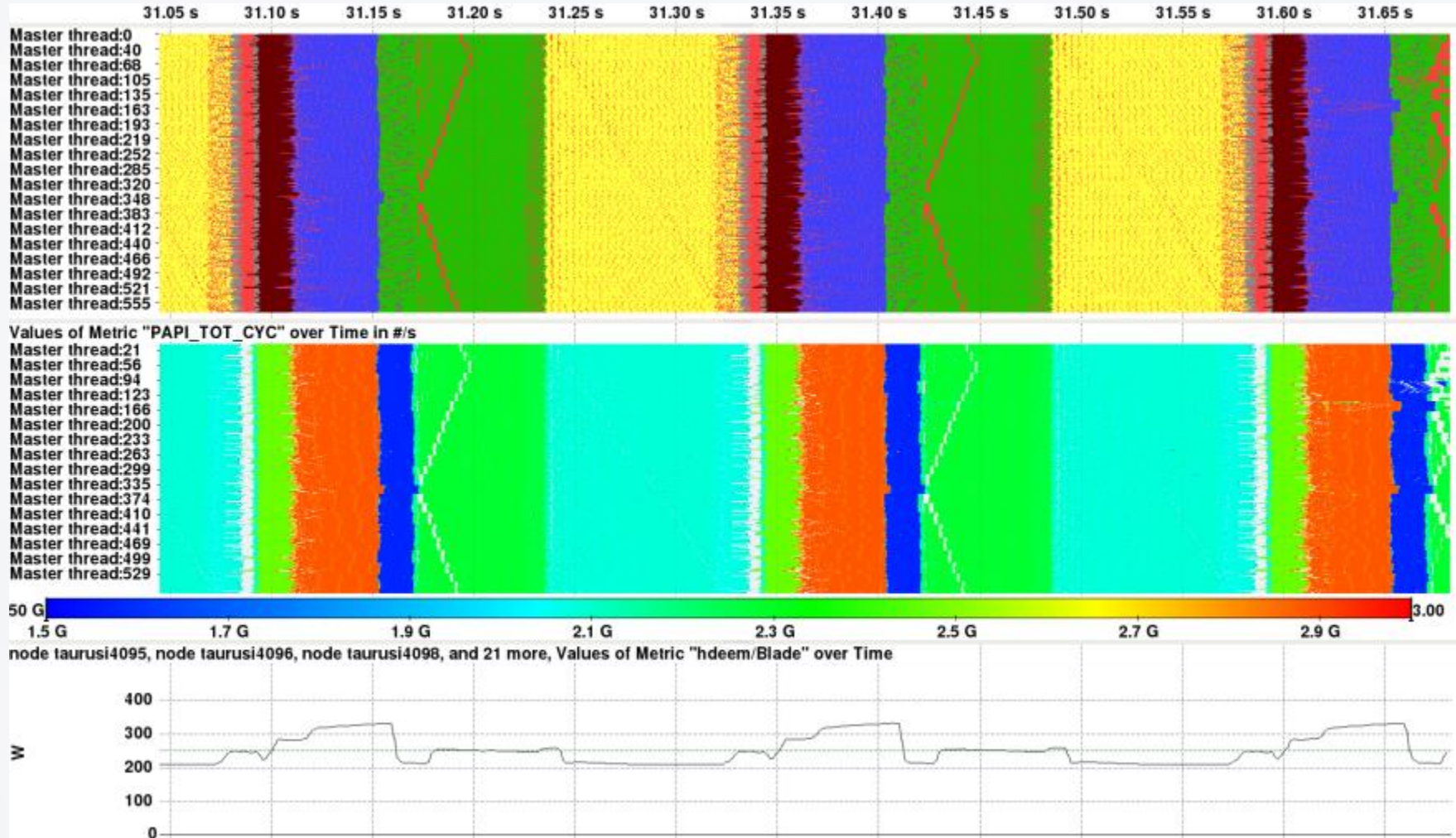

Per process calls to Substrate plugins

- Based on internal Substrate infrastructure

    – Management functions – related to Score-P internal processing

    – Event functions – related to program internals

- Only two non-optional functions

- Callbacks for plugins to get metadata for events from Score-P

- `export SCOREP_SUBSTRATE_PLUGINS=foo`

    `# load libscorep_substrate_foo`

Robert Schöne

# Substrate Plugins



Per location calls to Substrate plugins

- Calls for CPU locations

- All measurement events get location and timestamp

- Plugins can use callbacks to get metadata from handles (e.g., name of entered region, type of location, …)

- Support for storage of location specific data in Score-P (thread local storage)
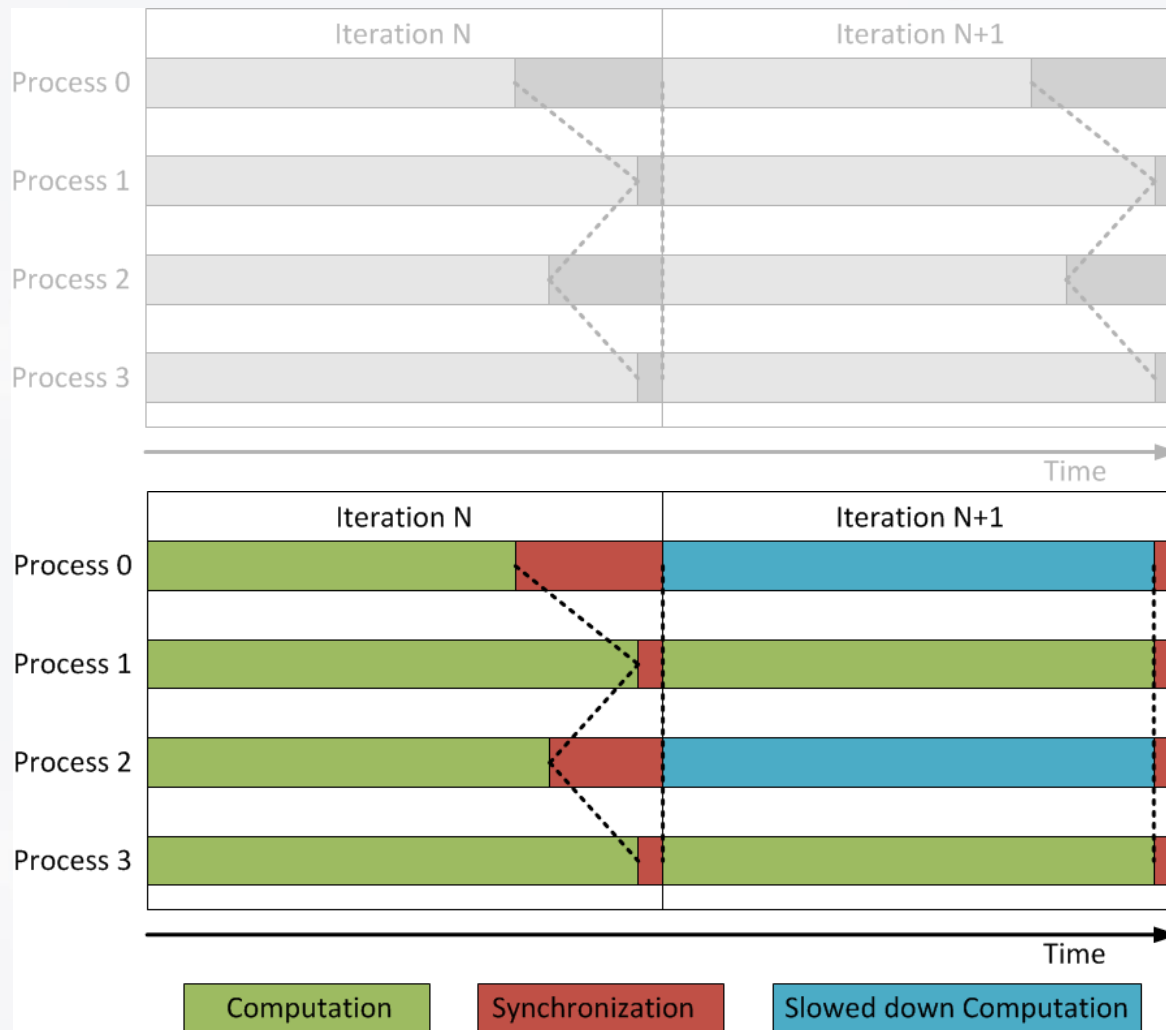
- Per region optimization using libadapt
  - As presented in [SM13]
  - Plugin reads config file and registers for enter/exit events to change HW/SW environment
  - DVFS/DCT/low-level hardware options

- Load balancing using DVFS
  - Like [RLdS+09], plus OpenMP
  - Successive execution of (compute region, synchronization region) pairs
  - Slow down computation to arrive just in time for synchronization
  - Plugin registers for enter/exit events

- Event Flow Graphs
  - Based on [AFL14]
  - Write event flow graph of exclusive regions as graphviz diagram
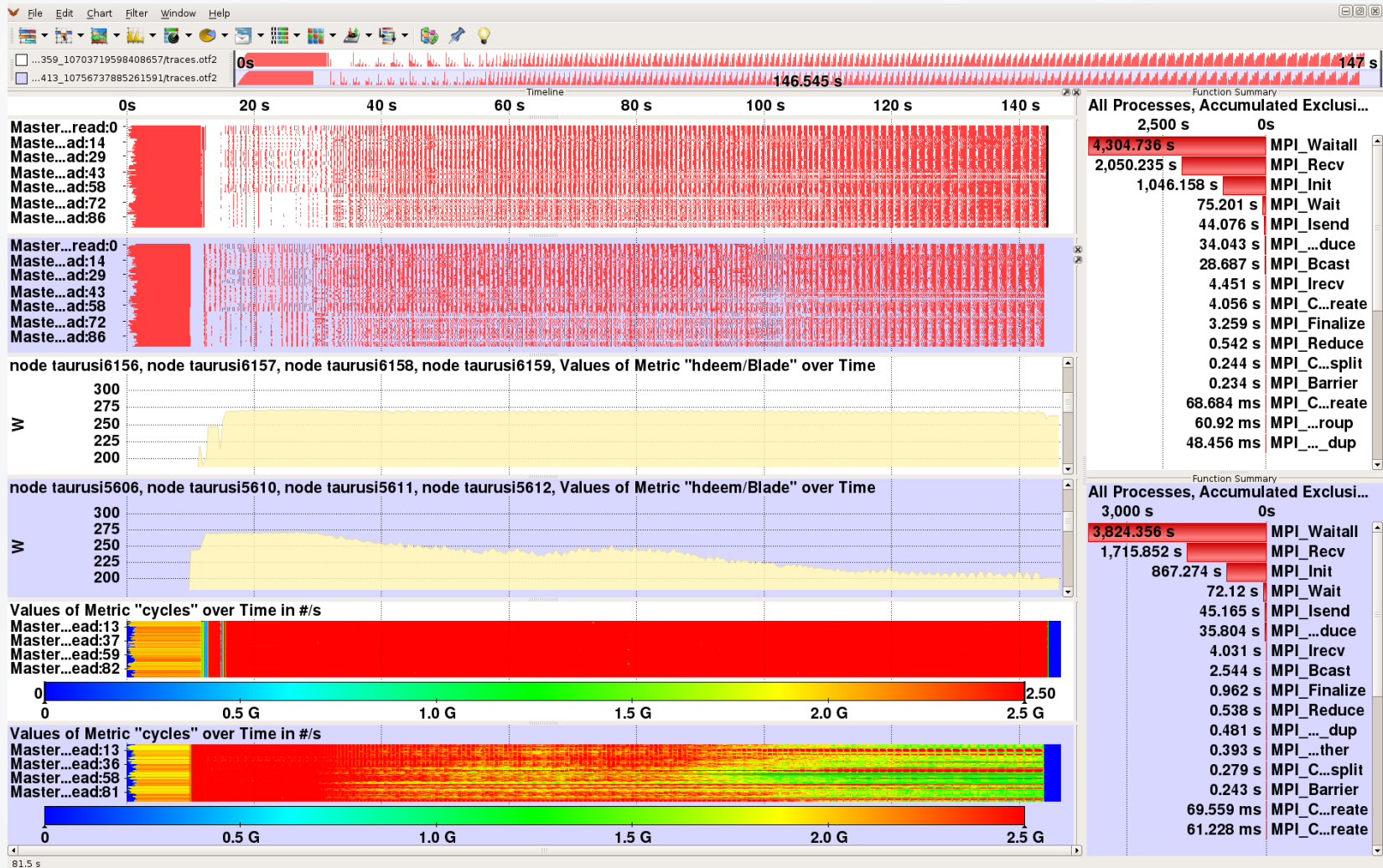
# Region Based Optimization



NPB-MPI benchmark BT, Class D, 576 ranks, optimized version: varying frequency and power consumption (HDEEM)

Robert Schöne

# Load Balancing Substrate
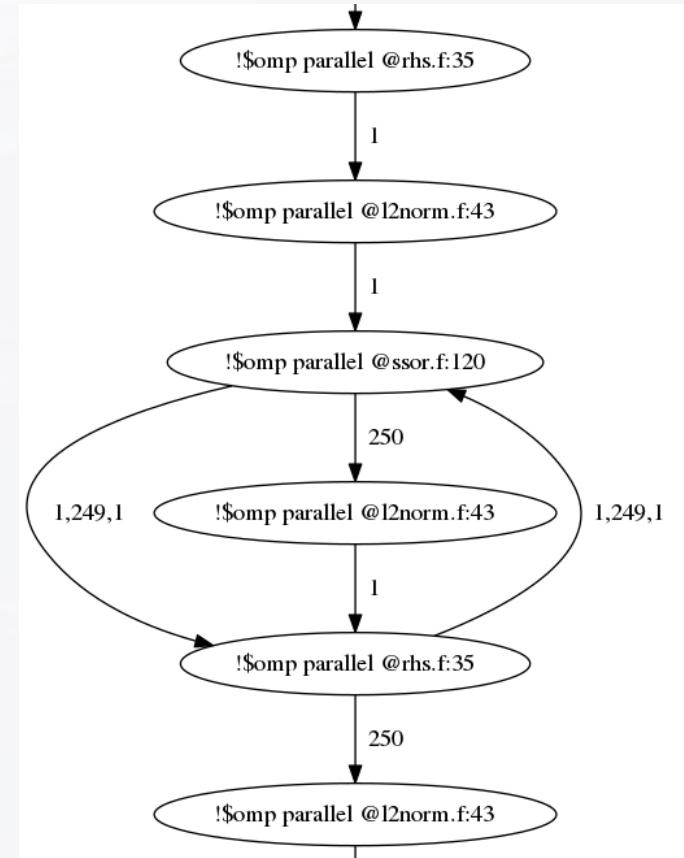


Computation | Synchronization | Slowed down Computation

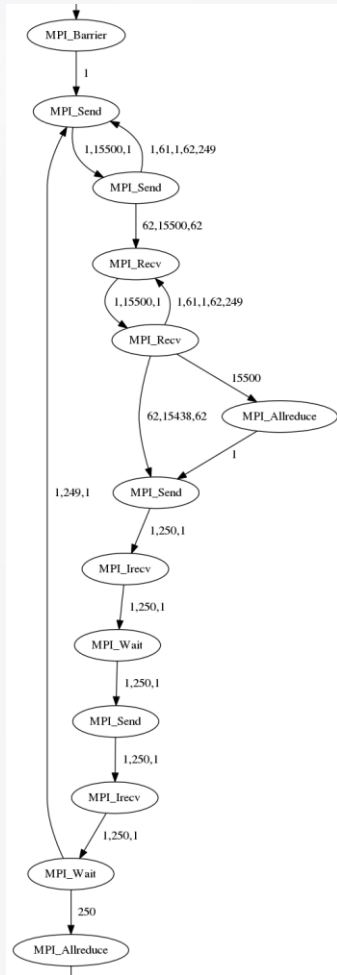# Load Balancing SPEC COSMO FD4 via DVFS



SPEC COSMO FD4 on 4 nodes / 96 ranks, MPI instrumentation, DVFS load balancing, HDEEM
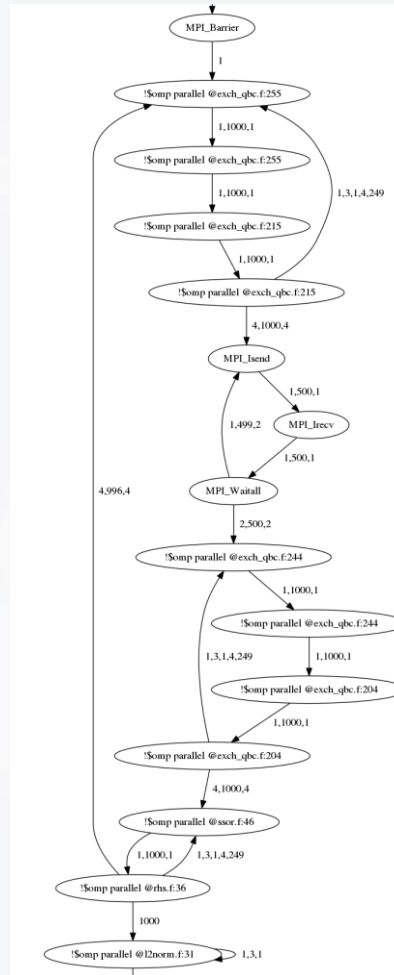
Robert Schöne

- Record exclusive functions, no nested calls

- E.g., OpenMP parallel regions

- Node = parallel region with specific stack state

- Edge = possible successor

- Node label: region name

- Edge label:

  - No label: Only transition from node

  - Single number N: on Nth iteration

  - Three numbers i,j,k: transition is taken ith, i+kth,…jth execution of predecessor node

  - …

# Event flow graphs (NPB lu CLASS A, core loop)



(a) MPI (rank 0)

(b) MZ (rank 0, thread 0)

(c) OMP (thread 0)

# Summary and Outlook

- Score-P as a common infrastructure for measurement, tuning, and debugging?

- Metric Plugins provide additional status information

- Substrate Plugins make use of Score-P's instrumentation frameworks and metrics to enhance functionality

Substrate Plugin Interface **not final**. Any suggestions?

Talk to me in the next coffee break or write an email!

Robert Schöne

- [RLdS+09] *Barry Rountree, David K. Lownenthal, Bronis R. de Supinski, Martin Schulz, Vincent W. Freeh, and Tyler Bletsch*. **Adagio: Making DVS Practical for Complex HPC Applications**. In Proceedings of the 23rd international conference on Supercomputing, pages 460–469. ACM, 2009. DOI: 10.1145/1542275.1542340.

- [STHI11] *Robert Schöne, Ronny Tschüter, Daniel Hackenberg, and Thomas Ilsche*. **The VampirTrace Plugin Counter Interface: Introduction and Examples**. In Euro-Par 2010 Parallel Processing Workshops, volume 6586 of Lecture Notes in Computer Science, pages 501–511. Springer-Verlag, 2011. DOI: 10.1007/978-3-642-21878-1_62.

- [HIS+13] *Daniel Hackenberg, Thomas Ilsche, Robert Schone, Daniel Molka, Martin Schmidt, and Wolfgang E. Nagel*. **Power Measurement Techniques on Standard Compute Nodes: A Quantitative Comparison**. In Performance Analysis of Systems and Software (ISPASS), 2013 IEEE International Symposium on, pages 194–204. IEEE, 2013.

- [SM13] *Robert Schöne and Daniel Molka*. **Integrating Performance Analysis and Energy Efficiency, Optimizations in a Unified Environment**. Computer Science - Research and Development, pages 1–9, 2013. DOI: 10.1007/s00450-013-0243-7.

- [AFL14] *Xavier Aguilar, Karl Fürlinger, and Erwin Laure*. **MPI Trace Compression Using Event Flow Graphs.** In Euro-Par 2014 Parallel Processing, pages 1--12. Springer International Publishing, Cham, 2014. DOI: 10.1007/978-3-319-09873-9_1

- [HIS+14] *Daniel Hackenberg, Thomas Ilsche, Joseph Schuchart, Robert Schöne, Wolfgang E Nagel, Marc Simon, and Yiannis Georgiou*. **HDEEM: High Definition Energy Efficiency Monitoring**. In Energy Efficient Supercomputing Workshop (E2SC), 2014, pages 1–10. IEEE, 2014. DOI: 10.1109/E2SC.2014.13.

# Backup Overhead Plugin

- Compare trunk, branch, branch with minimal registered plugin

- NPB OpenMP, BT Class A, 24 threads on Haswell dual-socket system

- many short regions, high overhead

- 5 measurements, use mean

- Average time increased by 1.5% if plugin registered, otherwise 0 overhead

- Sampling using perf record



trunk — bt.A.x 50%, scorep 22%, profile 17%, gomp 11%, plugin 0%, other 0%

branch — bt.A.x 50%, scorep 23%, profile 17%, gomp 10%, plugin 0%, other 0%

plugin — bt.A.x 49%, scorep 22%, profile 17%, gomp 10%, plugin 2%, other 0%